

CHEMOMETRICS FOR ANALYTICAL CHEMISTRY

Volume 2: PC-Aided Regression and Related Methods

MILAN MELOUN

Department of Analytical Chemistry, University of
Chemical Technology, Pardubice, Czech Republic

JIRI MILITKÝ

Department of Textile Materials, College of Mechanical and
Textile Design, Liberec, Czech Republic

MICHELE FORINA

Faculty of Pharmacy, University of Genoa, Italy

Translation Editors:

MARY M. MASSON and SUSSEN MATHEWS

Department of Chemistry, University of Aberdeen



ELLIS HORWOOD

NEW YORK LONDON TORONTO SYDNEY TOKYO SINGAPORE

First published 1994 by
Ellis Horwood Limited
Campus 400, Maylands Avenue
Hemel Hempstead
Hertfordshire, HP2 7EZ
A division of
Simon & Schuster International Group



© Ellis Horwood Limited 1994

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form, or by any means, electronic, mechanical, photocopying, recording or otherwise, without prior permission, in writing, from the publisher.

Printed and bound in Great Britain by
Hartnolls Limited, Bodmin, Cornwall

Library of Congress Cataloging-in-Publication Data

Available from the publisher

British Library Cataloguing in Publication Data

A catalogue record for this book is available from
the British Library

ISBN 0-13-123788-7

1 2 3 4 5 98 97 96 95 94

Table of contents

Preface	ix
Glossary	xi
6. Linear regression models	1
6.1 Formulation of a linear regression model	1
6.2 Conditions for the least-squares method	10
6.3 Statistical properties of the least-squares method	12
6.3.1 Construction of confidence intervals	20
6.3.2 Testing of hypotheses	24
6.3.2.1 Testing for multicollinearity	28
6.3.2.2 Test of significance of the intercept term	30
6.3.2.3 Simultaneous test of a composite hypothesis	33
6.3.2.4 Test of agreement of two linear models	36
6.3.2.5 Acceptance test for a proposed linear model	39
6.3.3 Comparison of regression lines	46
6.3.3.1 Test for homogeneity of intercepts	47
6.3.3.2 Test for homogeneity of slopes	48
6.3.3.3 Test for coincidence of regression lines	49
6.4 Numerical problems in the computer calculation of linear regression	52
6.4.1 The method of orthogonal functions	55
6.4.2 The method of rational ranks	57
6.5 Regression diagnostics	62
6.5.1 Exploratory regression analysis	62
6.5.2 Examination of data quality	64
6.5.2.1 Statistical analysis of residuals	64
6.5.2.2 Analysis of projection matrix elements	69
6.5.2.3 Plots for identification of influential points	72

6.5.2.4 Other characteristics of influential points	76
6.5.3 Examination of a proposed regression model	87
6.5.3.1 Partial regression leverage plots	87
6.5.3.2 Partial residual plots	89
6.5.3.3 Sign test for model specification	93
6.5.4 Examination of conditions for the least-squares method	94
6.5.4.1 Heteroscedasticity	94
6.5.4.2 Autocorrelation	95
6.5.4.3 Normality of errors	97
6.6 Procedures when conditions for least-squares are violated	98
6.6.1 Restrictions placed on the parameters	98
6.6.2 The method of generalized least squares (GLS)	102
6.6.2.1 Heteroscedasticity	104
6.6.2.2 Autocorrelation	110
6.6.3 Multicollinearity	116
6.6.4 Variables subject to random errors	121
6.6.5 Other error distributions of the dependent variable	126
6.6.5.1 The M-estimates method	126
6.6.5.2 The L_1 -approximation method	130
6.6.5.3 Robust estimates with bounded influence	133
6.7 Calibration	141
6.7.1 Types of calibration and calibration models	142
6.7.2 Calibration straight line	145
6.7.3 The precision of calibration	151
6.8 Procedure for linear regression analysis	155
6.9 Additional problems	157
References	175
7. Correlation	178
7.1 Correlation models	179
7.1.1 Correlation models for two random variables	179
7.1.2 The correlation model for many random variables	185
7.2 Correlation coefficients	194
7.2.1 Paired correlation coefficient	194
7.2.2 Partial correlation coefficients	200
7.2.3 Multiple correlation coefficient	202
7.2.4 Rank correlation	204
7.3 Procedure for correlation analysis	206
References	206
8. Nonlinear regression models	207
8.1 Formulation of a nonlinear regression model	209
8.2 Models of measurement errors	214
8.3 Formulation of the regression criterion	219
8.4 Geometry of nonlinear regression	225

8.5 Numerical procedure for parameter estimation	235
8.5.1 Non-derivative optimization procedures	236
8.5.1.1 Direct search methods	237
8.5.1.2 Simplex methods	238
8.5.1.3 Random optimization	244
8.5.1.4 Special procedures for the least-squares method	247
8.5.2 Derivative procedures for the least-squares method	250
8.5.2.1 Gauss–Newton methods	253
8.5.2.2 Marquardt-type methods	257
8.5.2.3 Dog-leg type procedures	259
8.5.3 Complications in nonlinear regression	260
8.5.3.1 Parameter estimability	260
8.5.3.2 Existence of a minimum of $U(\mathbf{b})$	262
8.5.3.3 Existence of local minima	263
8.5.3.4 Ill-conditioning of parameters	264
8.5.3.5 Small range of experimental data	264
8.5.4 Examination of the reliability of the regression algorithm	267
8.6 Statistical analysis of nonlinear regression	270
8.6.1 Degree of nonlinearity of a regression model	272
8.6.1.1 Bias of parameter estimates	273
8.6.1.2 Asymmetry of parameter estimates	276
8.6.2 Interval estimates of parameters	277
8.6.2.1 Confidence regions of parameters	277
8.6.2.2 Confidence intervals of parameters	282
8.6.2.3 Confidence intervals of prediction	285
8.6.3 Hypothesis tests about parameter estimates	286
8.6.4 Goodness-of-fit tests	288
8.6.4.1 Graphical analysis of residuals	289
8.6.4.2 Statistical analysis of residuals	289
8.6.4.3 Identification of influential points	292
8.7 Procedure for building and testing a nonlinear model	294
8.8 Additional problems	295
References	303
9. Interpolation and approximation	306
9.1 Classical interpolation procedures	307
9.1.1 Lagrange and Newton interpolation formulae	309
9.1.2 Hermite interpolation	315
9.1.3 Rational interpolation	316
9.2 Spline interpolation	319
9.2.1 Local Hermite interpolation	324
9.2.2 Cubic spline	328
9.3 Approximation of functions	336
9.4 Approximation of tabular data	340
9.4.1 Polynomial approximation	340

9.4.2 Piecewise regression	342
9.5 Numerical smoothing	349
9.5.1 Smoothing spline	350
9.5.2 Nonparametric regression	360
9.5.3 Digital filtration	362
9.6 Procedures for interpolation and approximation	372
9.7 Additional problems	373
References	374
10. Derivatives and integrals	375
10.1 Derivatives	376
10.1.1 Analytical derivatives	378
10.1.2 Numerical derivatives	378
10.2 Integrals	381
10.2.1 Analytical integration	381
10.2.2 Numerical integration	381
10.3 Procedure for numerical differentiation and integration	387
10.4 Additional problems	388
References	390
Appendix: CHEMSTAT	391
Index	393

Glossary

Akaike's information criterion (AIC). A model selection criterion to distinguish between models with different numbers of parameters.

analysis of residuals. Part of regression diagnostics, based on examining the discrepancy between the model and the observed data.

Andrews–Pregibon (AP) statistic. A diagnostic for finding influential points.

approximation. A set of techniques for replacing an unknown function by a simpler empirical model.

***a posteriori* class probability $p(c/x)$.** The probability that the object described by the vector x belongs to the class c .

***a priori* class probability $p(c)$.** Probability of the occurrence of the class c .

Atkinson influence statistic. A diagnostic for finding influential points.

augmented distance. In SIMCA, a distance that takes into account the residuals and the boundaries of the model.

autocorrelation. The internal correlation between members of a series of observations ordered in time or space.

autocorrelation function. A function of a time-dependent variable, defined for a stationary stochastic process.

Bayes classifier. A classification rule that minimizes the conditional average risk.

Bayes rule. The theorem that computes the *a posteriori* probability from the *a priori* probability and the conditional probability.

best linear unbiased estimator (BLUE). A linear and unbiased estimator with minimum variance.

biased estimator. An estimator with an expected value that is different from the true value of the estimated parameter.

calibration. A two-phase procedure. In the first phase (calibration model-building) the curve describing the relationship between the response y and the independent (explanatory) variable x is estimated. In the second phase (calibration), the value of the variable x^* corresponding to a measured response y^* is estimated.

canonical variables of LDA. Directions of maximum information, useful for

classification in LDA.

category. Same as class.

chi-squared test for goodness-of-fit. A widely used statistic that, when associated with discrete data, can test whether the observed frequencies deviate significantly from the theoretical frequencies.

class model. A mathematical model of a class plus the allowed dispersion due to errors.

class space. A space surrounding the class centre; when an object is in this space, it fits the class model.

class-modelling analysis. The determination of the characteristics of a class.

class-modelling method. A technique that computes a class model and a class space.

class. A population of objects with similar properties.

classification analysis. The determination of the class of one or more objects.

classification loss. The experimental measure of the conditional average risk.

classification method. A technique that produces a classification rule.

classification rate. Percentage of the objects in the training set correctly classified by a classification rule.

classification rule. A mathematical equation that assigns objects to one of several categories.

coefficient of determination. The square of the multiple correlation coefficient.

collinearity. An approximate linear relationship among the explanatory variables. This has a negative effect on the modelling power in regression, but does not necessarily reduce the goodness-of-fit.

conditional average risk. For a category, the expected loss for the classification in the category of an object of other categories.

conditional probability $p(\mathbf{x}/c)$. The probability that an object of the class c is described by the vector \mathbf{x} .

confidence interval. The interval between the upper and lower confidence limits of the estimate of a parameter (e.g., of location or spread). The theoretical parameter lies in this interval with selected probability.

confidence region. Region in which the vector of theoretical parameters lies, with a selected probability.

constant potential. A PFM where the objects give the same contribution to the computed probability density function.

Cook's influence statistic. A diagnostic based on analysis of residuals, for detection of influential points.

Coomans plot. A plot of the distance (or the square of distance) from their model of the objects in two categories; also a plot with the distance of the same objects from two class models, computed by different techniques.

correlation. A quantitative measure of the linear association between two or more random variables.

correlation coefficient (multiple). A measure of linear association between response y and a set of explanatory variables x_1, \dots, x_n . When it is close to zero, the variability of response cannot be explained by the explanatory variables. When it is close to one, a linear relationship exists between y and x_1, \dots, x_n .

correlation coefficient (paired). A scaled version of the covariance between two variables. When close to -1 or $+1$, it indicates strong correlation between variables, but not dependence.

correlation matrix. A symmetric matrix with nondiagonal elements that are paired correlation coefficients, and diagonal elements equal to one.

covariance. The first product moment of two variables about their mean values; the expectancy that the values measured for two objects deviate in a similar way from the true means.

covariance matrix. A symmetric matrix in which the values are the covariances.

critical distance. The distance of the class boundary from the centre of the category, in a class-modelling technique.

curve-fitting. Searching for a mathematical expression relating a set of observed values to a mathematical function or to another set of observed values.

Defrise correction. An improvement of the estimate of the Mahalanobis distance.

dependent variable. A term indicating a mathematical or statistical dependence of a variable on one or more other variables. In regression problems, it is often called the response variable or simply the response.

diagnostics of regression. See *regression diagnostics*.

digital filters. A technique for elimination of noise from a signal.

discriminant function. Difference between the discriminant scores for two categories.

discriminant scores. Scores measuring the degree of belonging of an object to a category in LDA.

discriminant power. The measure of the importance of a variable in the separation of two categories using SIMCA.

distance. (in class-modelling techniques). A measure of the fit to the mathematical model, centre of the model.

doubt interval. A non-decision interval around the delimiter used with classification rules.

error mean square. The expected square of the difference between the true and estimated value of either the regression coefficients or the response variable. It is equal to the error sum of squares divided by the number of degrees of freedom.

error sum of squares. The sum of squared differences between the observed and predicted response. It is a measure of goodness-of-fit.

evaluation set. The objects used to validate classification rules or class models.

extended range. In SIMCA, the range used to build the model.

Gauss-Newton method. An optimization method used for minimizing a quadratic nonlinear function.

generalized least-squares regression (GLS). A modification of ordinary least squares to deal with heteroscedasticity and correlated errors.

generalized linear model (GLM). A statistical model that consists of a random response variable and a set of predictor variables.

goodness-of-fit. A measure of how well a regression model accounts for the variance of the response variable.

goodness-of-prediction. A measure of how well a regression model estimates the value of the response variable, given a set of values for the predictor variables.

gradient method. An optimization method which, in the calculation of the set of parameters that minimizes a function, requires the evaluation of the derivatives of the function as well as the function values themselves.

gradient vector. A set of the first derivatives of a function f with respect to its parameters.

hat (projection) matrix. A matrix enabling orthogonal projection of an n -dimensional vector into m -dimensional space spanned by vectors corresponding to the explanatory variables.

Hermite interpolation. A special technique for interpolation of data and calculation of the first derivatives simultaneously.

Hessian matrix. The matrix of the second derivatives of a scalar valued function f with respect to its parameters.

heteroscedastic. Having unequal variances; the opposite of homoscedastic.

high-leverage point. See leverage point.

homoscedastic. Having equal variances; if the variance of one variable is the same for all values of the other, the distribution is said to be homoscedastic in the first variable.

homothetic. With the same value of the probability density function.

index-residual plot. The dependence of residuals on the index of measurements.

influence analysis. A method for detection of influential points.

influential points. Particular points that exert great influence on regression parameters and the results of regression analysis.

instrumental variable. An independent variable singled out in estimation of regression coefficients when the random error is correlated with the independent variables.

iteratively reweighted least squares. A technique for iterative computation of robust estimators based on special weighting in each iteration.

Jack-knife residuals. A type of residual for detection of outliers in data.

Jacobian matrix. A matrix containing the first derivatives of a vector.

joint distribution. A simultaneous distribution of a random vector.

KNN. A distance-based non-parametric classification method.

L_1 -estimator. Least absolute values estimator minimizing sum of absolute values of residuals.

L_2 -estimator. Least-squares estimator minimizing sum of squares of residuals.

L_∞ -estimator. An estimator minimizing maximal absolute residuals.

LDA. Linear statistical discriminant analysis, a classification method based on the multivariate normal distribution and the pooled covariance matrix.

least absolute residual regression. See L_1 -estimator.

least-squares estimator. See L_2 -estimator.

least-squares regression. Regression method where the model is linear in parameters and the regression coefficients are calculated with the least-squares estimator.

leave-more-out. A validation procedure, where the evaluation set is obtained in several steps; in each step some objects are omitted from the training set.

leave-one-out. A validation procedure, where the evaluation set is obtained in as many steps as there are objects in the category; at each step, one object is omitted from the training set.

leverage. The extent of the influence of an explanatory variable on the regression model. A measure of the leverage of a point is the corresponding diagonal element of the hat matrix.

linear estimator. Estimator where the predicted response variables can be expressed as a linear combination of the observed response variables.

linear regression model. A regression where the response variable is a linear function of the regression coefficients.

LM. Learning machines, a non-parametric non-probabilistic classification method.

loss matrix. A matrix of the loss incurred when an object is assigned to a false category.

M-estimator. A robust estimator based on the maximum likelihood estimator for a particular error distribution.

Mahalanobis distance. A distance weighted by the standard deviations and the correlations of the variables.

Marginal distribution. The probability distribution for a single variable in a multivariate problem.

maximum likelihood estimator (MLE). An estimator maximizing the likelihood function.

McCulloh–Meeter plot. A graphical diagnostic for discovering outliers and leverages.

mean absolute deviation. A robust scale estimator; the average of the absolute values of deviations.

mean covariance matrix. The average of the category covariance matrices.

mean squared error. The expected squared difference between the true and estimated value of either the regression coefficient or the response variable.

minimax strategy. If a strategy is selected from a group of admissible strategies as being the one which, on a basis of the expected loss, has the smallest maximum loss, this strategy is a minimax strategy.

misclassification probability. A measure of misclassification based on the *a posteriori* probabilities.

model (mathematical). A mathematical description of some part of reality.

model error. Or residual is the quantity remaining after some other quantity has been subtracted. In regression, it is the part of the response variable not described by the regression model, i.e. the difference between the observed and the predicted value of the response variable.

model parameter. A term used to denote an unknown quantity which may vary over a certain set of values. The values of the parameters are estimated to get the best possible fit for given data.

model sum-of-squares. The sum of squared differences between the estimated responses and the average response, or between the observed variable and the predicted response variable.

modelling power. A measure of the importance of a variable in a SIMCA model.

multicollinearity. An approximate linear relationship among the predictor variables. It causes high variance in the least-squares estimates of the regression coefficients, resulting in instability in the estimated values or even wrong signs.

multiple correlation. A quantitative measure of the linear relationship between several

random variables.

multiple correlation coefficient. A measure of the linear association between the observed response variable and the predicted response variable, or between the observed response variable and a linear combination of the predictor variables in a regression model.

multiple regression: (multivariate regression). A regression model where the response is a function of more than one predictor variable.

multiple regression analysis (MRA). A regression of more than one predictor variable.

natural histogram. A histogram without class intervals; the contribution of an object is positioned at the true value of the variable.

nonlinear least-squares regression (NLS). Techniques for estimating parameters in nonlinear regression models.

nonlinear regression model. A regression where the response variable is a nonlinear function of the regression coefficients.

nonlinear partial least-squares (NLPLS). An extension of the PLS regression model, which permits response to be modelled as a nonlinear function of the latent variables.

nonparametric regression. A technique for creating an approximate relationship, based on a weighted linear combination of responses.

normal range. In SIMCA, the range of the scores for principal components.

optimization. Finding the optimum value (minimum or maximum) of a function called the objective function, with respect to parameters.

parameter. See model parameter.

partial least-squares. A biased regression method that relates a set of predictor variables x to a set of response variables y . The least-squares regression is performed on a set of uncorrelated latent variables that are standard linear combinations of the original predictor variables.

PC. Principal (significant) components.

PFM. Potential functions method; a classification technique based on a potential function density computed as the sum of individual contributions of each object in the training set.

pooled covariance matrix. A weighted average of the category subgroup covariance matrices.

predicted value. A value calculated from a statistical model. When it is calculated by fitting an object to a parametric model, it is called the fitted value.

prediction. The process of finding the value of a variable based on a statistical model.

prediction rate. Percentage of the objects in the evaluation set correctly classified by a classification rule.

Pregibon's plot. A graphical diagnostic for finding outliers and leverages.

proportional sample. A sample where the number of objects in the categories is proportional to the *a priori* class probability.

QDA. Quadratic discriminant analysis; a classification method based on the multivariate normal probability density and on the class.

recursive residual. A type of residual for detection of influential points.

regression analysis. A collection of statistical methods that are used to model the relationships among measured or observed quantities.

regression coefficient. The coefficient of a predictor variable or independent variable in a regression model.

regression curve. A graphical presentation of a regression model.

regression diagnostic. A set of techniques used to detect and assess the degree of discrepancy between the model and the observed data.

regression model. A mathematical equation describing the relationship among explanatory and response variables.

regression parameter. See model parameter.

regression surface. In the case of multiple regression, the model is represented by a response surface.

repeated evaluation set. A validation procedure, where the training set is created many times, with a random extraction of objects from the sample.

residual. Error, difference between a datum and the corresponding predicted value.

residual analysis. Part of regression diagnostics.

residual mean square. The squared estimated model error, i.e. the expected squared difference between the true and estimated value of either the regression coefficients or the response variables.

residual plot. A scatter plot of the residuals *vs.* the independent variable.

residual sum of squares (RSS). The sum of the squared differences between the observed and predicted response.

ridge regression. A biased regression based on the assumption that large regression coefficients are likely to be caused by multicollinearity; it shrinks them toward zero by adding a small constant to each diagonal element of the covariance matrix.

robust bounded-influence regression (GM estimator). Limits the influence of leverages in a regression model by means of some weight function.

robust iteratively reweighted least-squares regression. Estimates of regression parameters found by minimizing a weighted criterion of least squares. The weights are calculated simultaneously with the estimates of standard deviation, in an iterative fashion.

robust least absolute residual regression (L_1 -regression). Estimates of regression parameters found by minimizing the sum of absolute residuals. This method safeguards against outliers in the response but not against outliers in the predictors.

scatter plot. A cartesian diagram showing the joint variation of two variables (e.g. x and y).

scedasticity: dispersion, especially as measured by variance. In a bivariate distribution, the graph of the variance one variable *vs.* the variance of the other is called a scedastic curve.

SIMCA. A class-modelling technique based on principal components.

simplex method. An optimization by a direct search method based on comparing the values of a function at the vertices of a simplex, and moving the simplex toward the optimum by an iterative procedure.

single evaluation set. A validation procedure in which the evaluation set is created only once.

single regression. See univariate regression.

SLDA. Stepwise linear discriminant analysis, a procedure for selection of features.

smoothing. The process of removing fluctuations from an ordered series of data so that the resulting trend is smooth, the main differences become regular, and higher order differences small.

smoothing coefficient. A parameter that, multiplied by the standard deviation of a variable, gives the smoothing parameter, in PFM.

smoothing parameter. A parameter determining the dispersion of the contribution of an object to the probability density function in PFM.

spline curves. These consist of polynomial segments smoothly joined. Cubic splines have continuous first and second derivatives at the joins (knots).

spline function. A segmented polynomial function of class C^m .

steepest descent. An optimization that minimizes a function by estimating the optimal values of parameters by a linear search method in the direction of negative gradient.

stepwise regression. A method which models the response variable as a function of only a selected subset of the predictor variables. It is a biased regression, because it is based on the assumption that not all predictor variables are relevant in the regression problem.

stochastic. Implies the presence of a random variable.

stochastic process. A process that incorporates an element of randomness.

test set. Objects of unknown class; the classification of these objects is the final aim of classification analysis.

training set. The objects used to compute classification rules or class models.

unequal covariance matrix classification (UNEQ). Modelling version of QDA.

univariate regression (single regression). A regression model where the response variable is a function of just one explanatory variable.

variable. A symbol (x , y etc) representing an unspecified member of a class of objects, numbers, etc.

variable, random. A quantity that varies with a given frequency distribution, so that values occur with specific probabilities.

variable potential. A PFM method where the objects give a different contribution to the computed probability density, according to their distance from the K nearest neighbour.

Williams plot. A graphical diagnostic for detection of influential points.

Wilks' lambda. A parameter measuring the separation between categories.

6

Linear regression models

6.1 FORMULATION OF THE LINEAR REGRESSION MODEL

In instrumental methods of chemical analysis, the instrument's response y (output variable) for selected values of the input variables \mathbf{x} is often measured. For example, an absorbance A (here, the output variable y) is measured on the scale of a spectrophotometer at

- (1) a selected value of wavelength λ (here, the first independent variable $x_{i,1}$),
- (2) a concentration of colour-forming solution c (here, the second independent variable $x_{i,2}$),
- (3) a value of adjusted pH of solution (here, the third independent variable $x_{i,3}$), and
- (4) in kinetic measurements, at an actual time (here, the fourth independent variable $x_{i,4}$).

This results in n observed values of y , measured at four kinds of selected values of independent variables, $m = 4$, written as $\{y_i, x_{ij}\}$, $i = 1, \dots, n$, and $j = 1, \dots, m$

$$\begin{bmatrix} y_1 \\ y_2 \\ \cdot \\ \cdot \\ \cdot \\ y_n \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ x_{n1} & x_{n2} & x_{n3} & x_{n4} \end{bmatrix}$$

In matrix notation, this is written as $\{y, \mathbf{X}\}$. Vector y has dimensions $(n \times 1)$ and matrix \mathbf{X} has dimensions $(n \times m)$.

The statistical analysis is intended to find a relationship between the response (output) variable y and the controllable (independent) variables \mathbf{x} . The type of function $y = f(\mathbf{x}, \boldsymbol{\beta})$ depends on the nature of both the variables y and \mathbf{x} . There are three possible scenarios.

- (a) Variables y and \mathbf{x} have no random errors. The function $y = f(\mathbf{x}, \boldsymbol{\beta})$ contains

a vector of unknown parameters β of dimension $(m \times 1)$. To estimate them, at least $n = m$ measurements y_i , $i = 1, \dots, n$, at adjusted values x_i are necessary to solve a set of n equations of the form

$$y_i = f(x_i, \beta) \quad (6.1)$$

with regard to unknown parameters β . The measured variables y_i are assumed to be measured completely precisely, without any experimental errors. The model function $f(x, \beta)$ is assumed to be correct and to correspond to data y . In the chemical laboratory, none of these assumptions is usually fulfilled.

(b) Variable y is subject to random errors, but variables x are controllable. This case is the *classical regression model*, for which the conditional mean of random variable y at a point x is given by

$$E(y/x) = f(x, \beta) \quad (6.2)$$

The method of estimation of parameters β depends on the distribution of random variable y . The *additive model of measurement errors* (Chapter 1) is assumed:

$$y_i = f(x_i, \beta) + \varepsilon_i \quad (6.3)$$

where ε_i is a random variable containing the measurement errors $\varepsilon_{M,i}$, and the model errors $\varepsilon_{T,i}$ coming from an approximate model which does not correspond to the true theoretical model $f_T(x_i, \beta)$. Decomposition of the total error ε_i into components $\varepsilon_{M,i}$

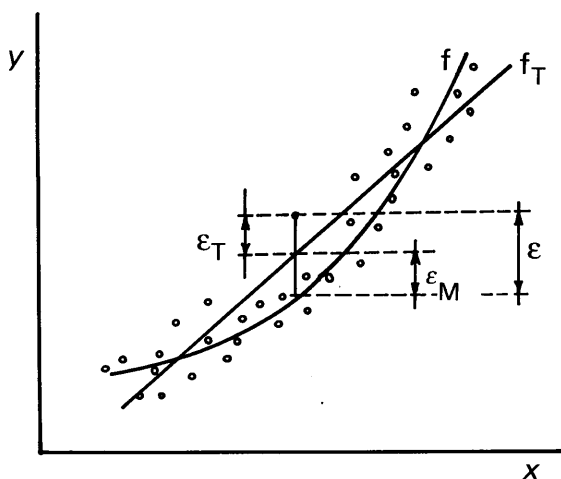


Fig. 6.1—Decomposition of the total error ε into two components, the measure error ε_M and the model error ε_T .

and $\varepsilon_{T,i}$ is illustrated in Fig. 6.1.

Treatment of chemical data by regression analysis involves first the choice of a linear regression model

$$E(y/x) = \sum_{j=1}^m \beta_j x_j \quad (6.4)$$

which either can be an approximation of the unknown theoretical function f_T (Fig. 6.1) or can be derived from a knowledge of the chemical system. In Eq. (6.4), instead of variables x_j , their functions which do not contain parameters β may be used. Parameter estimates of model (6.4) may be determined, on the assumption that Eq. (6.3) is valid, either by the *method of maximum likelihood* or the *method of least-squares*.

(c) Variables y, \mathbf{x} are a sample from the random vector (η, ζ^T) with $m + 1$ components. Regression is conditioned by the mean value (6.2) where \mathbf{x} represents an actual quantity from the random vector ξ . Unlike regression models, in these "correlation models" the regression function can be derived from a simultaneous probability density frequency function $p(y, \mathbf{x})$ and a conditional probability density function $p(y/\mathbf{x})$. The analysis of correlation models is discussed in Chapter 7.

For either correlation or for regression models, the same expressions are valid, although they differ significantly in meaning.

In this chapter, only linear regression models [Eq. (6.4)] are considered, i.e. models which may be written in the form

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \cdot & \cdot & \cdot & x_{1m} \\ x_{21} & x_{22} & \cdot & \cdot & \cdot & x_{2m} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ x_{n1} & x_{n2} & \cdot & \cdot & \cdot & x_{nm} \end{bmatrix} \cdot \begin{bmatrix} \beta_1 \\ \beta_2 \\ \cdot \\ \beta_m \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \cdot \\ \varepsilon_n \end{bmatrix} \quad (6.5a)$$

of dimensions

$$(n \times 1) \qquad (n \times m) \qquad (m \times 1) \quad (n \times 1)$$

In matrix notation, Eq. (6.5a) takes the simple form

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (6.5b)$$

Columns \mathbf{x}_j define geometrically the m -dimensional co-ordinate system or the hyperplane L in n -dimensional Euclidean space E^n . The vector \mathbf{y} does not have to lie in this hyperplane L , as shown in Fig. 6.2, which is for two independent variables ($m = 2$).

The vector $\mathbf{X}\boldsymbol{\beta}$ lies in hyperplane L and parameters $\boldsymbol{\beta}$ may be understood as the coefficients of proportionality of the individual components x_j of the co-ordinate system. The regression model is formed by their linear combination. Whatever regression criterion is used for linear regression models, the model function $\mathbf{X}\boldsymbol{\beta}$ and the theoretical model $\mathbf{X}\boldsymbol{\beta}$ will lie in an m -dimensional hyperplane L .

The least-squares method (LS) is the most frequently used method in regression analysis. The geometry of the least-squares method is very simple. For a linear regression (Fig. 6.2), the parameter estimates \mathbf{b} may be found by minimization of distance between the vector \mathbf{y} and the hyperplane L . This is equivalent to finding the minimal length of the residual vector

$$\hat{\mathbf{e}} = \mathbf{y} - \hat{\mathbf{y}}_p \quad (6.6)$$

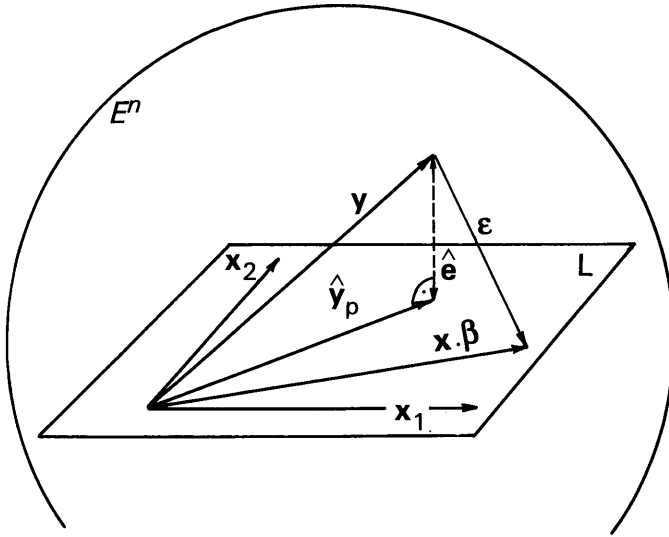


Fig. 6.2—Geometric illustration of a linear regression model for two independent variables.

where $\hat{y}_p = \mathbf{X}\mathbf{b}$ is the *prediction vector*. In Euclidean space the length of vector $\hat{\mathbf{e}}$ can be expressed by

$$D = \sqrt{\langle \hat{\mathbf{e}}, \hat{\mathbf{e}} \rangle} = \sqrt{\sum_{i=1}^n \hat{e}_i^2} \quad (6.7)$$

The symbol $\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^n x_i y_i$ means the scalar product of two vectors. The square of vector $\hat{\mathbf{e}}$ length is consistent with criterion $U(\mathbf{b})$ of the least-squares method $D^2 = U(\mathbf{b})$ so that the estimates of model parameters \mathbf{b} minimize the expression

$$U(\mathbf{b}) = D^2 = \sum_{i=1}^n \left[y_i - \sum_{j=1}^m x_{ij} b_j \right]^2 = \sum_{i=1}^n (y_i - \hat{y}_{p,i})^2 \quad (6.8)$$

Vectors $\hat{\mathbf{e}}$ and $\hat{\mathbf{y}}_p$ are illustrated on Fig. 6.2. The vector $\hat{\mathbf{y}}_p$ represents a perpendicular projection of vector \mathbf{y} onto hyperplane L . The vector $\hat{\mathbf{e}}$ for which a function D is minimal lies in $n - m$ dimensional hyperplane L^\perp that is perpendicular to the hyperplane L and is called the *residual vector*.

The residual vector $\hat{\mathbf{e}}$ is perpendicular to all columns of matrix \mathbf{X} and therefore all corresponding scalar products are zero

$$\langle \mathbf{x}_j, \hat{\mathbf{e}} \rangle = \sum_{i=1}^n x_{ij} \hat{e}_i = 0, \quad j = 1, \dots, m \quad (6.9)$$

This set of equations may be written in matrix notation as

$$\mathbf{X}^T \hat{\mathbf{e}} = 0 \quad (6.10)$$

Substituting $(y - \mathbf{Xb})$ for $\hat{\mathbf{e}}$ leads to a set of linear equations in the known vector \mathbf{b}

$$\mathbf{X}^T \mathbf{y} = \mathbf{X}^T \mathbf{X} \mathbf{b} \quad (6.10a)$$

The estimate \mathbf{b} which minimizes the distance \mathbf{D} is then

$$\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (6.11)$$

where \mathbf{A}^{-1} represents the inverse of matrix \mathbf{A} .

The perpendicular projection of \mathbf{y} into hyperplane L can be made by using projection matrix \mathbf{H} and may be expressed by

$$\hat{\mathbf{y}}_P = \mathbf{H} \mathbf{y} \quad (6.12)$$

By substitution from Eq. (6.11), Eq. (6.12) may be rewritten as

$$\hat{\mathbf{y}}_P = \mathbf{Xb} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (6.12a)$$

The projection matrix $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ has the property of projecting any vector \mathbf{V} into a plane L . When the vector \mathbf{V} already lies in plane L , $\mathbf{H}\mathbf{V} = \mathbf{V}$. However, when vector \mathbf{V} is perpendicular to plane L , $\mathbf{H}\mathbf{V} = \mathbf{0}$ where $\mathbf{0}$ is a zero vector.

The projection matrix \mathbf{P} for perpendicular projection into a hyperplane L^\perp that is orthogonal to hyperplane L is

$$\mathbf{P} = \mathbf{E} - \mathbf{H} \quad (6.13)$$

where \mathbf{E} is an $n \times n$ identity matrix. With the use of these two projection matrices the total decomposition of vector \mathbf{y} into two orthogonal components may be written as

$$\mathbf{y} = \mathbf{H}\mathbf{y} + \mathbf{P}\mathbf{y} = \hat{\mathbf{y}}_P + \hat{\mathbf{e}}$$

The geometric interpretation is that vector \mathbf{y} is decomposed into two mutually perpendicular vectors (Fig. 6.2).

The same expressions can be reached by an analytical minimization, i.e. by a differentiation of Eq. (6.8) and rearrangement.

Problem 6.1. *Parameter estimates of a calibration line*

Apply the expressions already derived to a model of a straight calibration line

$$E(y/x) = \beta_1 x + \beta_2$$

where y is the measured quantity and x is usually a concentration. Derive estimates b_1 , b_2 and the elements of the projection matrix \mathbf{H} .

Solution: For this case, we have

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} x_1 & x_2 & \dots & x_n \\ 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \cdot & \cdot \\ \cdot & \cdot \\ x_n & 1 \end{bmatrix} = \begin{bmatrix} \sum x_i^2 & \sum x_i \\ \sum x_i & n \end{bmatrix}$$

$$\mathbf{X}^T \mathbf{y} = \begin{bmatrix} x_1 & x_2 & \dots & x_n \\ 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \cdot \\ \cdot \\ y_n \end{bmatrix} = \begin{bmatrix} \sum x_i y_i \\ \sum y_i \end{bmatrix}$$

where all sums are considered for $i = 1$ to n . For determination of the inversion matrix $(\mathbf{X}^T \mathbf{X})^{-1}$ the method based on the adjugate matrix may be used

$$(\mathbf{X}^T \mathbf{X})^{-1} = (1/\det(\mathbf{X}^T \mathbf{X}) \cdot \text{adj}(\mathbf{X}^T \mathbf{X}))$$

where

$$\det(\mathbf{X}^T \mathbf{X}) = n \sum_{i=1}^n x_i^2 - \left[\sum_{i=1}^n x_i \right]^2$$

and

$$\text{adj}(\mathbf{X}^T \mathbf{X}) = \begin{bmatrix} n & -\sum x_i \\ -\sum x_i & \sum x_i^2 \end{bmatrix}$$

Recall that for matrix $\mathbf{A} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$, $\det(\mathbf{A}) = a \times d - c \times b$ and $\text{adj}(\mathbf{A}) = \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$. Substitution into Eq. (6.11) leads to the vector of parameter estimates \mathbf{b} , which is given by

$$\mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = \frac{1}{n \sum x_i^2 - \left[\sum x_i \right]^2} \begin{bmatrix} n & -\sum x_i \\ -\sum x_i & \sum x_i^2 \end{bmatrix} \begin{bmatrix} \sum x_i y_i \\ \sum y_i \end{bmatrix} \quad (6.14)$$

Multiplication yields estimates of the two parameters β_1, β_2 in the closed form

$$b_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{D} \quad (6.14a)$$

$$b_2 = \frac{\sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i}{D} \quad (6.14b)$$

where

$$D = n \sum x_i^2 - \left[\sum x_i \right]^2 \quad (6.14c)$$

Equation (6.12a) allows the diagonal elements of projection matrix \mathbf{H} to be determined

$$H_{jj} = \frac{1}{D} [x_j \quad 1] \begin{bmatrix} n & -\sum x_i \\ -\sum x_i & \sum x_i^2 \end{bmatrix} \begin{bmatrix} x_j \\ 1 \end{bmatrix}$$

$$= \frac{\sum x_i^2 + nx_j^2 - 2x_j \sum x_i}{D}, \quad j = 1, \dots, n$$

and for nondiagonal elements H_{jk}

$$H_{jk} = \frac{1}{D} [x_k \quad 1] \begin{bmatrix} n & -\sum x_i \\ -\sum x_i & \sum x_i^2 \end{bmatrix} \begin{bmatrix} x_j \\ 1 \end{bmatrix}$$

$$= \frac{nx_j x_k + \sum x_i^2 - (x_k + x_j) \sum x_i}{D}, \quad j, k = 1, \dots, n$$

Introduction of the arithmetic mean $\bar{x} = 1/n \sum_{i=1}^n x_i$ into the expressions for H_{jj} and H_{jk} gives

$$H_{jj} = \frac{1}{n} + \frac{n(x_j - \bar{x})^2}{D}$$

and

$$H_{jk} = \frac{1}{n} + \frac{n(x_j - \bar{x})(x_k - \bar{x})}{D}$$

Conclusion: With the use of simple matrix operations, estimates of the parameters of a straight line and of the elements of the projection matrix may be calculated.

Problem 6.2. *Geometric interpretation of a calibration line*

Derive expressions for estimates of parameters b_1 and b_2 for a calibration straight line and make a geometric representation. Use a perpendicular projection of vector \mathbf{y} into the plane defined by the columns of matrix \mathbf{X} and also make a geometric representation of the individual projections.

Solution: The model of a calibration straight line is expressed in matrix form by

$$\mathbf{y} = b_1^* \mathbf{x}_C + b_2^* \mathbf{J} + \hat{\mathbf{e}}$$

where \mathbf{x}_C is a centred variable with elements $x_{C,i} = x_i - \bar{x}$ representing the concentration or content of component and \mathbf{J} is an $(n \times 1)$ vector of ones. Parameters b_1^* and b_2^* refer to the model with centred variables. The advantage of the use of a centred variable instead of the original one is that the vectors \mathbf{x}_C and \mathbf{J} are orthogonal and their scalar product is equal to zero,

$$\langle \mathbf{x}_C, \mathbf{J} \rangle = \sum_{i=1}^n x_{C,i} \cdot 1 = \sum_{i=1}^n (x_i - \bar{x}) = 0$$

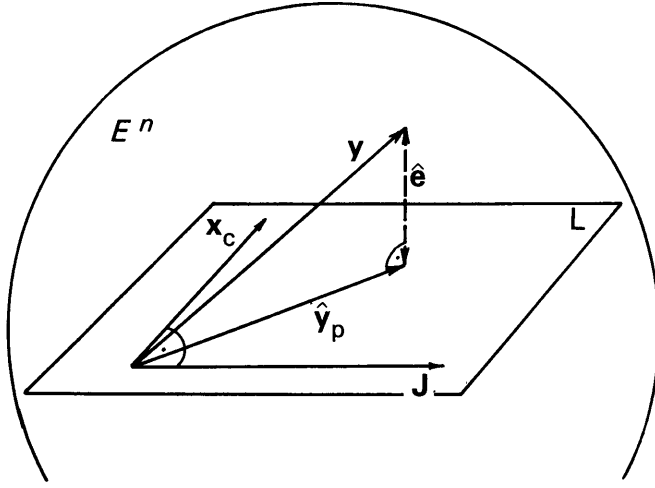


Fig. 6.3—Geometrical representation of a calibration straight line (when y is the only random variable) in Euclidean space E^n .

(see Fig. 6.3)

The perpendicular projection P_L (of vector y into a plane L) is [owing to the orthogonality of the components] equal to the sum of the projection P_J (of vector y on vector J) and the projection P_X (of vector y on a vector x_c),

$$P_L = P_J + P_X = b_1^* x_c + b_2^* J = \hat{y}_p$$

(see Fig. 6.4).

Projection P_X lies on vector x_c and vector $y - P_X$ is on the perpendicular to this vector. Then the corresponding scalar product must be zero.

$$\langle y - P_X, x_c \rangle = \langle y - b_1^* x_c, x_c \rangle = \sum_{i=1}^n (y_i - b_1^* x_{c,i}) x_{c,i} = 0$$

The estimate b_1^* will then be

$$b_1^* = \frac{\sum_{i=1}^n y_i (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \sum_{i=1}^n y_i w_i \quad (6.14d)$$

where w_i are weight coefficients. Because the vector $(y - P_J)$ is perpendicular to the vector J , the following equality is valid

$$\langle y - P_J, J \rangle = \sum_{i=1}^n (y_i - b_2^*) \cdot 1 = 0$$

and therefore the estimate b_2^* is equal to

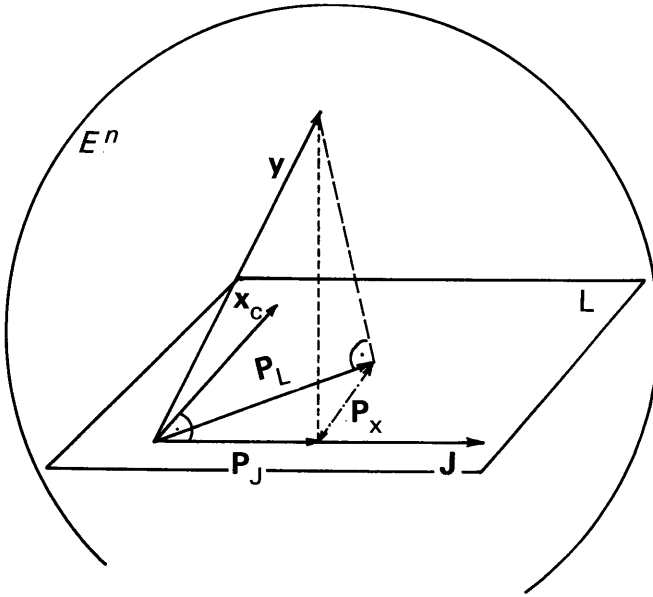


Fig. 6.4—Geometrical representation of the individual projections P_L , P_J and P_X .

$$b_2^* = \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

The estimates b_1 and b_2 for a model with a non-centred variable x , and the estimates b_1^* and b_2^* of a model with centred variable x_C are related in the following ways.

$$b_1 = b_1^*$$

$$b_2 = b_2^* - b_1 \bar{x} = \sum_{i=1}^n \left(\frac{1}{n} - \bar{x} w_i \right) y_i \quad (6.14e)$$

where $w_i = [x_i - \bar{x}] / [\sum_{i=1}^n (x_i - \bar{x})^2]$ are the weight coefficients of the individual values y_i when b_1^* is calculated from Eq. (6.14d). Equations (6.14a) and (6.14d) or (6.14b) and (6.14e) are equivalent. Equations (6.14d, e) show that the parameter estimates are the weighted linear combinations of all y_i . The magnitudes of the weight coefficients depend only on the location of experimental data. This important conclusion tells us that when a value of x_i is far from the mean \bar{x} , the weight w_i is large, so the point (x_i, y_i) has great “weight”, and is more significant in the estimate b_1 .

Conclusion: The estimates of straight line parameters may be found from geometric considerations. These estimates correspond to those found by the least-squares method.

6.2 CONDITIONS FOR THE LEAST-SQUARES METHOD

In determination of the statistical properties of random vectors $\hat{\mathbf{y}}_p$, $\hat{\mathbf{e}}$, and \mathbf{b} , there are some conditions necessary for the least-squares method (LS) to be valid [1].

- (1) The regression parameters $\boldsymbol{\beta}$ can have any value. In chemometric practice, however, there are some restrictions on the parameters, based on physical meaning.
- (2) The regression model is linear in the parameters, and an additive model for the measurement errors is valid [Eq. (6.5b)].
- (3) The matrix of non-random controllable values of the independent variables \mathbf{X} has a column rank equal to m . This means that the two columns \mathbf{x}_j , \mathbf{x}_k are not collinear (i.e. parallel) vectors. This is the same as saying that the matrix $\mathbf{X}^T\mathbf{X}$ is a symmetric regular invertible matrix with non-zero determinant. That is, plane L is m -dimensional, and vector $\mathbf{X} \cdot \mathbf{b}$ and the parameter estimates \mathbf{b} are unambiguously determined.
- (4) The mean value of the random errors ε_i is zero; $E(\varepsilon_i) = 0$. This is valid for all correlation models. It may happen that $E(\varepsilon_i) = K$, $i = 1, \dots, n$, which means that the model does not contain an intercept term. If an intercept term is used in such a model, it will be found that $E(\varepsilon'_i) = 0$ where $\varepsilon'_i = y_i - \hat{y}_{p,i} - K$.
- (5) The random errors ε_i have constant and finite variance, $E(\varepsilon_i^2) = \sigma^2$. The conditional variance $D(y/x) = \sigma^2$ is also constant and therefore the data are said to be **homoscedastic**.
- (6) The random errors ε_i are uncorrelated and therefore $\text{cov}(\varepsilon_i, \varepsilon_j) = E(\varepsilon_i, \varepsilon_j) = 0$. When the errors follow the normal distribution they are also independent. (This corresponds to independence of the measured variables y .)
- (7) The random errors ε_i have a normal distribution $N(0, \sigma^2)$. The vector \mathbf{y} then has a multivariate normal distribution with mean $\mathbf{X}\boldsymbol{\beta}$ and covariance matrix $\sigma^2\mathbf{E}$ where \mathbf{E} is the identity matrix.

When first six conditions are met, the parameter estimates \mathbf{b} found by minimization of a least-squares criterion are **best unbiased linear estimates** of the regression parameters $\boldsymbol{\beta}$ [2].

The term **best** estimates (\mathbf{b}) means that any linear combination of these estimates has the *smallest* variance of all linear unbiased estimates. That is, the variances of the individual estimates $D(b_j)$ are the smallest from all possible linear unbiased estimates (Gauss–Markov theorem).

It should be noted that there exist *biased estimates*, the variance of which can be smaller than the variance of estimates $D(b_j)$.

The term **unbiased** estimates means that $E(\boldsymbol{\beta} - \mathbf{b}) = 0$ and the mean value of an estimate vector $E(\mathbf{b})$ is equal to a vector of regression parameters $\boldsymbol{\beta}$.

The term **linear** estimates means that they can be written as a linear combination of measurements \mathbf{y} with weights Q_{ij} which depend only on the locations of variables \mathbf{x}_j , $j = 1, \dots, m$. If we write Eq. (6.11) $\mathbf{Q} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$ for the weight matrix, we can then say

$$b_j = \sum_{i=1}^n Q_{ij} y_i \quad (6.15)$$

Each estimate b_j is the weighted sum of all measurements. Also, the estimates \mathbf{b} have an asymptotic multivariate normal distribution with covariance matrix [2, 4].

$$D(\mathbf{b}) = \sigma^2(\mathbf{X}^T\mathbf{X})^{-1} \quad (6.16)$$

When condition (7) is valid, all estimates \mathbf{b} have a normal distribution, even for finite sample sizes n .

Problem 6.3. *Variance of parameter estimates for calibration line*

Derive the equations for calculation of the variance of estimates of parameters $D(b_1)$, $D(b_2)$ and the covariance $\text{cov}(b_1, b_2)$ for the calibration straight lines from Problem 6.1.

Solution: From Eq. (6.16) and the answer to Problem 6.1, we can write

$$D(b_1) = \frac{n\sigma^2}{D}$$

$$D(b_2) = \frac{\sigma^2 \sum_{i=1}^n x_i^2}{D}$$

and

$$\text{cov}(b_1, b_2) = \frac{-\sigma^2 \sum_{i=1}^n x_i}{D}$$

The correlation coefficient R_0 , expressing the correlation between estimates b_1 and b_2 is calculated from

$$R_0 = \frac{\text{cov}(b_1, b_2)}{\sqrt{D(b_1) \cdot D(b_2)}} = \frac{-\sum_{i=1}^n x_i}{\sqrt{n \sum_{i=1}^n x_i^2}}$$

For small values of n and positive values of x , the coefficient R_0 can be near to -1 . This means that the estimate of the slope b_1 and the estimate of the intercept b_2 in model $y = \beta_1 x + \beta_2$ are negatively correlated, and the corresponding correlation coefficient can reach very high value.

For the variances $D(b_1^*)$ and $D(b_2^*)$, the equations for estimates b_1^* and b_2^* may be used,

$$D(b_1^*) = \sum_{i=1}^n w_i^2 D(y_i) = \sigma^2 \sum_{i=1}^n w_i^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$D(b_2^*) = \frac{\sum_{i=1}^n D(y_i)}{n^2} = \frac{\sigma^2}{n}$$

Based on Eq. (6.14e), the expression may be written as

$$\begin{aligned} D(b_2) &= \sum_{i=1}^n \left(\frac{1}{n} - \bar{x}w_i \right)^2 D(y_i) \\ &= \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \end{aligned}$$

Conclusion: For a calibration straight line, the variance of the estimates of intercept and slope, and the correlation coefficient between the parameters may be calculated from the simple expressions derived. Estimates b_1 and b_2 for positive data $x_i > 0$, $i = 1, \dots, n$, are always negatively correlated.

6.3. STATISTICAL PROPERTIES OF THE LEAST-SQUARES METHOD

When conditions (1) to (7) for the least-squares method are met, some statistical properties of vectors \mathbf{b} , $\hat{\mathbf{y}}_p$ and $\hat{\mathbf{e}}$ may be utilized. As the projection matrix \mathbf{H} is non-random, for a **covariance matrix of prediction**, the following expression is valid

$$D(\hat{\mathbf{y}}_p) = \sigma^2 \mathbf{H} \quad (6.17)$$

and for a **covariance matrix of residuals** the expression

$$D(\hat{\mathbf{e}}) = \sigma^2 \mathbf{P} = \sigma^2 (\mathbf{E} - \mathbf{H}) \quad (6.18)$$

Both (6.17) and (6.18) are based on important properties of projection matrices, i.e. idempotence $\mathbf{H} = \mathbf{H}\mathbf{H}$ and symmetry $\mathbf{H} = \mathbf{H}^T$.

Variances of the parameter estimates \mathbf{b} are derived from Eq. (6.11) and given by Eq. (6.16). The residual sum of squares RSS, denoted also by $U(\mathbf{b})$, may be written as:

$$\text{RSS} = U(\mathbf{b}) = \hat{\mathbf{e}}^T \hat{\mathbf{e}} = \mathbf{y}^T (\mathbf{E} - \mathbf{H}) \mathbf{y} = \mathbf{y}^T (\mathbf{E} - \mathbf{H}) \mathbf{y} = \mathbf{e}^T \mathbf{P} \mathbf{e}$$

and the mean residual sum of squares sum is expressed as

$$E(\text{RSS}) = \sigma^2 \text{tr}(\mathbf{P}) = \sigma^2 (n - m) \quad (6.19)$$

where $\text{tr}(\mathbf{P})$ is a trace matrix \mathbf{P} . With reference to the idempotence and symmetry of the projection matrix \mathbf{P} , the trace of matrix \mathbf{P} is equal to its rank.

An unbiased estimate of the variance of errors σ^2 can be calculated with the use of the variance of the residuals

$$\hat{\sigma}^2 = \frac{U(\mathbf{b})}{n - m} = \frac{\hat{\mathbf{e}}^T \hat{\mathbf{e}}}{n - m} \quad (6.20)$$

Problem 6.4. *Estimation of the variances of prediction and residuals for a calibration line*

Derive expressions for calculation of an estimate of the prediction variance $D(\hat{y}_{p,i})$ and of an estimate of the residual variance $D(\hat{e}_i)$ for the calibration straight line from

Problem 6.1.

Solution: From Problem 6.1 and Eq. (6.17) we know that for an estimate of the prediction variance:

$$D(\hat{y}_{P,i}) = \sigma^2 \left[\frac{1}{n} + \frac{n(x_i - \bar{x})^2}{D} \right]$$

and from Eq. (6.18) for an estimate of the residual variance

$$D(\hat{e}_i) = \sigma^2 \left[\frac{n-1}{n} - \frac{n(x_i - \bar{x})^2}{D} \right]$$

Thus, the prediction variance and the residual variance are quadratic functions of the distance from \bar{x} . At the point $x_i = \bar{x}$, the prediction variance has a minimum and

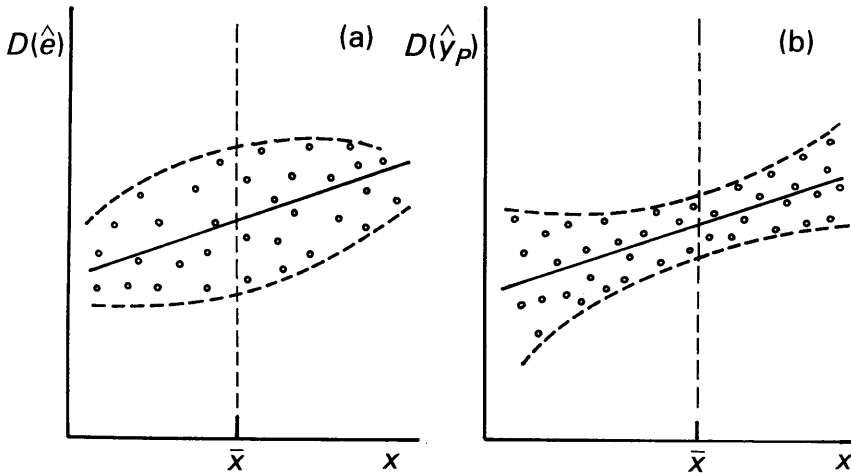


Fig. 6.5—Dependence of (a) the residual variance $D(\hat{e}_i)$, and (b) the prediction variance $D(\hat{y}_{P,i})$ on the independent variable x .

the residual variance a maximum (Fig. 6.5a, b).

Conclusion: At the point \bar{x} the prediction variance is minimal and the residual variance maximal. When a point x_i is far away from \bar{x} , the prediction is less precise but an estimate of the residual is more precise.

From Fig. 6.1 it is evident that a square with the length of vector \mathbf{y} is equal to the sum of squares of the lengths of vectors $\hat{\mathbf{y}}_P$ and $\hat{\mathbf{e}}$,

$$\mathbf{y}^T \mathbf{y} = \hat{\mathbf{y}}_P^T \hat{\mathbf{y}}_P + \hat{\mathbf{e}}^T \hat{\mathbf{e}} = \mathbf{y}^T \mathbf{H} \mathbf{y} + \mathbf{y}^T (\mathbf{E} - \mathbf{H}) \mathbf{y} \quad (6.21)$$

or in abbreviated notation

$$TSS = SS + RSS \quad (6.21a)$$

Equations (6.21) and (6.21a) may be understood to mean that the total sum of squares TSS may be decomposed into two components, the sum of squares SS caused by the

regression model and the unelucidated residual sum of squares RSS . The mean value of the sum-of-squares of a regression model is given by

$$E(SS) = E(\mathbf{y}^T \mathbf{H} \mathbf{y}) = \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} + m\sigma^2 \quad (6.22)$$

Instead of the quantities RSS and SS , their average values are often used. The mean regression sum of squares is defined by

$$MSS = SS/m \quad (6.23a)$$

and its expected value by

$$E(MSS) = \sigma^2 + \frac{1}{m} \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} \quad (6.23b)$$

The mean residual sum of squares is defined by

$$MRS = RSS/(n - m) \quad (6.24a)$$

and its expected value by

$$E(MRS) = \sigma^2 \quad (6.24b)$$

If $\boldsymbol{\beta} = 0$, i.e. all parameters of regression model are zero, then

- (a) SS is independent of RSS
- (b) $SS = \sigma^2 \chi_m^2$ where χ_m^2 is a random variable with the χ^2 -distribution with m degrees of freedom,
- (c) $RSS = \sigma^2 \chi_{n-m}^2$ where χ_{n-m}^2 is a random variable with the χ^2 -distribution with $(n - m)$ degrees of freedom.

On the basis of these three facts, the ratio

$$F = MSS/MRS \quad (6.25)$$

has the Fisher–Snedecor F -distribution with m and $(n - m)$ degrees of freedom.

Problem 6.5. *Decomposition of the total sum of squares for the model of a calibration line*

For the calibration straight line defined in Problem 6.1 try to decompose the total sum of squares [Eq. (6.21)].

Solution: The full expression for $RSS = \sum_{i=1}^n \hat{e}_i^2$ is

$$RSS = \sum_{i=1}^n \hat{e}_i (y_i - b_2 - b_1 x_i) = \sum_{i=1}^n \hat{e}_i y_i - b_2 \sum_{i=1}^n \hat{e}_i - b_1 \sum_{i=1}^n \hat{e}_i x_i$$

For models with an intercept, $\sum_{i=1}^n \hat{e}_i = 0$ always. The vector $\hat{\mathbf{e}}$ is perpendicular to vector \mathbf{x} so that $\sum_{i=1}^n \hat{e}_i x_i = 0$. Therefore, the last two terms in the sum are equal to zero and

$$RSS = \sum_{i=1}^n \hat{e}_i y_i = \sum_{i=1}^n [y_i - b_2 - b_1 x_i] y_i$$

$$= \sum_{i=1}^n y_i^2 - \left[b_2 \sum_{i=1}^n y_i + b_1 \sum_{i=1}^n x_i y_i \right] = TSS - SS$$

and the regression sum-of-squares is expressed by

$$SS = b_2 \sum_{i=1}^n y_i + b_1 \sum_{i=1}^n x_i y_i$$

Expressing SS directly from the definition we get

$$SS = \sum_{i=1}^n \hat{y}_{P,i}^2 = \sum_{i=1}^n (b_2 + b_1 x_i)^2 = nb_2^2 + 2b_1 b_2 \sum_{i=1}^n x_i + b_1^2 \sum_{i=1}^n x_i^2$$

The resulting equation can be used to express the expected value of SS . From the elemental properties of an expected value, we have

$$E(b_2^2) = \beta_2^2 + D(b_2)$$

$$E(b_1^2) = \beta_1^2 + D(b_1)$$

and

$$E(b_1 b_2) = \beta_1 \beta_2 + \text{cov}(b_1, b_2) = \beta_1 \beta_2 - \frac{\sigma^2 \sum_{i=1}^n x_i}{D}$$

From these three expressions the expected value of the regression model sum-of-squares is given by

$$E(SS) = 2\sigma^2 + \left[n\beta_2^2 + 2\beta_1 \beta_2 \sum_{i=1}^n x_i + \beta_1^2 \sum_{i=1}^n x_i^2 \right]$$

This can also be derived by straight substitution into Eq. (6.22).

Conclusion: From estimates b_1 and b_2 it is possible to calculate not only RSS but also SS .

When the model contains an intercept term we will use the linear combination of vectors

$$E(y/\mathbf{x}) = \beta_1 \mathbf{x}_1 + \dots + \beta_{m-1} \mathbf{x}_{m-1} + \beta_m \mathbf{J} \quad (6.26)$$

where $\mathbf{J} = (1, 1, \dots, 1)^T$ is the vector of all ones. On introducing centred variables

$$\mathbf{x}_{C,j} = \mathbf{x}_j - \mathbf{J}\bar{x}_j, \quad j = 1, \dots, m-1 \quad (6.27)$$

the scalar products become $\langle \mathbf{x}_{C,j}, \mathbf{J} \rangle = 0$. This means also that vectors $\mathbf{x}_{C,j}$ and \mathbf{J} are orthogonal. In Eq. (6.27) the symbol $\bar{x}_j = 1/n \sum_{i=1}^n x_{ij}$ means the arithmetic mean of the j th controllable independent variable. By using centred variables the regression model will be expressed in matrix form

$$\mathbf{y} = \mathbf{X}_C \boldsymbol{\beta}^* + \beta_m^* \mathbf{J} + \boldsymbol{\varepsilon} \quad (6.28)$$

where \mathbf{X}_C is a matrix of dimension $(n \times (m - 1))$ and $\boldsymbol{\beta}^*$ is vector of dimension $(m - 1) \times 1$. Because of the orthogonality of the two variables in Eq. (6.28), the estimates \mathbf{b}^* and b_m^* of parameters $\boldsymbol{\beta}^*$ and β_m^* may be determined independently, from a projection into a plane L_1 as defined by the columns of matrix \mathbf{X}_C or from projection on a vector \mathbf{J} (see Problem 6.2).

By projection into a plane L_1 we find

$$\mathbf{X}_C \mathbf{b}^* = \mathbf{X}_C (\mathbf{X}_C^T \mathbf{X}_C)^{-1} \mathbf{X}_C^T \mathbf{y} \quad (6.29a)$$

or

$$\mathbf{b}^* = (\mathbf{X}_C^T \mathbf{X}_C)^{-1} \mathbf{X}_C^T \mathbf{y} \quad (6.29b)$$

By using projection onto a vector \mathbf{J} we find

$$\mathbf{J} b_m^* = \mathbf{J} (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \mathbf{y} \quad (6.30a)$$

or

$$b_m^* = (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \mathbf{y} = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y} \quad (6.30b)$$

With the use of Eq. (6.30), the expression for an estimate of intercept term variance may be derived:

$$D(b_m) = \sigma^2 (\mathbf{J}^T \mathbf{J})^{-1} = \frac{\sigma^2}{n} \quad (6.31)$$

By introducing a centred dependent variable

$$\mathbf{y}_C = \mathbf{y} - \mathbf{J} \bar{y} \quad (6.32)$$

the regression model (6.28) is transformed into a model without an intercept term

$$\mathbf{y}_C = \mathbf{X}_C \boldsymbol{\beta}_C + \boldsymbol{\varepsilon} \quad (6.33)$$

For this model, the total sum of squared deviations from the average TSC may be decomposed into the sum of squared deviations from the regression model and the sum of squared residuals RSC (equal to RSS). Then

$$\mathbf{y}_C^T \mathbf{y}_C = \hat{\mathbf{y}}_{PC}^T \hat{\mathbf{y}}_{PC} + \hat{\mathbf{e}}^T \hat{\mathbf{e}} \quad (6.34)$$

or

$$TSC = SSC + RSC \quad (6.35)$$

Decomposition of the total variations from an average $(y_i - \bar{y})$ into a part explained

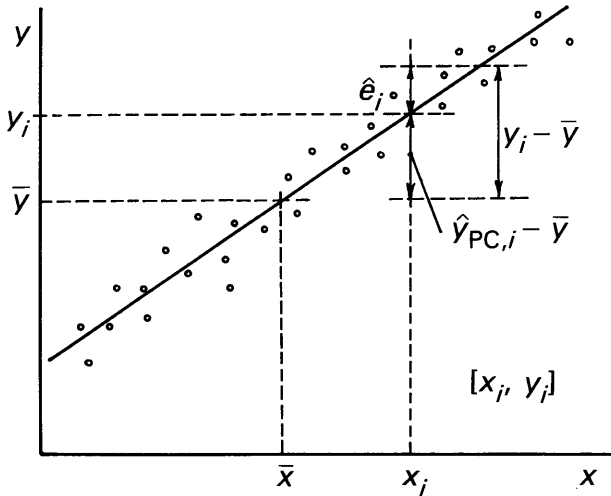


Fig. 6.6—Decomposition of the total deviation from the mean ($y_i - \bar{y}$) into an explained part ($\hat{y}_{PC,i} - \bar{y}$) and a part not explained by the regression model, \hat{e}_i .

($\hat{y}_{PC,i} - \bar{y}$) and a part not explained \hat{e}_i by a regression model is illustrated in Fig. 6.6 for a regression straight line.

We know that

$$TSC = TSS - n\bar{y}^2 \quad (6.36)$$

and also

$$SSC = \mathbf{y}_C^T \mathbf{H}_C \mathbf{y}_C = \mathbf{b}^T \mathbf{X}^T \mathbf{X} \mathbf{b} - n\bar{y}^2 \quad (6.37)$$

where \mathbf{b} is a vector of dimension $(m \times 1)$ containing an intercept term. The cosine of the angle between vectors \mathbf{y}_C and $\hat{\mathbf{y}}_{PC}$ can be calculated by trigonometry (Fig. 6.7):

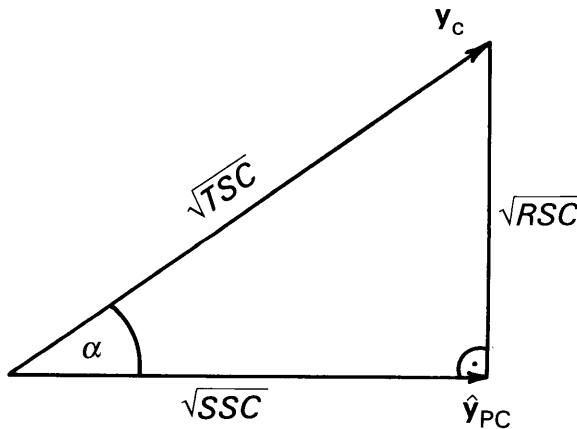


Fig. 6.7—Geometry of the determination of $\cos \alpha$.

$$\cos \alpha = \sqrt{\frac{SSC}{TSC}} = \sqrt{1 - \frac{RSC}{TSC}} = \sqrt{1 - \frac{\sum_{i=1}^n \hat{e}_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (6.38)$$

The quantity $\cos \alpha$ is numerically equal to the value of the multiple correlation coefficient expressed for correlation models.

The square of the correlation coefficient R is called the determination coefficient R^2 . For regression models a quantity $R = \cos \alpha$ is interpreted as the measure of relative difference between regression model $M_1: \hat{y}_p = \mathbf{X}\mathbf{b}$ and model $M_0: \hat{y}_p = \mathbf{J}\bar{y}$. If R tends to zero, the regression model has all parameters except the intercept term equal to zero and model M_0 is valid. This means that in practical chemometric calculations the quantity R cannot be used as a measure of linearity, even for investigation of the quality of a regression model.

From Eq. (6.25) we can calculate the ratio

$$F_R = \frac{(TSC - RSC)(n - m)}{RSC(m - 1)} = \frac{\hat{R}^2(n - m)}{(1 - \hat{R}^2)(m - 1)} \quad (6.39)$$

which has the Fisher–Snedecor F -distribution with $(m - 1)$ and $(n - m)$ degrees of freedom. In Eq. (6.39) the quantity \hat{R}^2 means an estimate of the determination coefficient calculated from $R = \cos \alpha$ and the use of estimates \mathbf{b} . With the use of F_R (6.39) the null hypothesis $H_0: \beta_C = 0$ may be tested; this is equivalent to the hypothesis $H_0: R^2 = 0$. A test of significance of multiple correlation coefficient is the same as a test of significance of all regression coefficients except the intercept term.

Problem 6.6. *Investigation of abrasion resistance and the composition of rubber*

The dependence of the abrasion resistance of rubber, y , on the content of silica filler x_1 and a binding substance x_2 was studied. Whereas the filler increases abrasion resistance, the binding substance also increased its resistance efficiency. Estimate parameters β_1 , β_2 and β_3 of this proposed linear model

$$E(y/x) = \beta_1 x_1 + \beta_2 x_2 + \beta_3$$

and test the statistical significance of the correlation coefficient.

Data: $n = 11$, $m = 3$

y	83	113	92	82	100	96	98	95	80	100	92
x_1^*	1	1	-1	-1	0	0	0	0	0	1.5	-1.5
x_2^*	-1	1	1	-1	0	0	0	1.5	-1.5	0	0

Variables x_1^* and x_2^* are transformed from the raw variables x_1 and x_2 as follows:

$$x_1 = 6.7x_1^* + 50$$

$$x_2 = 2x_2^* + 4.$$

Solution: Since the data are the result of a planned experiment, the vectors \mathbf{J} , \mathbf{x}_1^* and \mathbf{x}_2^* are orthogonal. The matrix $\mathbf{X}^T\mathbf{X}$ is then diagonal, with the form

$$\mathbf{X}^T\mathbf{X} = \begin{bmatrix} \sum x_{i1}^{*2} & 0 & 0 \\ 0 & \sum x_{i2}^{*2} & 0 \\ 0 & 0 & n \end{bmatrix} = \begin{bmatrix} 8.5 & 0 & 0 \\ 0 & 8.5 & 0 \\ 0 & 0 & 11 \end{bmatrix}$$

The vector $\mathbf{X}^T\mathbf{y}$ has the components

$$\mathbf{X}^T\mathbf{y} = \begin{bmatrix} \sum y_i x_{i1}^* \\ \sum y_i x_{i2}^* \\ \sum y_i \end{bmatrix} = \begin{bmatrix} 34 \\ 62.5 \\ 1031 \end{bmatrix}$$

The following estimates are calculated from Eq. (6.11).

$$b_1^* = \frac{\sum_{i=1}^n y_i x_{i1}^*}{\sum_{i=1}^n x_{i1}^{*2}} = 4$$

$$b_2^* = \frac{\sum_{i=1}^n y_i x_{i2}^*}{\sum_{i=1}^n x_{i2}^{*2}} = 7.3529$$

$$b_3^* = \frac{\sum_{i=1}^n y_i}{n} = 93.7273$$

where the stars denote estimates in the transformed variables model. The regression model in the raw variables has the form

$$\begin{aligned} \hat{y}_P &= \left[\frac{x_1 - 50}{6.7} \right] 4 + \left[\frac{x_2 - 4}{2} \right] 7.3529 + 93.7273 \\ &= 0.597x_1 + 3.6765x_2 + 49.1708 \end{aligned}$$

The corresponding residuals sum of squares RSC is given by

$$RSC = \sum_{i=1}^n (y_i - \hat{y}_{P,i})^2 = 326.6$$

and the estimate of residual standard deviation ($n = 11$, $m = 3$)

$$\hat{\sigma} = \sqrt{\frac{RSC}{n-m}} = 6.39$$

Because of the diagonality of the matrix $\mathbf{X}^T\mathbf{X}$ we can say that

$$D(\mathbf{b}^*) = \hat{\sigma}^2 \begin{bmatrix} 1/8.5 & 0 & 0 \\ 0 & 1/8.5 & 0 \\ 0 & 0 & 1/11 \end{bmatrix} = \begin{bmatrix} 4.804 & 0 & 0 \\ 0 & 4.804 & 0 \\ 0 & 0 & 3.712 \end{bmatrix}$$

For estimates corresponding to the raw variables, then

$$D(b_1) = \frac{D(b_1^*)}{6.7^2} = 0.107$$

$$D(b_2) = \frac{D(b_2^*)}{2^2} = 1.201$$

$$D(b_3) = \left[\frac{50}{6.7}\right]^2 D(b_1^*) + \left[\frac{4}{2}\right]^2 D(b_2^*) + D(b_3^*) = 290.47$$

The total sum of deviations from the mean TSC is

$$TSC = \sum_{i=1}^n y_i^2 - n\bar{y} = 922.1818$$

Introducing this into Eq. (6.38) leads to an estimate of the coefficient of determination:

$$\hat{R}^2 = 1 - \frac{326.6}{922.1818} = 0.6458$$

From Eq. (6.39), the test criterion F_R is

$$F_R = \frac{(922.1818 - 326.6) \times 8}{326.6 \times 8} = 7.2943$$

The value F_R is higher than the corresponding quantile of the Fisher-Snedecor distribution $F_{0.95}(2, 8) = 4.46$, so at the significance level $\alpha = 0.05$, the coefficient of determination is considered to be significantly different from zero.

Conclusion: By using planned experimental data, all the columns of matrix \mathbf{X} are mutually orthogonal, and with use of a suitable transformation of variables, the statistical characteristics of the linear model may be calculated.

6.3.1 Construction of confidence intervals

When parameter estimates \mathbf{b} are determined, it is necessary to remember that \mathbf{b} represents the point estimates of parameters $\boldsymbol{\beta}$. These estimates are random quantities and in practice they are less important than the **confidence intervals** in which the true (theoretical) value of parameter $\boldsymbol{\beta}$ lies with some selected probability $(1 - \alpha)$. As for univariate data samples, the significance level is usually chosen with $\alpha = 0.05$ or 0.01 .

These levels correspond to the 95% or the 99% confidence intervals.

Confidence intervals are constructed on the assumption that a random quantity $(n - m)\hat{\sigma}^2/\sigma^2$ has the χ^2 -distribution with $(n - m)$ degrees of freedom, and a random quantity $(\mathbf{b} - \boldsymbol{\beta})^T \mathbf{X}^T \mathbf{X} (\mathbf{b} - \boldsymbol{\beta})/\sigma^2$ has a χ^2 -distribution with m degrees of freedom. The corrected ratio of these quantities has a Fisher distribution with m and $n - m$ degrees of freedom. The bounds of the $100(1 - \alpha)\%$ confidence region are described by

$$(\mathbf{b} - \boldsymbol{\beta})^T \mathbf{X}^T \mathbf{X} (\mathbf{b} - \boldsymbol{\beta}) = m\hat{\sigma}^2 F_{1-\alpha}(m, n - m) \quad (6.40)$$

where $F_{1-\alpha}(m, n - m)$ is the $(1 - \alpha)$ quantile of the Fisher-Snedecor F -distribution with m and $(n - m)$ degrees of freedom. Because the matrix $\mathbf{X}^T \mathbf{X}$ is regular, Eq. (6.40) defines a hyperellipsoid with axes oriented in the directions of the eigenvectors \mathbf{V}_j of the matrix $(\mathbf{X}^T \mathbf{X})^{-1}$. The lengths of the individual half-axes are equal to $p\sqrt{\lambda_j}$ where λ_j are eigenvalues of the matrix $(\mathbf{X}^T \mathbf{X})^{-1}$ and coefficient p is defined by

$$p^2 = m\hat{\sigma}^2 F_{1-\alpha}(m, n - m) \quad (6.41)$$

Neglecting any correlation between parameter estimates, from Eq. (6.40) the $100(1 - \alpha)\%$ simple confidence interval for parameter β_j has the form

$$b_j - t_{1-\alpha/2}(n - m)\hat{\sigma}\sqrt{c_{jj}} \leq \beta_j \leq b_j + t_{1-\alpha/2}(n - m)\hat{\sigma}\sqrt{c_{jj}} \quad (6.42)$$

where c_{jj} is the j th diagonal element of the matrix $(\mathbf{X}^T \mathbf{X})^{-1}$ and $t_{1-\alpha/2}(n - m)$ is the $(1 - \alpha/2)$ quantile of the Student distribution with $(n - m)$ degrees of freedom. Simple confidence intervals are, however, too narrow for correlated estimates \mathbf{b} . Therefore we will define the *extreme confidence intervals* as the extremes on the confidence ellipsoid given by

$$b_j - p\sqrt{c_{jj}} \leq \beta_j \leq b_j + p\sqrt{c_{jj}} \quad (6.43)$$

In some cases the confidence ellipsoid is created for q regression parameters on the assumption that they are the last q components of vector $\boldsymbol{\beta}$. Then, for the $100(1 - \alpha)\%$ boundary confidence ellipsoid

$$(\mathbf{b}_2 - \boldsymbol{\beta}_2)^T \mathbf{D}_2^{-1} (\mathbf{b}_2 - \boldsymbol{\beta}_2) = q\hat{\sigma}^2 F_{1-\alpha}(q, n - m) \quad (6.44)$$

where the matrix \mathbf{D}_2 of dimension $(q \times q)$ is formed from the matrix $(\mathbf{X}^T \mathbf{X})^{-1}$ by leaving out the first $(m - q)$ columns and $(m - q)$ rows. The symbol $\boldsymbol{\beta}_2$ denotes the vector of the last q components of vector $\boldsymbol{\beta}$, and vector \mathbf{b}_2 is defined analogously.

Similarly, the confidence interval for a prediction $\hat{y}_{P,i}$ for point $\mathbf{x}_0 = (x_{01}, \dots, x_{0m})^T$ can be calculated. For $100(1 - \alpha)\%$ confidence interval of prediction we can write [4]

$$\mathbf{x}_0^T \mathbf{b} - t_{1-\alpha/2}(n - m)\hat{\sigma}_{P,0} \leq \mathbf{x}_0^T \boldsymbol{\beta} \leq \mathbf{x}_0^T \mathbf{b} + t_{1-\alpha/2}(n - m)\hat{\sigma}_{P,0} \quad (6.45)$$

where $\hat{\sigma}_{P,0}^2$ is the variance of prediction for which [Eq. (6.17)] the following expression may be used

$$D(\hat{y}_{P,0}) \approx \hat{\sigma}_{P,0}^2 = \hat{\sigma}^2 \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0 \quad (6.46)$$

The relationships between the limits of the confidence interval of prediction (6.46) and \mathbf{x}_0 form the $100(1 - \alpha)\%$ confidence band. The band is narrowest at the centre of gravity of the controllable (independent) variables, $x_{0j} = \bar{x}_j$.

When the confidence bands for all possible values of vector $\mathbf{x} = (x_1, \dots, x_m)^T$ are to be calculated, the Scheffe method should be used. With probability $(1 - \alpha)$ the theoretical value $\mathbf{x}^T \boldsymbol{\beta}$ lies in an interval

$$\mathbf{x}^T \mathbf{b} \pm \sqrt{m F_{1-\alpha}(m, n-m) \hat{\sigma}^2 \mathbf{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}} \quad (6.45a)$$

Confidence bands constructed from Eq. (6.45a) are called Working–Hottelling bands [4] and have the same properties as those constructed on the basis of Eq. (6.45).

Problem 6.7. Validation of a new analytical method

Validate a new analytical method by comparison of results (y) with results obtained by a classical standard method (x) for a set of parallel determinations. If both methods lead to same results, the dependence of y on x is linear [$y = \beta_1 x + \beta_2$] with zero intercept $\beta_2 = 0$ and unit slope $\beta_1 = 1$. Estimate the parameters b_1 and b_2 and construct the 95% confidence interval of intercept and slope, and the 95% confidence interval of prediction for a sample with $x_0 = \bar{x}$.

Data: amount of reagent in mg determined by new (y) and standard (x) methods, $n = 24$, $m = 2$

x	40.2	43.8	47.6	50.7	56.8	81.3	83.3	97.1	102.5	118.7
y	48.9	39.1	42.6	56.9	70.3	71.5	97.6	99.9	105.2	102.3

129.4	184.8	287.5	295.4	420.3	421.3	427.9	566.1	608.5	640.7
106.8	162.9	234.0	303.4	388.8	391.1	369.3	611.6	580.2	643.3

692.8	705.2	714.4	881.4
596.6	612.6	633.5	669.8

Solution: The least-squares estimates of slope is $b_1 = 0.868(\pm 0.030)$ and intercept $b_2 = 14.73(\pm 12.61)$ (\pm standard deviations) were computed. The determination coefficient $\hat{R}^2 = 0.974$ and the estimate of standard deviation of residuals $\hat{\sigma} = 39.54$. From Eq. (6.42)

$$b_2 - t_{1-\alpha/2}(22) \sqrt{D(b_2)} \leq \beta_2 \leq b_2 + t_{1-\alpha/2}(22) \sqrt{D(b_2)}$$

whence

$$14.73 - 2.08 \times 12.61 \leq \beta_2 \leq 14.73 + 2.08 \times 12.61$$

and

$$-11.499 \leq \beta_2 \leq 40.959$$

Since this confidence interval includes zero, the intercept β_2 is not significantly different from zero. The confidence interval for the slope is

$$0.868 - 2.08 \times 0.0302 \leq \beta_1 \leq 0.868 + 2.08 \times 0.0302$$

or

$$0.805 \leq \beta_1 \leq 0.930$$

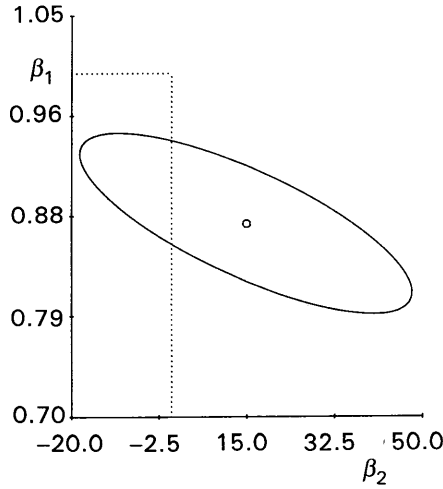


Fig. 6.8—Construction of the 95% confidence ellipse and the point $\beta_1 = 1$ and $\beta_2 = 0$.

Because this interval does not include 1.0, the slope can not be considered to be equal to one. Figure 6.8 demonstrates the 95% confidence ellipse for parameters β_1 and β_2 .

For a regression straight line and $x_0 = \bar{x} = 320.7$, according to Eq. (6.46) we may write

$$D(\hat{y}_P) = \hat{s}_{P,0}^2 = \hat{\sigma}^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{D} \right] = \frac{\hat{\sigma}^2}{n}$$

Then introducing numbers into Eq. (6.45) leads to

$$\begin{aligned} 14.73 + 0.868 \times 320.7 - \frac{2.08 \times 39.54}{\sqrt{24}} &\leq \mathbf{x}_0^T \boldsymbol{\beta} \leq 14.73 \\ + 0.868 \times 320.7 + \frac{2.08 \times 39.54}{\sqrt{24}} & \end{aligned}$$

whence

$$276.89 \leq \mathbf{x}_0^T \boldsymbol{\beta} \leq 309.89$$

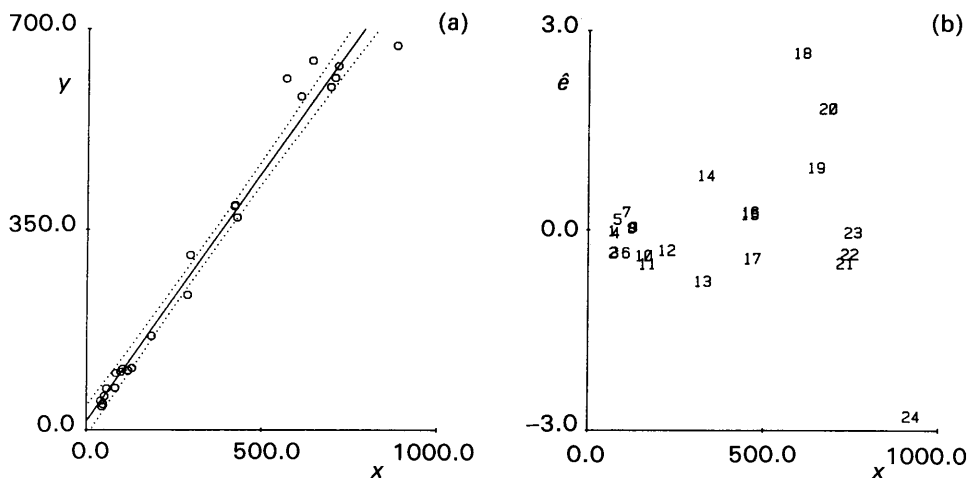


Fig. 6.9—(a) Construction of the 95% confidence bands and (b) dependence of residuals $\hat{e} = f(x)$ on variable x .

Despite the small variance of prediction in the point \bar{x} , the 95% confidence interval is rather broad. Figure 6.9 shows construction of the 95% confidence bands of the calculated regression straight line together with experimental points.

Conclusion: The confidence intervals of the intercept and the slope indicate that the intercept of regression straight line can be considered to be equal to zero, but the slope significantly differs from unity. Thus the results of the new analytical method differ from those obtained by the standard method by a multiplicative constant.

6.3.2 Testing of hypotheses

Tests for significance of parameters are closely connected with the construction of confidence intervals. To test the null hypothesis $H_0: \beta = \beta_0$ where β_0 is the vector of known numbers, against the alternative $H_A: \beta \neq \beta_0$, the test criterion based on Eq. (6.40) may be expressed as

$$F = \frac{(\mathbf{b} - \beta_0)^T \mathbf{X}^T \mathbf{X} (\mathbf{b} - \beta_0)}{m \hat{\sigma}^2} \quad (6.47)$$

which has approximately the Fisher–Snedecor F -distribution with m and $(n - m)$ degrees of freedom. If H_0 is valid. The criteria defined by Eqs. (6.25) and (6.39) are special cases of the test statistic F as defined by Eq. (6.47).

For a test of the simple hypothesis $H_0: \beta_j = \beta_{j,0}$ against the alternative $H_A: \beta_j \neq \beta_{j,0}$, the test criterion based on Eq. (6.42) is

$$T_j = \frac{|b_j - \beta_{j,0}|}{\hat{\sigma} \sqrt{c_{jj}}} \quad (6.48)$$

which has approximately the Student t -distribution with $n - m$ degrees of freedom when H_0 is valid.

Most regression programs perform a Fisher–Snedecor test of significance of the determination coefficient Eq. (6.39) and a Student t -test on the significance of the individual parameters β_j calculated from Eq. (6.48) with $\beta_{j,0} = 0$. The F -test also determines simultaneous significance of all components of vector β except an absolute member. There are four cases:

- (1) The F -test is not significant and all t -tests also are not significant. Then the model is considered to be unsuitable because it does not explain the variability of y .
- (2) The F -test and all t -tests are significant. Then the model is considered to be suitable to express the variability of y . It does not mean, however, that the model is correct and acceptable.
- (3) The F -test is significant but one or more t -test is not significant. Then the model is considered to be suitable, but the controllable variables x_j for which the parameters β_j are not significantly different from zero are rejected.
- (4) The F -test is significant but all t -tests for parameters β indicate that all controllable variables are insignificant. Paradoxically, this shows that although the model as a whole is suitable, none of controllable variables is significant. This may result from multicollinearity (Section 6.3.21).

It should be noted that a model is considered to be significant when the form of the model is $f(\mathbf{x}, \boldsymbol{\beta}) = \mathbf{X}\boldsymbol{\beta}$ but not $f(\mathbf{x}, \boldsymbol{\beta}) = J\bar{y}$.

Problem 6.8. *Disadvantages of classical statistical analysis for linear regression*

Anscombe [5] published test data for four simulated samples of size $n = 11$. Test the statistical significance of parameters β_1 and β_2 , i.e. $H_0: \beta_1 = 0$ against $H_A: \beta_1 \neq 0$ and $H_0: \beta_2 = 0$ against $H_A: \beta_2 \neq 0$, and compare the results of tests with graphical analysis of residuals.

Data: To check the efficiency of statistical algorithms for linear regression models, the following test data samples are often used. The four data samples have the same statistical characteristics $b_1 = 0.5$, $b_2 = 3.0$, $D(b_1) = 0.0139$ and $D(b_2) = 1.2656$.

Table 6.1. Test data for linear regression

Data sample	A		B		C		D	
Variable	x	y	y	y	x	y	x	y
Point								
1	10	8.04	9.14	7.46	8	6.58		
2	8	6.95	8.14	6.77	8	5.76		
3	13	7.58	8.74	12.74	8	7.71		
4	9	8.81	8.77	7.11	8	8.84		
5	11	8.33	9.26	7.81	8	8.47		
6	14	9.96	8.10	8.84	8	7.04		
7	6	7.24	6.13	6.08	8	5.25		
8	4	4.26	3.10	5.39	19	12.50		
9	12	10.84	9.13	8.15	8	5.56		
10	7	4.82	7.26	6.42	8	7.91		
11	5	5.68	4.74	5.73	8	6.89		

Solution: When the linear regression model is $E(y/x) = \beta_1 x + \beta_2$, all four sets of data lead to the same parameter estimates i.e. $b_1 = 0.5$ and $b_2 = 3.0$, with parameter variances $D(b_1) = 0.0139$ and $D(b_2) = 1.2656$, and the test criteria are $T_1 = 2.667$ and $T_2 = 4.241$. The test criterion F_R (6.39) has the same value $F_R = 17.97$, the determination coefficient $\hat{R}^2 = 0.66$ and the residual standard deviation $\hat{\sigma} = 1.237$. This leads to the conclusion that both regression parameters β_1 and β_2 are significantly different from zero.

Since the quantile of the F -distribution $F_{0.95}(1, 9) = 5.117$ is less than the calculated F_R , the determination coefficient differs from zero. It would seem that for all four data samples the linear regression model fits quite well. Figures 6.10 to 6.13 show

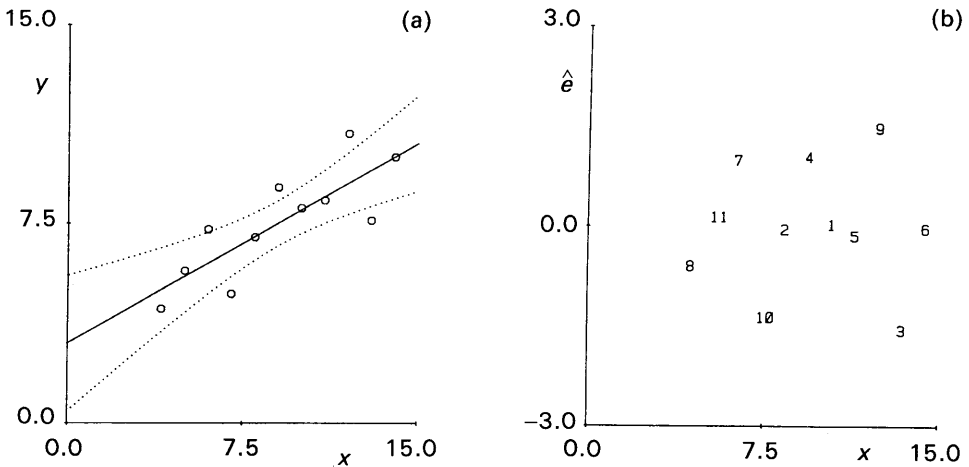


Fig. 6.10—(a) Linear regression model $\hat{y}_p = f(x)$ for data sample A, and (b) dependence of residuals $\hat{e} = f(x)$ on variable x .

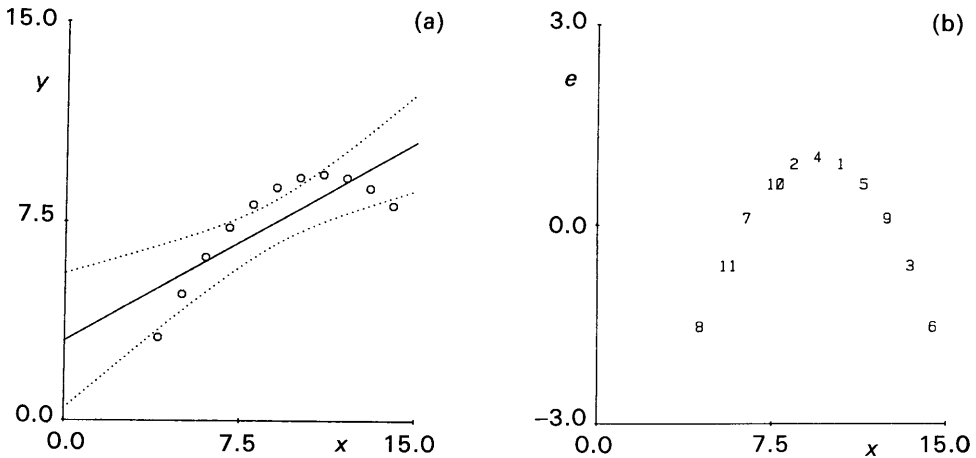


Fig. 6.11—(a) Linear regression model $\hat{y}_p = f(x)$ for data sample B, and (b) dependence of residuals $\hat{e} = f(x)$ on variable x .

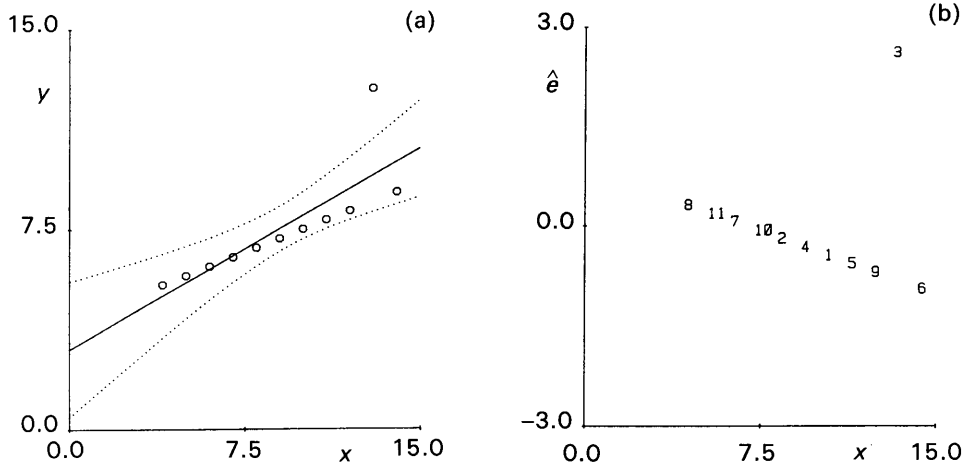


Fig. 6.12—(a) Linear regression model $\hat{y}_p = f(x)$ for data sample C, and (b) dependence of residuals $\hat{e} = f(x)$ on variable x .

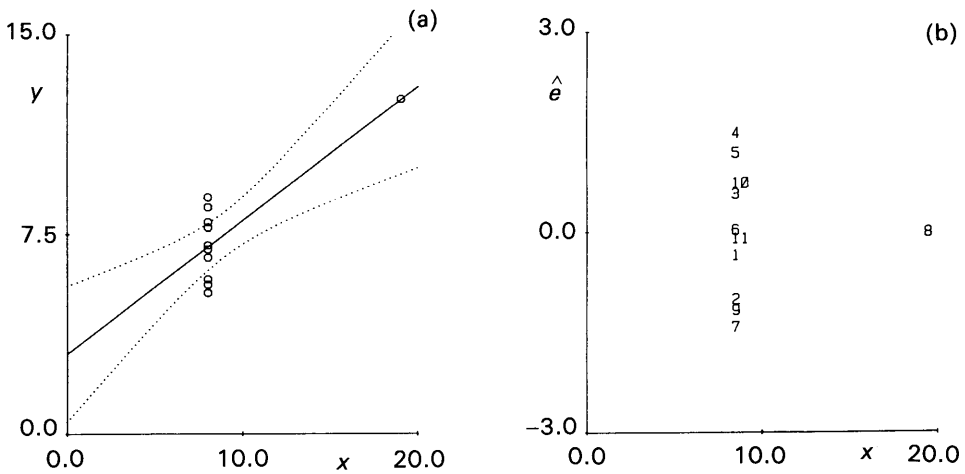


Fig. 6.13—(a) Linear regression model $\hat{y}_p = f(x)$ for data sample D, and (b) dependence of residuals $\hat{e} = f(x)$ on variable x .

that only data sample A is well characterized by a linear model. A good approximation is also reached for sample C, where just one outlier prevents the data from corresponding to a linear model.

It may be rather surprising to a user of linear regression analysis that the statistical characteristics do not indicate here either the nonlinear trend of sample B (Fig. 6.11) or the silly data of sample D (Fig. 6.13).

Conclusion: Classical regression analysis may suggest models that do not correspond at all to the data set. Any model must be confirmed by graphical examination of residuals.

6.3.2.1 Testing for multicollinearity

Paradoxical cases where the F -test is significant and all t -tests are not significant may result from strong multicollinearity among columns of matrix \mathbf{X} . In correlation models, this corresponds to a situation when there are high values of paired correlation coefficients between controllable variables. Multicollinearity may be recognized by a finding that vectors \mathbf{x}_j and \mathbf{x}_k , $j \neq k$, (which represents columns of matrix \mathbf{X}) are approximately parallel.

In the presence of multicollinearity, it is not possible to find the influence of the individual controllable variables x_j . Multicollinearity may exist in models which fit experimental reality quite well. Here, RSC has a small value and the predictions $\hat{y}_{P,i}$ are quite close to the experimental values y_i . Multicollinearity also appears in polynomial models and data which come from unplanned experiments.

Multicollinearity can be removed, for example, by selecting the location of experimental points such that the columns of matrix \mathbf{X} will be mutually orthogonal, i.e. their scalar product will be zero.

$$\langle \mathbf{x}_j, \mathbf{x}_k \rangle = \sum_{i=1}^n x_{ij}x_{ik} = 0 \quad \text{for } j \neq k$$

If all columns of matrix \mathbf{X} are mutually orthogonal, the matrix $\mathbf{X}^T\mathbf{X}$ is diagonal and a solution of Eq. (6.11) can be expressed in the form

$$b_j = \frac{\sum_{i=1}^n x_{ij}y_i}{\sum_{i=1}^n x_{ij}^2}, \quad j = 1, \dots, m$$

and the test criterion F_R is

$$F_R = \frac{\sum_{j=1}^{m-1} T_j^2}{m-1} = T_S$$

where T_S is an average value of all test statistics T_j^2 defined by Eq. (6.48) for $\beta_{0j} = 0$. It is supposed here that β_m is the intercept term.

To examine the suitability of a proposed linear model with regard to possible multicollinearity, Scott [8] uses a test criterion M_T

$$M_T = \frac{\frac{F_R}{T_S} - 1}{\frac{F_R}{T_S} + 1} \quad (6.49)$$

and the following rules for identification of multicollinearity.

- (a) If $M_T > 0.8$ the model is not suitable because of multicollinearity, so a model correction is necessary.
- (b) If $0.33 \leq M_T \leq 0.8$ the model is poor because of multicollinearity, so some

- model correction is recommended.
- (c) If $M_T < 0.33$, the model has little problem from multicollinearity, so no model correction is necessary.

The M_T criterion is useful in cases where it is necessary to discover all controllable variables which significantly affect the variability of the dependent variable y . When data are approximated by an empirical model, for example, by a polynomial, the M_T values need not be considered.

Problem 6.9. *Approximation of an absorption spectrum by a polynomial*
Find a model which describes the dependence of the molar absorptivity ϵ on wave-length λ , $\epsilon = f(\lambda)$. Use a polynomial of the second degree $E(\epsilon/\lambda) = \beta_1 + \beta_2 \lambda^2 + \beta_3 \lambda^2$.
Data: $n = 15$, $m = 3$

ϵ , mol ⁻¹ . dm ³ . cm ⁻¹	3	3.4	4.3	5	6	6.8	8.1	9.2
λ , nm	460	470	480	490	500	510	520	530

10.7	11.6	12.9	13.6	14.6	15.3	15.5
540	550	560	570	580	590	600

Solution: Table 6.2 lists the numerical values of the parameter estimates b_1 , b_2 and b_3 with their standard deviations and test statistics T_j for $\beta_j = 0$. Since the test criterion $F_R = 696$ is greater than the corresponding quantile of the Fisher–Snedecor F -distribution $F_{0.95}(2, 12) = 3.885$, the proposed model is statistically significant. In contrast, the quantile of the Student t -distribution $t_{0.975}(12) = 2.2$ is greater than both T_1 and T_2 , therefore both parameters β_1 and β_2 are insignificant. The test criterion $M_T = 0.989$ (Eq. (6.49)) indicates very strong multicollinearity in the model.

Table 6.2. Parameter estimates and their statistical characteristics

j	Parameter	Estimate b	$\sqrt{D(b_j)}$	T_j
1	β_1	-43.93	19.38	-2.267
2	β_2	0.1018	0.0735	1.386
3	β_3	-2.51×10^{-6}	6.923×10^{-5}	-0.0361

Conclusion: In polynomial models, the significance of individual terms of the equation can not be judged from the result of the Student t -test alone.

Statistical tests are rather insensitive to small deviations of the error distribution from normality. However, in cases of strong non-normality or heteroscedasticity it is necessary to make a correction to the number of degrees of freedom for determination of the quantiles of the Fisher–Snedecor and Student distributions.

6.3.2.2 Test of significance of the intercept term

In chemometrics practice, it is important to examine the significance of the intercept term β_m by testing the null hypothesis $H_0: \beta_m = 0$ against the alternative $H_A: \beta_m \neq 0$. An intercept term always exists in correlation models. In regression models the intercept term ensures a zero sum of residuals $\sum_{i=1}^n \hat{e}_i = 0$.

In programs for regression analysis the following difficulties exist regarding the intercept term:

- The intercept term β_m always exists for centred data.
- Because the value $\bar{y} = 0$ is used in its calculation, the determination coefficient (6.38) for models without an intercept, \hat{R}_B^2 will be significantly higher than \hat{R}^2 for models with an intercept. The residual sum of squares for a model without an intercept, RSC_B , is always higher than or equal to the residual square sum for a model with an intercept, RSC .

Good programs allow calculation for a model with or without an intercept term, and correctly evaluate the determination coefficient because they do not substitute $\bar{y} = 0$. The difficulty can be avoided by introduction of a fictional point (x_{n+1}, y_{n+1}) with the co-ordinates [9]:

$$x_{n+1,j} = n^* \bar{x}_j, \quad j = 1, \dots, m-1$$

$$y_{n+1} = n^* \bar{y}$$

where $n^* = 1 + \sqrt{n+1}$ and \bar{x}_j and \bar{y} are arithmetic means calculated from the data. With the use of this extended data set, the classical linear regression with an intercept term leads to the same results as the regression without an intercept term for the original set of n data points.

The influence of an intercept term on regression model may be understood by considering the location of point (x_{n+1}, y_{n+1}) with regard to other points. When this point is an outlier, the model without an intercept term is not suitable. The significance of an intercept term may be evaluated by use of Jack-knife residuals (6.97). If the point (x_{n+1}, y_{n+1}) is far from other points, the data are not in a suitable range to allow testing for presence of an intercept term. The significance of an intercept term may be also examined by the test statistic T_j (6.48), with $\beta_{oj} = 0$.

Problem 6.10. A Lambert–Beer Law calibration line

Estimate the parameters of the calibration line for the Lambert–Beer law for the dependence of a measured absorbance A on the concentration c_i of *p*-nitroaniline. Does the straight line pass through the origin?

Data: $n = 6$, $m = 2$, $d = 1.000 \text{ cm}$

$c, 10^5 \cdot \text{mol} \cdot \text{dm}^{-3}$	1.98	2.58	3.42	4.43	5.51	6.58
A	0.293	0.374	0.500	0.642	0.804	0.963

Solution: If the Lambert–Beer law holds, the straight line $A = \varepsilon_M d c$, where ε_M is the

molar absorptivity, d is the cell length and c the molar concentration, passes through the origin. We will add the fictional point with co-ordinates $A_{n+1} = (1 + 7)\bar{A} = 2.172$ and $c_{n+1} = (1 + 7)\bar{c} = 14.88$. For this extended data set the model found was $A = (0.146 \pm 0.0003)c + (2.703 \pm 0.206) \times 10^{-5}$. The test statistics $T_1 = 490$ and $T_2 = 0.013$ show that the straight line goes through the origin.

To check whether it is correct to neglect the intercept term, the extended data set

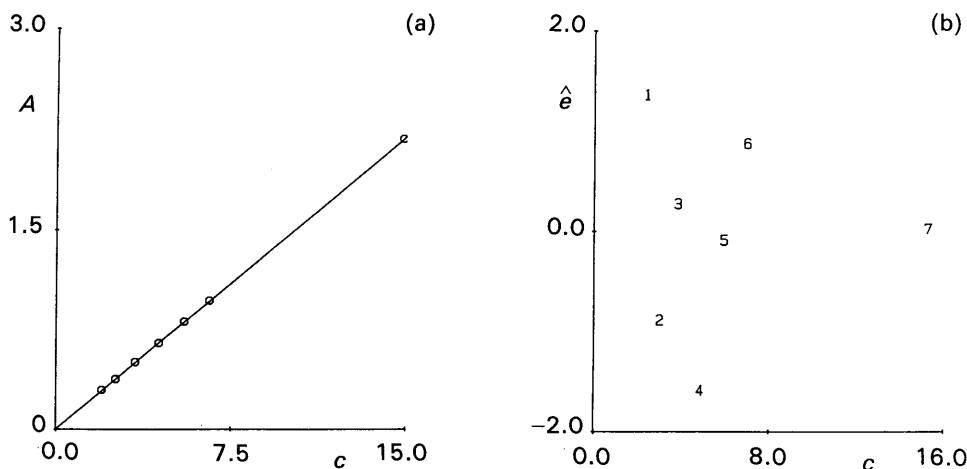


Fig. 6.14—Test of significance of the intercept term in a Lambert–Beer law calibration: (a) the line for $\hat{A} = \epsilon c$ and the fictional point, and (b) the dependence of the residuals \hat{e} on variable c .

is plotted in Fig. 6.14. It is obvious that the point (x_{n+1}, y_{n+1}) does not differ from the others, and the Jack-knife residual (Eq. 6.97) $\hat{e}_{j,n+1} = 0.037$ also indicates that the intercept is insignificant.

Problem 6.11. *Estimates of the parameters of a calibration line that passes through the origin*

Demonstrate a procedure for parameter estimation in the case of a calibration straight line which must pass through the origin.

Solution: The regression model $E(y/x) = \beta_1 x$ is shown in Fig. 6.15. The vector $\hat{\mathbf{y}}_p$ is the perpendicular projection of vector \mathbf{y} on vector \mathbf{x} . The estimate b_1 may be calculated by the simple expression from Problem 6.1, but we use here an analytical derivation of the least-squares criterion $U(\beta) = \sum_{i=1}^n (y_i - \beta_1 x_i)^2$,

$$\frac{\delta U(\beta)}{\delta \beta_1} = -2 \sum_{i=1}^n (y_i - \beta_1 x_i) x_i = 0$$

and rewriting

$$b_1 = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} = \sum_{i=1}^n v_i y_i$$

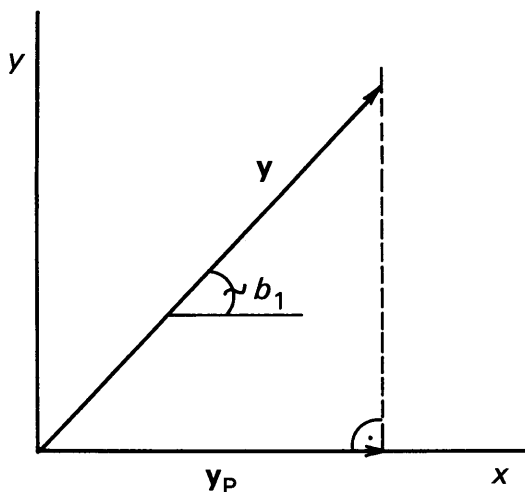


Fig. 6.15—Geometrical illustration of the regression model $y = \beta_1 x$, and the projection of vector y on vector x .

The variance of this estimate is calculated from

$$D(b_1) = \sum_{i=1}^n v_i^2 D(y_i) = \frac{\sigma^2}{\sum_{i=1}^n x_i^2}$$

The residual sum of squares RSC is given by

$$RSC = \sum_{i=1}^n y_i^2 - b_1 \sum_{i=1}^n x_i y_i$$

and the theoretical sum of squares SS is

$$SS = b_1 \sum_{i=1}^n x_i y_i = b_1^2 \sum_{i=1}^n x_i^2$$

For the determination coefficient R^2 it holds

$$R^2 = 1 - \frac{\sum_{i=1}^n y_i^2 - b_1 \sum_{i=1}^n x_i y_i}{\sum_{i=1}^n y_i^2 - n\bar{y}^2} = \frac{b_1 \sum_{i=1}^n y_i x_i - n\bar{y}^2}{\sum_{i=1}^n y_i^2 - n\bar{y}^2}$$

The projection matrix \mathbf{H} contains elements

$$H_{jj} = \frac{x_j^2}{\sum_{i=1}^n x_i^2} \quad \text{and} \quad H_{jk} = \frac{x_j x_k}{\sum_{i=1}^n x_i^2}, \quad j, k = 1, \dots, n.$$

From Eq. (6.17) we have

$$D(\hat{y}_{P,i}) = \frac{\hat{\sigma}^2 x_j^2}{\sum_{i=1}^n x_i^2}$$

and the confidence interval for a model of regression straight line is calculated from Eq. (6.45)

$$\begin{aligned} x_0 b_1 - t_{1-\alpha/2}(n-1) \frac{\hat{\sigma} \cdot x_0}{\sqrt{\sum_{i=1}^n x_i^2}} &\leq x_0 \beta_1 \\ &\leq x_0 b_1 + t_{1-\alpha/2}(n-1) \frac{\hat{\sigma} \cdot x_0}{\sqrt{\sum_{i=1}^n x_i^2}} \end{aligned}$$

The end-points of the two confidence intervals are straight lines going through an origin, whereas for a general model of a regression straight line they are parabolic curves.

6.3.2.3 Simultaneous test of a composite hypothesis

The likelihood ratio test (Section 8.6.2) may be used for testing general parametric hypotheses. In a case where the null hypothesis $H_0: \beta_2 = 0$ is to be tested against the alternative $H_A: \beta_2 \neq 0$, where β_2 represents the last q elements of the vector β , the regression model is expressed in the divided form

$$y = [X_1 \quad X_2] \cdot \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} + \varepsilon = X_1 \beta_1 + X_2 \beta_2 + \varepsilon$$

where X_1 is the matrix of dimension $[n \times (m - q)]$ containing those controllable variables with regression coefficients that are not included in a test vector β_2 . Similarly, X_2 is the matrix of dimension $(n \times q)$ containing those controllable variables with regression coefficients that are included in a test vector β_2 .

When the hypothesis H_0 is valid, it is evident that

$$\hat{y}_{P,1} = X_1 b_1$$

where

$$b_1 = (X_1^T X_1)^{-1} X_1^T y$$

and the corresponding residual sum of squares RSC_1 is

$$RSC_1 = (y - \hat{y}_{P,1})^T (y - \hat{y}_{P,1})$$

When the hypothesis H_A is valid, we have

$$\hat{y}_P = Xb$$

where

$$\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

and the corresponding residual sum of squares RSC is

$$RSC = (\mathbf{y} - \hat{\mathbf{y}}_P)^T (\mathbf{y} - \hat{\mathbf{y}}_P)$$

The difference $(RSC_1 - RSC)$ corresponds to an increase in the residual sum of squares caused by validity of the null hypothesis H_0 . The test criterion has the form

$$F_1 = \frac{(RSC_1 - RSC)(n - m)}{RSCq}$$

which if the H_0 hypothesis is valid, has the Fisher–Snedecor F -distribution with q and $(n - m)$ degrees of freedom.

A mistake often made in the application of linear regression in chemical laboratories is a false approach to a choice of test criteria. Instead of the test criterion F_1 , the individual test statistics T_j from Eq. (6.48) are calculated, and on their basis, the significance of a composite hypothesis $H_0: \beta_2 = \beta_{2,0}$ against $H_A: \beta_2 \neq \beta_{2,0}$, is tested. Here $\beta_{2,0}$ is the vector of known parameters.

For tests of composite hypotheses, the test statistic F_1 should be used, where RSC_1 is the residual sum of squares for the model

$$\hat{\mathbf{y}}_{P,i} = \mathbf{X}_1 \mathbf{b}_1 + \mathbf{X}_2 \beta_{2,0}$$

where \mathbf{b}_1 is the estimate of parameters β_1 on the assumption that the restriction $\beta_2 = \beta_{2,0}$ is valid.

Problem 6.12. *Simultaneous test of a composite hypothesis for a Lambert–Beer law model*

For the data from Problem 6.10, test the composite null hypothesis $H_0: \beta_2 = 0, \beta_1 = 0.148$ against $H_A: \beta_2 \neq 0, \beta_1 \neq 0.148$. The false approach would be two separate tests of two null hypotheses, $H_0: \beta_2 = 0$ and $H_0: \beta_1 = 0.148$.

Solution: On substitution into Eq. (6.48), we obtain

$$T_2 = \frac{|1.461 \times 10^{-4} - 0|}{0.00398} = 0.037$$

$$T_1 = \frac{|0.1459 - 0.148|}{0.000908} = 2.314$$

Because T_1 and T_2 are less than the quantile of the Student t -distribution, $t_{0.975}(4) = 2.7764$, both tests lead to a conclusion that $H_0: \beta_2 = 0, \beta_1 = 0.148$ should be accepted. This conclusion is, however, *false*.

The more rigorous approach uses a simultaneous test of the composite hypothesis $H_0: \beta_2 = 0$ and $\beta_1 = 0.148$.

The procedure starts with a calculation of $RSC = 5.12 \times 10^{-5}$ for estimates $b_1 = 0.1459$ and $b_2 = 1.461 \times 10^{-4}$. Then, $RSC_1 = 5.3476 \times 10^{-4}$ for parameters $\beta_{2,0} = 0$ and $\beta_{1,0} = 0.148$ is calculated. From Eq. (6.50), the test criterion F_1 is

$$F_1 = \frac{(5.347 \times 10^{-4} - 5.12 \times 10^{-5}) \times 4}{5.12 \times 10^{-5} \times 2} = 18.89$$

Because the quantile of Fisher–Snedecor F -distribution is $F_{0.95}(2, 4) = 6.944$, the null hypothesis $H_0: \beta_2 = 0$ and $\beta_1 = 0.148$ cannot be accepted. The result of this F -test

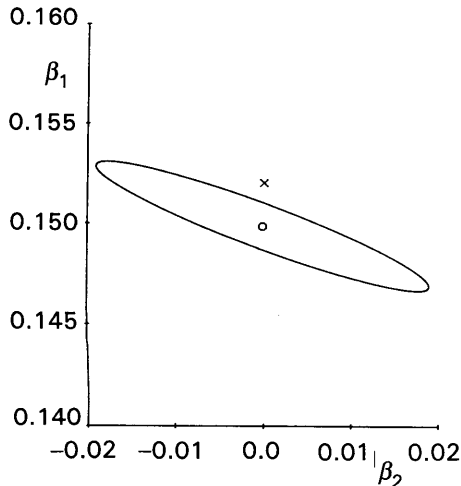


Fig. 6.16—The 95% confidence interval for the parameters β_1 and β_2 . The point $\beta_2 = 0$, $\beta_1 = 0.148$ is marked by a cross.

is not in agreement with conclusion of the previous t -tests. Figure 6.16 shows the 95% confidence ellipse of parameters β_1 and β_2 , and the point $\beta_{1,0} = 0.148$ and $\beta_{2,0} = 0$ marked by a cross. This point lies outside the 95% confidence interval of the two parameters.

Conclusion: It may be concluded that a simultaneous test of the composite hypothesis cannot be replaced by tests of two separate hypotheses. Thus, testing of individual parameters in a vector β_0 can lead to quite false conclusions.

Problem 6.13. *Validation of a new analytical method by a simultaneous test of a composite hypothesis*

Try to test a composite hypothesis $H_0: \beta_2 = 0$ and $\beta_1 = 1$ in Problem 6.7 against the alternative $H_A: \beta_2 \neq 0$ and $\beta_1 \neq 1$.

Data: from Problem 6.7

Solution: From the results of Problem 6.7, we have $RSC = 3440$, and when we set $\beta_{1,0} = 1$ and $\beta_{2,0} = 0$, we obtain $RSC_1 = 8221$. On substitution into Eq. (6.50), we find

$$F_1 = \frac{(8220 - 3440) \times 22}{3440 \times 2} = 15.28$$

which is greater than the quantile of the Fisher–Snedecor F -distribution

$F_{0.95}(2, 22) = 3.44$, so the null hypothesis H_0 cannot be accepted. This conclusion is also in agreement with the partial t -tests and confidence intervals of the two

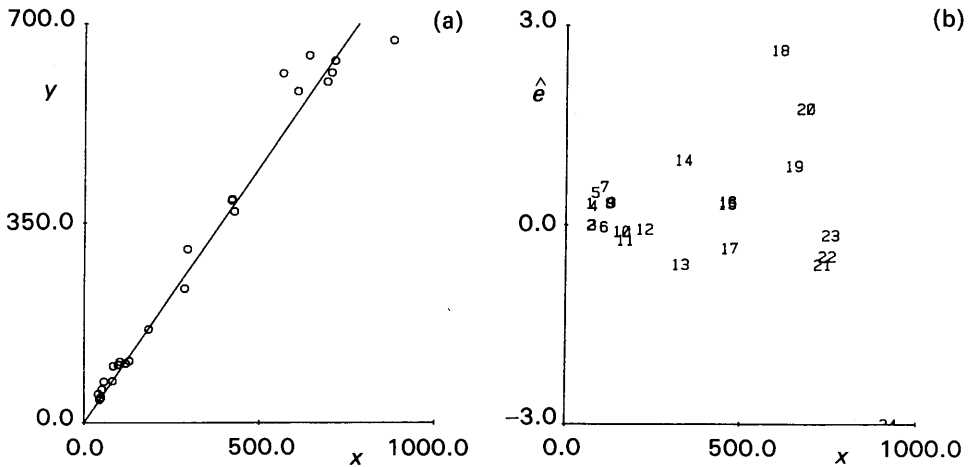


Fig. 6.17—(a) Linear regression model of validation of a new analytical method $\hat{y}_p = x$, and (b) dependence of the residuals on x .

parameters. Figure 6.17 shows the regression straight line $\hat{y}_p = x$, with experimental points and a graphical analysis of residuals.

Conclusion: A simultaneous test of the composite hypothesis ($H_0: \beta_2 = 0$ and $\beta_1 = 1$) confirmed that a new analytical method is not in agreement with the results of a standard one.

6.3.2.4 Test of agreement of two linear models

The test of a composite hypothesis just described may be re-arranged to allow for testing of agreement of parameters in two linear models

$$y_1 = X_1\beta_1 + \varepsilon_1 \quad (6.51a)$$

$$y_2 = X_2\beta_2 + \varepsilon_2 \quad (6.51b)$$

where X_1 is a matrix of dimension $(n_1 \times m)$, y_1 is a vector of dimension $(n_1 \times 1)$, X_2 is a matrix of dimension $(n_2 \times m)$, and y_2 is a vector of dimension $(n_2 \times 1)$. RSC_1 is the residual sum of squares corresponding to model (6.51a), RSC_2 is the residual sum of squares corresponding to model (6.51b) and RSC is the residual sum of squares corresponding to the composite model:

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \times \beta + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \end{bmatrix} \quad (6.52)$$

We use the Chow test of the null hypothesis $H_0: \beta_1 = \beta_2$ against the alternative $H_A: \beta_1 \neq \beta_2$, based on the test criterion

$$F_C = \frac{(RSC - RSC_1 - RSC_2)(n - 2m)}{(RSC_1 + RSC_2) \times m} \quad (6.53)$$

where $n = n_1 + n_2$. If the variances of the two samples are the same ($\sigma_1^2 = \sigma_2^2$, homoscedasticity), the test criterion F_C has the Fisher–Snedecor F -distribution with m and $(n - 2m)$ degrees of freedom.

When the variances of the two samples are not the same ($\sigma_1^2 \neq \sigma_2^2$, heteroscedasticity), the Fisher–Snedecor F -distribution may be used with m and r degrees of freedom where

$$r = \frac{[(n_1 - m)\sigma_1^2 + (n_2 - m)\sigma_2^2]^2}{(n_1 - m)\sigma_1^4 + (n_2 - m)\sigma_2^4} \quad (6.54)$$

A more accurate version of this equation is given in [6].

Problem 6.14. *Comparison of measurement results from two laboratories*

Determination of the free energy ΔG of the vapour of boron oxide as a function of temperature T was carried out in two laboratories [11]. Compare the results and test whether the values measured in the two laboratories can be considered to be the same.

Data: $n = 6$, $m = 2$,

T, K	$-\Delta G, \text{kcal/mol}$	
	Lab A	Lab B
1409	34.9	34.9
1441	34.6	33.8
1457	31.9	33.4
1492	33.1	32.4
1569	30.1	30.3
1610	29.3	29.1

Solution: If a linear regression model is valid for both data samples, the models are

$$E(\Delta G/T) = \beta_{1,A}T + \beta_{2,A}$$

$$E(\Delta G/T) = \beta_{1,B}T + \beta_{2,B}$$

We will test the null hypothesis $H_0: \beta_A = \beta_B$ against the alternative $H_A: \beta_A \neq \beta_B$, where $\beta_A = (\beta_{1,A}, \beta_{2,A})^T$ and $\beta_B = (\beta_{1,B}, \beta_{2,B})^T$. We use the Chow test:

Laboratory A: $b_{1,A} = -0.02768(\pm 0.00525)$

$$b_{2,A} = 73.73(\pm 7.865)$$

$$\hat{\sigma} = 0.916$$

$$RSC_A = 3.358$$

Laboratory B: $b_{1,B} = -0.02776(\pm 0.000157)$

$$b_{2,B} = 73.82(\pm 0.235)$$

$$\hat{\sigma} = 0.0274$$

$$RSC_B = 0.002992$$

Laboratory A + B: $b_{1,A+B} = -0.02772(\pm 0.00235)$

$$b_{2,A+B} = 73.77(\pm 3.521)$$

$$\hat{\sigma}_{A+B} = 0.58$$

$$RSC_{A+B} = 3.364$$

The standard deviations of the parameters are given in brackets. Substitution into Eq. (6.53) for $RSC = RSC_{A+B}$, $RSC_1 = RSC_A$ and $RSC_2 = RSC_B$, leads to

$$F_C = \frac{(3.364 - 3.358 - 0.002992)(12 - 4)}{(3.358 + 0.002992) \times 2} = 0.0036$$

Because the variances of the samples differ, we calculate the degrees of freedom r from Eq. (6.54).

$$r = \frac{[4 \times 0.916^2 + 4 \times 0.0274^2]^2}{4 \times 0.916^4 + 4 \times 0.0274^4} = 4.007 \div 4$$

The quantile of the Fisher-Snedecor F -distribution $F_{0.95}(2, 4) = 6.94$ is greater than

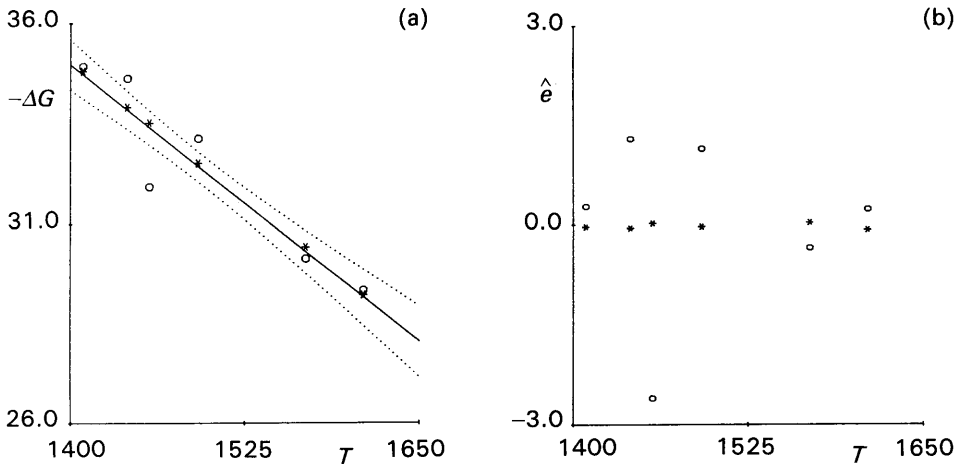


Fig. 6.18—(a) Regression straight line for data $(-\Delta G, T)$ for two laboratories. Data for laboratory A (denoted by circle) have more spread; (b) dependence of residuals: lab A (denoted by circle) and lab B (denoted by star).

F_C , so $H_0: \beta_A = \beta_B$ is accepted. Figure 6.18 shows a graphical interpretation of the laboratory measurements. The data from lab A show much more spread.

Conclusion: On the basis of a Chow test, it may be concluded that results from the two laboratories can be considered to be the same. The data from lab A are less precise.

Problem 6.15. Comparison of two calibration straight lines

Two insulin samples, A and B, are compared according to their ability to decrease the level of blood sugar. The sample A was injected into 11 randomly chosen rats and sample B into 9 rats. The decrease of blood sugar level was determined. Compare the efficiency of the two insulin samples.

Data: Insulin A: x is the amount in μl of insulin A, and y is a decrease in sugar level

in blood, $n_A = 11$, $m_A = 2$,

x	120	160	200	240	280	320	360	400	440	480	500
y	17	26	30	27	45	47	48	63	60	69	69

Insulin B: the same for insulin B, $n_B = 9$, $m_B = 2$

x	169	200	240	280	320	360	400	440	480
y	9	18	17	25	39	45	47	57	61

Solution: If the two insulin types have the same effect, the two regression straight lines will not be significantly different. To test the agreement between the lines we use the test criterion F_C (6.53). The statistical characteristics are:

Insulin A: $\hat{y}_A = 1.808(\pm 3.504) + 0.1369(\pm 0.0103)x$
 $RSC_A = 159.6$
 $\hat{\sigma} = 4.211$

Insulin B: $\hat{y}_B = -18.67(\pm 3.535) + 0.1688(\pm 0.0105)x$
 $RSC_B = 74.25$
 $\hat{\sigma} = 3.26$

Insulin A + B: $\hat{y}_{A+B} = -6.397(\pm 4.45) + 0.1481(\pm 0.0132)x$
 $RSC_{A+B} = 821.1$

Then from Eq. (6.53) we find:

$$F_C = \frac{(821.1 - 159.6 - 74.25)(20 - 4)}{(159.6 + 74.25) \times 2} = 467.73$$

which is a greater value than the quantile of the Fischer-Snedecor F -distribution $F_{0.95}(2, 16) = 3.63$, so that the null hypothesis H_0 is rejected. From Fig. 6.19 it is evident that, although the straight lines have similar slopes, they differ in intercept.

Conclusion: The insulin samples have significantly different activity.

6.3.2.5 Acceptance test for a proposed linear model

Utts [7] has introduced a test of acceptance of a proposed linear regression model $f(\mathbf{x}, \boldsymbol{\beta}) = \mathbf{X}\boldsymbol{\beta}$, based on a ratio of the residual sum of squares.

If the regression model $f(\mathbf{x}, \boldsymbol{\beta})$ is non-linear there exists a group of points n_1 to which a linear model will fit, as shown in Fig. 6.20.

Let us denote by RSC_1 the residual sum of squares corresponding to a linear regression of n_1 points, and by RSC the residual sum of squares corresponding to a linear regression of all n points. Utts criterion for model acceptance is formulated as

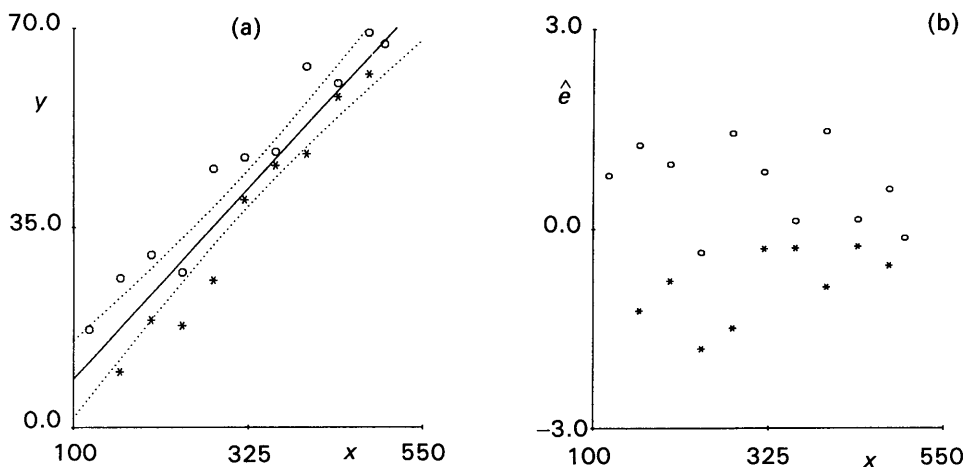


Fig. 6.19—(a) Regression straight line for insulin data A and B and (b) dependence of residuals \hat{e} on variable x . \circ —sample A, *—sample B.

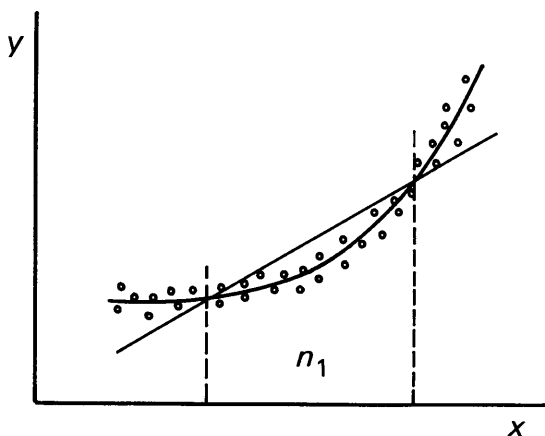


Fig. 6.20—Principle of the Utts test.

$$F_{\mu} = \frac{(RSC - RSC_1) \times (n_1 - m)}{RSC_1 \times (n - n_1)} \quad (6.55)$$

which for the hypothesis H_0 : “the linear regression model is valid” has the Fisher-Snedecor F -distribution with $(n - n_1)$ and $(n_1 - m)$ degrees of freedom. Utts recommends choosing $n_1 \approx n/2$ and selecting the points that give the smallest values of the diagonal elements $H_{ii} = \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i$ of the projection matrix \mathbf{H} . Such selected points lie close to the centre of gravity of the controllable independent variables. If the calculated F_{μ} is smaller than the corresponding quantile of F -distribution, the linear regression model can be accepted.

Another group of tests of acceptance of a proposed regression model [7] is based on an application of extended regression model. The proposed linear model is usually

extended to include higher powers of the independent variables and eventually also their interactions. If original linear model was correct, the parameters for higher members of the new model will be statistically insignificant.

For a model of a regression straight line, the assumption of linearity may be checked by a test of significance of parameter β_2 in an extended model

$$E(y/x) = \beta_1 x + \beta_2 x^2 + \beta_3 \quad (6.55a)$$

Here, the t -test could be used, but as any multicollinearity would change the results of the t -test, the Fisher–Snedecor F -test is applied. Let RSC_Q be the residual sum of squares for a quadratic model and RSC_L the residual sum of squares for a linear model. The test criterion of linearity will have the form

$$F_L = \frac{(RSC_L - RSC_Q)(n - 3)}{RSC_Q \times 1} \quad (6.56)$$

If the null hypothesis $H_0: \beta_2 = 0$ is valid, the F_L criterion has the Fisher–Snedecor F -distribution with 1 and $(n - 3)$ degrees of freedom. The statistic $\sqrt{F_L}$ has the Student distribution with $(n - 3)$ degrees of freedom.

Instead of the various test criteria for testing linearity, the statistical characteristics for comparison of different models may also be used. One of these is the *mean quadratic error of prediction* defined by:

$$MEP = \frac{\sum_{i=1}^n (y_i - \mathbf{x}_i^T \mathbf{b}_{(i)})^2}{n} \quad (6.57a)$$

where $\mathbf{b}_{(i)}$ is the estimate of regression parameters when all points except the i th one were used and \mathbf{x}_i is i th row of matrix \mathbf{X} . The statistic MEP uses a prediction $\hat{y}_{p,i}$ from an estimate constructed without including the i th point. Another mathematical expression for MEP is:

$$MEP = \sum_{i=1}^n \frac{\hat{e}_i^2}{(1 - H_{ii})^2 n} \quad (6.57b)$$

For large sample sizes n the element H_{ii} tends to zero ($H_{ii} \approx 0$) and then

$$MEP = RSC/n \quad (6.57c)$$

If MEP is used instead of RSC in the equation for the determination coefficient (6.38), the resulting statistic \hat{R}_p^2 is called the *predicted determination coefficient*

$$\hat{R}_p^2 = 1 - \frac{n \times MEP}{\sum_{i=1}^n y_i^2 - n \times \bar{y}^2} \quad (6.58)$$

Another statistical characteristic has quite general use and is derived from information theory and entropy [12] and known as the *Akaike information criterion*, AIC

$$AIC = n \ln\left(\frac{RSC}{n}\right) + 2m \quad (6.59)$$

The most suitable model is the one which gives the lowest value of the Akaike information criterion.

Problem 6.16. *Selection from three polynomial models*

For the data from Problem 6.9, do a regression analysis, and test whether the data sample should be fitted by a polynomial of the third or fifth degree.

Data: Problem 6.9

Solution: Table 6.3 lists the statistical characteristics MEP , \hat{R}_p^2 , \hat{R}^2 and AIC for the hypotheses that the regression model is expressed by a polynomial of the second, third and fifth degree.

Table 6.3. Selection of best model according to the statistics MEP , \hat{R}_p^2 , \hat{R}^2 and AIC

Polynomial degree	MEP	\hat{R}_p^2	\hat{R}^2	AIC
2	0.3502	0.9905	0.9915	-21.65
3	0.0283	0.9992	0.9992	-56.02
5	0.0613	0.9985	0.9997	-55.04

Because of the good precision of the experimental data, the statistics MEP and \hat{R}_p^2 do not indicate that the polynomial of third degree is the most suitable. Only statistic AIC indicates that the best model is the polynomial of third degree:

$$\hat{y}_P = 860.2(\pm 85.17) - 5.057(\pm 0.485)x + 0.00977(\pm 0.00092)x^2 - 6.146 \times 10^{-6}(\pm 5.78 \times 10^{-7})x^3$$

and estimates of all three parameters are statistically significant. In the case of the polynomial of fifth degree, all the parameter estimates except β_3 are statistically insignificant, as a consequence of multicollinearity.

Figure 6.21 shows the curve fitting of the data by a polynomial of the second degree and Fig. 6.22 by a polynomial of the third degree. The numerical statistics \hat{R}^2 and \hat{R}_p^2 are not able to distinguish between these two polynomials, but the graphical analysis of residuals is a more efficient tool for deciding among several plausible models. *Conclusion:* There are cases when statistical characteristics fail, and graphical examination of residuals gives more satisfactory results for model specification.

Problem 6.17. *Examination of linearity of four samples of test data*

Examine the linearity of the four data samples from Problem 6.8. These four samples have the same values of statistical characteristics, but only sample A may be considered as an acceptable straight line.

Data: Problem 6.8

Solution: To test the linearity of the data, the linear model $E(y/x) = \beta_1 x + \beta_2$ is compared with the quadratic model $E(y/x) = \beta_1 x + \beta_2 x^2 + \beta_3$. The table on p. 43 lists the statistical characteristics RSC , MEP , \hat{R}_p^2 , AIC , T_2 for a test of $H_0: \beta_2 = 0$

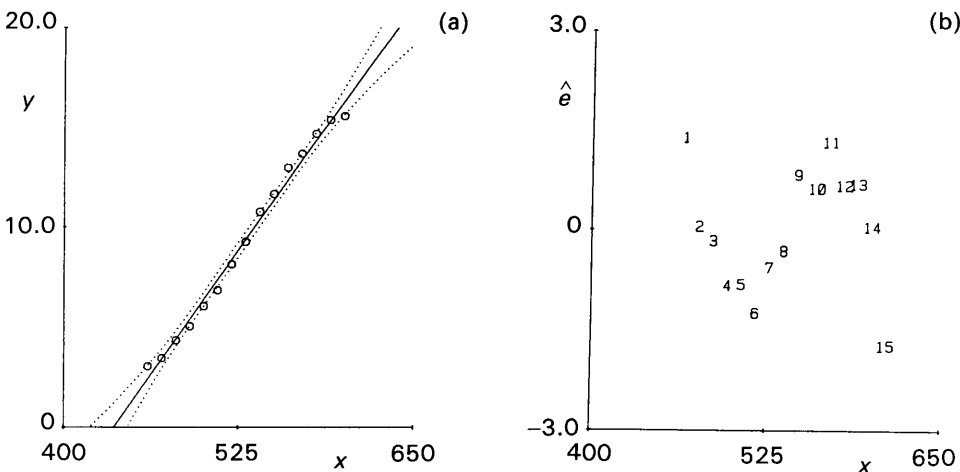


Fig. 6.21—(a) Curve fitting of data from Problem 6.9 by a polynomial of the second degree, and (b) graphical examination of residuals.

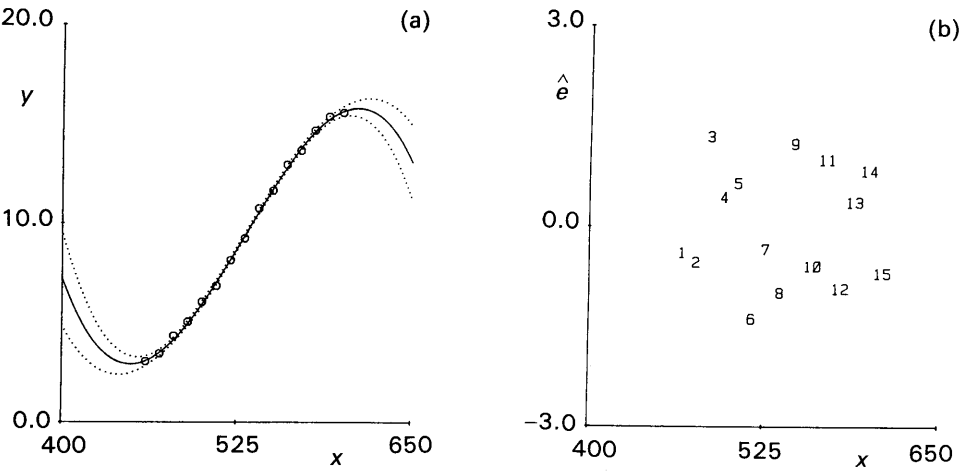


Fig. 6.22—(a) Curve fitting of data from Problem 6.9 by a polynomial of the third degree, and (b) graphical examination of residuals.

in the quadratic model and F_L . Nonlinearity is clearly indicated by F_L and by AIC . MEP and \hat{R}_p^2 show that there is an improvement of fit with the use of quadratic model for sample B, which exhibits a nonlinear curve.

	Proposed linear model				Proposed quadratic model				Tests	
	RSC	MEP	\hat{R}_p^2	AIC	RSC	MEP	\hat{R}_p^2	AIC	T_2	F_L
A	13.76	1.871	0.708	6.46	12.9	1.955	0.976	31.74	0.72*	0.53*
B	13.76	2.204	0.642	6.46	2.23×10^{-5}	3.11×10^{-6}	1	-114	2219	4.9×10^1
C	13.76	2.147	0.653	6.46	13.0	3.107	0.961	31.82	0.68*	0.467*
D	13.76	9×10^{99}	0.767	6.46	13.8	9×10^{99}	-	32.43	4	0*

Conclusion: For examination of linearity, the F_L test criterion for comparing the linear with the quadratic model seems to be the most reliable. Other statistical characteristics, i.e. AIC , MEP and \hat{R}_p^2 , can be used, but the rejection of the nonlinearity assumption does not lead to automatic acceptance of the linear model (see samples C and D).

There are many criteria for examination of linearity in regression models. Suitability of a proposed model can be easily checked if information about the measurement variance or the errors variance, σ^2 , is available. When a proposed model is correct, RSC will be approximately equal to $(n - m)\sigma^2$. When the model is incorrect, the residual sum of squares RSC can be decomposed into the sum of squares SSE corresponding to "pure" errors and the sum of squares SSL corresponding to the poor choice of model, the so-called lack of fit:

$$RSC = SSE + SSL \quad (6.59a)$$

To estimate a variance σ^2 that is independent of a proposed linear regression model, the method of repeated measurements can be used. At M different values of vector \mathbf{x}_i , $i = 1, \dots, M$, there are always n_i repeated measurements y_{ij} , $i = 1, \dots, M$, $j = 1, \dots, n_i$. The regression model is then expressed by

$$y_{ij} = \sum_{k=1}^m \beta_k x_{ik} + \varepsilon_{ij}, \quad j = 1, \dots, n_i \quad (6.60)$$

This model represents a linear regression model for $n = \sum_{i=1}^M n_i$ measurements. The matrix \mathbf{X} has dimension $(n \times m)$ and the vector \mathbf{y} has dimension $(n \times 1)$. On substitution into Eq. (6.11) the estimate \mathbf{b} is obtained and the residual sum of squares is:

$$RSC = \mathbf{y}^T \mathbf{y} - \mathbf{b}^T \mathbf{X} \mathbf{y} = \sum_{i=1}^n \hat{\varepsilon}_i^2 \quad (6.61)$$

An independent estimate of the sum of squares due to pure errors SSE is calculated from

$$SSE = \sum_{i=1}^M \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2 \quad (6.62)$$

where $\bar{y}_i = [\sum_{j=1}^{n_i} y_{ij}] / n_i$ is the arithmetic mean of repeated y_{ij} values for a given \mathbf{x}_i . The sum of squares corresponding to the lack of fit is $SSL = RSC - SSE$. The independent estimate of variance is here $SSE / (n - M)$ and the mean value of SSL is $SSL / (M - m)$. The test criterion

$$F_N = \frac{(RSC - SSE)(n - M)}{SSE(M - m)} \quad (6.63)$$

has, for a correct model, the Fisher-Snedecor F -distribution with $(M - m)$ and $(n - M)$ degrees of freedom. When $F_N < F_{1-\alpha}(M - m, n - M)$, the proposed model is correct and the sum of squares SSL is not significantly different from zero. The

quantity $SSE/(n - M)$ represents an unbiased estimate of the variance σ^2 , whether the proposed model is correct or not.

Problem 6.18. *Examination of a kinetic model of recombination of the Bromocresol Green anion*

The kinetic constant has been determined for the recombination reaction of the anion of Bromocresol Green (BCG) with a proton in a solution of glycerol and water [13]. For different concentrations x of BCG the reciprocal values of the relaxation times y_{ij} , $i = 1, \dots, n$ were measured. The kinetic model

$$y_{ij} = -k_D + k_R \times x_i + \varepsilon_{ij}$$

is proposed, where k_D is the kinetic constant of recombination, and k_R is the kinetic constant of dissociation. Determine the two kinetic constants and examine the proposed model.

Data: $n = 24$, $n_i = 2$, $i = 1, \dots, M$, $M = 12$, $m = 2$

x_i [10^{-6} mol dm $^{-3}$]	7.98	8.96	10.37	12.08	16.81	24.22	29.5
y_{i1} [10^6 sec $^{-1}$]	0.44	0.36	0.37	0.43	0.79	0.8	1.04
y_{i2} [10^6 sec $^{-1}$]	0.35	0.40	0.46	0.51	0.59	0.9	1.04

36.75	37.69	65.32	87.32	145.5
1.04	1.22	2.06	2.32	3.70
0.94	1.27	1.73	2.08	3.66

Solution: With the use of the least-squares method, the estimates $\hat{k}_D = 0.232 (\pm 0.033)$ and $\hat{k}_R = 0.0238 (\pm 0.0006)$ were obtained (standard deviations in brackets). The determination coefficient $\hat{R}^2 = 0.987$ and the test criterion $F_R = 1616$ (Eq. (6.39)) prove that the linear model is valid.

The residual sum of squares $RSC = 0.285$, the residual standard deviation $\hat{\sigma} = 0.1139$ and an independent estimate of sum of squares corresponding to pure errors SSE (6.62) is for $M = 12$ and $n_i = 2$ equal to $SSE = 0.1274$. The independent estimate of the residual standard deviation is then $\hat{\sigma}_I = \sqrt{0.1274/(24 - 12)} = 0.103$. On substituting into Eq. (6.63), we find the test criterion is

$$F_N = \frac{(0.285 - 0.1274)(24 - 12)}{0.1274(12 - 2)} = 1.484$$

As the quantile $F_{0.95}(10, 12) = 2.75$ is greater than F_N , the proposed kinetic model is correct and the residual sum of squares corresponding to lack of fit,

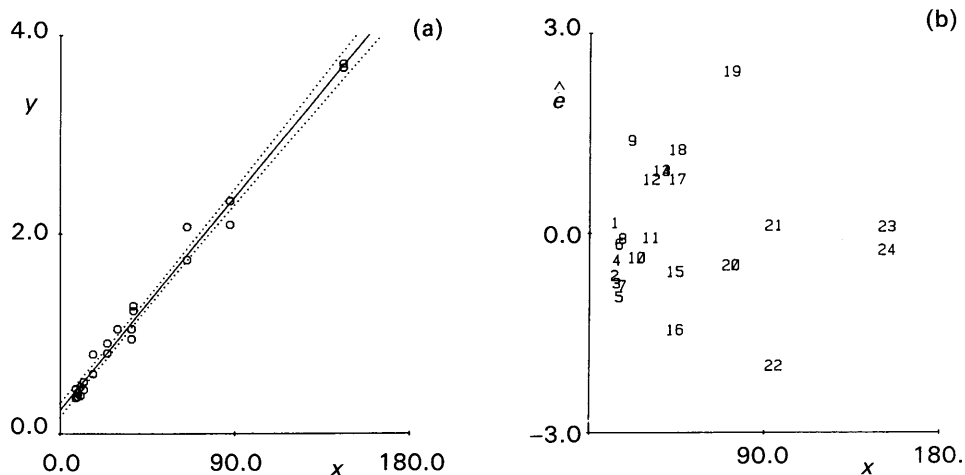


Fig. 6.23—(a) Degree of fit of straight line for proposed kinetic model, and (b) graphical examination of residuals.

$SSL = 0.285 - 0.1274 = 0.1575$ may be considered as insignificant. Figure 6.23 shows the regression model with the 95% confidence interval.

Conclusion: When replicate results are available, it is relatively easy to test the validity of a regression model.

When replicate measurements for all i are not available, the data may be divided into groups of approximately the same x values [14].

6.3.3 Comparison of regression lines

Often in chemometrics we need to compare M proposed regression models

$$y_{ij} = \beta_{2j} + \beta_{1j}x_{ij} + \varepsilon_{ij}, \quad \begin{matrix} j = 1, \dots, M \\ i = 1, \dots, n_j \end{matrix} \quad (6.64)$$

for M groups of experimental data $[(x_{ij}, y_{ij}), i = 1, \dots, n_j], j = 1, \dots, M$. Typical examples are Lambert–Beer law calibration lines, i.e. the dependence of absorbance on concentration at M different wavelengths. We want to know:

- if the regression lines have the same intercept;
- if the regression lines have the same slope;
- if the regression lines are identical.

The first step of the statistical analysis is always estimation of the parameters b_{2j} , b_{1j} and $\hat{\sigma}_j^2$ for each set of data, individually, by the least-squares method.

The second step involves examination of homoscedasticity, i.e. constancy of variance $\hat{\sigma}_j^2$, because testing hypotheses (a), (b) and (c) requires constant and identical variance in all groups.

The Bartlett test for homoscedasticity is a commonly used test. In this test, we compare M independent variance estimates $\hat{\sigma}_j^2$, $j = 1, \dots, M$, with v_j degrees of freedom. The null hypothesis $H_0: \sigma_j^2 = \sigma^2, j = 1, \dots, M$, is tested. For models of a regression straight line, the degrees of freedom are $v_j = n_j - 2$. We define:

$$V = \sum_{j=1}^M v_j \quad (6.65a)$$

$$\hat{\sigma}_C^2 = \frac{\sum_{j=1}^M v_j \times \hat{\sigma}_j^2}{V} \quad (6.65b)$$

and

$$L = 1 + \frac{\sum_{j=1}^M v_j^{-1} - V^{-1}}{3M - 3} \quad (6.66)$$

The test criterion of the Bartlett test of homoscedasticity is given by

$$B = \frac{V \times \ln \hat{\sigma}_C^2 - \sum_{j=1}^M v_j \times \ln \hat{\sigma}_j^2}{L} \quad (6.67)$$

which, if the null hypothesis is valid, has the χ^2 -distribution with $(M - 1)$ degrees of freedom. If $B < \chi_{1-\alpha}^2(M - 1)$ [where $\chi_{1-\alpha}^2(M - 1)$ is the $100(1 - \alpha)\%$ quantile of the χ^2 distribution], the null hypothesis is accepted and the estimate of constant variance σ^2 is called the **pooled variance** $\hat{\sigma}_C^2$ (Eq. 6.65b). The Bartlett test, is, however, sensitive to deviations of the residuals from normality.

To compare two groups of points, $M = 2$, the identity of two variances $H_0: \sigma_1^2 = \sigma_2^2$ may be tested by the test criterion

$$F_2 = \frac{\max(\hat{\sigma}_1^2, \hat{\sigma}_2^2)}{\min(\hat{\sigma}_1^2, \hat{\sigma}_2^2)} \quad (6.68)$$

which, if the null hypothesis is valid, has the Fisher–Snedecor F -distribution with $(n_1 - 2)$ and $(n_2 - 2)$ degrees of freedom when $\hat{\sigma}_1^2 > \hat{\sigma}_2^2$. Generally, the degrees of freedom used in calculation of $\hat{\sigma}_i^2$, $i = 1, 2$ are also used here.

6.3.3.1 Test for homogeneity of intercepts

When the null hypothesis $H_0: \beta_{21} = \beta_{22} = \dots = \beta_{2j} = \dots = \beta_{2M} = \beta_{2C}$ is valid, the pooled estimate of the overall intercept β_{2C} as a weighted combination of the estimates of the individual intercepts b_{2j} may be obtained from

$$b_{2C} = \frac{\sum_{j=1}^M w_{Bj} \times b_{2j}}{\sum_{j=1}^M w_{Bj}} \quad (6.69)$$

where the j th weight coefficient w_{Bj} corresponding to the estimate of the j th straight line is given by

$$w_{Bj} = \frac{n_j \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2}{\sum_{i=1}^{n_j} x_{ij}^2} \quad (6.70)$$

For testing the estimate of errors, variance σ^2 is calculated from the variance of individual parameter estimates b_{2j} around their weighted average b_{2c} and from a combination of the variability of all points around the regression line inside the individual data groups. The test criterion is

$$F_1 = \frac{\sum_{j=1}^M w_{B,j}(b_{2j} - b_{2c})^2 / (M - 1)}{\sum_{j=1}^M \sum_{i=1}^{n_j} \hat{e}_{ij}^2 / (n - 2M)} \quad (6.71)$$

where $n = \sum_{j=1}^M n_j$. When the null hypothesis H_0 is valid the test criterion F_1 has the Fisher-Snedecor distribution with $(M - 1)$ and $(n - 2M)$ degrees of freedom. The residuals \hat{e}_{ij} are calculated for the individual regression lines. We can write

$$\sum_{j=1}^M \sum_{i=1}^{n_j} \hat{e}_{ij}^2 = \sum_{j=1}^M RSC_j$$

where RSC_j is the residual sum of squares for the j th group.

When $F_1 < F_{1-\alpha}(M - 1, n - 2M)$, then all straight lines have, at significance level α , the same intercept; and its estimate is given by Eq. (6.69). The variance of this intercept is calculated from

$$D(b_{2c}) = \frac{\hat{\sigma}^2}{\sum_{j=1}^M w_{Bj}} = \frac{\sum_{j=1}^M \sum_{i=1}^{n_j} \frac{\hat{e}_{ij}^2}{n - 2M}}{\sum_{j=1}^M w_{Bj}} \quad (6.72)$$

The intercept estimate has an asymptotically normal distribution and represents an unbiased estimate of parameter β_{2c} .

6.3.3.2 Test for homogeneity of slopes

The test of homogeneity of slopes is known as a test of parallelism of regression straight lines. If the null hypothesis $H_0: \beta_{11} = \beta_{12} = \dots = \beta_{1j} = \dots = \beta_{1M} = \beta_{1c}$ is valid, the pooled estimate of overall (common) slope β_{1c} as a weighted combination of individual slope estimates b_{1j} may be calculated from

$$b_{1c} = \frac{\sum_{j=1}^M w_{Sj} b_{1j}}{\sum_{j=1}^M w_{Sj}} \quad (6.73)$$

where

$$W_{sj} = \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2 \quad (6.74)$$

As in the test for homogeneity of intercepts, a test criterion may be derived

$$F_s = \frac{\sum_{j=1}^M w_{sj}(b_{1j} - b_{1c})^2 / (M - 1)}{\sum_{j=1}^M \sum_{i=1}^{n_j} \hat{e}_{ij}^2 / (n - 2M)} \quad (6.75)$$

which, when the null hypothesis is valid, has the Fisher–Snedecor F -distribution with $(M - 1)$ and $(n - 2M)$ degrees of freedom. When $F_s < F_{1-\alpha}(M - 1, n - 2M)$, all the regression straight lines can, at significance level α , be considered as parallel. The best estimate of overall slope is b_{1c} , from Eq. (6.73), and its variance estimate may be calculated from

$$D(b_{1c}) = \frac{\sum_{j=1}^M \sum_{i=1}^{n_j} \hat{e}_{ij}^2 / (n - 2M)}{\sum_{j=1}^M w_{sj}} \quad (6.76)$$

When the null hypothesis H_0 is valid, the slope estimate b_{1c} has a normal distribution and represents an unbiased estimate of parameter β_{1c} .

6.3.3.3 Test for coincidence of regression lines

The test for coincidence of regression lines $H_0: \beta_{2j} = \beta_{2c}, \beta_{1j}, \beta_{1c}, j = 1, \dots, M$ is a combination of the two previous tests F_1 and F_s (Sections 6.3.3.1 and 6.3.3.2). The test compares two residual sums of squares, RSC_K with RSC_C . RSC_K was obtained after fitting all M groups of data by a single common straight line with estimates b_{1K} and b_{2K} and RSC_C is calculated from the individual groups of data separately, $RSC_C = \sum_{j=1}^M RSC_j$. The test criterion is

$$F_A = \frac{(RSC_K - RSC_C) / (2M - 2)}{RSC_C / (n - 2M)} \quad (6.77)$$

When the null hypothesis H_0 is valid, the test criterion F_A has the Fisher–Snedecor F -distribution with $(2M - 2)$ and $(n - 2M)$ degrees of freedom. When $F_A < F_{1-\alpha}[(2M - 2), (n - 2M)]$, then all regression straight lines may be considered as identical, with slope b_{1K} and intercept b_{2K} . Individual groups of data are then collected into a single common sample of size n . When the null hypothesis is not accepted, it is usually possible to find subgroups of data which are homogeneous enough.

Problem 6.19. Comparison of three methods for determination of gold

The classical method for determination of gold in jewellery alloys is rather tedious and time-consuming. Three variants of the X-ray fluorescence (XRF) method for gold have been proposed. Ten samples were measured by the classical method (values x)

and the new methods (values y_1 , y_2 and y_3). Test whether the results from the three variants differ from one another and if these results are in agreement with those measured by the classical method.

Data: $M = 4$, $n_j = 10$

i	x	y_1	y_2	y_3
1	59.0	63.8	72.2	64.3
2	59.0	62.6	71.1	63.1
3	59.0	64.0	69.9	64.3
4	59.3	58.9	63.5	59.2
5	62.7	60.2	64.8	60.5
6	64.8	62.9	70.8	63.4
7	72.3	70.0	68.7	70.0
8	75.2	77.4	76.7	77.4
9	75.2	78.7	85.2	79.1
10	75.2	80.7	78.5	80.6

Solution: The parameter estimates b_{2j} , b_{1j} , RSC_j and the residual standard deviations for all three variants of determination of gold are listed in Table 6.4. All three variants have the same sample size $n_j = 10$ and have a common x -axis. Therefore the weights w_{Bj} and w_{Sj} are independent of the j value, so that:

$$w_{Bj} = \frac{10 \sum_{i=1}^{10} (x_i - \bar{x})^2}{\sum_{i=1}^{10} x_i^2} = 0.1124$$

and

$$w_{Sj} = \sum_{i=1}^{10} (x_i - \bar{x})^2 = 497.59$$

Before testing the agreement of residuals, the variances must be examined. On introducing numbers into Eqs. (6.64)–(6.67), we find $V = 24$, $\hat{\sigma}_c^2 = 15.804$, $L = 1.0556$ and $B = 4.9087$. As the quantile $\chi_{0.95}^2(2) = 5.99$ is greater than B , the null hypothesis about homoscedasticity is accepted. The variance for the second X-ray fluorescence method is significantly greater than for the first and third one.

Table 6.4. Regression analysis of gold determination by three X-ray fluorescence methods.

Variant j	b_{2j}	b_{1j}	RSC_j	$\hat{\sigma}_j$
1	1.087	1.01	92.98	3.409
2	31.55	0.61	192.8	4.91
3	2.74	0.989	93.44	3.418
1 + 2 + 3	11.79	0.871	540.4	4.393

To test the homogeneity of intercepts, the pooled estimate of the intercept in Eq. (6.69), $b_{2c} = 11.79$, is calculated, and

$$\sum_{j=1}^3 RSC_j = 379.22$$

The test criterion F_1 is calculated from Eq. (6.71):

$$F_1 = \frac{65.9692/2}{379.22/24} = 2.088$$

Since the quantile $F_{0.95}(2, 24) = 3.4$ is greater than F_1 , the three intercepts can be considered to be identical and equal to the estimate $b_{2c} = 11.792$ with variance $D(b_{2c}) = 35.151$ (Eq. (6.72)).

If the three methods are not systematically biased, the null hypothesis $H_0: \beta_{2c} = 0$ should be valid. From Eq. (6.48), $T_2 = 1.989$ is smaller than the quantile $t_{0.975}(8) = 2.3$, so the null hypothesis H_0 is accepted, and all three variants do not lead to systematically biased results.

To test the homogeneity of the slopes, the pooled estimate of the slope is calculated (6.73), $b_{1c} = 0.870$, and from Eq. (6.75) the test criterion $F_S = 1.596$, which is a lower value than the quantile $F_{0.95}(2, 24) = 3.4$. Therefore all slopes can be considered to be identical and equal to $b_{1c} = 0.870$ with variance [Eq. (6.76)] $D(b_{1c}) = 0.0079$.

When all three variants of the XRF method (y_1, y_2, y_3) give the same results as the standard method (x), the null hypothesis $H_0: \beta_{1c} = 1$ should be valid. From Eq. (6.48), $T_1 = 1.459$, which is smaller than the quantile $t_{0.975}(8) = 2.3$, so the null hypothesis is accepted and all results y_1, y_2, y_3 are identical with x .

In the test for identity of the three regression lines by Eq. (6.77), the test criterion

$$F_A = \frac{(540.4 - 379.22)/4}{379.22/24} = 2.55$$

is smaller than the quantile $F_{0.95}(4, 24) = 2.776$, so all three XRF methods can be considered to be identical.

To test the homoscedasticity of the three variants, three repetitive measurements at $x = 59.0$ and $x = 75.2$ are used. Table 6.5 lists independent estimates of variance, calculated from these repetitive measurements.

Table 6.5. Variance estimates calculated from repeated measurements, and the F -test criterion [Eq. (6.65)]

Variant	$\hat{\sigma}_1^2$ ($x = 59.0$)	$\hat{\sigma}_2^2$ ($x = 75.2$)	$\hat{\sigma}_K^2$ (both x)	$F_2 = \frac{\hat{\sigma}_1^2}{\hat{\sigma}^2}$	$F_j = \frac{\hat{\sigma}_j^2}{\hat{\sigma}_K^2}$
1	0.573	2.763	1.668	4.882	6.97
2	1.323	20.06	10.692	15.163	1.59
3	0.480	2.564	1.522	5.341	7.67

The null hypothesis $H_0: \sigma_1^2 = \sigma_2^2$ is tested by the Fisher–Snedecor F -test, and Table 6.5 lists the F_2 values. The second null hypothesis $H_0: \sigma_j^2 = \sigma_K^2$ says that any j th variance is the same as the variance of all measurements. When F_2 values are

compared with the quantile $F_{0.95}(2, 2) = 19$, the null hypothesis H_0 (i.e. an assumption about homoscedasticity) is accepted. When F_j values are compared with the quantile $F_{0.95}(24, 4) = 5.77$, it may be concluded that the residual variances of the first and the third variants are higher than the variance of measurement. The conclusion is here affected by the fact that the replicate measurements were carried out only for two different levels of x .

Conclusion: Examination of three XRF variants proved that they do not differ significantly and lead to the same results as the standard method. The first variant seems to be the most precise one.

6.4 NUMERICAL PROBLEMS IN THE COMPUTER CALCULATION OF LINEAR REGRESSION

The determination of parameter estimates of a linear model [Eq. (6.11)] seems to be a simple task. When subprograms for matrix operations are available in a package of algorithms, the formal solution of Eq. (6.11) is quite easy. Some difficulties arise when the matrix $\mathbf{X}^T\mathbf{X}$ appears to be singular, from the point of view of the machine precision and the algorithm. In some cases, especially with polynomial models, the parameter estimates may be without physical meaning. The regression curve goes quite close to experimental points but oscillates among them (for polynomials of higher degree) or is systematically shifted.

The reasons for numerical difficulties in the computer evaluation of parameter estimates \mathbf{b} are as follows:

- (1) Neglect of the limited precision of computer in building the matrix $\mathbf{X}^T\mathbf{X}$.
- (2) Inconvenient procedures for matrix inversion or solving the set of linear equations.
- (3) Multicollinearity leading to the ill-conditioning of matrix $\mathbf{X}^T\mathbf{X}$.
- (4) Linear dependence of some columns of matrix $\mathbf{X}^T\mathbf{X}$, leading to its non-invertability because of a singularity.

Good linear-regression programs overcome these difficulties and always give correct solutions. Among the most effective programs are algorithms which do not build matrix $\mathbf{X}^T\mathbf{X}$ but instead solve the overdetermined set of n linear equations of m unknowns $\mathbf{y} = \mathbf{X} \cdot \mathbf{b}$. For example, the algorithm SVD (Singular value decomposition [16]) works even on computer with poor data precision.

Problem 6.20. Examination of the quality of a regression algorithm (LS)

Many test examples are available for examining the quality and effectiveness of linear regression algorithms. An example suitable for numerical control of quality of regression programs comes from the linear model $E(y/x) = \beta_1 x_1 + \beta_2 x_2$. Calculate the estimates b_1 and b_2 by the least-squares method.

Data: The numerical constant ε in algorithms examination is selected such that the condition $\varepsilon < 10^{-(d+1)/2}$ is fulfilled, where d is the number of valid digits used in the actual computer.

i	y	x_1	x_2
1	3	1	1
2	ε	ε	0
3	2ε	0	ε

Solution: (a) The analytical approach:

$$\mathbf{X}^T\mathbf{X} = \begin{bmatrix} 1 + \varepsilon^2 & 1 \\ 1 & 1 + \varepsilon^2 \end{bmatrix}$$

$$\mathbf{X}^T\mathbf{y} = \begin{bmatrix} 3 + \varepsilon^2 \\ 3 + 2\varepsilon^2 \end{bmatrix}$$

According to Problem 6.1 the inversion matrix is determined as

$$(\mathbf{X}^T\mathbf{X})^{-1} = \frac{1}{(1 + \varepsilon^2)^2 - 1} \begin{bmatrix} 1 + \varepsilon^2 & -1 \\ -1 & 1 + \varepsilon^2 \end{bmatrix}$$

and on substituting into Eq. (6.11), the estimates of the parameters are found to be

$$b_1 = \frac{(3 + \varepsilon^2)(1 + \varepsilon^2)}{(1 + \varepsilon^2)^2 - 1} - \frac{3 + 2\varepsilon^2}{(1 + \varepsilon^2)^2 - 1} = 1$$

$$b_2 = \frac{-(3 + \varepsilon^2)}{(1 + \varepsilon^2)^2 - 1} + \frac{(1 + \varepsilon^2)(3 + 2\varepsilon^2)}{(1 + \varepsilon^2)^2 - 1} = 2$$

The estimates b_1 and b_2 do not depend on the magnitude of ε , and moreover $RSC = 0$.

(b) *The numerical approach, by computer:* If the condition $\varepsilon < 10^{-(d+1)/2}$ is valid, then $1 + \varepsilon^2 = 1$. All elements of the matrix $\mathbf{X}^T\mathbf{X}$ will be ones, and its inversion will be impossible because $\det(\mathbf{X}^T\mathbf{X}) = 0$. If the computer works with a precision of 11 digits, the choice of $\varepsilon = 10^{-6}$ will cause the computation to fail.

Conclusion: Because of limited precision of computers and possible ill-conditioning of the normal equations $\mathbf{X}^T\mathbf{X}$, even simple tasks may cause numerical difficulties.

To make statistical analysis easier, most programs work with matrix $\mathbf{X}^T\mathbf{X}$. To avoid difficulties with large differences its centred or normalized version is used. The variables are expressed as deviations from the arithmetic mean:

$$x_{Cij} = x_{ij} - \bar{x}_j$$

and

$$y_{Ci} = y_i - \bar{y}.$$

The resulting centred variables then form the elements of the matrix \mathbf{X}_C or the vector \mathbf{y}_C . Centring all variables results in the intercept term cancelling out.

The matrix $(\mathbf{X}_C^T \mathbf{X}_C)^{-1}$ is a submatrix of $(\mathbf{X}^T \mathbf{X})^{-1}$, so that after its inversion all elements of matrix $(\mathbf{X}^T \mathbf{X})^{-1}$ may be calculated. These elements are necessary for the statistical analysis. Then

$$(\mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} c_0 & \mathbf{c}^T \\ \mathbf{c} & (\mathbf{X}_C^T \mathbf{X}_C)^{-1} \end{bmatrix}$$

If we define the averages vector $\mathbf{x}_p = (\bar{x}_1, \dots, \bar{x}_{m-1})^T$, we may write

$$c_0 = \frac{1}{n} - \mathbf{c}^T \mathbf{x}_p$$

and

$$\mathbf{c}^T = (-1) \mathbf{x}_p^T (\mathbf{X}_C^T \mathbf{X}_C)^{-1}$$

With normalized variables, the standard deviations $\hat{\sigma}(x_j)$ and $\hat{\sigma}(y)$ are used. If we introduce normalized variables

$$Z_{ij} = \frac{x_{Cij}}{\hat{\sigma}(x_j) \sqrt{n-1}}$$

$$q_i = \frac{y_{Ci}}{\hat{\sigma}(y) \sqrt{n-1}}$$

the matrix $\mathbf{R} = \mathbf{Z}^T \mathbf{Z}$ is formally identical with the correlation matrix of controllable variables and the vector $\mathbf{r} = \mathbf{Z}^T \mathbf{q}$ formally contains the correlation coefficients of all controllable variables with the response variable. In correlation models there are real correlation coefficients. In the least-squares method, the parameter estimates $\mathbf{b}_N = \mathbf{R}^{-1} \mathbf{r}$ are found, for which

$$b_{Nj} = b_j \frac{\hat{\sigma}(x_j)}{\hat{\sigma}(y)}$$

The advantage of normalized variables is that the elements of matrix \mathbf{R} are the numbers in the interval from -1 to $+1$. A disadvantage is the possible distortion of the calculated matrix \mathbf{R} , e.g. by use of the "for pocket-calculator" modified expression:

$$\sum_{i=1}^n x_i^2 - n\bar{x}^2 \quad \text{instead of} \quad \sum_{i=1}^n (x_i - \bar{x})^2.$$

With the modified expression, the final result is the difference between two large numbers, with a result close to zero. The limited precision of computer or calculator can result in a sum of squares of residuals from the mean with a value that is zero or even negative.

From many techniques of numerical solution of the least-squares problem, we select here the following two cases:

(a) the *method of orthogonal functions*, which is simple and convenient for polynomial

models; and

- (b) the *method of rational ranks*, used in the program CHEMSTAT. Another algorithm is described by Lawson and Hanson [18].

6.4.1 The method of orthogonal functions

Orthogonal functions are frequently used because they result in considerable simplification of the statistical analysis. For the linear regression model

$$E(y/\mathbf{x}) = \sum_{j=1}^m \beta_j f_j(\mathbf{x}) \quad (6.78)$$

where $f_j(\mathbf{x})$ are any functions of the input variables \mathbf{x} which do not contain regression parameters. Equation (6.11) is used to estimate the parameters β . For this case, Eq. (6.11) contains the matrix \mathbf{F} of dimension $(n \times m)$ with elements $f_j(x_i)$, $j = 1, \dots, m$, $i = 1, \dots, n$, instead of matrix \mathbf{X} . For further analysis it is convenient for the matrix $\mathbf{F}^T \mathbf{F}$ to be diagonal. Therefore the scalar products of all pairs of columns of \mathbf{F} must be equal to zero. For $\mathbf{F}^T \mathbf{F}$ to be diagonal,

$$\sum_{i=1}^n f_j(\mathbf{x}_i) \times f_k(\mathbf{x}_i) \begin{cases} \rightarrow 0 & \text{for } k \neq j \\ \rightarrow \sum_{i=1}^n f_j^2(\mathbf{x}_i) & \text{for } k = j \end{cases} \quad (6.79)$$

From Eq. (6.79) it follows that the diagonality of $\mathbf{F}^T \mathbf{F}$ may be achieved

(a) by the adjustment of values \mathbf{x}_i for a given function f_j , $i = 1, \dots, n$. This is the case in designed experiments;

(b) by the special choice of functions f_j for given locations \mathbf{x}_i , $i = 1, \dots, n$. That is, the construction of orthogonal functions $g_j(\mathbf{x})$ from the original ones $f_j(\mathbf{x}_i)$.

Orthogonal functions are generated by use of the recurrent relation

$$g_j(\mathbf{x}) = f_j(\mathbf{x}) + \sum_{L=j-1}^1 Q_{jL} g_L(\mathbf{x}) \quad (6.80)$$

The coefficients Q_{jL} may be found from the conditions of orthogonality. A set of j linear equations is formed, and each equation contains just one unknown Q_{jL} for which:

$$Q_{jL} = \frac{- \sum_{i=1}^n f_j(\mathbf{x}_i) \times g_L(\mathbf{x}_i)}{\sum_{i=1}^n g_L^2(\mathbf{x}_i)} \quad (6.81)$$

With the orthogonal functions $g_j(\mathbf{x})$ generated from the original ones $f_j(\mathbf{x})$ by Eqs. (6.80) and (6.81), the linear regression model may be expressed in the form

$$E(y/\mathbf{x}) = \sum_{j=1}^m c_j g_j(\mathbf{x}) \quad (6.82)$$

Since $g_j(\mathbf{x})$ are orthogonal, the estimates of parameters c_j may be obtained by straight substitution into

$$c_j = \frac{\sum_{i=1}^n g_j(\mathbf{x}_i) \times y_i}{\sum_{i=1}^n g_j^2(\mathbf{x}_i)} \quad j = 1, \dots, m \quad (6.83)$$

The variances of the parameter estimates are given by

$$D(c_j) = \frac{\sigma^2}{\sum_{i=1}^n g_j^2(\mathbf{x}_i)} \quad (6.84)$$

Thus, when the orthogonal functions are known, the data evaluation for a linear regression requires only substitution into simple expressions.

Problem 6.21. *Examination of the quality of the orthogonal functions method*

Estimate parameters β_1 and β_2 in Problem 6.20 by the method of orthogonal functions.

Data: from Problem 6.20

Solution: The original functions are $f_1(x) = x_1$, $f_2(x) = x_2$. On substituting into Eq. (6.80) we get

$$g_1(x) = x_1$$

$$g_2(x) = x_2 + Q_{21}$$

$$g_1(x) = x_2 + Q_{21}x_1$$

$$g_2(x) = x_2 + Q_{21}g_1(x) = x_2 + Q_{21}x_1$$

For functions $g_1(x)$ and $g_2(x)$ to be orthogonal for the given values x_i , $i = 1, \dots, n$, their scalar product must be equal to zero, that is

$$\sum_{i=1}^n g_1(x_i)g_2(x_i) = \sum_{i=1}^n x_{1i}(x_{2i} + Q_{21}x_{1i}) = 0$$

From this equation, the term Q_{21} is given by

$$Q_{21} = \frac{-\sum_{i=1}^n x_{1i}x_{2i}}{\sum_{i=1}^n x_{1i}^2}$$

For the data in the problem

$$\sum_{i=1}^n x_{1i}x_{2i} = 1 + 0 + 0 = 1$$

$$\sum_{i=1}^n x_{1i}^2 = 1 + \varepsilon^2$$

and therefore

$$Q_{21} = \frac{-1}{1 + \varepsilon^2}$$

Then, on substituting into Eq. (6.83), we have

$$c_1 = \frac{\sum_{i=1}^n x_{1i} y_i}{\sum_{i=1}^n x_{1i}^2} = \frac{3 + \varepsilon^2}{1 + \varepsilon^2}$$

and

$$c_2 = \frac{\sum_{i=1}^n [x_{2i} - x_{1i}/(1 + \varepsilon^2)] \times y_i}{\sum_{i=1}^n [x_{2i} - x_{1i}/(1 + \varepsilon^2)]^2} = \frac{4 + 6\varepsilon^2 + 2\varepsilon^4}{2 + 3\varepsilon^2 + \varepsilon^4} = 2$$

The regression model of type (6.82) then has the form

$$E(y/x) = \frac{3 + \varepsilon^2}{1 + \varepsilon^2} g_1(x) + 2g_2(x)$$

After substituting for $g_1(x)$ and $g_2(x)$, we obtain

$$E(y/x) = \frac{3 + \varepsilon^2}{1 + \varepsilon^2} x_1 + 2 \left(x_2 - \frac{x_1}{1 + \varepsilon^2} \right) = 1x_1 + 2x_2$$

Thus, the estimates of the parameters are $b_1 = 1$ and $b_2 = 2$.

Conclusion: The method of orthogonal functions can find estimates of parameters for linear models. The quality of the parameter estimates achieved by computer is determined by the precision of calculation of the individual orthogonal functions.

This example demonstrates that the application of orthogonal functions is quite simple. The errors caused by limited computer precision accumulate according to Eq. (6.80) and increase with the number of orthogonal functions m . The use of orthogonal functions is nearly equivalent to the use of adequate methods for the matrix inversion.

The advantage of orthogonal functions is that when some functions $g_e(x)$ are omitted, the coefficients c_j for the remaining functions $g_j(x)$ will be unchanged. The method may be used to search for the optimum combination of polynomial terms (degrees of polynomial model). The disadvantage of orthogonal functions is the rather complicated manipulation. It is useful to know that all types of orthogonal polynomials may be expressed by three-term recurrent expressions [19].

6.4.2 The method of rational ranks

To detect ill-conditioning of $\mathbf{X}^T \mathbf{X}$ or the \mathbf{R} , the matrices are decomposed into eigenvalues and eigenvectors. Since the matrix \mathbf{R} is symmetrical it may be expressed by eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_m$, and corresponding eigenvectors $\mathbf{P}_j, j = 1, \dots, m$, in the form of the sum

$$\mathbf{R} = \sum_{j=1}^m \lambda_j \mathbf{P}_j \mathbf{P}_j^T \quad (6.85)$$

The inverse matrix \mathbf{R}^{-1} may be expressed in the form

$$\mathbf{R}^{-1} = \sum_{j=1}^m \lambda_j^{-1} \mathbf{P}_j \mathbf{P}_j^T \quad (6.86)$$

With the use of Eq. (6.86) Eq. (6.11) can be rewritten in the form

$$\mathbf{b}_N = \sum_{j=\omega}^m [\lambda_j^{-1} \mathbf{P}_j \mathbf{P}_j^T] \mathbf{r} \quad (6.87)$$

The covariance matrix of normalized estimates \mathbf{b}_N may be rewritten in form

$$D(\mathbf{b}_N) = \hat{\sigma}_N^2 \sum_{j=\omega}^m \lambda_j^{-1} \mathbf{P}_j \mathbf{P}_j^T \quad (6.88)$$

In the case of the least-squares in Eqs. (6.86) and (6.88) the parameter ω is set to $\omega = 1$. From both equations it follows that when the eigenvalues λ_j are small, the estimates \mathbf{b}_N and their variances are rather high. According to the magnitude of the eigenvalues λ_j , regression problems can be divided into three groups:

- (1) All eigenvalues are significantly higher than zero. The use of the least-squares method does not cause any problems.
- (2) Some eigenvalues are close to zero. This is a typical example of multicollinearity when some common methods fail.
- (3) Some eigenvalues are equal to zero. Then the matrix $\mathbf{X}^T \mathbf{X}$ or \mathbf{R} is singular and cannot be inverted.

The only way of avoiding difficulties with groups (2) and (3) is the use of the method of rational ranks. Here, the terms (or parts of them) with small values of eigenvalues λ_j are neglected [20]. The criterion for omitting terms corresponding to small eigenvalues has the form

$$\text{abs} \left[\frac{\sum_{j=1}^{\omega} \lambda_j}{\sum_{j=1}^m \lambda_j} \right] = P \quad (6.89)$$

where P is the chosen precision (usually 10^{-5}). The value ω determines the lower limit from which, in Eqs. (6.87) and (6.88), the summation is carried out.

Let us define

$$W = \sum_{j=1}^{\omega} \lambda_j$$

and

$$E = \sum_{j=1}^m \lambda_j$$

When the condition

$$\frac{W}{E} > P$$

is valid i.e. the value ω is not an integer, the summation is made from $\omega - 1$ and the eigenvalue $\lambda_{\omega-1}$ is “weighted” by the factor

$$u = \frac{W - EP}{\lambda_{\omega}} \quad (6.90)$$

Therefore, the length of estimates $\|\mathbf{b}_N\|$ with their variances may be continuously decreased as a function of increasing precision P . However, it is followed by an increase of the estimate bias and a decrease in the multiple correlation coefficient. The bias of estimates is here caused by neglecting terms in Eqs. (6.87) and (6.88) at $\omega > 1$.

It has been proposed [20] that the squared bias

$$h_V^2(\mathbf{b}_N) = [\boldsymbol{\beta} - E(\mathbf{b})]^2$$

achieved by the method of rational ranks is equal to

$$h_V^2(\mathbf{b}_N) = \boldsymbol{\beta}_N^T \left[\sum_{j=1}^{\omega} \mathbf{P}_j \mathbf{P}_j^T \right] \boldsymbol{\beta}_N \quad (6.91)$$

The optimum magnitude of P may be determined by finding a minimum of the mean quadratic error of prediction MEP , Eq. (6.57). In program VLR the user chooses the value of precision P , or it takes the default value $P = 10^{-32}$.

Problem 6.22. Examination of the method of rational ranks

Determine the parameter estimates b_1 and b_2 for the data from problem 6.20 by the method of rational ranks, with $P = 5 \times 10^{-9}$ and $\varepsilon = 10^{-4}$.

Data: from Problem 6.20

Solution: By decomposition of matrix $\mathbf{X}^T \mathbf{X}$ into eigenvalues and eigenvectors it is found that

$$\lambda_1 = \varepsilon^2$$

and

$$\lambda_2 = 2 + \varepsilon^2$$

and that

$$\mathbf{P}_1^T = (-\sqrt{0.5}, \sqrt{0.5})$$

$$\mathbf{P}_2^T = (\sqrt{0.5}, -\sqrt{0.5})$$

Because

$$\frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{10^{-8}}{2 + 2 \times 10^{-8}} = 5.0 \times 10^{-9}$$

$\omega = 1$ and in Eq. (6.87) only the second term is used. Because the decomposition of the matrix $\mathbf{X}^T\mathbf{X}$ (and not of \mathbf{R}) was made the parameter estimates from Eq. (6.87) are calculated

$$\mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = \frac{1}{2 + \varepsilon^2} \begin{bmatrix} \sqrt{0.5} \\ \sqrt{0.5} \end{bmatrix} [\sqrt{0.5} \quad \sqrt{0.5}] \begin{bmatrix} 3 + \varepsilon^2 \\ 3 + 2\varepsilon^2 \end{bmatrix}$$

$$= \begin{bmatrix} 1.5(2 + \varepsilon^2)/(2 + \varepsilon^2) \\ 1.5(2 + \varepsilon^2)/(2 + \varepsilon^2) \end{bmatrix} = \begin{bmatrix} 1.5 \\ 1.5 \end{bmatrix}$$

The fact was used here that $\mathbf{X}^T\mathbf{y} = [3 + \varepsilon^2 \quad 3 + 2\varepsilon^2]^T$. The parameter estimates found differ from the true values $\beta_1 = 1$ and $\beta_2 = 2$, but the residual sum of squares $RSC = 0.5\varepsilon^2 = 5 \times 10^{-9}$ is rather small. By using the classical least-squares method, and a computer working with precision of 7 valid digits the problem is not at all solved. It is evident that, even with smaller values of ε , these biased estimates remain the same, and “numerical underflowing” may occur when ε^2 is smaller than the computer precision.

Conclusion: The method of rational ranks enables biased parameter estimates to be found, and for singular or ill-conditioned matrices $\mathbf{X}^T\mathbf{X}$ are more suitable than estimates found by the least-squares method, which are always unbiased.

For the ill-conditioned matrix $\mathbf{X}^T\mathbf{X}$, the biased parameter estimates are shorter and smaller than least-squares estimates. They are “more precise” because they have smaller parameter variances. Moreover, these estimates exist even for a singular $\mathbf{X}^T\mathbf{X}$ matrix, when the least-squares method always fails.

Problem 6.23. *Approximation of a convex increasing function by a polynomial*

Many problems in chemometrics concern approximation of instrumental data of convex (or concave) increasing (or decreasing) values by a polynomial of any degree, so that the polynomial represents the shape of a data curve. Use the method of rational ranks for approximation of convex increasing data. For approximation, choose a polynomial of the sixth degree $E(y/x) = \sum_{j=1}^6 b_j x^j + b_7$. Calculate the value of the dependent variable y_0 at the origin i.e. parameter estimate β_7 .

Data: $n = 10$

x	25	35	45	55	65	75	85	95	105	115
y	150	160	170	190	210	230	270	310	370	450

Solution: Table 6.6 lists parameter estimates found by the classical least-squares method (LS), with $P = 10^{-30}$, and by the method of rational ranks (RV), with $P = 3.5 \times 10^{-4}$ for which the statistic MEP was smallest.

Figure 6.24a shows the regression model, with the 95% confidence interval for the LS method and Fig. 6.24b the regression model for the RV method, for $P = 3.5 \times 10^{-4}$. From the figures and Table 6.6 it is evident that the parameter estimates for the LS method do not match the data very well.

Table 6.6. The parameter estimates for different P values

Method	P	MEP	b_7	b_1	b_2	b_3
LS	10^{-30}	160.8	195.5	-5.92	0.258	-4.9×10^{-3}
RV	3.5×10^{-4}	8.59	134.7	0.35	0.0092	3.2×10^{-5}

b_4	b_5	b_6
5.33×10^{-5}	-2.9×10^{-7}	6.98×10^{-10}
-5.3×10^{-8}	3.9×10^{-9}	4.5×10^{-11}

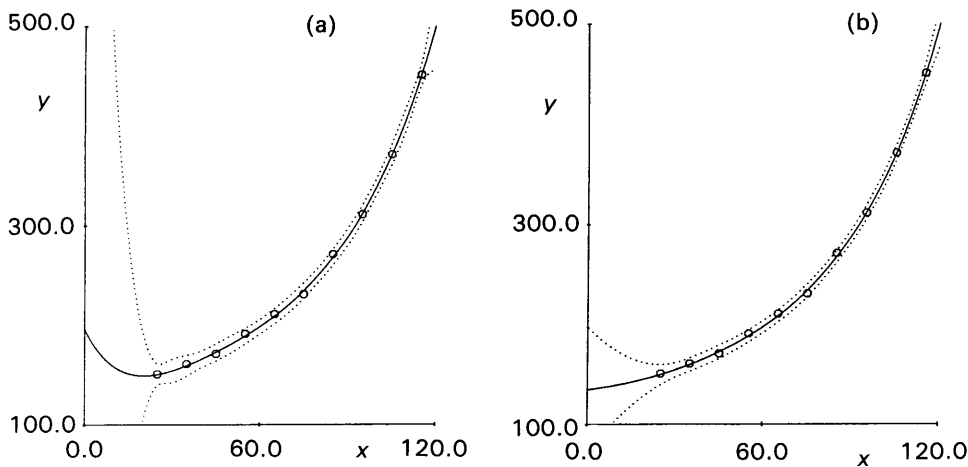


Fig. 6.24—The curve fitting of the model proposed, with the 95% confidence interval and the experimental data: (a) by the LS method ($P = 10^{-30}$), and (b) by the RV method ($P = 3.5 \times 10^{-4}$).

The estimate of parameter β_7 is larger than the value of y_1 , so the proposed model equation has a minimum between the origin and the point (x_1, y_1) . The confidence intervals are rather broad, and do not permit prediction of y outside the measurement interval.

The estimates determined by the rational ranks method are biased. The parameter b_7 is, however, smaller than y_1 and the confidence intervals show some possibility of prediction, even outside the measurement interval. As the polynomial degree was known and had physical meaning, no corrections based on the statistical analysis will be attempted.

Conclusion: The method of rational ranks allows us to find parameter estimates which give a curve shape corresponding to the data trend, with no extra extremes or inflections. With classical LS, such problems can be solved only with the introduction of some restrictions on parameters. Moreover, it is found that, for practical purposes, the biased estimates are not unsatisfactory.

6.5 REGRESSION DIAGNOSTICS

In linear and nonlinear regression analysis, the method of least-squares is often used. This method, however, does not ensure that the model is fully acceptable from statistical and physical point of view. A source of problems may be found in components of a regression triplet [data, model and method of estimation]. The least-squares method provides accurate estimates only when all assumptions about data and about a regression model are fulfilled. When some assumptions are not fulfilled, the least-squares method is inconvenient. Regression diagnostics represent the procedures for identification of

- (a) the data quality for a proposed model,
- (b) the model quality for a given set of data, and
- (c) fulfilment of all least-squares assumptions.

In the literature [21] the term regression diagnostics refers to methods for identification of influential points and multicollinearity. Atkinson [22] also includes as part of regression diagnostics, methods for proposing an actual regression model, perhaps with use of transformation of variable(s). Weisberg [23] includes as regression diagnostics

- (1) the examination of all assumptions for parameter estimation,
- (2) the statistical analysis of parameters, i.e. testing of the model
- (3) the identification of influential points, i.e. critical examination of data.

In this book we understand by regression diagnostics

- (1) methods of exploratory data analysis of individual variables (Chapter 2),
- (2) methods for analysis of influential points and
- (3) methods for identification of violations of the conditions for least-squares (Section 6.2).

The main difference between the use of regression diagnostics and classical statistical tests is that there is no necessity for an alternative hypothesis, but all types of deviations from an ideal regression triplet are discovered. Our concept of exploratory regression analysis is based on the fact that "the computer user knows more about the data than the computer". The personal computer serves us as an efficient tool for interactive diagnosis of data, model, and estimation method. The procedure of model building with the help of a personal computer involves interactive co-operation between the user and the computer program. Therefore, formal models that do not have physical meaning should not be proposed and analysed.

6.5.1 Exploratory regression analysis

Methods of exploratory data analysis have been described in Chapter 2. In exploratory regression analysis we will use these methods for (a) determination of statistical peculiarities of individual variables or residuals, (b) examination of assumptions regarding the distribution of variables and residuals.

In some cases, simply plotting the measured variable y_i against an index i may

uncover a latent variable, often related to time or order of measurement [25].

The first view into the relationship between individual variables comes from an x - y scatter plot of y against x . Some information about multicollinearity can be obtained by plotting pairs of controllable variables, x_j against x_k , $j \neq k$. An approximately linear dependence indicates strong multicollinearity. However, a plot of the response y against variable x_j , $j = 1, \dots, m$, may suggest nonlinearity of a model which is, in fact, of linear nature.

Problem 6.24. *Danger of false conclusions from inappropriate application of the scatter plot*

Draw the scatter plot of response y against variable x_1 and y against x_2 for a linear regression model $y_i = 10 - 6x_{1i} + 0.5x_{2i}$, $i = 1, \dots, 10$. Choose values of independent variable $x_{1i} = i - 5$ and $x_{2i} = x_{1i}^2$. Draw conclusions from these two plots.

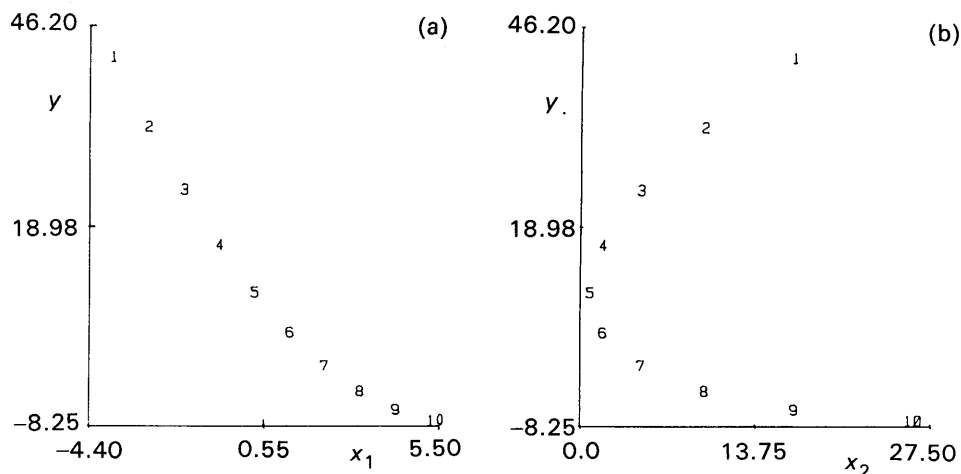


Fig. 6.25—The scatter plots for (a) response y as a function of variable x_1 , and (b) response y as a function of x_2 .

Solution: Figure 6.25 shows two scatter plots indicating quite strong nonlinearity. This could lead to a hypothesis that y is a nonlinear function of variables x_1 and x_2 . The apparent nonlinear nature is caused by a particular choice of controllable values x_{2i} and by quite strong multicollinearity between the independent variables x_1 and x_2 . **Conclusion:** For multivariate cases, scatter plots of y against x_j often are not helpful in identification of the regression model type.

Data normality is examined by the quantile–quantile (Q–Q) plot (Chapter 2). The principle methods of exploratory regression analysis include the determination of data range, data variability and the presence of outliers. All the graphical diagnostic tools described in Chapter 2 may be applied here. The EDA techniques allow identification of situations where

- (1) the range of measured data is too restricted,

- (2) the proposed model is false because there are some latent variables,
- (3) multicollinearity exists,
- (4) data do not have the normal distribution when the controllable variables are random numbers.

6.5.2 Examination of data quality

Data quality has a strong influence on any proposed regression model. Examination of data quality involves detection of the **influential points** (IP), which cause many problems in regression analysis by shifting the parameter estimates or increasing the variance of the parameters. Influential points may be classified into three groups:

- (a) **Gross errors**, caused by outliers in the measured variable or by the leverage points (extremes in the controllable variables).
- (b) **Golden points** are special chosen points which have been very precisely measured to extend the prediction capability of model.
- (c) **Latently influential points** are the consequence of a poor regression model.

Influential points may instead be classified according to data location:

- (a) **Outliers** differ from the other points in value on the y-axis;
- (b) **Leverage points** differ from the other points in values on the x-axis or in a combination of these quantities (in the case of multicollinearity).

There are also points, however, which are outliers and leverage points together. Outliers are identified by examination of the residuals. Leverage points are found from the diagonal elements H_{ii} of the projection hat matrix.

6.5.2.1 Statistical analysis of residuals

(1) Classical residuals

Residuals \hat{e}_i are defined by the expression

$$\hat{e}_i = y_i - \mathbf{x}_i \mathbf{b}$$

where \mathbf{x}_i is the i th row of matrix \mathbf{X} . Classical analysis is based on the assumption that residuals are estimates of errors ε_i . With the use of residuals the properties of errors are examined. This assumption, however, is not quite correct and can sometimes lead to false results. The false assumptions about residuals include the following:

- (a) the distribution of residuals is the same as the error distribution and the statistical properties of the residuals are identical with those of the errors, and
- (b) if the residual value is large, a large effect is caused by the corresponding point, so the point should be excluded from the data.

Let us note the differences between errors ε_i and residuals \hat{e}_i . The geometric illustration (Fig. 6.2) shows that the residuals \hat{e}_i are *not independent* even when the errors ε_i are independent. The residuals \hat{e}_i are a projection of vector \mathbf{y} into a subspace of dimension $(n - m)$. By using projection matrix \mathbf{P} we can write

$$\hat{\mathbf{e}} = \mathbf{P}\mathbf{y} = \mathbf{P}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \mathbf{P}\boldsymbol{\varepsilon} = (\mathbf{E} - \mathbf{H})\boldsymbol{\varepsilon} \quad (6.92)$$

To rearrange Eq. (6.92), we use the fact the vector $\mathbf{X}\boldsymbol{\beta}$ lies in the plane perpendicular to the projection plane, so that a zero vector results. For the i th residual

$$\hat{e}_i = (1 - H_{ii})y_i - \sum_{j \neq i}^n H_{ij}y_j = (1 - H_{ii})\varepsilon_i - \sum_{j \neq i}^n H_{ij}\varepsilon_j \quad (6.93)$$

Each residual \hat{e}_i is a linear combination of all errors ε_i . The distribution of residuals depends on

- (a) the error distribution,
- (b) the elements of the projection matrix \mathbf{H} ,
- (c) the sample size n .

Because the residual \hat{e}_i represents a sum of random quantities with bounded variance, the supernormality effect appears for small sample sizes. Even when the errors ε do not have a normal distribution, the distribution of residuals is close to normal. In small samples, the elements of the projection matrix \mathbf{H} are large and the main role of an actual point is to influence the sum of terms $H_{ij}\varepsilon_j$. The distribution of this sum is closer to a normal one than the distribution of errors ε . For large sample sizes, where $1/n \approx 0$, we find that $\hat{e}_i \approx \varepsilon_i$ and analysis of the residual distribution gives direct information about the distribution of errors.

Equation (6.18) may be used to calculate the residual variance

$$D(\hat{e}_i) = (1 - H_{ii})\hat{\sigma}^2 \quad (6.94)$$

The variance of residuals $D(\hat{e}_i)$ is not constant even when the variance of errors is constant. According to Eq. (6.18), the paired correlation coefficient r_{ij} between two residuals e_i and e_j is given by

$$r_{ij} = \frac{-H_{ij}}{\sqrt{(1 - H_{ii})(1 - H_{jj})}} \quad (6.95)$$

which shows that residuals are correlated even when errors ε_i and ε_j are independent.

For strong leverage points (extremes), the diagonal elements $H_{ii} \rightarrow 1$ while non-diagonal elements $H_{ij} \approx 0$. From Eq. (6.93), it may be concluded that an equation $\hat{e}_i = 0$ is valid, whatever the magnitude of y_i . The residuals do not always indicate correctly some strongly deviant values.

When a regression analysis is carried out by the least-squares method, for a model with an intercept term it is true that

$$\frac{\sum_{i=1}^n \hat{e}_i}{n} = 0$$

which corresponds to saying that the mean value of errors is equal to zero, $E(\varepsilon) = 0$.

Classical residuals are always associated with non-constant variance; they sum to be more normal and may not indicate strongly deviant points. The common practice

of chemometrics programs for statistical analysis of residuals is to use for examination some statistical characteristics of residuals such as the mean, the variance, the skewness and the kurtosis.

It should be particularly noted that in the case of small sample sizes the estimates of skewness and kurtosis are rather distorted and cannot reliably indicate the correctness of a proposed model.

Problem 6.25. *Inappropriate application of some simple statistics for residual analysis*

To illustrate the overestimated approach to examination of the reliability of parameter estimates, the following residual statistics are calculated: the mean of absolute values of residuals $|\bar{e}|$, the estimate of the standard deviation of residuals $\hat{\sigma}(\hat{e})$, the estimate of residual skewness $\hat{g}_1(\hat{e})$, the estimate of residual kurtosis $\hat{g}_2(\hat{e})$. Calculate these statistics for the four data samples from Problem 6.8.

Data: from Problem 6.8

Solution: The numerical values of some statistical characteristics for the four data samples from Problem 6.8 are listed in Table 6.7.

Table 6.7. Statistical characteristics for the four data samples of Problem 6.8

Sample	$ \bar{e} $	$\hat{\sigma}(\hat{e})$	$\hat{g}_1(\hat{e})$	$\hat{g}_2(\hat{e})$	$L(\hat{e})$
A	0.837	1.237	0.13	2.24	0.29
B	0.967	1.237	0.63	1.94	1.25
C	0.716	1.237	-2.28	6.57	15.39
D	0.903	1.237	-0.011	1.93	0.52

The last column of Table 6.7 lists values of the criterion of the Jarque–Berra test, which combines both skewness and kurtosis (Section 6.5.4). When $L(\hat{e}) > \chi^2_{1-\alpha}(2) = 5.99$, the normality of the data distribution is not proved.

From Table 6.7 it is evident that the statistics $|\bar{e}|$ and $\hat{\sigma}(\hat{e})$ do not indicate the model quality. Only sample C, which has one strong outlier, exhibits significant deviation in skewness and kurtosis, and the Jarque–Berra test does not prove the normality of the sample distribution. Neither the nonlinear dependence (sample B) nor the spurious data with one strong outlier (sample D) are correctly detected by the four statistics of residuals.

Conclusion: The statistical characteristics $|\bar{e}|$, $\hat{\sigma}(\hat{e})$, $\hat{g}_1(\hat{e})$, $\hat{g}_2(\hat{e})$ often do not give a correct indication of the quality of a model.

(2) *Normalized residuals*

In chemometrics the normalized residuals \hat{e}_{Ni} defined by

$$\hat{e}_{Ni} = \frac{\hat{e}_i}{\hat{\sigma}}$$

are often recommended. It is often assumed that these residuals are normally distributed quantities with zero mean and variance equal to 1, $\hat{e}_{Ni} \sim N(0, 1)$. When normalized residuals are used, the rule of 3σ is classically recommended: quantities

with \hat{e}_{Ni} of magnitude greater than $\pm 3\sigma$ are classified as the outliers. For a normal distribution, only 0.3% of all values lie outside the interval $\bar{x} \pm 3\hat{\sigma}$. Such assumptions about normalized residuals are misleading.

From Eq. (6.94) it is obvious that the variance $D(\hat{e}_{Ni}) = (1 - H_{ii})$ is not constant, and also not equal to one. For strong leverage points, $\hat{e}_i \approx 0$, so application of $\pm 3\sigma$ rule could lead to exclusion of correct points but retention of erroneous values.

Problem 6.26. *Inappropriate application of normalized residuals for identification of influential points*

For a dependence $y_i = x_i^2$, $i = 1, \dots, 12$, the data contain one gross error, the number x_{12} is replaced by y_{12} . Estimate parameters and with the use of normalized residuals try to locate the false leverage point 12.

Data:

i	1	2	3	4	5	6	7	8	9	10	11	12
x_i	1	2	3	4	5	6	7	8	9	10	11	144
y_i	1	4	9	16	25	36	49	64	81	100	121	12

Solution: By the least-squares method for the model $E(y/x) = \beta_1 x^2$, the parameter estimate was found $b_1 = 0.000671$. Point 12 has $\hat{e}_{N12} = -0.03$ and point 1 has the highest value, $\hat{e}_{N1} = 1.94$. With the use of the standardized residual, Eq. (6.96), the maximum for the 12th point was indicated, $\hat{e}_{S12} = -3.16$.

Conclusion: Normalized residuals are not able to indicate leverage points. Such points may be discovered, for example, by standardized residuals.

(3) *Standardized residuals*

The standardized residuals \hat{e}_{Si} , defined by

$$\hat{e}_{Si} = \frac{\hat{e}_i}{\hat{\sigma} \sqrt{1 - H_{ii}}} \quad (6.96)$$

exhibit constant variance. The statistical properties of standardized residuals are the same as those of classical residuals. Standardized residuals are, apart from the multiplicative constant $1/\sqrt{n - m}$, equal to cosine θ_i the angle between the vector $\hat{\mathbf{e}}$ and the vector \mathbf{i}^\perp (which is the projection of the i th column of matrix \mathbf{E} onto hyperplane L^\perp) so that

$$\cos \theta_i = \frac{\hat{e}_{Si}}{\sqrt{n - m}}$$

The maximum value of \hat{e}_{Si} is bounded by $\sqrt{n - m}$. The variable $\hat{e}_{Si}^2/(n - m)$ has the beta distribution $\text{Be} [0.5, (n - m - 1)/2]$.

(4) *Jack-knife residuals*

If, in Eq. (6.96), instead of the standard deviation we use the estimate of standard deviation $\hat{\sigma}_{(-i)}$ obtained by leaving out the i th point, we obtain the Jack-knife or fully

studentized residuals \hat{e}_{ji} ,

$$\hat{e}_{ji} = \hat{e}_{si} \sqrt{\frac{n-m-1}{n-m-\hat{e}_{si}^2}} = \sqrt{n-m} \times \cotg \theta_i \quad (6.97)$$

which, with an assumption of normality of errors, have the Student distribution with $(n-m-1)$ degrees of freedom. Jack-knife residuals correspond to the criterion of a t -test of the null hypothesis $H_0: C = 0$ in the model of a simple shift

$$y = \mathbf{X}\boldsymbol{\beta} + \mathbf{C}\mathbf{i} + \varepsilon \quad (6.98)$$

where \mathbf{i} is the identity vector with the i th element equal to one and other elements equal to zero. The model [Eq. (6.98)] expresses the case of an outlier where C is directly equal to the value of deviation, but also the case of a leverage point $C = \mathbf{d}_i^T \boldsymbol{\beta}$ where \mathbf{d}_i is the vector of the deviation of the individual x -components of the i th point. Jack-knife residuals are often used instead of classical residuals \hat{e}_i for identification of outliers. In the case of leverage points, these residuals do not give a reliable indication.

(5) Predicted residuals

The estimate of parameter C in Eq. (6.98) is represented by the predicted residuals

$$\hat{e}_{Pi} = y_i - \mathbf{x}_i \mathbf{b}_{(i)} = \frac{\hat{e}_i}{1 - H_{ii}} \quad (6.99)$$

where $\mathbf{b}_{(i)}$ is a vector of the parameter estimates obtained by the least-squares method with the i th point omitted. Predicted residuals sensitively monitor the magnitude of shift C .

(6) Recursive residuals

All the residuals already mentioned are correlated. To find uncorrelated residuals, the recursive least-squares method can be used. The resulting recursive residuals are very useful diagnostically as they allow identification of any instability in a model, for example, instability in time. Recursive residuals are defined

$$\hat{e}_{Ri} = 0, \quad i = 1, \dots, m \quad (6.100a)$$

$$\hat{e}_{Ri} = \frac{(y_i - \mathbf{x}_i \mathbf{b}_{i-1})}{\sqrt{1 + \mathbf{x}_i (\mathbf{X}_{i-1}^T \mathbf{X}_{i-1})^{-1} \mathbf{x}_i^T}}, \quad i = m+1, \dots, n \quad (6.100b)$$

where \mathbf{b}_{i-1} are estimates obtained from the first $(i-1)$ points. The matrix \mathbf{X}_{i-1} contains the first $(i-1)$ rows of matrix \mathbf{X} . These recursive residuals are independent and have constant variance. They are often used in normality tests or in tests of stability of regression coefficients.

Problem 6.27. Identification of influential points by various types of residuals

The outlier in sample C and the leverage point in sample D from Problem 6.8 can be identified by use of some types of residuals.

Data: from Problem 6.8

Solution: The classical residuals \hat{e}_i , standardized residuals \hat{e}_{Si} and Jack-knife residuals \hat{e}_{Ji} are used to detect the outlier in sample C ($x_3 = 13, y_3 = 12.74$) and one leverage point ($x_8 = 19, y_8 = 12.5$). The results are shown in Table 6.8. The outlier in sample C is most effectively detected by its Jack-knife residual. The leverage point (8 in sample D) is not detected by any residual.

Table 6.8. Various types of residuals for samples C and D

Sample	i	x_i	y_i	\hat{e}_i	\hat{e}_{Si}	\hat{e}_{Ji}
C	3	13	12.74	3.24	3	1203.54
D	8	19	12.5	0	0	0

Conclusion: Neither standardized nor Jack-knife residuals are always suitable for identification of influential points.

The various types of residuals differ in suitability for diagnostic purposes.

- (1) The standardized residuals \hat{e}_{Si} serve for identification of heteroscedasticity.
- (2) The Jack-knife residuals \hat{e}_{Ji} or the predicted residuals \hat{e}_{Pi} are suitable for identification of outliers.
- (3) Recursive residuals \hat{e}_{Ri} are used for identification of autocorrelation.

For analysis of residuals a variety of plots are used. Three principal types of plots can indicate inaccuracy of a proposed model, some trends, heteroscedasticity or influential points in data.

Plot type I (the index sequence plot) is a plot of residuals \hat{e}_i against the index i .

Plot type II (the plot against the independent variables) is a plot of residuals \hat{e}_i vs. the independent variable x_j .

Plot type III (the plot against the prediction) is represented by a plot of residuals \hat{e}_i against the predicted value \hat{y}_i .

Figure 6.26 shows possible graph shapes which can occur in plots of residuals. If the graph shape is a random pattern (Fig. 6.26a), the least-squares assumption is correct. Some systematic pattern indicates that the approach is incorrect in some way. A sector pattern in graph types I, II and III indicates heteroscedasticity in data (Fig. 6.26b). A band pattern in graph types I and II indicates some error in calculation or absence of x_j in model (type II). The band pattern may be also caused by outlying points or in type III by a missing intercept term in the regression model.

It should be noted that the plot of \hat{e}_i against the dependent variable y_i is not recommended, because the two quantities are strongly correlated. The smaller the correlation coefficient, the more linear is this plot.

A nonlinear pattern in all three graph types I, II and III indicates that the model proposed is incorrect.

6.5.2.2 Analysis of projection matrix elements

Analysis of elements of the projection hat matrix plays an important role in regression diagnostics because the diagonal elements of this matrix

$$H_{ii} = \mathbf{x}_i(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{x}_i^T$$

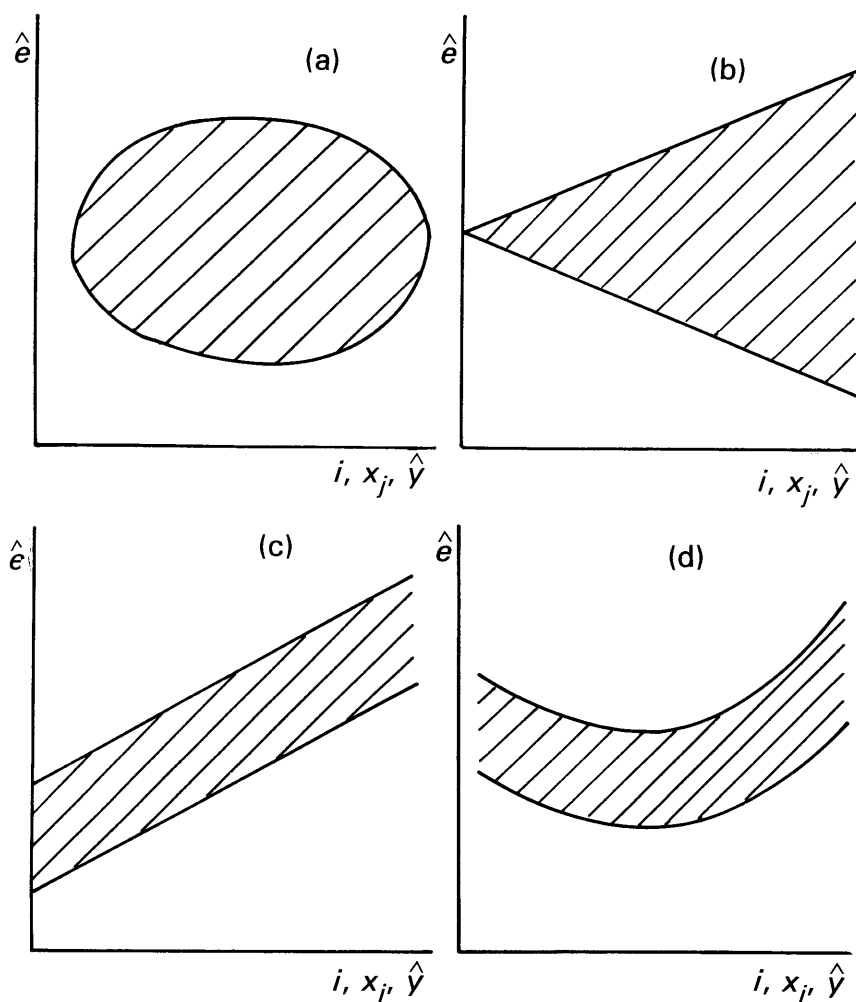


Fig. 6.26—Possible shapes of residual plots: (a) random pattern shape, (b) sector pattern shape, (c) band shape, (d) nonlinear curved band shape.

indicate the presence of leverage points which are not detected by analysis of residuals. Diagonal elements (denoted in literature as “leverage”) have some properties which come from the symmetry and idempotency of matrix \mathbf{H} . Among the properties of matrix \mathbf{H} are:

- (1) The condition for the diagonal elements of a projection matrix is $0 \leq H_{ii} \leq 1$ and for nondiagonal elements $-1 \leq H_{ij} \leq 1$. When a model also contains an intercept term and the rank of matrix \mathbf{X} is m , another condition for diagonal elements is valid, $1/n \leq H_{ii} \leq 1/c$, where C is the number of replicate measurements at each value of the controllable variable.
- (2) For a model with an intercept term and the full rank of matrix, \mathbf{X} :

$$\sum_{i=1}^n H_{ii} = m$$

$$\sum_{i=1}^n H_{ij} = 1$$

The mean value of the diagonal element is $H_{ii} = m/n$.

- (3) From the idempotency of matrix \mathbf{H} it follows that

$$H_{ii} = H_{ii}^2 + \sum_{j \neq i}^n H_{ij}^2 = \sum_{j=1}^n H_{ij}^2$$

From this equation two important properties of diagonal elements H_{ii} follow:

- (a) If the diagonal elements are close to zero, $H_{ii} \rightarrow 0$, all nondiagonal elements are also close to zero, $H_{ij} \rightarrow 0$, for $j = 1, \dots, n$;
 - (b) If the diagonal elements are close to 1, $H_{ii} \rightarrow 1$, all nondiagonal elements are close to zero, $H_{ij} \rightarrow 0$, for $j = 1, \dots, n$.
- (4) If the matrix \mathbf{X} comes from the multivariate normal distribution, the quantity

$$F = (n - m)[H_{ii} - 1/n][(1 - H_{ii})(m - 1)]$$

has the Fisher-Snedecor distribution $F(m - 1, n - m)$.

- (5) The larger the diagonal elements H_{ii} , the more the i th point of prediction \hat{y}_i is affected. If the H_{ii} elements are close to 1 ($H_{ii} \rightarrow 1$, and $\hat{y}_i = y_i$) then all of the variability in x_i is explained by the regression model.
- (6) The diagonal elements $H_{ii} = \delta \hat{y}_i / \delta y_i$ express the sensitivity of the prediction \hat{y}_i to any change in variable y_i . A zero value, $H_{ii} = 0$, indicates a point which has no influence on prediction.
- (7) The diagonal elements H_{ii} are a nondecreasing function of the controllable variables m , and a nonincreasing function of the number of points n .
- (8) The further point x_i lies from the centre of gravity of all points, the more it is likely to be a leverage point, and the more the value of diagonal elements H_{ii} will increase.
- (9) If the controllable variables \mathbf{x} have the normal distribution, for large sample sizes n ($nH_{ii} - 1$) has approximately the $\chi_m^2(2)$ distribution.

For more complex analysis, it is useful to form the extension of matrix \mathbf{X} by a vector \mathbf{y} to give matrix $\mathbf{X}^* = (\mathbf{X} | \mathbf{y})$. This matrix corresponds to the projection matrix

$$\mathbf{H}^* = \mathbf{H} + \frac{\hat{\mathbf{e}}\hat{\mathbf{e}}^T}{\hat{\mathbf{e}}^T\hat{\mathbf{e}}} \quad (6.101)$$

Since the matrix \mathbf{H}^* contains information about all variables it can be used as the total measure of influential points. Diagonal elements of this matrix are given by

$$H_{ii}^* = H_{ii} + \frac{\hat{e}_i^2}{(n - m)\delta^2} \quad (6.102)$$

To look at elements of the projection matrix, the *index graph* of H_{ii} elements against the index i is used.

Problem 6.28. *Identification of influential points from elements of the projection matrix*

The outlying point in sample C and the leverage point in sample D from Problem 6.8 may be used to test the identification of influential points by elements H_{ii} and H_{ii}^* of projection matrix \mathbf{H} .

Data: from Problem 6.8

Solution: The calculated diagonal elements H_{ii} and H_{ii}^* of the projection matrix \mathbf{H} are listed in Table 6.9

Table 6.9. Elements of the projection matrix H_{ii} and the extended projection matrix H_{ii}^* for samples C and D

Sample	x_i	y_i	H_{ii}	H_{ii}^*
C	13	12.75	0.236	1
D	19	12.5	1	1

The diagonal elements of the extended projection matrix indicate a strong influential point in both samples. The leverage point in sample D is indicated even by the diagonal element H_{ii} of the original projection matrix.

Conclusion: The diagonal elements of an extended projection matrix are useful for detecting outlier and leverage points in data. The leverage point was not detected by any type of residuals (Problem 6.27).

6.5.2.3 Plots for identification of influential points

For identification of different types of influential points, various types of residuals are combined with the diagonal elements of the projection matrix \mathbf{H} .

(1) *Graph of predicted residuals (GPR)*

(x -axis: the predicted residuals \hat{e}_{pi} ; y -axis: the classical residuals \hat{e}_i)

This graph is one of the simplest graphs. The leverage points are easy detected by their location as they lie outside the line $y = x$, and they are located quite far from this line. The outliers are located on the line $y = x$ but far from its central pattern (Fig. 6.27).

(2) *Williams graph (WG)*

(x -axis: the diagonal elements H_{ii} ; y -axis: the Jack-knife residuals \hat{e}_{ji})

In this graph two boundary lines are drawn. The first line is for outliers, $y = t_{0.95}(n - m - 1)$ and the second line is for leverage points, $x = 2m/n$ (Fig. 6.28). Denote that $t_{0.95}(n - m - 1)$ is the 95% quantile of the Student distribution with $(n - m - 1)$ degrees of freedom.

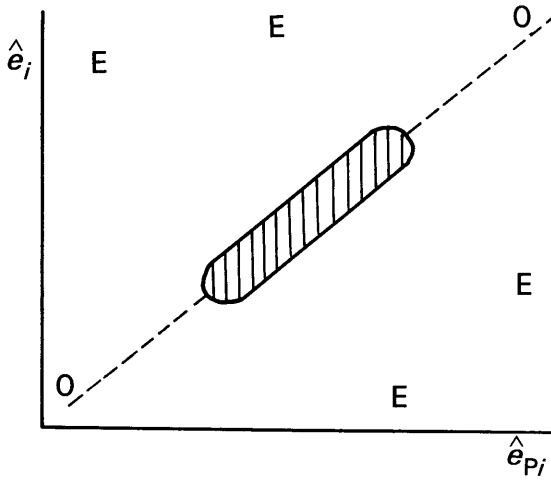


Fig. 6.27—Graph of predicted residuals (GPR): E is a high leverage point and O is an outlier.

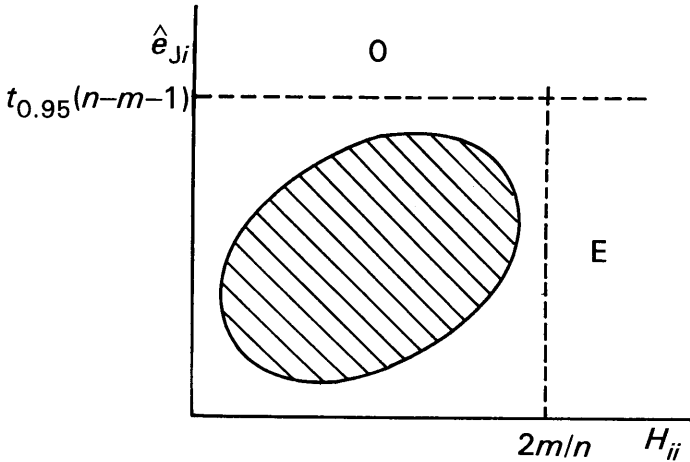


Fig. 6.28—Williams graph (WG): E is the leverage point and O is the outlier.

(3) Pregibon graph (PG)

(x-axis: the diagonal elements H_{ii} , y-axis: the normalized residuals \hat{e}_{Ni}^2)
 Since the expression $E(H_{ii} + \hat{e}_{Ni}^2) = (m + 1)/n$ is valid for this graph, two different constraining lines can be drawn,

$$y = -x + 2(m + 1)/n$$

and

$$y = -x + 3(m + 1)/n.$$

To distinguish among influential points the following rules are used:

- (a) a point is strongly influential if it is located above the upper line;

- (b) a point is influential if it is located between the two lines. The influential point can be either an outlier or a leverage point.

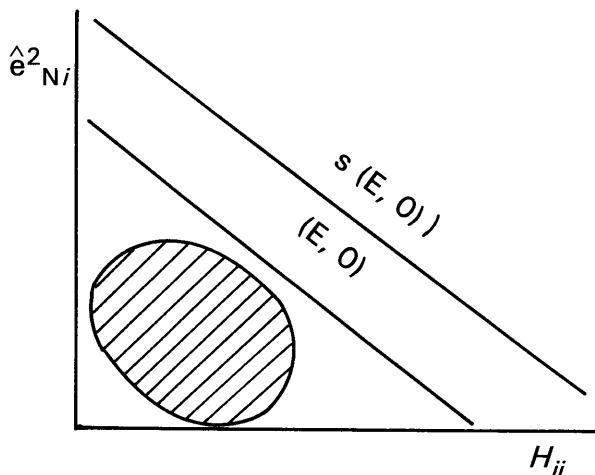


Fig. 6.29—Pregibon graph (PG): (E, O) are influential points, and $s(E, O)$ are strongly influential points.

(4) *McCulloch and Meeter graph (MMG)*

(x-axis: $\ln [H_{ii}/(m(1 - H_{ii}))]$; y-axis: the standardized residuals \hat{e}_{Si}^2)

In this plot the solid line drawn represents the locus of points with identical influence, with slope -1 . The 90% confidence line is defined by

$$y = -x - \ln F_{0.9}(n - m, m)$$

The boundary line for leverage points is defined as

$$x = \ln [2/(n - 2m)]$$

The boundary line for outliers is defined by

$$y = \ln [(n - m) \times (t_{0.95}^2(n - m))]$$

where $t_{0.95}(n - m)$ is the 95% quantile of the Student distribution with $(n - m - 1)$ degrees of freedom.

(5) *Index graph (IG)*

(x-axis: the index i ; y-axis: the residuals \hat{e}_i , \hat{e}_{Si} , \hat{e}_{Ni} , \hat{e}_{Pi} , \hat{e}_{Ji} , \hat{e}_{Ri} , or the diagonal elements H_{ii} or H_{ii}^* , or estimates b_i)

The x-axis always contains the order index i , but the y-axis can be a residual or the diagonal elements of the projection matrix. Sometimes also the parameter estimates b_i are on this axis.

(6) *Rankit graph (Q-Q plot)*

(x-axis: the quantile of the standardized normal distribution u_{Pi} ; y-axis: the

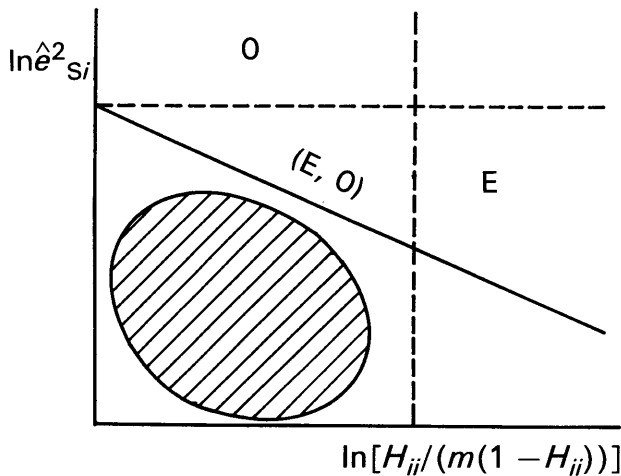


Fig. 6.30—McCulloch and Meeter graph (MMG): E is a leverage point, O is an outlier and (E, O) is an influential point.

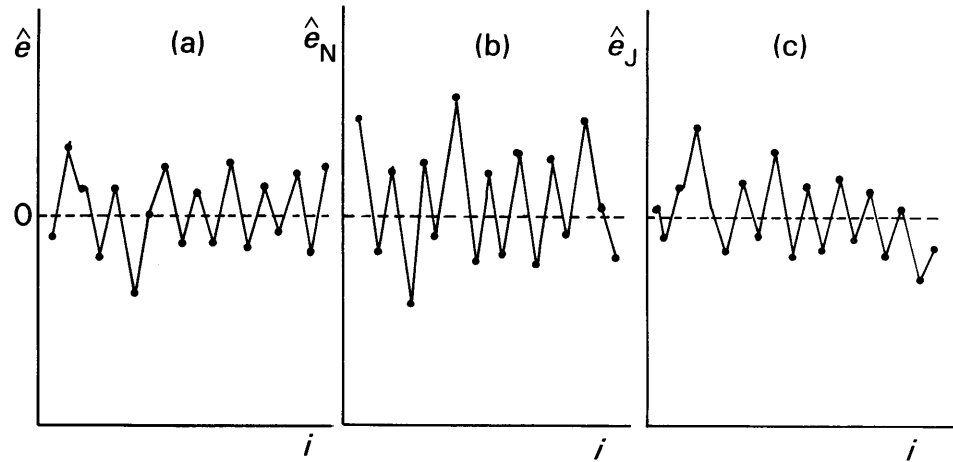


Fig. 6.31—Various types of index graph:
(a) \hat{e}_i vs. i , (b) \hat{e}_{Ni} vs. i , (c) \hat{e}_{Ji} vs. i ,

ordered residuals $\hat{e}_{(i)}$, $\hat{e}_{S(i)}$, $\hat{e}_{N(i)}$, $\hat{e}_{P(i)}$, $\hat{e}_{J(i)}$, $\hat{e}_{R(i)}$
On the x-axis are quantiles of the standardized normal distribution u_{p_i} for $P_i = i/(n + 1)$ and on the y-axis the order statistics of the residuals, i.e. increasing ordered values of various types of residuals.

Problem 6.29. Identification of influential points by graphical analysis of residuals and examination of elements of projection matrix
Use a variety of methods of graphical analysis of residuals and elements of projection matrix to identify influential points in the data from Problem 6.7. Compare the efficiency of the various graphical tools for detecting outliers and leverage points.

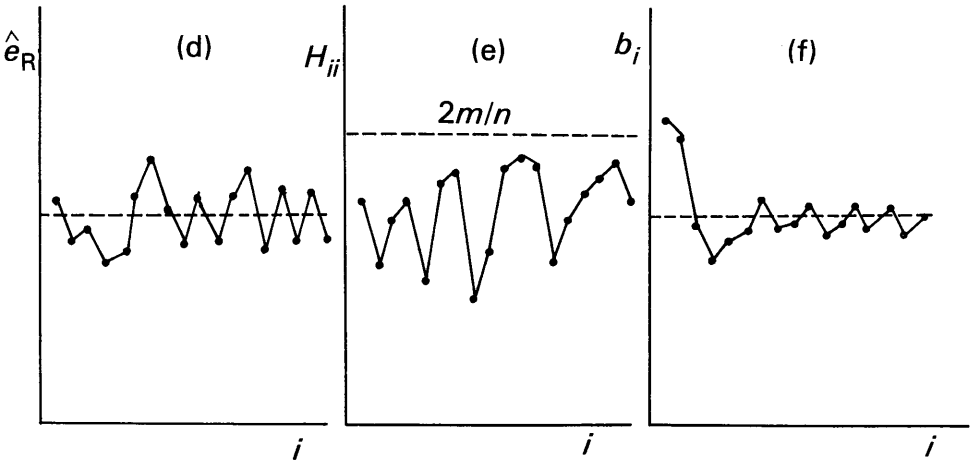


Fig. 6.31—Continued. (d) \hat{e}_{Ri} vs. i , (e) H_{ii} vs. i , (f) b_i vs. i .

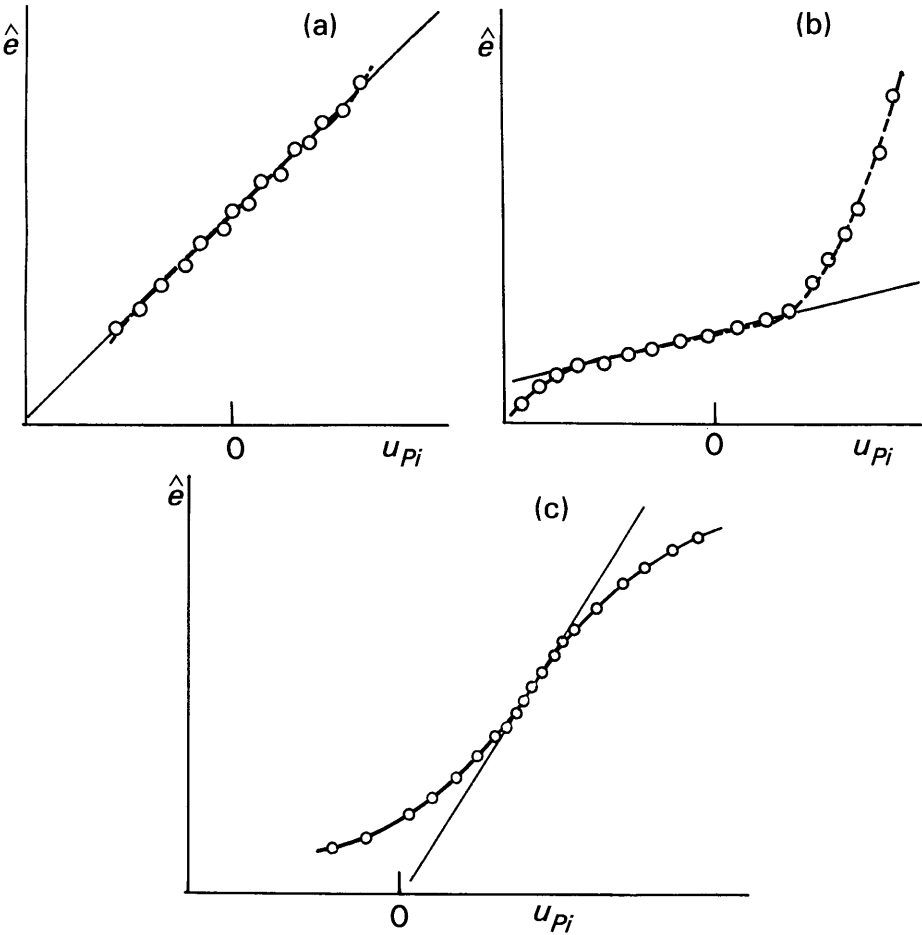


Fig. 6.32—Possible variations of the rankit (Q-Q) graph of residuals.

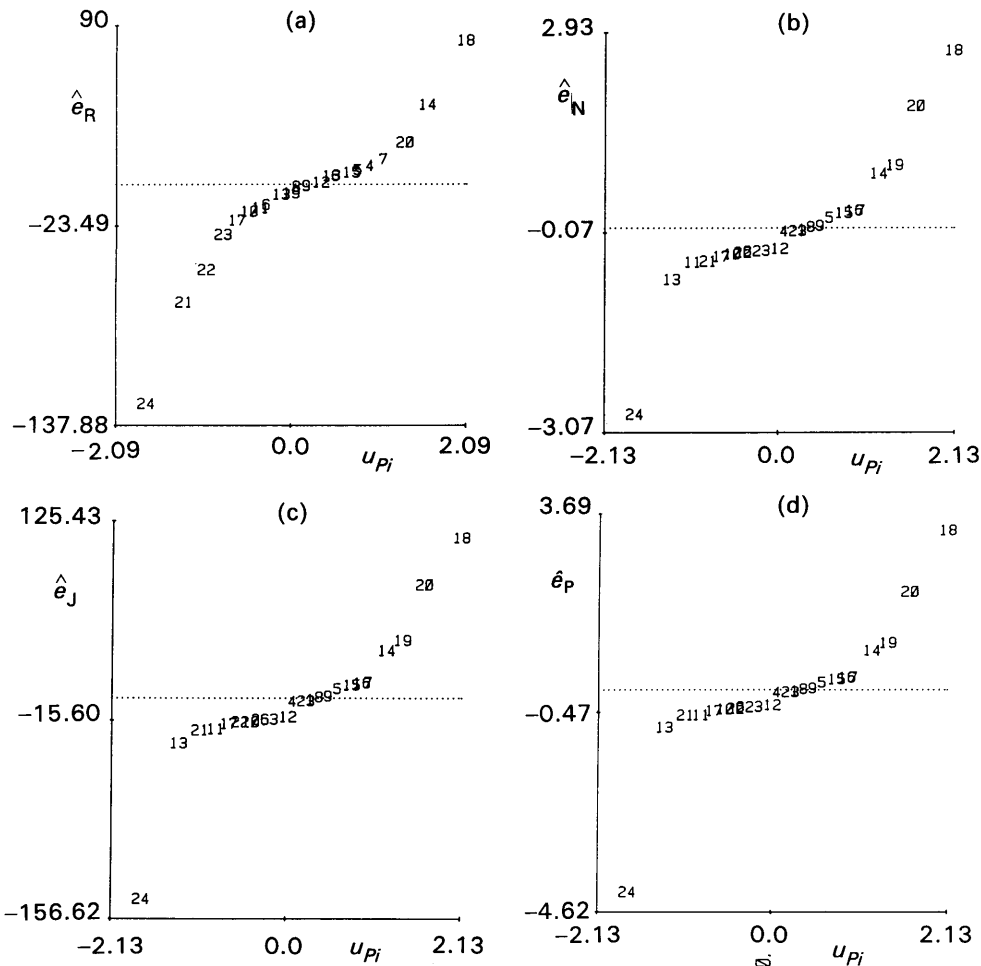


Fig. 6.33—Four types of rankit graph: (a) recursive residuals, (b) normalized residuals, (c) Jack-knife residuals, and (d) predicted residuals.

Data: from Problem 6.7

Solution: Figure 6.33 shows four types of rankit graph and Fig. 6.34 shows four types of index graph.

The graph of the third type of residuals \hat{e}_i against prediction \hat{y}_{P_i} indicates that the model proposed is incorrect; some trends and heteroscedasticity in data are also seen (Fig. 6.35).

For identification of influential points we can use the information from four graphs.

The graph of predicted residuals (GPR) on Fig. 6.36a discovers three outliers, points 24, 20 and 18, which, although they lie on the axis $y = x$, are rather far from other points.

The Williams graph (WG) on Fig. 6.36b also discovers three outliers, points 18, 20 and 24.

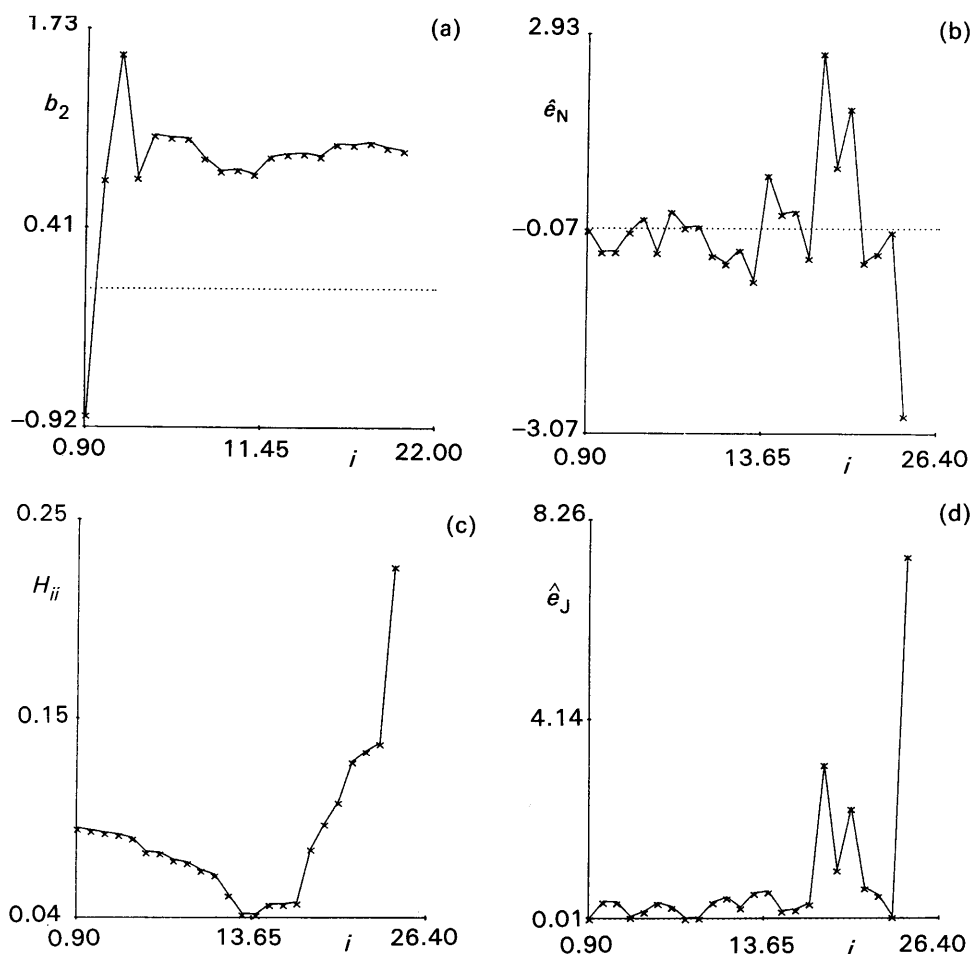


Fig. 6.34—Four types of index graph: (a) recursive estimates of slope b_{2i} (b) normalized residuals, (c) elements of projection matrix, and (d) Jack-knife residuals.

The Pregibon graph (PG) in Fig. 6.36c shows two outlying points above the upper limiting line, detecting that these two points are strongly influential. Point 20, lying between the two parallel limiting lines, is the only influential point. This point can be either an outlier or a leverage point.

The McCulloh and Meeter graph (MMG) in Fig. 6.36d includes the line with slope -1 , connecting points which are of the same influence, and two boundary lines. This plot also discovers three strongly influential points, 24, 20 and 18.

Conclusion: Graphical methods of analysis of residuals are rather simple and illustrative. They enable quick identification of influential points, i.e. here points 18, 20 and 24. To distinguish between outliers and leverage points, some boundary lines must be constructed.

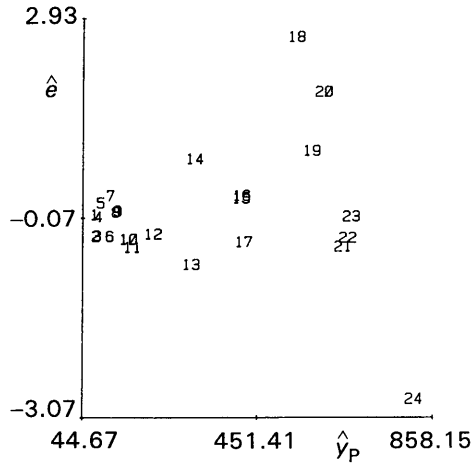


Fig. 6.35—Graph of type III, of classical residuals \hat{e}_i against prediction, \hat{y}_{Pi} .

6.5.2.4 Other characteristics of influential points

In the classification of influential points, it is important to remember that they can affect the various regression characteristics differently. Points affecting the prediction \hat{y}_i , for example, may not affect the parameter variance. The degree of influence of individual points can be classified according to the characteristics that are affected. For identification of influential points, there are many additional diagnostics which may be divided according to two principal approaches.

The first is based on the examination of changes which occur when certain points are omitted.

The second approach concerns the validity of the linear regression model (6.5b) when the variance of errors is abnormal. For the i th error, the normal distribution $N(0, \sigma^2/w_i)$ is valid, but the other errors ε_j , $j \neq i$, have the normal distribution of constant variance σ , i.e. $N(0, \sigma^2)$. The weight parameter lies in the interval $0 < w_i < 1$. This second approach leads to the *model of inflated variance*.

For $w_i = 1$ this assumption leads to the classical least-squares method. If we write $b(w_i)$ for the parameter estimate calculated according to Eq. (6.5b) when the variance of the i th error is just equal to σ^2/w_i , then the following expression is valid

$$\mathbf{b}(1) - \mathbf{b}(w_i) = \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i (1 - w_i) \hat{e}_i}{1 - (1 - w_i) H_{ii}} \quad (6.103)$$

where \mathbf{x}_i is the i th row of matrix \mathbf{X} which contains x components of the i th point.

For $w_i = 0$, Eq. (6.103) leads to

$$\mathbf{b}(1) - \mathbf{b}(0) = \mathbf{b} - \mathbf{b}_{(i)}$$

where $\mathbf{b}_{(i)}$ is the estimate reached by the least-squares method by using all points except the i th one. Leaving out the i th point is therefore the same as the case when this point has unbounded infinite variance.

To express the sensitivity of parameter estimates to the perturbation parameter w_i ,

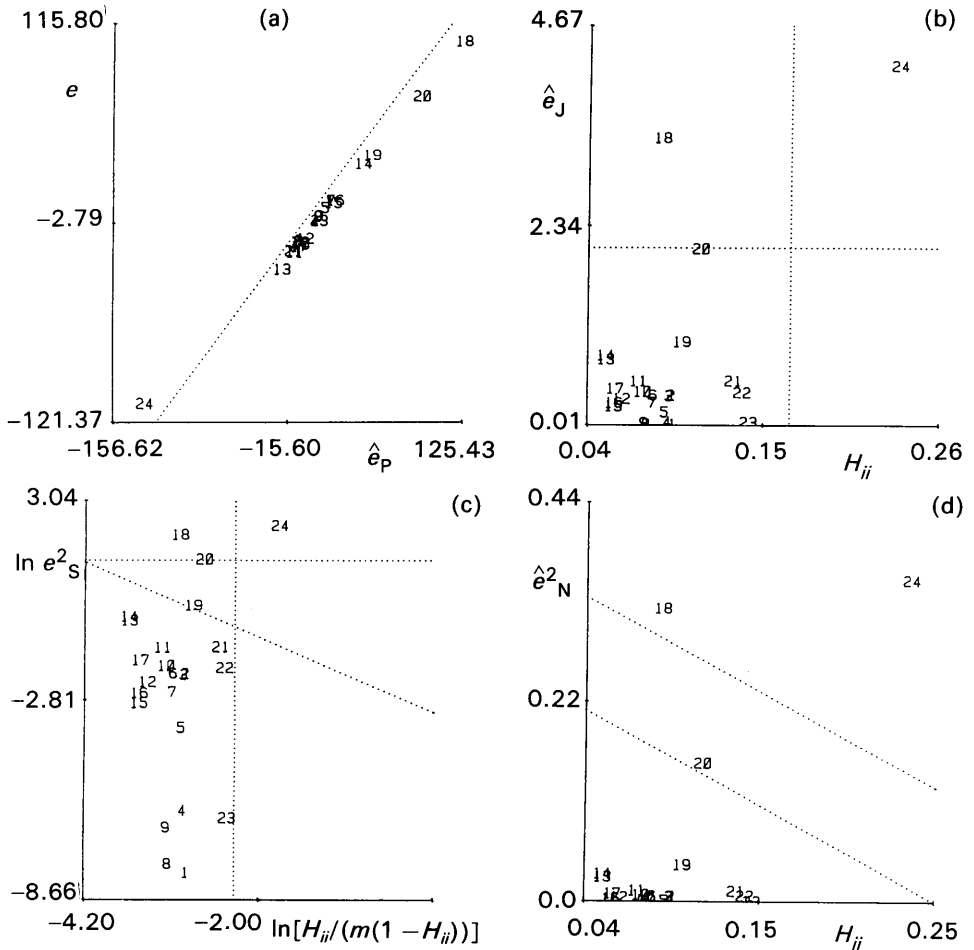


Fig. 6.36—Graphs for identification of outliers and leverage points: (a) GPR, (b) WG, (c) PG, and (d) MMG.

the sensitivity function $\delta \mathbf{b}(w_i)/\delta w_i$ can be used.

$$\frac{\delta \mathbf{b}(w_i)}{\delta w_i} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \hat{e}_i \frac{s_A + (1 - w_i)H_{ii}}{s_A^2} \quad (6.104)$$

where $s_A = 1 - (1 - w_i)H_{ii}$. The following types of sensitivity function of parameter estimates are possible.

(1) *The Jack-knife influence function*

The sensitivity function of parameter estimates defined by Eq. (6.104) at the value $w_i = 0$, is given by

$$\left. \frac{\delta \mathbf{b}(w_i)}{\delta w_i} \right|_{w_i=0} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \frac{\hat{e}_i}{(1 - H_{ii})^2} = \frac{JC_i}{n-1} \quad (6.105)$$

The term JC_i is the Jack-knife influence function. It is related to the sensitivity function of parameter estimates for the case when the i th point is omitted, because $\mathbf{b}(0) = \mathbf{b}_{(i)}$.

(2) *The empirical influence function*

The sensitivity function of parameter estimates [Eq. (6.104)] at the value of $w = 1$ is given by

$$\left. \frac{\delta \mathbf{b}(w_i)}{\delta w_i} \right|_{w_i=1} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \hat{e}_i = \frac{EC_i}{n-1} \quad (6.106)$$

The term EC_i is the empirical influence function. It is related to the sensitivity function of parameter estimates at the location of parameter estimates, \mathbf{b} , by the least-squares method.

(3) *The sample influence function SC_i*

The sample influence function is proportional to the change in the vector of parameter estimates when the i th point is left out. With the use of Eq. (6.103) we can write

$$SC_i = n(\mathbf{b} - \mathbf{b}_{(i)}) = n(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \frac{\hat{e}_i}{1 - H_{ii}} \quad (6.107)$$

All three influence functions differ only in a single term $(1 - H_{ii})$ so they are not identically sensitive to the presence of leverage points, for which $H_{ii} \rightarrow 1$. The disadvantage of all these influence functions is the fact that they are m -dimensional vectors. Their components define the influence of the i th point on the estimate of the j th parameter. Therefore, normalization of these vectors is used [26] to obtain scalar measures corresponding to distances which express the relative influence of the given point on all parameter estimates.

A popular scalar measure of the relative influence of the i th point on all parameter estimates is the *Cook distance* D_i . This is derived by normalization of the sample influence function SC_i . The resulting D_i has the form

$$D_i = \frac{(\mathbf{b} - \mathbf{b}_{(i)})^T \mathbf{X}^T \mathbf{X} (\mathbf{b} - \mathbf{b}_{(i)})}{m \times \hat{\sigma}^2} = \frac{(\hat{\mathbf{y}} - \hat{\mathbf{y}}_{(i)})^T (\hat{\mathbf{y}} - \hat{\mathbf{y}}_{(i)})}{m \times \hat{\sigma}^2} = \frac{\hat{e}_{Si}}{m} \times \frac{H_{ii}}{1 - H_{ii}} \quad (6.108)$$

The Cook distance is related to the confidence ellipsoid of the estimates, and it also permits comparison with the quantiles of the Fisher-Snedecor F -distribution. However, the shift of estimates appears here when the i th point is left out. It is approximately true that when $D_i > 1$, the shift is greater than the 50% confidence region, so the relevant point is rather influential. Another interpretation of the Cook distance D_i is based on the Euclidean distance between the prediction vector $\hat{\mathbf{y}}$ estimated by the least-squares method and the prediction vector $\hat{\mathbf{y}}_{(i)}$ estimated by the least-squares method when the i th point is left out. The Cook distance D_i expresses the influence of the i th point on the parameter estimate \mathbf{b} only.

When the i th point does not affect the parameter estimates \mathbf{b} significantly, the value of the Cook distance D_i is low. Such a point, however, can strongly affect the estimate of the residual variance $\hat{\sigma}^2$.

The relative changes in the parameter estimates caused by leaving out the i th point may be expressed by the standardized deviations of the j th estimate b_j of that parameter estimate $b_{(i)j}$ which has been obtained by leaving out the i th point. The corresponding diagnostic is defined by

$$DS_{ij} = \frac{b_j - b_{(i)j}}{\hat{\sigma}_{(i)} \sqrt{C_{ii}}} \quad (6.109)$$

where C_{ii} is the diagonal element of matrix $\mathbf{X}^T \mathbf{X}$. The influence of the i th point on the estimate of the j th regression parameter is significant when $DS > 2/\sqrt{n}$.

Problem 6.30. *The change in the estimate of the slope and intercept of a calibration straight line, caused by an outlier*

Determine the change of estimate value for the slope and intercept of the regression straight line $E(y/x) = \beta_1(x - \bar{x}) + \beta_2^*$ in the presence of one outlier.

Solution: From Eq. (6.103) for $w_i = 0$, the expression for the change in parameters $\Delta = \mathbf{b}_i - \mathbf{b}_{(i)}$ is given by

$$\Delta = \begin{bmatrix} b_1 - b_{(i)1} \\ b_2 - b_{(i)2} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n (x_i - \bar{x})^{-2} & 0 \\ 0 & n^{-1} \end{bmatrix} \begin{bmatrix} x_i - \bar{x} \\ 1 \end{bmatrix} \hat{e}_i \cdot \left[1 - n^{-1} - \frac{(x_i - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]^{-1}$$

For the slope change $\Delta_1 = b_1 - b_{(i)1}$

$$\Delta_1 = \frac{n \times \hat{e}_i (x_i - \bar{x})}{(n-1) \sum_{j=1}^n (x_j - \bar{x})^2 - n(x_i - \bar{x})^2}$$

and for the intercept change $\Delta_2 = b_2^* - b_{(i)2}^*$

$$\Delta_2 = \frac{\hat{e}_i}{(n-1) \sum_{j=1}^n (x_j - \bar{x})^2 - n(x_i - \bar{x})^2}$$

From these expressions it may be concluded that for $x_i = \bar{x}$, $\Delta_1 = 0$ regardless of the magnitude y_i . The slope of the regression straight line will not change whether the point located at $x_i = \bar{x}$ is an outlier or not. The estimate of the intercept will change, however, in dependence on the magnitude of \hat{e}_i .

Conclusion: The points of a calibration straight line located on the x -axis far from the mean \bar{x} have the most significant effect on the slope. A point having a negligible

effect on the slope may have a strong influence on the intercept estimate.

To express the sensitivity of distance measures to influential points, the *Atkinson distance* A_i is used

$$A_i = DF_i \sqrt{\frac{n-m}{m}} = |\hat{e}_{ji}| \times \sqrt{\frac{n-m}{m} \times \frac{H_{ii}}{1-H_{ii}}} \quad (6.110)$$

which is also convenient for graphical interpretation. With designed experiments, usually $H_{ii} = m/n$, and the Atkinson distance A_i is numerically equal to the Jack-knife residual \hat{e}_{ji} .

By normalizing the sample influence function and using the variance estimate $\sigma_{(i)}^2$ obtained from estimates $\mathbf{b}_{(i)}$, we obtain the characteristic DF_i defined by

$$DF_i^2 = \frac{(\hat{y}_i - \hat{y}_{(i)})^2}{\hat{\sigma}_{(i)}^2 \times H_{ii}} = \hat{e}_{ji}^2 \frac{H_{ii}}{1-H_{ii}} \quad (6.111)$$

The i th point is considered to be significantly influential when $DF_i > 2\sqrt{m/n}$. The characteristic DF_i was recommended by Belsey, Kuh and Welsch [21] as the basic diagnostic characterizing the influence of individual points on prediction \hat{y} . The term $H_{ii}/(1-H_{ii})$ in Eqs. (6.108)–(6.111) is equal to the ratio of variances $D(\hat{y}_i)/D(\hat{e}_i)$, and gives a measure of the sensitivity of regression to the location of the i th point.

There are many regression diagnostics indicating influential points which are based on the approach of leaving out the i th point. In addition to DS_{ij} and DF_i , several other characteristics may be useful [26].

The *Anders–Pregibon diagnostic* AP_i expresses the influence of the i th point on the volume of the confidence ellipsoid

$$AP_i = \frac{\det(\mathbf{X}_{(i)}^* \mathbf{T} \mathbf{X}_{(i)}^*)}{\det(\mathbf{X}^* \mathbf{T} \mathbf{X}^*)} \quad (6.112)$$

where $\mathbf{X}^* = (\mathbf{X} | \mathbf{y})$ is the matrix extended by the vector \mathbf{y} . The diagnostic AP_i is related to the elements of the extended projection matrix \mathbf{H}^* by the expression

$$AP_i = 1 - H_{ii} - \hat{e}_{Ni}^2 = 1 - H_{ii}^* \quad (6.113)$$

A point is considered to be influential if:

$$H_{ii}^* = 1 - AP_i > \frac{2(m+1)}{n}$$

To unify some of the expressions for identification of influential points, Cook and Weisberg [24] have recommended the use of a general diagnostic called the *likelihood distance* LD_i defined by

$$LD_i = 2[L(\hat{\theta}) - L(\hat{\theta}_{(i)})] \quad (6.114)$$

where $L(\hat{\theta})$ is the maximum of the logarithm of the likelihood function when all points are used and $L(\hat{\theta}_{(i)})$ is corresponding value when the i th point is omitted. The parametric vector θ contains either the parameter \mathbf{b} or the variance estimate $\hat{\sigma}^2$. For strongly influential points

$$LD_i > \chi^2_{1-\alpha}(m+1)$$

where $\chi^2_{1-\alpha}(m+1)$ is the quantile of the χ^2 distribution.

With the use of different variants of LD_i it is possible to examine the influence of the i th point on the parameter estimates or on the variance estimate or on both [26].

(a) To examine the influence of individual points on the parameter estimates \mathbf{b} the likelihood distance $LD_i(\mathbf{b})$ is expressed by

$$LD_i(\mathbf{b}) = n \times \ln \left[\frac{d_i \times H_{ii}}{1 - H_{ii}} + 1 \right] \quad (6.115)$$

where $d_i = \hat{e}_{Si}^2/(n-m)$

(b) To examine the influence of individual points on the residual variance estimate, the likelihood distance $LD_i(\hat{\sigma}^2)$ has the form

$$LD_i(\hat{\sigma}^2) = n \times \ln \left[\frac{n}{n-1} \right] + n \ln(1 - d_i) + \frac{d_i(n-1)}{1 - d_i} - 1 \quad (6.116)$$

(c) To examine the influence of individual points on the parameters \mathbf{b} and variance $\hat{\sigma}^2$ together, the likelihood distance $LD_i(\mathbf{b}, \hat{\sigma}^2)$ has the form

$$LD_i(\mathbf{b}, \hat{\sigma}^2) = n \times \ln \left[\frac{n}{n-1} \right] + n \ln(1 - d_i) + \frac{(n-1)d_i}{(1 - d_i)(1 - H_{ii})} - 1 \quad (6.117)$$

Investigation of the three variants of the likelihood distance leads to the following conclusions [26].

- (a) The diagnostic $LD_i(\mathbf{b})$ is a monotonic function of the Cook distance D_i (6.108) and has no advantage over the diagnostic D_i .
- (b) The diagnostic $LD_i(\hat{\sigma}^2)$ does not depend on H_{ii} and therefore it is not affected by high leverage points.
- (c) The diagnostic $LD_i(\mathbf{b}, \hat{\sigma}^2)$ expresses the influence of individual points on \mathbf{b} and $\hat{\sigma}^2$. It is more useful than both diagnostics A_i and DF_i , especially for models without intercept [26]. It seems to be enough to examine just this diagnostic. Generally the LD_i measures are not quite universal, and for estimation of influential points, many diagnostics must be combined.

Another test for influential points [27] is based on the influence of individual points on the sum of mean quadratic errors of estimates, of the mean quadratic errors of prediction and on the integral mean quadratic error of prediction. To test the influence of the i th point on all these characteristics, the Jack-knife residual \hat{e}_{ji} may be used as test criterion. This is suitable either for models of simple shift, Eq. (6.98), or models of inflated variance $D(\varepsilon_i) = \sigma^2/w_i$.

When more points are examined simultaneously for the model of simple shift, the validity of condition

$$\hat{e}_{ji}^2 \leq F_{1-\alpha/n}(1, n-m-1, 0.5) \quad (6.118)$$

means that no influential points are present in data. Here, $F_{1-\alpha/n}(1, n-m-1, 0.5)$ means the $100(1-\alpha/n)\%$ quantile of the non-central F -distribution with non-centrality parameter 0.5 and 1, and $(n-m-1)$ degrees of freedom. For the model of inflated variance, analogously the validity of the condition

$$\hat{e}_{ji}^2 \leq 2 \times F_{1-\alpha/n}(1, n - m - 1) \quad (6.119)$$

means that influential points are absent. Here $F_{1-\alpha/n}(1, n - m - 1)$ means the $100(1 - \alpha/n)\%$ quantile of the central F -distribution with 1 and $(n - m - 1)$ degrees of freedom. On the basis of these two tests, an approximate rule may be formulated: strongly influential points have squared Jack-knife residuals \hat{e}_{ji}^2 greater than 10.

Problem 6.31. *Comparison of various diagnostics for identification of influential points*

For the outlier from sample C and for the high leverage point from sample D (Problem 6.8) calculate the following five diagnostics: DF_i , D_i , $LD_i(\mathbf{b})$, $LD_i(\hat{\sigma}^2)$ and $LD_i(\mathbf{b}, \hat{\sigma}^2)$.

Data: from Problem 6.8.

Solution: The calculated diagnostics for identification of influential points are listed in Table 6.10. It can be seen that the leverage point in sample D leads to the indefinite relation 0/0 for D_i and DF_i and a computer interpreted it as zero. Even the characteristics LD_i do not indicate the leverage point in sample D.

Table 6.10. Comparison of five diagnostics for identification of influential points

Sample	x_i	y_i	D_i	DF_i	$LD_i(\mathbf{b})$	$LD_i(\hat{\sigma}^2)$	$LD_i(\mathbf{b}, \hat{\sigma}^2)$
C	13	12.75	1.39	670	2.97	1.81×10^6	2.37×10^6
D	19	12.5	0	0	0	4.84×10^{-2}	4.84×10^{-2}

Conclusion: If the influential points are leverage points, then $\hat{e}_i = 0$ and $H_{ii} = 1$. Detection of these points depends on calculation of indefinite relations by a computer.

To test for influential points, some diagnostic graphs may be used:

- The index graph (IG) shows the characteristics of influential points as a function of index i of the point. These graphs may also be plotted for elements of the projection matrix, H_{ii} , etc.
- The L-R graph introduced by Gray [28] has on the y -axis the squared residuals $\hat{e}_{Ni}^2 = \hat{e}_i^2 / \text{RSC}$ and on the x -axis the elements H_{ii} . All the points will lie under the hypotenuse of the triangle with a 90° angle in the origin of the two axes and the hypotenuse defined by the limiting equality $H_{ii} + \hat{e}_{Ni}^2 = 1$.

Most of the characteristics of influential points may be expressed in the form

$$K(m, n) \times f(H_{ii}, \hat{e}_{Ni}^2)$$

where $K(m, n)$ is a constant depending only on m and n . Therefore the characteristic DF_i [Eq. (6.111)] can be rewritten as

$$DF_i = \sqrt{n - m - 1} \times \sqrt{\frac{H_{ii} \hat{e}_{Ni}^2}{(1 - H_{ii})(1 - H_{ii} - \hat{e}_{Ni}^2)}} \quad (6.120)$$

In the L-R graph, contours of the same critical influence are plotted, and the locations of individual points are compared with them. It may be determined from Eq. (6.120) that for the characteristics DF_i , the contours are hyperbolic, as described by the equation

$$y = \frac{(2x - x^2 - 1)}{x(K - 1) - 1}$$

where $K = n(n - m - 1)/(c^2m)$ and c is a constant. For $c = 2$, the constant K corresponds to the limit $2/\sqrt{m/n}$. The constant c is usually equal to 2, 4 or 8. L-R graphs for other characteristics of influential points may also be drawn.

Problem 6.32. *Examination of influential points in the validation of a new analytical method*

Use the L-R graph for DF_i to examine the influential points in Problem 6.7 (validation of a new analytical method by a comparison with a standard one).

Data: from Problem 6.7

Solution: Figure 6.37a shows the L-R graph for DF_i . This indicates that points 18, 20 and 24 are strongly influential. With these three points omitted, the regression equation $y = 9.413 (\pm 5.67) + 0.876 (\pm 0.016)x$ is estimated, with determination coefficient $\hat{R}^2 = 0.994$. The standard deviations of parameter estimates are given in brackets. When these results are compared with those of Problem 6.7, it may be concluded that elimination of influential points will not significantly affect the parameter estimates, but does affect their variances.

The standard deviation of the residuals for the original model with $n = 24$ points, $\hat{\sigma} = 39.54$, decreased on elimination of points 18, 20 and 24, to the value $\hat{\sigma} = 17.24$. Omitting three influential points caused a significant decrease in the quadratic error of prediction from $\text{MEP} = 1942$ to $\text{MEP} = 333.6$.

Figure 6.37b shows the regression model with the 95% confidence bands. If these are compared with those on Fig. 6.9, the confidence band can be seen to have narrowed.

Conclusion: The L-R graph permits easy identification of influential points. Elimin-

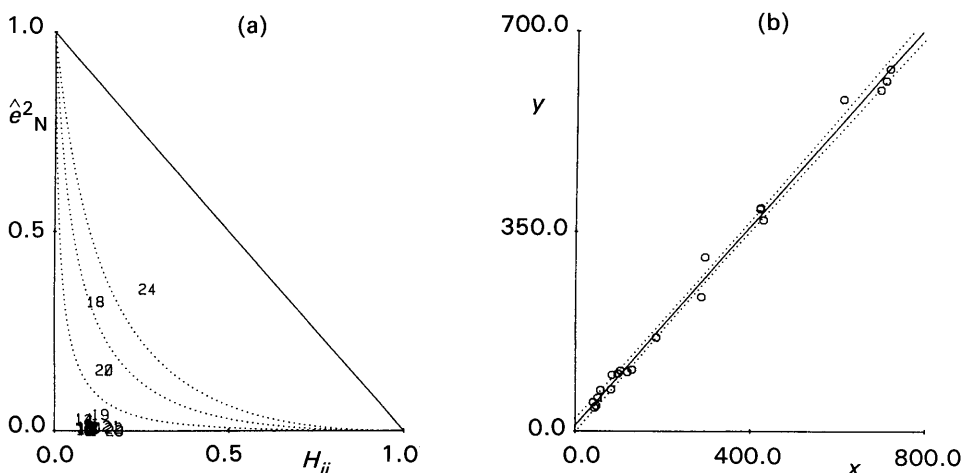


Fig. 6.37—(a) The L-R graph for the diagnostic DF_i , and (b) comparing a new analytical method with the standard one, with three influential points omitted.

ating influential points causes an improvement in the interval estimates, and this also affects the results of statistical tests.

6.5.3 Examination of a proposed regression model

The quality of a proposed model can be considered in case of one controllable variable x directly from the scatter plot of y vs. x . In the case of more controllable variables, scatter plots can falsely indicate nonlinearity in a linear model (cf. Problem 6.24). There are many various plots for considering y on x_j but we limit the choice here to (a) partial regression leverage plots, and (b) partial residual plots. Both plots are augmented here by the graph of residual \hat{e} vs. prediction \hat{y} , which can indicate a false model when the points form a nonlinear pattern.

6.5.3.1 Partial regression leverage plots

Belsey [21] named these graphs partial regression leverage plots (PRL plots) and considers them as the basic computer tools for interactive analysis of regression models. They permit classification of the quality of a regression model proposed and also indicate the presence of an influential point and lack of fulfillment of the assumptions of the classical least-squares method. They show the dependence between y and a selected controllable variable x_j when the other controllable variables forming columns in the matrix $\mathbf{X}_{(j)}$ are kept constant. By the symbol $\mathbf{X}_{(j)}$ we mean a matrix formed by leaving out the j th column \mathbf{x}_j .

To discuss the properties of these plots, we assume the regression model (6.5b) expressed in the form

$$\mathbf{y} = \mathbf{X}_{(j)}\boldsymbol{\beta}^* + \mathbf{x}_j c + \boldsymbol{\varepsilon} \quad (6.121)$$

where $\boldsymbol{\beta}^*$ is of dimension $(m-1) \times 1$ and c is the regression parameter of the j th variable. On projecting both sides of Eq. (6.121) into a space orthogonal to the space spanned by the columns of matrix $\mathbf{X}_{(j)}$, we obtain

$$\mathbf{P}_{(j)}\mathbf{y} = \mathbf{P}_{(j)}\mathbf{x}_j c + \mathbf{P}_{(j)}\boldsymbol{\varepsilon} \quad (6.122)$$

In Eq. (6.122), the product $\mathbf{P}_{(j)}\mathbf{X}_{(j)}$ is equal to zero. The projection matrix $\mathbf{P}_{(j)} = \mathbf{E} - \mathbf{H}_{(j)}$ leads to a projection into the space of the residuals. From Eq. (6.92) it follows that

- (a) the term $\hat{\mathbf{v}}_j = \mathbf{P}_{(j)}\mathbf{x}_j$ is the residual vector of regression of one variable \mathbf{x}_j on the other variables which form columns of the matrix $\mathbf{X}_{(j)}$;
- (b) the term $\hat{\mathbf{u}}_j = \mathbf{P}_{(j)}\mathbf{y}$ is the residual vector of regression of variable \mathbf{y} on other variables which form columns of the matrix $\mathbf{X}_{(j)}$.

The mean value $E(\hat{\mathbf{u}}_j)$ is then given by

$$E(\hat{\mathbf{u}}_j) = cE(\hat{\mathbf{v}}_j) \quad (6.123)$$

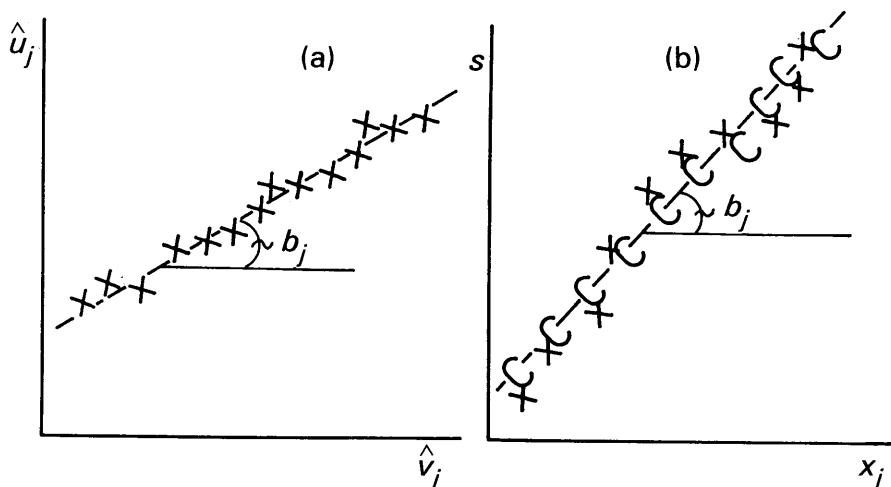


Fig. 6.38—(a) A partial regression leverage plot, and (b) a partial residual plot.

The dependences of \hat{u}_j on \hat{v}_j form the *partial regression leverage plots*.

When Eq. (6.121) is valid, Eq. (6.123) is linear with zero intercept. The slope estimate obtained by the least-squares method is calculated from

$$\hat{c} = \frac{\hat{\mathbf{u}}_j^T \hat{\mathbf{u}}_j}{\hat{\mathbf{v}}_j^T \hat{\mathbf{v}}_j} = \frac{\mathbf{x}_j^T \mathbf{P}_{(j)} \mathbf{y}}{\mathbf{x}_j^T \mathbf{P}_{(j)} \mathbf{x}_j} \quad (6.124)$$

After some rearrangements it may be shown that the slope estimate \hat{c} is identical with the estimate b_j determined by the classical least-squares method for *unpartitioned model* $E(y/x) = \mathbf{X}\beta$. Moreover, an important equality

$$\hat{e} = \hat{\mathbf{u}}_j - \hat{\mathbf{v}}_j \hat{c} \quad (6.125)$$

shows how the residual \hat{e} from the least-squares method (6.92) is connected with the partial residuals \hat{u}_j and \hat{v}_j .

The partial regression leverage plots have the following properties.

- The slope \hat{c} in the PRL plot is identical with the estimate b_j in an unpartitioned model and the intercept is equal to zero. This linear dependence is valid only when the proposed model [Eq. (6.121)] is correct.
- The correlation coefficient between \hat{v}_j and \hat{u}_j corresponds to the partial correlation coefficient $R_{yx_j}(\mathbf{x})$.
- Residuals corresponding to a regression straight line in the PRL plot are identical with residuals for an unpartitioned model.
- In the PRL plot the influential points stand out, and also any violation of the assumptions for the least-squares method, for example, about homoscedasticity.

Partial regression leverage plots have also some disadvantages.

- On the x -axis, the co-ordinates \hat{v}_j are not in the original scale of variable x_j . If

there is, for example, some scatter in the residuals of functions x_j , the PRL plots may not indicate it.

- (b) If individual controllable variables (in the columns of matrix \mathbf{X}) are strongly correlated, the PRL plot may not indicate correctly the nonlinearity, so a false hypothesis of the model may be proposed (6.121).

The partial regression leverage plots are in the standard output of the regression module of CHEMSTAT, because they correctly indicate various types of influential points.

Problem 6.33. *Application of partial regression leverage plots*

Construct PRL plots for a linear regression model with the simulated data from Problem 6.24.

Data: Generated from Problem 6.24

Solution: The PRL plots for variables x_1 and x_2 are shown in Fig. 6.39. The linear course and the zero residuals show that the data are in accord with a linear function of x_1 and x_2 . The strong multicollinearity between the variables does not influence their course.

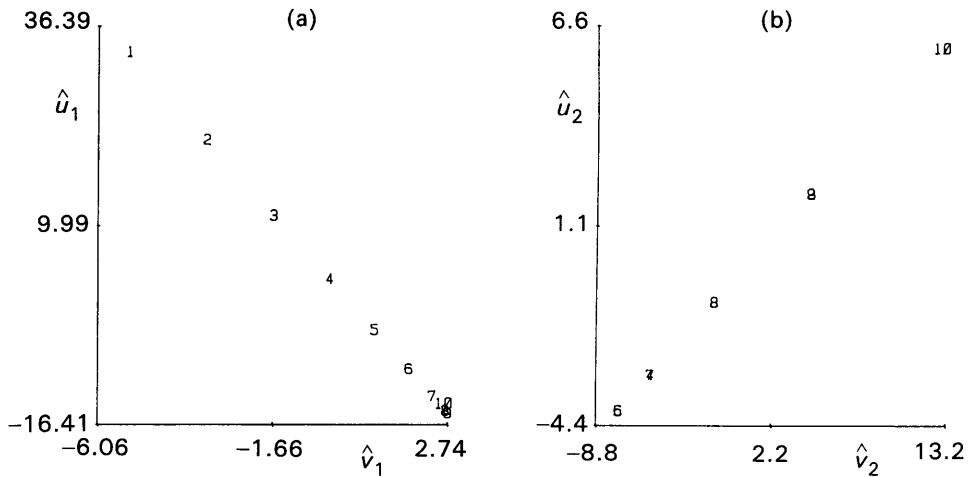


Fig. 6.39—The partial regression leverage plots for (a) variable x_1 , and (b) variable x_2 .

Conclusion: The linearity of all partial regression leverage plots proves the correctness of a proposed regression model. The quality of estimates may be classified according to the spread of points around the regression straight line in partial regression leverage plots.

6.5.3.2 Partial residual plots

Partial residual plots are also termed “component + residual” plots. Rewriting Eq. (6.125) in the form

$$\hat{\mathbf{u}}_j = \hat{\mathbf{e}} + b_j(\mathbf{E} - \mathbf{H}_{(j)})\mathbf{x}_j \quad (6.126)$$

gives the partial regression leverage plot expressed as a dependence of $\hat{e} + b_j(\mathbf{E} - \mathbf{H}_{(j)})\mathbf{x}_j$ on $(\mathbf{E} - \mathbf{H}_{(j)})\mathbf{x}_j$. The *partial residual plot* is the special case $\mathbf{H}_{(j)} = 0$. It is, in fact, a dependence of partial residuals s on variable \mathbf{x}_j . For variable s we have

$$s = \hat{e} + \mathbf{b}\mathbf{x}_j = y - \sum_{k \neq j}^m \mathbf{x}_k b_k \quad (6.127)$$

where \mathbf{x}_k is the k th column of matrix \mathbf{X} .

When the regression model contains an intercept, the modified partial residuals may be used

$$s_i^* = \hat{e}_i + (x_{ij} - \bar{x}_j)b_j + \bar{y} \quad (6.128)$$

where \bar{x}_j , \bar{y} are the arithmetic averages of variables x_j and y .

In "component + residual" plots a deterministic component is plotted separately.

$$c_{ij} = (x_{ij} - \bar{x}_j)b_j, \quad i = 1, \dots, n \quad (6.129)$$

which is usually marked on a plot by the letter "C". The partial residual $s_i = c_{ij} + \hat{e}_i$, $i = 1 \dots, n$, are in this plot marked by crosses. If \mathbf{x}_j is orthogonal to all the columns of matrix $\mathbf{X}_{(j)}$, then $\hat{\mathbf{v}}_j = \mathbf{x}_j$ and the partial regression leverage plot would be identical to the partial residual plot. The partial residual plots provide rather different information from the partial regression leverage plots. Partial residual plots have the following properties:

- the slope of s vs. \mathbf{x}_j is equal to \mathbf{b}_j and the intercept is zero. The linear dependence shows the suitability of proposed variable \mathbf{x}_j in the model;
- the residuals of these regression lines are the residuals \hat{e}_i for the unpartitioned model;
- if the angle between \mathbf{x}_j and some columns of matrix $\mathbf{X}_{(j)}$ is small (multicollinearity) the partial residual plot has falsely small scatter around the regression line $\mathbf{b}_j\mathbf{x}_j$, and the effect of influential points is suppressed.

Partial residual plots are recommended for indication of different types of nonlinearity in the case of a poorly proposed regression model.

Problem 6.34. Building partial residual plots

Construct the partial residual plots for the linear regression model and simulated data from Problem 6.24.

Data: from Problem 6.24

Solution: The partial residual plots for variables x_1 and x_2 are drawn in Fig. 6.40.

The linear course together with the zero residuals \hat{e}_i show again the linearity with respect to variables x_1 and x_2 . Since the x -axes are not transformed, the magnitude of slopes may be considered. This magnitude should correspond to parameter estimate b_j in the regression model.

Conclusion: The linearity in all partial residual plots shows the correctness of the regression model proposed. It is recommended to combine examination of the partial regression leverage plots with the partial residual plots.

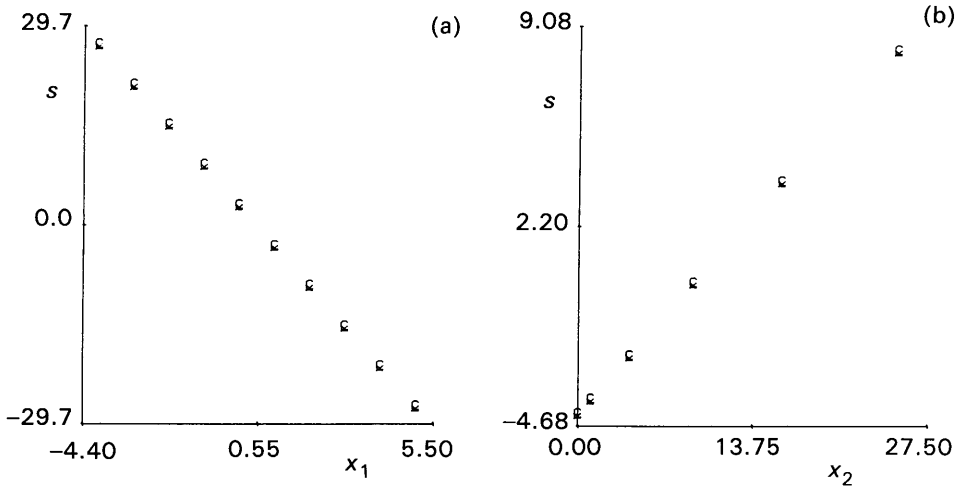


Fig. 6.40—Partial residual plots for (a) variable x_1 , and (b) variable x_2 .

Problem 6.35. *Examination of the model for the relationship between rubber composition and its abrasion resistance*

Make a graphical examination of a proposed linear model expressing the relationship between the composition of rubber and its abrasion resistance (Problem 6.6). Identify any influential points.

Data: from Problem 6.6

Solution: Figure 6.41 shows the partial regression leverage plots and Fig. 6.42 the partial residual plots for x_1 and x_2 .

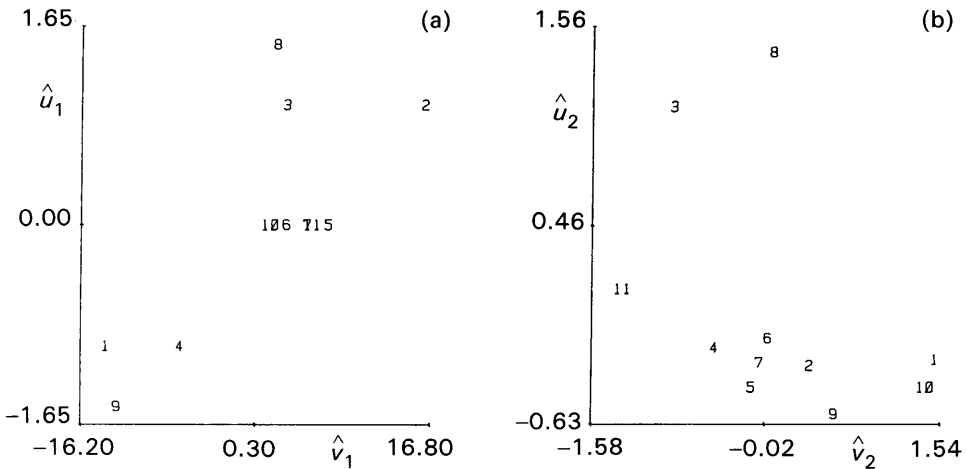


Fig. 6.41—Partial regression leverage plot for (a) the variable x_1 , and (b) the variable x_2 .

Because of orthogonality of the two variables, the plots in Figs. 6.41 and 6.42 are nearly the same. The variable x_1 is not significantly affected, as in both plots 6.41a

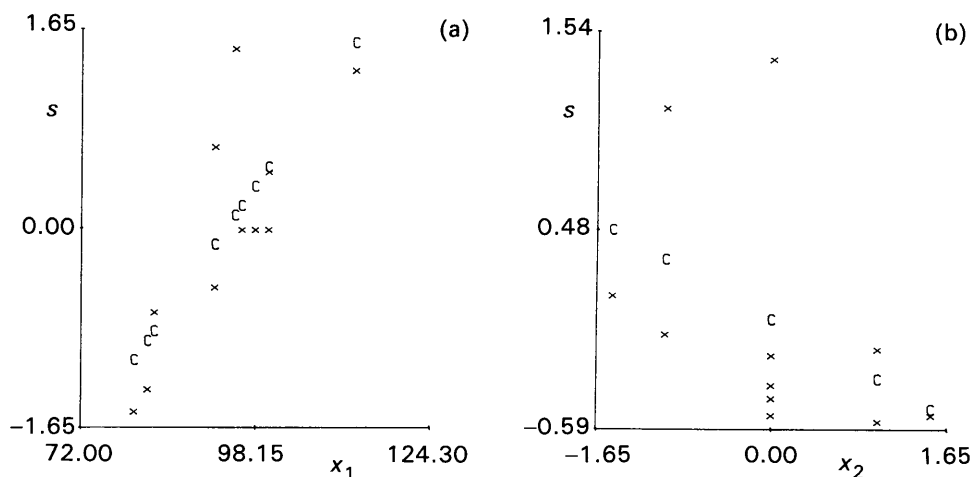


Fig. 6.42—Partial residual plot for (a) the variable x_1 , and (b) the variable x_2 .

and 6.42a the points form a random pattern. Variable x_2 shows a distinct trend which may result either from nonlinearity or from outliers in the data, and particularly point 8. In Fig. 6.43 the plot of residuals \hat{e}_i vs. the prediction \hat{y} shows a random pattern of points, proving that the proposed model is suitable, despite two points, 8 and 2, seeming to be outliers.

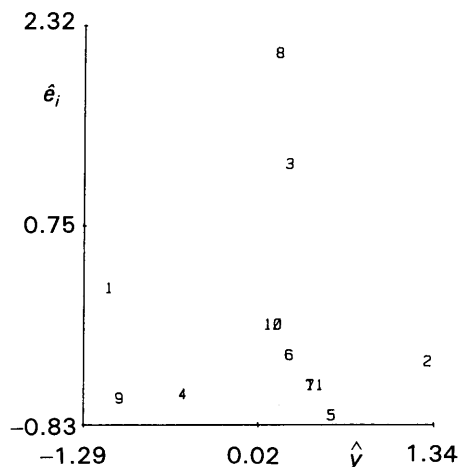


Fig. 6.43—Plot of residuals \hat{e}_i vs. prediction \hat{y} .

Both graphs in Fig. 6.44 show significant influence from point 8 and also from points 3 and 1. The values of the Jack-knife residuals \hat{e}_{ji} do not indicate strongly influential points, because the maximum value \hat{e}_{ji} for $i = 8$ is $\hat{e}_{j8} = -2.404$.

Conclusion: For a small sample size it is not possible to consider whether the model has been correctly proposed. From the graphs, it can be concluded that the data

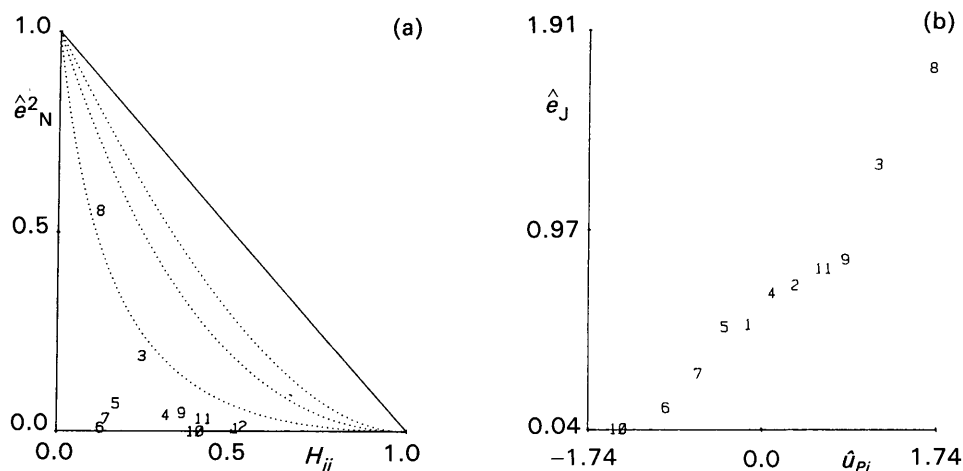


Fig. 6.44—(a) The L-R graph for DF_b , and (b) the Q-Q plot for the Anderson distance A_i .

may contain the outliers. Repetition of experiments proved that the linear model is not correct.

6.5.3.3 Sign test for model specification

To check a proposed regression model with reference to the data, all tests of specification (linearity) from Section 6.3 may be applied. A simple test based on the residuals \hat{e}_i is the *sign test*. Incorrectness of a proposed model causes non-randomness of residuals, and this non-randomness may be tested by a sign test. The number of sequences n_U of the same sign of residuals is estimated, e.g. for residuals $-1, -1, 1, -1, 1, 2, 1$ the number of sequences is equal to 4, $\hat{n}_U = 4$. Then the number of residuals with positive sign (n_+) and negative sign (n_-) is determined. For medium sample sizes the theoretical number of sequences n_t and its variance D_t are defined by

$$n_t = 1 + \frac{2n_+n_-}{n_+ + n_-} \approx 1 + \frac{n}{2} \quad (6.130)$$

$$D_t = \frac{2n_+n_-(2n_+n_- - n_+ + n_-)}{(n_+ + n_-)^2(n_+ + n_- - 1)} \approx \frac{n}{4}$$

When $n_U < n_t - 2\sqrt{D_t}$, there is a trend in the residuals and the model is incorrect.

Problem 6.36. Examination of a proposed model by the sign test

For samples A, B, C and D from Problem 6.8, test for correctness of the proposed model of a regression straight line.

Data: from Problem 6.8

Solution: Table 6.11 lists the numbers of sequences for samples A, B, C and D.

From the table it is evident that small values of \hat{n}_U (less than $n_t - 2\sqrt{D_t} \approx 4.84$) for samples B and C correspond to non-randomness of residuals and also the incorrectness of the proposed straight line model.

Table 6.11. The number of sequences for samples A, B, C and D

Data sample	A	B	C	D
The number of sequences \hat{n}_U	7	3	4	7

Conclusion: The sign test can test for non-randomness of residuals, caused either by a false model (sample B) or by outliers (sample C).

6.5.4 Examination of conditions for the least-squares method

The violation of the basic conditions for the least-squares method is discussed in Section 6.6. In this section the graphical diagnostics for indication of heteroscedasticity, autocorrelation and non-normality of errors ε are described.

6.5.4.1 Heteroscedasticity

Heteroscedasticity often appears in instrumental data measured in the chemical laboratory. The variance of measurement is usually an increasing function of variable y because the relative precision of the measurement is constant. This type of heteroscedasticity may be detected by the plot of \hat{e}_i^2 vs. \hat{y}_i , which gives a pattern of typically linear or nonlinear shape. If the measurement variance is dependent on x_j the plot \hat{e}_i^2 vs. x_{ij} leads also to a linear or nonlinear pattern. The heteroscedasticity may be detected by plots of residuals or by partial regression leverage plots.

Identification of heteroscedasticity in data is based on the idea that the variance of a measured quantity at the i th point is an exponential function of the variable $\mathbf{x}_i\boldsymbol{\beta}$ of the type

$$\sigma_i^2 = \sigma^2 \exp(\lambda \mathbf{x}_i \boldsymbol{\beta})$$

where \mathbf{x}_i is the i th row of matrix \mathbf{X} . The test for homoscedasticity is carried out by checking the null hypothesis $H_0: \lambda = 0$. Cook and Weisberg [30] introduced the test criterion

$$S_f = \frac{\left[\sum_{i=1}^n (\hat{y}_i - \bar{y}_p) \hat{e}_i^2 \right]^2}{2\sigma^4 \sum_{i=1}^n (\hat{y}_i - \bar{y}_p)^2} \quad (6.131)$$

where $\bar{y}_p = \left(\sum_{i=1}^n \hat{y}_i \right) / n$. When the null hypothesis is valid, the test statistic S_f has approximately the $\chi^2(1)$ distribution with one degree of freedom.

The corresponding diagnostic plot has the squares of standardized residuals \hat{e}_{Si}^2 on the y -axis and $(1 - H_{ii})\hat{y}_i$ on the x -axis. If heteroscedasticity is not present in the data, a random pattern of points appears. When heteroscedasticity is present, a wedge-shaped pattern appears and most of the points are located in this part of plot.

Problem 6.37. Examination of the homoscedasticity assumption in validation of a new analytical method

Figure 6.9 shows that for larger values of x the variance of points round the regression line increases. Are the data from Problem 6.7 homoscedastic or heteroscedastic?

Data: from Problem 6.7

Solution: The test criterion $S_f = 119.45$ [Eq. (6.131)] has a higher value than the quantile $\chi_{0.975}^2(1) = 5.02$ and the null hypothesis $H_0: \lambda = 0$ is rejected. The data exhibit heteroscedasticity. The plot of \hat{e}_i^2 against \hat{y}_i Fig. 6.45a shows a recognizable systematic trend indicating heteroscedasticity. Points 18, 20 and 24 are influential points. The typical wedge-shaped pattern of points in the plot of \hat{e}_{si}^2 against $(1 - H_{ii})\hat{y}_i$ in Fig. 6.45b also proves heteroscedasticity. When points 18, 20 and 24 are left out, the test criterion $S_f = 2.75$ is lower than $\chi_{0.975}^2(1) = 5.02$ and heteroscedasticity is not proved.

Conclusion: Plots indicating heteroscedasticity can also detect whether the heteroscedasticity is caused by the presence of influential points.

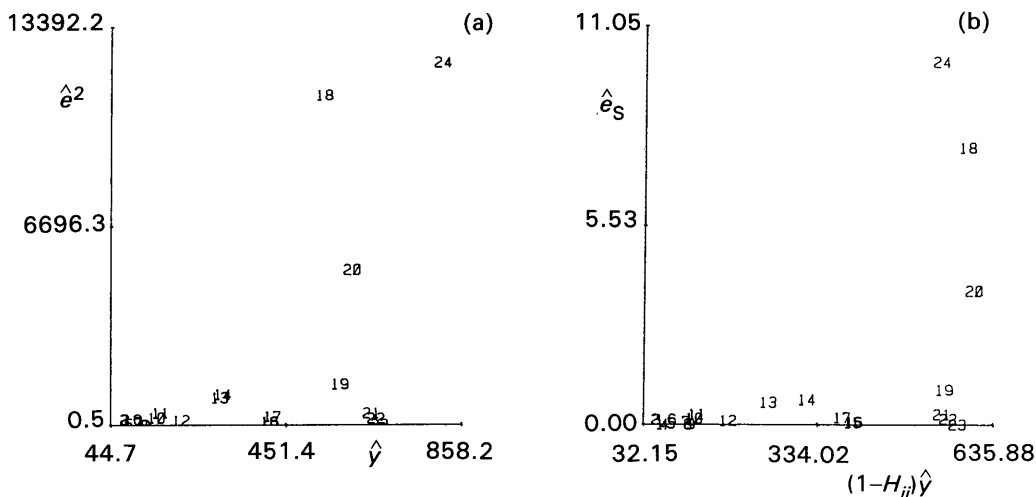


Fig. 6.45—Plots for testing for heteroscedasticity (a) \hat{e}_i^2 vs. \hat{y}_i , and (b) \hat{e}_{si}^2 vs. $(1-H_{ii})\hat{y}_i$.

6.5.4.2 Autocorrelation

When data are a time series, the errors ε are not independent but are correlated with one another. We will discuss only the most frequent case of autocorrelation of the first order, the *autoregressive process of the first order* AR(1), described by the expression

$$\varepsilon_i = \rho_1 \varepsilon_{i-1} + u_i \quad (6.132)$$

where $u_i \approx N(0, \sigma^2)$ is an independent, random variable with constant variance and $\rho_1 \leq 1$ is the *autocorrelation coefficient of the first order*. For $\rho_1 = 1$, Eq. (6.132) defines a case of cumulative errors, which appears quite often in chemometrics. When the model $X\beta$ does not contain all the significant variables and is falsely proposed,

the mean values of the residuals correspond to an AR(1) process, with a positive autocorrelation coefficient of the first order, ρ_1 . Tests of autocorrelation can be understood as tests of accuracy of a proposed model, with reference to the number of controllable variables. From Eq. (6.132) it may be concluded that for an AR(1) process, the dependence of ε_i on ε_{i-1} is linear, with slope ρ_1 . To test for autocorrelation, the graph of $\hat{\varepsilon}_i$ against $\hat{\varepsilon}_{i-1}$ is plotted, and an approximately linear trend proves significant autocorrelation.

Classical residuals are, however, correlated even in cases when the errors ε_i are not correlated. For small sample sizes, this may lead to a false finding of linearity of the dependence $\hat{\varepsilon}_i$ vs. $\hat{\varepsilon}_{i-1}$. The use of recursive residuals $\hat{\varepsilon}_{Ri}$ is more convenient.

Problem 6.38. Autocorrelation test for kinetic data

Kinetic data for inversion of a saccharide in 1M HCL at 30° C were measured. Find out whether the autocorrelation effect in the data is caused by the method of taking samples. Use graphical tests.

Data: x is time in minutes; y is the logarithm of the fraction of saccharide remaining unreacted in the reaction mixture, multiplied by 10.

x	0	10	20	30	40	50	60	70	80
y	1	0.954	0.895	0.843	0.791	0.735	0.685	0.628	0.581

Solution: The plot of $\hat{\varepsilon}_i$ vs. $\hat{\varepsilon}_{i-1}$ in Fig. 6.46a and the plot of $\hat{\varepsilon}_{Ri}$ vs. $\hat{\varepsilon}_{Ri-1}$ in Fig. 6.46b show significant negative autocorrelation. Since both plots are similar in nature, classical residuals may be used.

Conclusion: Plots for examination of autocorrelation of residuals allow the sign and

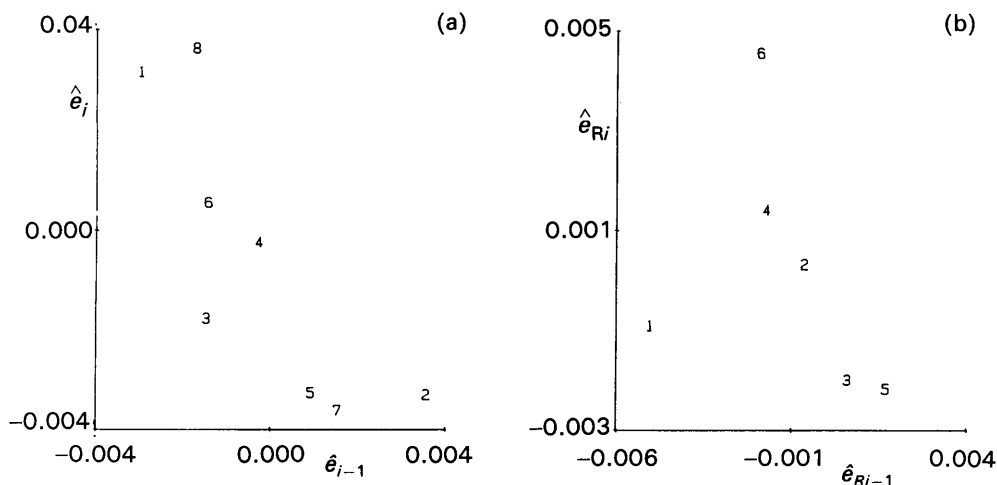


Fig. 6.46—Plots for examination of autocorrelation of the first order (a) $\hat{\varepsilon}_i$ vs. $\hat{\varepsilon}_{i-1}$, and (b) $\hat{\varepsilon}_{Ri}$ vs. $\hat{\varepsilon}_{Ri-1}$.

the magnitude of the autocorrelation coefficient of the first order ρ_1 to be estimated. Here, the classical \hat{e}_i and recursive residuals \hat{e}_{Ri} give similar results.

6.5.4.3 Normality of errors

The normality of errors is examined by a Q-Q plot containing the order statistics of classical residuals $\hat{e}_{(i)}$ in dependence on the quantile of the normalized normal distribution u_{Pi} for $P_i = i/(n+1)$. Since small samples exhibit a supernormality effect, independent recursive residuals \hat{e}_{Ri} are used instead of classical residuals, because this effect then does not exist.

To test the normality of residuals, some tests from Chapter 3 may be used. The most convenient test seems to be the Jarque-Berra test [32] which is based on the criterion

$$L(\hat{e}) = n \left[\frac{\hat{u}_3^2}{6\hat{u}_2^3} + \frac{(\hat{u}_4/\hat{u}_2^2) - 3}{24} \right] + n \left[\frac{3\hat{u}_1^2}{2\hat{u}_2} - \frac{\hat{u}_3 \times \hat{u}_1}{\hat{u}_2^2} \right] \quad (6.133)$$

where the symbol \hat{u}_j denotes the j th general moment of the sample residuals, and is defined by

$$\hat{u}_j = \frac{1}{n} \sum_{i=1}^n \hat{e}_i^j \quad (6.133a)$$

When the errors have a normal distribution, the test statistic $L(\hat{e})$ has asymptotically the $\chi^2_{1-\alpha}(2)$ distribution. When $L(\hat{e}) > \chi^2_{0.95}(2) = 5.99$, the null hypothesis H_0 about the error normality is rejected. In this test, the supernormality effect of small samples may again disturb statistical testing.

For linear models with an intercept term, $E(\hat{e}_i) = \hat{u}_1 = 0$ and the Jarque-Berra criterion can be simplified to the form

$$L(\hat{e}) = n \left[\frac{\hat{g}_1}{6} + \frac{(\hat{g}_2 - 3)^2}{24} \right] \quad (6.134)$$

where $\hat{g}_1 = \hat{u}_3^2/\hat{u}_2^3$ and $\hat{g}_2 = \hat{u}_4/\hat{u}_2^2$. This procedure however, is not convenient for small samples, because of the supernormality effect, and moreover the distribution of $L(\hat{e})$ differs from the asymptotic $\chi^2_{1-\alpha}(2)$. For small samples it is more convenient to determine the distribution of $L(\hat{e})$ from a simulation calculation for the given matrix X . As $L(\hat{e})$ is independent of the error variance σ^2 , the errors ε_i may be generated from the normalized normal distribution $N(0, 1)$.

Problem 6.39. Examination of normality for four samples

For samples A, B, C and D in Problem 6.8, use the Jarque-Berra test criterion $L(\hat{e})$ to test for normality.

Data: from Problem 6.8

Solution: For four samples the Jarque-Berra test criterion $L(\hat{e})$ is the last column in Table 6.7 (Problem 6.25). This test disproved normality only for sample C. The other samples A, B and D exhibit normality of residuals.

Conclusion: Tests of normality of residuals do not prove incorrectness of a proposed regression model or unsuitability of data. When normality is not proved, the presence of outliers is often the cause; the Q-Q plot is then recommended for detection of influential points.

6.6 PROCEDURES WHEN CONDITIONS FOR LEAST-SQUARES ARE VIOLATED

In Section 6.2, seven conditions were mentioned which must be met if the least-squares method is to give the best unbiased linear estimates of parameters. The construction of confidence intervals and hypothesis tests also depend on these conditions being satisfied. In the chemical laboratory some of the conditions, however, are not met. In this chapter we give our attention to regression procedures when

- (1) some restrictions are placed on the parameters;
- (2) the covariance matrix of errors is not diagonal and data do not exhibit the same variance;
- (3) the matrix $\mathbf{X}^T\mathbf{X}$ is ill-conditioned because of multicollinearity;
- (4) the distribution of data is not normal and some influential points exist in data;
- (5) the independent variables x are also subject to random errors.

The most important diagnostic procedures for identification of violations of the least-squares conditions are described in Section 6.5. This section gives a modified procedure for parameter estimation and some special tests.

6.6.1 Restrictions placed on the parameters

In many chemometrics problems some restrictions are placed on parameters because of their physical meaning and chemical interpretation. Positive values, for example, are often requested for most chemical parameters. The regression procedure with a restriction depends on whether the restrictions are *precise* (deterministic as they are fixed numbers) or *statistical* (they are random numbers). The restrictions can be stated in form of equalities, or inequalities when they concern restricted intervals.

The most frequent request in chemometrics problems is that the regression line should fit the data and also pass through the origin. This last request can be fulfilled by omitting the intercept term. We will discuss cases when the parameter restrictions are given as an equality, and when parameters should be numerically greater than a given limit.

A restriction in the form of an equality

This group of restrictions includes the following requests about parameters,

- (a) some parameters should reach specified values;
- (b) some parameters should have a specified mutual ratio;
- (c) the sums or differences of some parameters should be equal to a given number; and
- (d) the regression model should also fit certain points specified by co-ordinates.

To satisfy these four requests, the condition of linearity may be formulated as

$$\begin{array}{rcl}
P_{11}\beta_1 + P_{12}\beta_2 + \dots + P_{1m}\beta_m & = & p_1 \\
P_{21}\beta_1 + P_{22}\beta_2 + \dots + P_{2m}\beta_m & = & p_2 \\
\vdots & & \vdots \\
P_{k1}\beta_1 + P_{k2}\beta_2 + \dots + P_{km}\beta_m & = & p_k
\end{array} \tag{6.135}$$

or written in a matrix, notation

$$\mathbf{P}\boldsymbol{\beta} = \mathbf{p} \tag{6.135a}$$

where \mathbf{P} is the matrix of dimension $(k \times m)$ of known coefficients and \mathbf{p} is the vector of dimension $(k \times 1)$ of known components, estimated on the basis of the requested restrictions. The mathematical condition for solution is that the rank of matrix \mathbf{P} should be equal to k and also $k < m$. This means that rows of matrix \mathbf{P} are linearly independent.

To estimate parameters \mathbf{b}_R which fulfil the condition of the minimum of the least-squares method with a restriction (6.135a), the technique of Lagrange multipliers is used. This method involves minimization of the conditioned sum of squares

$$U_R = (\mathbf{y} - \mathbf{X}\mathbf{b}_R)^T(\mathbf{y} - \mathbf{X}\mathbf{b}_R) + \lambda^T(\mathbf{p} - \mathbf{P}\mathbf{b}_R) \tag{6.136}$$

where λ is the vector of Lagrange multipliers of dimension $(k \times 1)$; its estimate is also sought. The method is called *the conditioned least-squares method* (CLS). As in the classical least-squares method (LS), the estimates \mathbf{b}_R and $\hat{\lambda}$ may be found with the use of the first derivative of function U_R according to these parameters,

$$\frac{\delta U_R}{\delta \mathbf{b}} = -2\mathbf{X}^T\mathbf{y} + 2\mathbf{X}^T\mathbf{X}\mathbf{b}_R - \mathbf{P}^T\hat{\lambda} = 0 \tag{6.137a}$$

$$\frac{\delta U_R}{\delta \lambda} = \mathbf{p} - \mathbf{P}\mathbf{b}_R = 0 \tag{6.137b}$$

This equation defines $(n + k)$ linear equations according to parameters $\hat{\lambda}$ and \mathbf{b}_R . After rewriting we obtain the estimate $\hat{\lambda}$ in the form

$$\hat{\lambda} = 2[\mathbf{P}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{P}^T]^{-1}(\mathbf{p} - \mathbf{P}\mathbf{b}) \tag{6.138}$$

where $\mathbf{b} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$ is the parameter estimate found by classical least-squares. The estimate of restricted parameters is calculated from

$$\mathbf{b}_R = \mathbf{b} - (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{P}^T[\mathbf{P}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{P}^T]^{-1}(\mathbf{P}\mathbf{b} - \mathbf{p}) \tag{6.139}$$

When a given parameter restriction is valid

- (a) the estimate \mathbf{b}_R is unbiased;
- (b) its covariance matrix is given by

$$D(\mathbf{b}_R) = D(\mathbf{b})[\mathbf{E} - \mathbf{P}^T\mathbf{S}\mathbf{P}(\mathbf{X}^T\mathbf{X})^{-1}] \tag{6.140}$$

where $D(\mathbf{b})$ is a covariance matrix for estimates \mathbf{b} by the classical least-squares method (6.16), \mathbf{E} is the unit matrix and $\mathbf{S} = [\mathbf{P}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{P}^T]^{-1}$; and

- (c) the unbiased estimate of the residual variance σ^2 is calculated from

$$\hat{\sigma}_R^2 = \frac{(\mathbf{y} - \mathbf{X}\mathbf{b}_R)^T(\mathbf{y} - \mathbf{X}\mathbf{b}_R)}{n - m + k} \quad (6.141)$$

If the errors have a normal distribution, confidence intervals and tests of significance may be constructed as in Section 6.3. From Eq. (6.140) it may be concluded that the variance $D(b_{Rj})$ for restricted parameter estimates are always smaller than for $D(b_j)$. The main task here is to check the validity of Eq. (6.135a), i.e. $H_0: \mathbf{P}\boldsymbol{\beta} - \mathbf{p} = 0$ against $H_A: \mathbf{P}\boldsymbol{\beta} - \mathbf{p} \neq 0$. The hypothesis H_0 may be tested by the classical Fisher–Snedecor F -test with the test criterion

$$F_o = \frac{(\mathbf{P}\mathbf{b} - \mathbf{p})^T \mathbf{S}(\mathbf{P}\mathbf{b} - \mathbf{p})/k}{(\mathbf{y} - \mathbf{X}\mathbf{b})^T(\mathbf{y} - \mathbf{X}\mathbf{b})/(n - m)} \quad (6.142)$$

which, if H_0 is valid, has the Fisher–Snedecor F -distribution with k and $(n - m)$ degrees of freedom. When $F_o > F_{0.95}(k, n - m)$ the parameter restrictions are not suitable for the given data, and the estimate \mathbf{b}_R is biased. For a small bias, the variances often decrease, so that the estimate \mathbf{b}_R seems to be better than the estimate \mathbf{b} [33].

From the computational point of view, it is more convenient to express the test criterion (6.142) in terms of the residual sum of squares (Section 6.3). If we write

$$RSC = (\mathbf{y} - \mathbf{X}\mathbf{b})^T(\mathbf{y} - \mathbf{X}\mathbf{b})$$

and

$$RSC_o = (\mathbf{y} - \mathbf{X}\mathbf{b}_R)^T(\mathbf{y} - \mathbf{X}\mathbf{b}_R)$$

the Eq. (6.142) may be expressed in an equivalent form as

$$F_o = \frac{(RSC_o - RSC)/k}{RSC/(n - m)} \quad (6.143)$$

To use Eq. (6.143) both estimates \mathbf{b} and \mathbf{b}_R must be calculated and RSC and RSC_o evaluated.

Problem 6.40. *The conditional least-squares method in the case of a single parameter restriction*

Derive equations for estimation of parameters \mathbf{b}_R in a case where only one restriction is given:

$$P_1\beta_1 + \dots + P_m\beta_m = p \text{ or } \mathbf{P}\boldsymbol{\beta} = p \text{ where } \mathbf{P} \text{ is the row vector.}$$

Solution: Since the matrix \mathbf{S} contains just one element we will speak about the scalar S . The matrix $\mathbf{M} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{P}^T$ becomes the column vector \mathbf{M} . Let us introduce the matrix \mathbf{C} :

$$\mathbf{C} = (\mathbf{X}^T\mathbf{X})^{-1}$$

with elements C_{jk} . Then

$$S = \sum_{j=1}^m \sum_{k=1}^m P_j C_{jk} P_k$$

and the vector \mathbf{M} has the following elements

$$M_j = \sum_{k=1}^m C_{jk} P_k$$

Then, from Eq. (6.139) we have

$$\mathbf{b}_R = \mathbf{b} - \mathbf{M} \sum_{j=1}^m (P_j b_j - p)/S$$

Since the matrix \mathbf{S} becomes a scalar S , the expressions for the covariance matrix of estimates and the test criterion F_0 may also be simplified.

Problem 6.41. *Finding the relationship between the surface under a chromatographic peak and the ethanol concentration, by linear regression with restriction*

The relationship between the concentration of ethanol in water (x) and the corresponding area of the chromatographic peak (y) was examined, and the model $E(y/x) = \beta_1 x + \beta_2 x^2 + \beta_3$ was proposed. To have physical meaning, this curve should go through the two points with co-ordinates (0, 0) and (100, 100), corresponding to limits, the first for pure water and the second for pure ethanol. Estimate the model parameters and test whether the restrictions correspond to the given data set.

Data: x is the volume percentage of ethanol in water and y is the relative area of the chromatographic peak as a percentage.

x	10	20	30	40	50	60	70	80	90
y	8.16	15.9	22.7	31.5	39.8	49.4	59.7	70.6	83.6

Solution: Because the first restriction requires the regression curve to go through the origin, the intercept term β_3 should be equal to zero, $\beta_3 = 0$. The second restriction leads to the equation

$$100 = \beta_1 \times 100 + \beta_2 \times 100^2$$

which can be rewritten as

$$1 = \beta_1 + \beta_2 \times 100$$

On simplifying this equation by elimination of β_1 , we obtain the following equation (which is linear with respect to β_2)

$$E(y/x) = (1 - 100\beta_2)x + \beta_2 x^2 = x + \beta_2(x^2 - 100x)$$

By using the least-squares criterion and the analytical derivative we obtain

$$b_{R2} = \frac{\sum_{i=1}^n (y_i - x_i)(x_i^2 - 100x_i)}{\sum_{i=1}^n (x_i^2 - 100x_i)^2}$$

and find that $b_{R2} = 4.1199 \times 10^{-3}$ and $RSC_0 = 31.82$. The classical least-squares estimates for the regression model without restriction has the form

$$y = 2.724(\pm 0.648) + 0.557(\pm 0.0297)x + 3.726(\pm 0.29) \times 10^{-3}x^2$$

The standard deviations of the parameter estimates are given in brackets. The corresponding value of RSC is 1.555, and

$$F_0 = \frac{(31.82 - 1.555)/2}{1.555/(9 - 3)} = 58.39$$

Since F_0 is significantly greater than $F_{0.95}(2, 6) = 5.14$, the given restrictions do not correspond to the data. The two models are compared in Fig. 6.47.

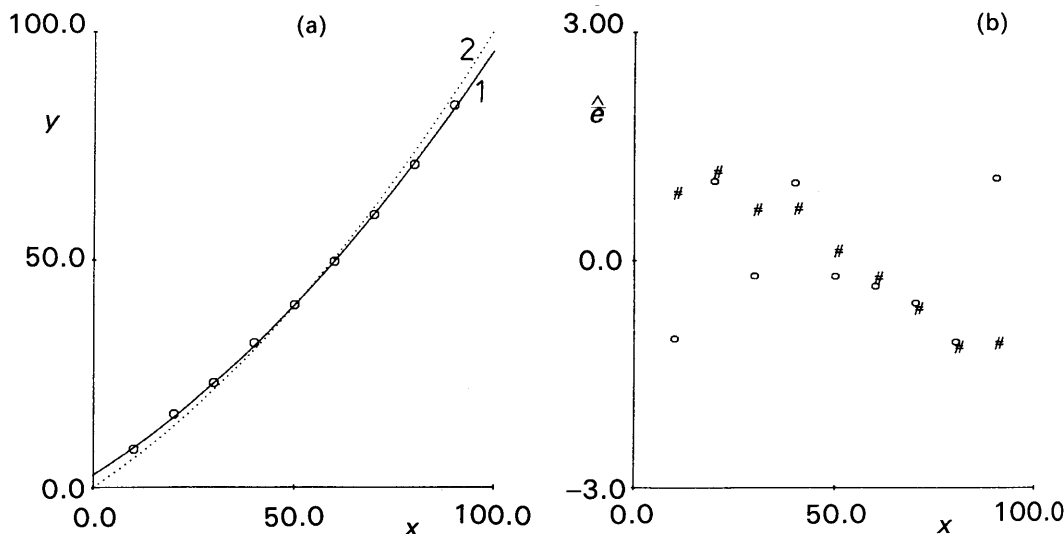


Fig. 6.47—(a) Regression model without restrictions (curve 1) and with restrictions (curve 2), fitted through the experimental points. The restriction requests the curve to go through two points, (0, 0) and (100, 100); (b) the graphical analysis of residuals: \circ model without and $\#$ model with restrictions.

Conclusion: In some cases, the restrictions given enable parameters to be derived in such a way that the method of Lagrange multipliers is not required. The statistical criterion F_0 examines whether the data are in agreement with the given restrictions.

6.6.2 The method of generalized least-squares (GLS)

In the analysis of instrumental data in the chemical laboratory, it is often found that the errors are often not independent or that they do not exhibit the same variance. The covariance matrix of errors $D(\epsilon) = C_\epsilon$ is then not equal to $\sigma^2 \times E$ and the more generalized relationship should be used

$$D(\epsilon) = C_\epsilon = \sigma^2 K \quad (6.144)$$

As the matrix K is, apart from the multiplicative constant σ^2 , the same as the

covariance matrix, it should be the case that it is symmetric and $K_{ij} = K_{ji}$. Moreover, the inequality

$$|K_{ij}| \leq \sqrt{K_{ii} K_{jj}}$$

which makes certain that matrix \mathbf{K} is positive definite. These properties allow the inversion matrix \mathbf{K}^{-1} to be expressed as the product of weight matrices \mathbf{V} in the form

$$\mathbf{K}^{-1} = \mathbf{V}^T \mathbf{V}$$

When all other conditions for the least-squares method are met, the parameter estimates may be obtained by minimization of the *generalized least-squares criterion* in the form

$$U(\mathbf{b}) = (\mathbf{y} - \mathbf{X}\mathbf{b}_G)^T \mathbf{K}^{-1} (\mathbf{y} - \mathbf{X}\mathbf{b}_G) \quad (6.145)$$

By using an analytical minimization of Eq. (6.145), the expression for the estimate \mathbf{b}_G may be derived in the form

$$\mathbf{b}_G = (\mathbf{X}^T \mathbf{K}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{K}^{-1} \mathbf{y} \quad (6.146)$$

where the index G denotes characteristics of the *method of generalized least-squares* (GLS). The estimate \mathbf{b}_G is called the *Aitken estimate*. When the weight matrix \mathbf{V} is introduced into Eq. (6.146) the resulting estimate by the GLS method is

$$\mathbf{b}_G = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{w} \quad (6.147)$$

where $\mathbf{Z} = \mathbf{V}\mathbf{X}$ and $\mathbf{w} = \mathbf{V}\mathbf{y}$. Equation (6.147) shows how the parameter estimate \mathbf{b}_G by the GLS method can easily be transformed into the estimate \mathbf{b} by the LS method with the use of simple multiplication by the weight matrices. When we multiply the Eq. (6.5b) by the weight matrix \mathbf{V} , we obtain

$$\mathbf{V}\mathbf{y} = \mathbf{V}\mathbf{X}\boldsymbol{\beta} + \mathbf{V}\boldsymbol{\varepsilon} \quad (6.148a)$$

or

$$\mathbf{w} = \mathbf{Z}\boldsymbol{\beta} + \mathbf{V}\boldsymbol{\varepsilon} \quad (6.148b)$$

The mean value $E(\mathbf{V}\boldsymbol{\varepsilon})$ is equal to zero, and the covariance matrix is calculated from

$$D(\mathbf{V}\boldsymbol{\varepsilon}) = \mathbf{V}^T \mathbf{V} \mathbf{K} \sigma^2 = \mathbf{E} \sigma^2 \quad (6.149)$$

This means that the transformed errors $\mathbf{V}\boldsymbol{\varepsilon}$ already satisfy the conditions for the classical LS method, and Eq. (6.147) may be used. With the use of variables \mathbf{Z} and \mathbf{w} , the expressions valid for the LS method, including the interval estimates and hypothesis tests, may be used. For example, the covariance matrix of estimates \mathbf{b}_G may be written, by using Eq. (6.16), as

$$D(\mathbf{b}_G) = \sigma^2 (\mathbf{Z}^T \mathbf{Z})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{K}^{-1} \mathbf{X})^{-1} \quad (6.150)$$

The estimates \mathbf{b}_G are unbiased and best in the class of all linear unbiased estimates. When the errors $\boldsymbol{\varepsilon}$ have the normal distribution $N(0, \sigma^2 \mathbf{K})$, the estimates \mathbf{b}_G are

normally distributed with mean value $E(\mathbf{b}_G) = \boldsymbol{\beta}$ and the covariance matrix is defined by Eq. (6.150). The estimate of the residual variance $\hat{\sigma}_G^2$ is calculated from

$$\hat{\sigma}_G^2 = \frac{(\mathbf{y} - \mathbf{X}\mathbf{b}_G)^T \mathbf{K}^{-1} (\mathbf{y} - \mathbf{X}\mathbf{b}_G)}{n - m} \quad (6.151)$$

When the classical LS method is used for a case when Eq. (6.144) is valid, instead of the more correct approach by the GLS method it is true that,

- (a) the parameter estimates \mathbf{b} remain unbiased;
- (b) the covariance matrix $D(\mathbf{b})$ does not correspond to the correct covariance matrix $D(\mathbf{b}_G)$, so that the estimates are not already best as they have greater variances;
- (c) the estimate of the residual variance $\hat{\sigma}^2$ will be biased.

For these reasons the interval estimates and statistical tests will give quite false results.

A special case of the GLS method is the method of weighted least-squares (WLS). The matrix \mathbf{K} is diagonal, condition 6 is valid and errors ε are independent in this case.

If K_{ii} are the diagonal elements of matrix \mathbf{K} , we can write for the diagonal element of a matrix \mathbf{V} that

$$V_{ii} = \sqrt{1/K_{ii}}$$

When this is introduced into Eq. (6.145), we obtain:

$$U(\mathbf{b}) = \sum_{i=1}^n K_{ii}^{-1} \left(y_i - \sum_{j=1}^m x_{ij} b_j \right)^2 = \sum_{i=1}^n \left[y_i V_{ii} - \sum_{j=1}^m V_{ii} x_{ij} b_j \right]^2 \quad (6.152)$$

When all variables are multiplied by the corresponding weights, the same conditions apply as for the classical LS method. Linear regression programs based on the least-squares method can readily be extended for weighted least-squares.

6.6.2.1 Heteroscedasticity

Heteroscedasticity in data means that condition 5, about constancy of variance, is violated. In chemometrics problems, nonconstancy of variance in measured data is common.

The operation of some instruments is such that there is a constant relative error, so that

$$\sigma_i^2 = \sigma_0^2 [E(y/x_i)]^2 \quad (6.153)$$

In other cases, the variance σ_i^2 may be estimated from error propagation for all operations, chemicals, glassware, procedures, etc. When the values of individual variances σ_i^2 can be exactly specified, then setting $V_{ii} = 1/\sigma_i$ permits application of the WLS method.

A large group of chemometrics problems in which there is heteroscedasticity in the data arises from the use of transformed variables. The variable y is often transformed to give a linear relationship (the linearization method). When the original

values y_i have constant variance $D(y_i) = \sigma^2$, nonlinear transformation $g(y_i)$ causes the variance $D(g(y_i))$ to be non-constant, and for the first approximation the following expression will hold

$$\sigma_i^2 = D(g(y_i)) = \left[\frac{\delta g(y_i)}{\delta y_i} \right]^2 D(y_i) \quad (6.154)$$

The variances can be equalized by introduction of weights given by

$$V_{ii} = \left[\frac{\delta g(y)}{\delta y} \right]_{y=y_i}^{-1} \quad (6.155)$$

with the use of the WLS method. The method of weighted least-squares with weights V_{ii} defined by Eq. (6.155) is called *quasilinear regression*. It is generally true that each transformation distorts the error distribution, so it is always better to use the method of nonlinear regression.

Problem 6.41. *Examination of the dependence of the solubility of Na_2SO_3 on temperature*

The dependence of the solubility of Na_2SO_3 (y) on temperature (x) may be described by the empirical expression $y = \exp(\beta_1 + \beta_2 x)$. Estimate parameters β_1 and β_2 with the use of the linearization, quasilinearization and non-linear methods of least-squares.

Data: solubility y , %, and temperature x , °C.

x	0	10	20	30	40	50	60	70	80
y	33.5	37.0	41.2	46.1	50.0	52.0	56.3	64.3	69.9

Solution: The expression $y = \exp(\beta_1 + \beta_2 x)$ can be transformed into linear form $\ln y = \beta_1 + \beta_2 x$. The weights V_{ii} will eliminate the heteroscedasticity and can be expressed by $V_{ii} = y_i$ [Eq. (6.155)]. Table 6.12 lists the parameter estimates by the three least-squares methods.

Table 6.12. Comparison of parameter estimates found by the LS, WLS and NLS methods

Method	Transformation	b_1	b_2	RSC*
LS	linearization	3.532	8812	10.22
WLS	quasilinearization	3.535	8756	10.13
NLS	nonlinear regression	3.537	8720	10.11

*in transformed variables

The estimates achieved by the quasilinearization (WLS) are in quite good agreement with those found by nonlinear regression (NLS). The precision of prediction should be considered in the original and not in the transformed variables.

Conclusion: Application of statistical weights V_{ii} from Eq. (6.155) in the WLS method increases the accuracy of parameter estimates.

To solve chemometrics problems with heteroscedasticity in the data, the procedure is usually as follows:

- (1) Identification of the presence of heteroscedasticity in the data with the use of tests from section 6.5.4.1.
- (2) Identification of the actual type of heteroscedasticity, which determines the effect of the errors variance on the variables of the regression model.
- (3) Determine parametric estimates for the known type of heteroscedasticity.

Step 1: identification of heteroscedasticity

Instead of a sample diagnostic test (Section 6.5.4) or the various plots, there are also *nonconstructive tests*, which do not require knowledge of the heteroscedasticity model, and *constructive tests*, which require knowledge of the heteroscedasticity model.

A common test is the *test of residual trend*, which has the test criterion

$$D = \sum_{i=1}^n [P(|\hat{e}_i|) - i]^2 \quad (6.156)$$

where $P(|\hat{e}_i|)$ stands for the order of absolute value of the i th residual. This criterion D is connected with the Spearman correlation coefficient ρ_s by the expression

$$\rho_s = 1 - \frac{6D}{n^3 - n} \quad (6.157)$$

The heteroscedasticity test therefore becomes the test of a null hypothesis $H_0: \rho_s = 0$ (i.e. homoscedasticity) against an alternative $H_A: \rho_s \neq 0$ (i.e. heteroscedasticity). For larger sample sizes, $n > 10$, another test-criterion can be also used

$$t_s = \sqrt{\frac{\hat{\rho}_s^2(n-2)}{1 - \hat{\rho}_s^2}} \quad (6.158)$$

which, when the null hypothesis is valid (i.e. homoscedasticity) has the Student t -distribution with $n - 2$ degrees of freedom.

The Szroeter test requires the data to be rearranged in ascending order of variance, $\sigma_{i-1}^2 \leq \sigma_i^2$, $i = 2, \dots, n$; in order to examine the values of the variable which is a monotonic function of the variances. If Eq. (6.153) is valid, the ordering is made according to the magnitude of the y -values (or prediction \hat{y}_i , respectively) in ascending order. The null hypothesis of homoscedasticity $H_0: \sigma_i^2 = \sigma_{i-1}^2$, $i = 2, \dots, n$, is tested against the alternative $H_A: \sigma_i^2 > \sigma_{i-1}^2$. The test criterion of the Szroeter test is defined by

$$Q_T = \sqrt{\frac{6n}{n^2 - 1}} \left(Q - \frac{n+1}{2} \right) \quad (6.159)$$

where

$$Q = \frac{\sum_{k=1}^n k \hat{e}_k^2}{\sum_{k=1}^n \hat{e}_k^2}$$

and residuals \hat{e}_k correspond to ordered data. The statistic Q_T has, asymptotically, the standardized normal distribution $N(0, 1)$. When $Q_T > 1.645$, heteroscedasticity is proved at the significance level $\alpha = 0.05$.

The constructive tests are based on the known model of heteroscedasticity and on significance tests.

Step 2: identification of the type of heteroscedasticity

When the matrix $\sigma^2 \mathbf{K}$ is not known, it is necessary to estimate its diagonal elements, which correspond to variances σ_i^2 . For large sample sizes, the variance estimates σ_i^2 can be replaced by the squared residuals \hat{e}_i^2 obtained by the classical LS method. This procedure is usually used in seeking parametric models of heteroscedasticity.

Horn [35] suggested application of so-called AUE estimates of variances, defined by

$$\hat{\sigma}_i^2 = \frac{\hat{e}_i^2}{1 - H_{ii}}$$

The estimates of variance $\hat{\sigma}_i^2$ may be used directly in the method of weighted least-squares (WLS) where $V_{ii} = 1/\hat{\sigma}_i$, or for examination of various types of heteroscedasticity. Three principal models of heteroscedasticity are distinguished.

(a) The *multiplicative model of heteroscedasticity* is expressed by

$$\sigma_i^2 = \sigma_0^2 \exp(\delta x_{ij}) \quad (6.160a)$$

or

$$\sigma_i^2 = \sigma_0^2 |x_{ij}|^\delta \quad (6.160b)$$

where δ is a parameter.

Instead of variable x_j in these two models, the theoretical value $E(y/x_i) = \eta_i$ may be used. The multiplicative model is valid when the dependence of $\ln(\hat{e}_i^2)$ on x_{ij} or \hat{y}_i is approximately linear. The significance test of the slope δ here corresponds to the test for the multiplicative model of heteroscedasticity.

(b) The *additive model of heteroscedasticity* is expressed by

$$\sigma_i^2 = \sigma_0^2 (1 + \delta x_{ij})^2 \quad (6.161)$$

where instead of x_j , $E(y/x_i) = \eta_i$ may be used. The additive model is valid for cases when the dependence of $|\hat{e}_i|$ on x_{ij} or \hat{y}_i is approximately linear. The significance test of the slope δ corresponds to the test for the additive model of heteroscedasticity.

(c) The *mixed model of heteroscedasticity* is expressed by

$$\sigma_i^2 = \delta_0 + \delta_1 x_{ij} \quad (6.162)$$

where instead of x_j , $E(y/x_i) = \eta_i$ may be used. The mixed model is valid for cases

when the dependence of \hat{e}_i^2 on x_{ij} or \hat{y}_i is approximately linear. The significance test of this slope δ_1 shows the presence of the mixed type of heteroscedasticity.

In chemometrics practice, the most common model seems to be the model of constant relative error (6.153) which corresponds to the multiplicative type of heteroscedasticity.

Step 3: Estimation of parameters

There are many methods of parameter estimation for linear models with heteroscedasticity in the data, but we restrict ourselves to the simplest one, i.e. the method of weighted least-squares (WLS), which is possible with most linear regression programs. The general procedure consists of the following steps:

- (1) Estimation of parameters β in the linear model by the classical LS method and estimation of the residuals \hat{e}_i .
- (2) Estimation of the parameters of the chosen heteroscedasticity type, with $\hat{\sigma}_i^2 = \hat{e}_i^2$.
- (3) Estimation of weights from $V_{ii} = 1/\hat{\sigma}_i^2$ where $\hat{\sigma}_i^2$ is the estimate of the standard deviation determined from the parametric model of heteroscedasticity.

The main problem of parameter estimation for heteroscedastic models lies in the transformation of the squared residuals \hat{e}_i^2 . This problem can be solved partly by use of quasilinear regression.

For a case defined by Eq. (6.153) the weights may be chosen such that $V_{ii} = 1/|y_i|$, or with the use of predicted values calculated by the classical least-squares method, $V_{ii} = 1/|\hat{y}_i|$.

Problem 6.42. Tests for heteroscedasticity in the validation of a new analytical method

Data from Problem 6.7, on the validation of a new analytical method by comparison with a standard one, were examined in Problem 6.37 and heteroscedasticity was proved. For the multiplicative model of heteroscedasticity [Eq. (6.153)], estimate the unknown parameters by the weighted LS method.

Data: from Problem 6.7

Solution: Data are examined for two assumptions:

- (a) Assumption of constant relative error of measurement.

With the use of the WLS method and weight $V_{ii} = 1/|y_i|$ the regression equation is found to be $y = 8.23 (\pm 5.177) + 0.879 (\pm 0.0249)x$, with determination coefficient $\hat{R}^2 = 0.983$ and quadratic error of prediction $MEP = 20560$.

Figure 6.48b illustrates "reverse" heteroscedasticity with regard to variable x : the variance decreases with increasing values of x .

- (b) Assumption of multiplicative heteroscedasticity.

The results of Problems 6.7 and 6.37 suggest that the multiplicative model of heteroscedasticity is applicable. Figure 6.49 shows the plot of $\ln \hat{e}_i^2$ vs. x_i , with the straight line $\ln \hat{e}_i^2 = 3.239 + 0.005098x$.

For parameter estimation, the WLS method was used with weights

$$V_{ii} = 1/\sqrt{\exp(3.239 + 0.005098x_i)}$$

and the regression equation was estimated as:

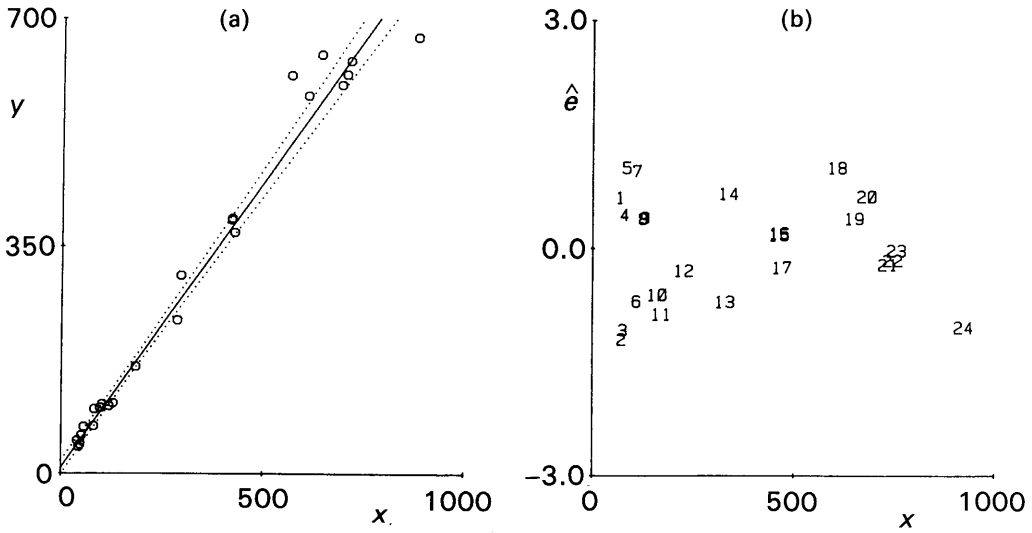


Fig. 6.48—(a) Regression model with 95% confidence interval of prediction, and (b) the graphical examination of residuals \hat{e} . The weight $V_{ii} = 1/|y_i|$ is used.

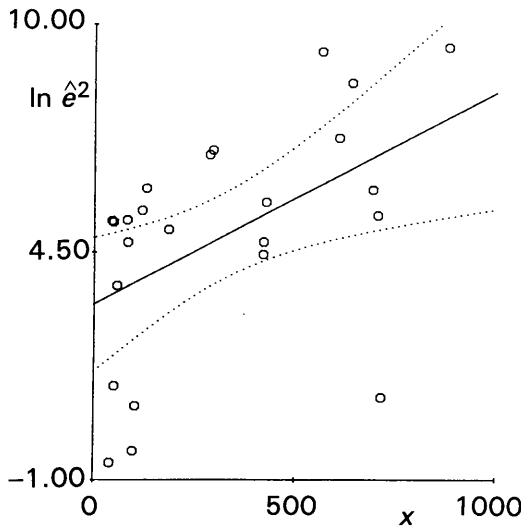


Fig. 6.49—The plot of $\ln \hat{e}_i^2$ vs. x_i indicates a multiplicative model of heteroscedasticity.

$$y = 7.937 (\pm 6.898) + 0.895 (\pm 0.0259)x$$

with determination coefficient $\hat{R}^2 = 0.982$ and the mean quadratic error $MEP = 1410$.

In Fig. 6.50, the residuals form a random pattern, and therefore the heteroscedasticity has been removed.

Conclusion: It was found that application of weights according to $V_{ii} = 1/|y_i|$ is not always the best solution. It is better to determine the actual type of heteroscedasticity.

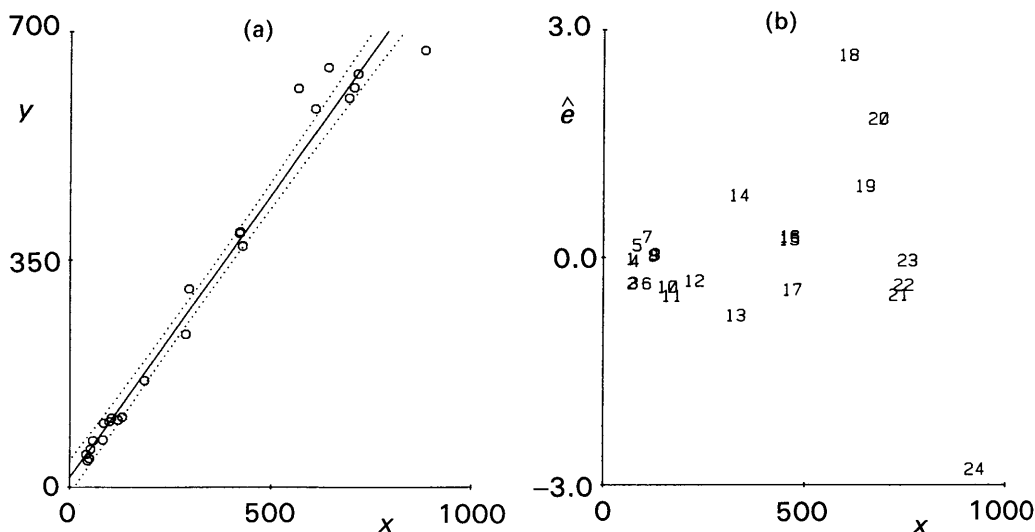


Fig. 6.50—(a) Regression model with the 95% confidence intervals and (b) graphical examination of residuals. Weights V_{ii} were found from a multiplicative model of heteroscedasticity.

6.6.2.2 Autocorrelation

Autocorrelation in data represents a violation of condition 6 for least-squares methods, concerning *independence of measurement errors*. Autocorrelation may be found in chemometrics problems involving data concerned with time dependencies, for example, the data from the kinetics of a reaction. The covariance matrix of errors C_e contains off-diagonal elements.

In chemometrics problems, we are often faced with cumulative errors which are the consequence of a sampling technique used when all experiments are carried out on a single solution. For example, investigation of the kinetics of a chemical reaction is performed by measurement of the concentration of initial substances or resulting reaction products in a single experiment. The *process error* ε_t at time t is, in an ideal case, given by

$$\varepsilon_t = \sum_{j=1}^t u_j \quad (6.163)$$

where u_j are independent random variables of the normal distribution $N(0, \sigma^2)$. This equation shows that the process error ε_t is a sum of all the random effects which have affected the process throughout the experiment. The model [Eq. (6.163)] is a special case of the *autoregressive model of the first order* AR(1), for which Eq. (6.132) is valid. Other eventualities leading to a non-diagonal matrix C_e rarely appear in chemometrics.

In the case of model AR(1), Eq. (6.132) may be expressed in matrix form as

$$\varepsilon = A_1 u \quad (6.164)$$

where A_1 is the lower triangular matrix

$$\mathbf{A}_1 = \begin{bmatrix} 1 & 0 & 0 & . & . & . & 0 \\ \rho_1 & 1 & 0 & . & . & . & 0 \\ \rho_1^2 & \rho_1 & 1 & . & . & . & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \rho_1^{n-1} & \rho_1^{n-2} & \rho_1^{n-3} & . & . & . & 1 \end{bmatrix} \quad (6.165)$$

and moreover it is valid that $E(\varepsilon) = 0$ and the variance of the i th error is given by

$$D(\varepsilon_t) = \sigma^2 \sum_{i=0}^{t-1} \rho_1^{2i} \approx \frac{\sigma^2}{1 - \rho_1^2} \quad (6.166)$$

The last term in this equation is valid on the assumption that t has a sufficiently high value, or that the autoregressive process started at $t = -\infty$.

For an autoregressive process of the first order, simple expressions for covariance of errors, $E(\varepsilon_t \varepsilon_s)$ may be found and the stationary covariance matrix of errors formed

$$\mathbf{C}_\varepsilon = \frac{\sigma^2}{1 - \rho_1^2} \begin{bmatrix} 1 & \rho_1 & \rho_1^2 & . & . & . & \rho_1^{n-1} \\ \rho_1 & 1 & \rho_1 & . & . & . & \rho_1^{n-2} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \rho_1^{n-1} & \rho_1^{n-2} & \rho_1^{n-3} & . & . & . & 1 \end{bmatrix} \quad (6.167)$$

with a general element $C_{ij} = \rho_1^{|i-j|}$. The inverse matrix has the simple structure

$$\mathbf{C}_\varepsilon^{-1} = \frac{1}{\sigma^2} \begin{bmatrix} 1 & -\rho_1 & 0 & 0 & . & 0 \\ -\rho_1 & 1 + \rho_1^2 & -\rho_1 & 0 & . & 0 \\ 0 & -\rho_1 & 1 + \rho_1^2 & -\rho_1 & . & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \quad (6.168)$$

For the autoregressive model of the first order, the covariance matrix of errors may be determined by substitution into Eq. (8.21a) (Chapter 8). In calculation of the inverse matrix $\mathbf{C}_\varepsilon^{-1}$ the matrix \mathbf{A}_1^{-1} should be known

$$\mathbf{A}_1^{-1} = \begin{bmatrix} 1 & 0 & 0 & . & . & . & 0 \\ -\rho_1 & 1 & 0 & . & . & . & 0 \\ 0 & -\rho_1 & 1 & . & . & . & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & . & . & . & 1 \end{bmatrix}$$

This matrix is composed from unit diagonal elements and one underdiagonal band of identical elements $-\rho_1$. The corresponding matrix $\mathbf{C}_\varepsilon^{-1}$ differs from the matrix of Eq. (6.168) only in that the first element on the main diagonal is not equal to 1 but to $(1 + \rho_1^2)$. For the case of cumulative errors in Eq. (6.169), $\rho_1 = 1$. When ρ_1 is unknown, the GLS method with a weight matrix $\mathbf{V} = \mathbf{A}_1^{-1}$ is used.

Let us derive the equation for the transformed vector \mathbf{w} and the matrix \mathbf{Z} used in Eq. (6.147). By a straight multiplication we determine that $w_1 = y_1$ and $w_j = y_j - \rho_1 y_{j-1}$ for $j = 2, \dots, n$. The first row in the matrix \mathbf{Z} is \mathbf{x}_1 . The general

element Z_{ij} of other rows of this matrix is given by

$$Z_{ij} = x_{ij} - \rho_1 x_{i-1,j}, \quad i = 2, \dots, n \\ j = 1, \dots, m \quad (6.169)$$

When the first experimental point is omitted and the first row in matrix \mathbf{Z} is neglected, we obtain the Cochrane–Orcutt estimate, which corresponds to the minimum of the LS criterion for the first differences

$$U_0 = \sum_{i=2}^n [y_i - \rho_1 y_{i-1} - (\mathbf{x}_i - \rho_1 \mathbf{x}_{i-1}) \mathbf{b}_Z]^2 \quad (6.170)$$

where \mathbf{x}_i stands for the i th row in matrix \mathbf{X} . However, it is better to use all experimental points, and the regression criterion

$$U_K = U_0 + (y_1 - \mathbf{x}_1 \mathbf{b}_Z)^2 \quad (6.171)$$

When the autocorrelation coefficient ρ_1 is not known it can be estimated by

$$\hat{\rho}_1 = \frac{\sum_{i=2}^n \hat{e}_i \times \hat{e}_{i-1}}{\sum_{i=2}^n \hat{e}_{i-1}^2} \quad (6.172)$$

Equation (6.172) represents the slope of the regression straight line of the plot of \hat{e}_i vs. \hat{e}_{i-1} , estimated by the classical LS method. As the residuals do not have constant variance, it is more convenient, in Eq. (6.172), to use the standardized residuals $\hat{e}_{si} = \hat{e}_i / \sqrt{1 - H_{ii}}$. By substituting for ρ_1 from Eq. (6.172), U_0 or U_K can be minimized and the estimates \mathbf{b}_Z of parameters $\boldsymbol{\beta}$ may be found. These estimates are biased, because the estimate of the autocorrelation coefficient of the first order, $\hat{\rho}_1$, from Eq. (6.172) is not ideal for small sample sizes. A significant improvement can be achieved by iterative refinement, as follows:

- (1) For a given $\hat{\rho}_1$ the estimates \mathbf{b}_Z and residuals \hat{e} are evaluated;
- (2) With the use of the residuals \hat{e} , estimate $\hat{\rho}_1$ is refined, then step (1) is repeated.

The iteration process terminates when the estimates for $\hat{\rho}_1$ in two successive steps do not differ.

It is permissible to make a simultaneous search for both estimates \mathbf{b}_Z and ρ_1 by minimization of U_K by nonlinear regression, because of ρ_1 , even though the model is linear in $\boldsymbol{\beta}$.

Problem 6.43. *Estimates of the parameters of a regression straight line in data with cumulative errors*

Write expressions for the parameter estimates of the calibration straight line $E(y/x) = \beta_1 x + \beta_2$ when the experimental arrangement produces data with cumulative errors.

Solution: Because the case of cumulative errors is a special case of model AR(1) for $\rho_1 = 1$, we start with the more general solution which is valid for any ρ_1 . For a regression straight line, $E(y/x) = \beta_1 x + \beta_2$ the regression criterion U_K is

$$U_K = (y_1 - b_{1Z}x_1 - b_{2Z})^2 + \sum_{i=2}^n [(y_i - \rho_1 y_{i-1}) - b_{1Z}(x_i - \rho_1 x_{i-1}) - b_{2Z}(1 - \rho_1)]^2.$$

From both derivatives $\delta U_K / \delta b_{1Z}$ and $\delta U_K / \delta b_{2Z}$, estimates may be found which minimize U_K . Then a set of two linear equations is formed

$$\begin{bmatrix} y_1 x_1 + \sum_{i=2}^n (y_i - \rho_1 y_{i-1})(x_i - \rho_1 x_{i-1}) \\ y_1 + (1 - \rho_1) \sum_{i=2}^n (y_i - \rho_1 y_{i-1}) \end{bmatrix} = \begin{bmatrix} x_1^2 + \sum_{i=2}^n (x_i - \rho_1 x_{i-1})^2 & x_1 + (1 - \rho_1) \sum_{i=2}^n (x_i - \rho_1 x_{i-1}) \\ x_1 + (1 - \rho_1) \sum_{i=2}^n (x_i - \rho_1 x_{i-1}) & 1 + (1 - \rho_1)^2 \end{bmatrix} \cdot \begin{bmatrix} b_{1Z} \\ b_{2Z} \end{bmatrix}$$

from which the estimates b_{1Z} and b_{2Z} are calculated. For a case of cumulative errors, when $\rho = 1$, the formulation is simpler

$$\begin{bmatrix} y_1 x_1 + \sum_{i=2}^n (y_i - y_{i-1})(x_i - x_{i-1}) \\ y_1 \end{bmatrix} = \begin{bmatrix} x_1^2 + \sum_{i=2}^n (x_i - x_{i-1})^2 & x_1 \\ x_1 & 1 \end{bmatrix} \begin{bmatrix} b_{1Z} \\ b_{2Z} \end{bmatrix}$$

From this, the following estimates are calculated

$$b_{2Z} = y_1 - b_{1Z} \times x_1 \quad (6.174)$$

and

$$b_{1Z} = \frac{\sum_{i=2}^n (y_i - y_{i-1})(x_i - x_{i-1})}{\sum_{i=2}^n (x_i - x_{i-1})^2} \quad (6.175)$$

Equation (6.175) corresponds to the minimum of U_K , and also of the simplified criterion U_0 of the LS method in first differences. From the set of normal equations, the estimates of variance may be derived

$$D(b_{1Z}) = \frac{\sigma^2}{\sum_{i=2}^n (x_i - x_{i-1})^2}$$

and

$$D(b_{2Z}) = \sigma^2 \left[1 + \frac{x_1^2}{\sum_{i=2}^n (x_i - x_{i-1})^2} \right]$$

From Eq. (6.175), we can also prove that for the case of constant difference between the location of experimental points $\Delta = x_i - x_{i-1}$, $i = 2, \dots, n$, a simple expression can be derived for the slope

$$b_{1Z} = \frac{y_n - y_1}{\Delta \times (n - 1)} \quad (6.176)$$

Conclusion: For the case of cumulative errors, the estimates of the calibration straight line parameters and their variances can be found from Eqs. (6.174) and (6.175). The estimate of a slope is not affected by the use of the simple criterion of the LS method.

Problem 6.44. *Parameters of the kinetics of sugar inversion, with consideration of various errors in the data*

Problem 6.38 presented kinetic data for the inversion of sugar. Because samples were taken from a single sugar solution after different time intervals, it can be expected that the data contain cumulative errors. Estimate parameters β_1 and β_2 of the regression straight line $E(y/x) = \beta_1 \times x + \beta_2$ with the use of (a) the classical LS method, (b) the generalized LS method for the case of cumulative errors, and (c) the generalized LS method for the AR(1) model of errors.

Data: from Problem 6.38

Solution: (a) The LS method.

By use of the classical LS method, the regression equation found was

$$y = 1.002(\pm 0.0017) - 0.005303(\pm 0.0000357)x$$

with the residual standard deviation $\hat{\sigma} = 0.00276$. The sign test confirmed a trend in residuals. The number of sequences $n_U = 8$ is significantly higher than the expected mean value $E(n_U) \approx 4$ for independent residuals. The estimate of the autocorrelation coefficient $\hat{\rho}_1 = -0.715$ shows that the assumption of cumulative errors is not quite correct. The value $\hat{\rho}_1$ is strongly affected by the small number of data points.

(b) The GLS method for cumulative errors.

Substitution into Eqs. (6.173) and (6.174) and the corresponding expressions for variances yields the regression equation

$$y = 1.000(\pm 2.19 \times 10^{-5}) - 0.005238(\pm 0.000166)x$$

with residual standard deviation $\hat{\sigma} = 0.00469$.

(c) The GLS method for the AR(1) model of errors.

With the use of estimates from Eq. (6.173) the best estimate of ρ_1 was refined iteratively, to $\hat{\rho}_1 = -0.864$. The calculated regression equation is then

$$y = 1.003(\pm 6.8 \times 10^{-4}) - 0.00532(\pm 1.49 \times 10^{-5})x$$

with residual standard deviation $\hat{\sigma} = 0.00189$.

Conclusion: The large number of sequences of residuals n_U in comparison with the mean value $E(n_U) \approx n/2$ shows that the model has structure AR(1). With a model of cumulative errors, it may happen that the results are much worse than those from the LS method. The general expression (6.173), with iterative refinement of $\hat{\rho}_1$, is more convenient to use. This method gave decreased residual standard deviation and variances of parameter estimates.

Many varied tests may be used to test the significance of the autocorrelation coefficient ρ_1 . The Wald test is a simple one which examines the null hypothesis $H_0: \rho_1 = 0$ against the alternative one $H_A: \rho_1 \neq 0$ by using the test criterion

$$W_a = \frac{n\hat{\rho}_1^2}{1 - \hat{\rho}_1^2} \quad (6.177)$$

When H_0 is valid, the test statistic W_a has approximately the $\chi^2(1)$ distribution with one degree of freedom.

The Durbin-Watson test is based on the test criterion

$$D_w = \frac{\sum_{i=2}^n (\hat{e}_i - \hat{e}_{i-1})^2}{\sum_{i=1}^n \hat{e}_i^2} \quad (6.178)$$

and $D_w \approx 2 - 2\hat{\rho}_1$. The range of rejection of a null hypothesis $H_0: \rho_1 = 0$ depends not only on the selected significance level α but also on the location of experimental points \mathbf{x}_i . For positive autocorrelation, $0 \leq D_w \leq 2$, while for negative autocorrelation $2 < D_w \leq 4$. If $D_w \approx 2$ then the autocorrelation coefficient is not significant. If $D_w \approx 0$ or $D_w \approx 4$, respectively, the null hypothesis H_0 is rejected and ρ_1 is significantly different from zero. In statistical tables both critical limits, the lower, d_L , and the upper, d_U , for a given significance level α and number of controllable variables m , may be found. When $\rho_1 > 0$, for $D_w > d_U$ the null hypothesis H_0 is accepted, and for $D_w < d_L$ it is rejected. When $d_L \leq D_w \leq d_U$, the test is not conclusive. When the value of the autocorrelation coefficient ρ_1 is very high, the proposed regression model may be false, and a significant variable may have been excluded from the model.

Problem 6.45. *Examination of the autocorrelation coefficient for the kinetics of inversion of sugar*

Examine the significance of the autocorrelation coefficient ρ_1 . In Problem 6.44, a value $\hat{\rho}_1 = -0.715$ was found for its estimate.

Data: $\hat{\rho}_1 = -0.715$, $n = 9$.

Solution: For the Wald test (6.177), the test criterion is

$$W_a = \frac{9(-0.715^2)}{(1 - 0.715^2)} = 9.413$$

Since W_a is greater than the quantile $\chi_{0.95}^2(1) = 3.84$, the null hypothesis $H_0: \rho_1 = 0$ is rejected and the autocorrelation coefficient ρ_1 may be considered to be significantly different from zero.

Conclusion: Examination of the autocorrelation coefficient confirmed the conclusion of Problem 6.44. By using the iterative method of refining the autocorrelation coefficient, the refined estimate is $\hat{\rho}_1 = -0.864$.

6.6.3 Multicollinearity

Multicollinearity does not mean a violation of the conditions for the least-squares methods. It concerns an assumption about positive definite matrix $\mathbf{X}^T\mathbf{X}$ and therefore the solution of Eq. (6.11).

According to Section 6.1, we understand the columns of matrix \mathbf{X} as the column vectors which define the hyperplane L in n -dimensional Euclidean space E^n (Fig. 6.2). According to the angle θ_{jk} between two vectors \mathbf{x}_j and \mathbf{x}_k (or between columns of matrix \mathbf{X}) two limiting cases may be distinguished:

- (1) Orthogonality is found when the cosine of angle θ_{jk} is zero

$$\cos \theta_{jk} = \frac{\langle \mathbf{x}_j, \mathbf{x}_k \rangle}{\|\mathbf{x}_j\| \times \|\mathbf{x}_k\|} \quad (6.179)$$

and also the scalar product $\langle \mathbf{x}_j, \mathbf{x}_k \rangle = 0$ where the symbol $\|\mathbf{x}_j\| = \sqrt{\langle \mathbf{x}_j, \mathbf{x}_j \rangle}$ means the length of vector \mathbf{x}_j . If all the columns of matrix \mathbf{X} are mutually orthogonal, then the matrix $\mathbf{X}^T\mathbf{X}$ is diagonal and the regression analysis simplifies (Section 6.4.1).

(2) Collinearity is found when the cosine of angle θ_{jk} is equal to 1, $\cos \theta_{jk} = 1$, because the angle between vectors \mathbf{x}_j and \mathbf{x}_k is zero, $\theta_{jk} = 0$, and the two vectors \mathbf{x}_j and \mathbf{x}_k are parallel and linearly dependent, and the following expression holds for them

$$c_j \mathbf{x}_j + c_k \mathbf{x}_k = 0 \quad (6.180)$$

where c_j and c_k are nonzero constants. When Eq. (6.180) holds for q pairs of columns of matrix \mathbf{X} , its rank is equal to $m - q$ and the matrix $\mathbf{X}^T\mathbf{X}$ is singular.

Equation (6.180) can be valid for more vectors yet, when one of the columns \mathbf{x}_j is the result of a linear combination of several other columns. This situation is called *perfect multicollinearity*. The term multicollinearity, however, can include other cases when some columns of matrix \mathbf{X} have nearly zero angle and are therefore approximately linearly dependent.

$$\sum_{j=1}^m c_j \mathbf{x}_j = \boldsymbol{\delta} \quad (6.181)$$

where $\boldsymbol{\delta}$ is the vector with components near to zero, and the vector \mathbf{c} with elements c_j is nonzero, $\|\mathbf{c}\| \gg \|\boldsymbol{\delta}\|$. The multicollinearity causes ill-conditioning of the matrix $\mathbf{X}^T\mathbf{X}$, and this has two consequences:

- (a) the determinant of matrix $\mathbf{X}^T\mathbf{X}$ is close to zero;
- (b) some eigenvalues of matrix $\mathbf{X}^T\mathbf{X}$ are close to zero.

Multicollinearity causes many difficulties in inversion of matrix $\mathbf{X}^T\mathbf{X}$ and also numerical errors, depending on the machine precision of the computer used. As well as numerical difficulties, multicollinearity causes statistical difficulties. From Eqs. (6.88) and (6.87), it is evident that for λ values near to zero, the parameter estimates and their variance will be abnormally high. The special difficulties are caused by the sensitivity of the parameter estimates \mathbf{b} to small changes in data, such as adding another point to the data.

Figure 6.51 shows a geometric interpretation of the LS method for two nearly collinear controllable variables. Figure 6.51a shows vector \mathbf{y} projected into a segment of angle θ_{12} , and Fig. 6.51b shows a case when a small change in vector \mathbf{y} causes its perpendicular projection to lie out of this segment.

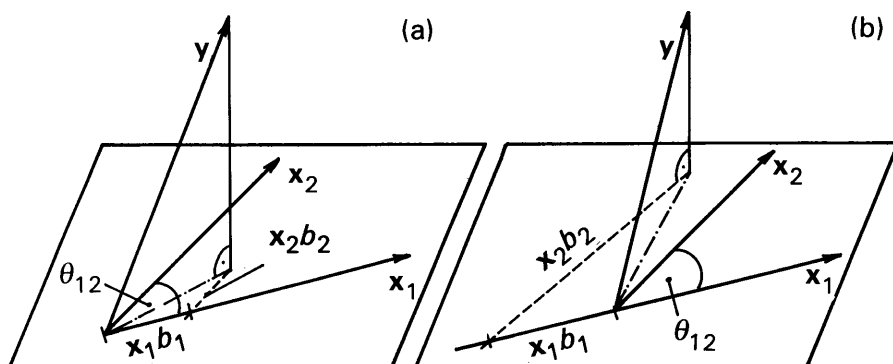


Fig. 6.51—Geometric interpretation of the sensitivity of the estimates in the case of multicollinearity: (a) the estimates b_1 and b_2 are small and positive, and (b) the estimate b_1 is negative while b_2 is large and positive.

Multicollinearity causes the following statistical difficulties:

- Non-stability of estimates* is caused by great sensitivity of parameter estimates to small changes in the data. The estimates often have the wrong sign, and this damages their physical interpretation.
- Large variances* $D(b_j)$ of individual estimates cause t -tests to indicate that parameter β_j (cf. Section 6.3) is statistically insignificant.
- Strong correlation* between elements of the estimates vector \mathbf{b} means that they cannot be interpreted separately.

On the other hand, in cases of multicollinearity the determination coefficient is always high and the regression model may fit the data quite well. For data approximation and data smoothing by regression, multicollinearity does not cause difficulties apart from numerical ones related to the ill-conditioning of matrix $\mathbf{X}^T\mathbf{X}$.

When data are measured according to an experimental design the problem of multicollinearity is removed. The plan of designed experiments leads to orthogonality of the columns of matrix \mathbf{X} .

When data are not measured according to a designed experiment multicollinearity always exists to some extent. However, strong multicollinearity causes the parameter estimates and hypotheses tests to be affected more by linear connections between the

columns of matrix \mathbf{X} than by the regression model itself. In the chemical laboratory, the values of the controllable variables may be adjusted freely, so the problem of multicollinearity may be avoided by an appropriate data measurement step.

With reference to multicollinearity in data, we can identify three cases of interest.

- (a) The *over-estimated regression model* contains too many controllable variables expressing the same basic factors. An example is a structure – properties model in which properties of substances are described by various measurable changeable structures.
- (b) *Inappropriate location of experimental points* causes multicollinearity to form “artificially” because of the choice of location of points. Often the values of significantly important variables oscillate in a small range and seem to be nearly constant, and they are collinear with the vector corresponding to the intercept term.
- (c) *Physical constraints in model or data* refers to limits on the values of the controllable variables derived from the chemistry of the system. An example is an investigation of multicomponent mixtures where the controllable variables are represented by the content of each component. Because the sum of all relative concentrations should be equal to 100%, in a q -component mixture there will be $(q - 1)$ independent components. In a model, only $(q - 1)$ variables are assumed: for a two-component mixture there is only one variable, for a three-component mixture, only two, etc. Similar restriction may apply to stoichiometric ratios, etc.

From knowledge about the controllable variables, and their significance and restrictions, multicollinearity can be completely removed from the data. In the case of polynomial models, the multicollinearity is defined by the model structure. If the experimental strategy cannot be changed, other techniques for decreasing the influence of multicollinearity should be used, despite the fact that the parameter estimates are then biased, as in the case of the method of rational ranks (section 6.4.2).

Multicollinearity can be detected from scatter plots for \mathbf{x}_j and \mathbf{x}_k when the approximate linear dependence proves the strong multicollinearity. The multicollinearity may be exposed or masked by the presence of influential points and especially by high leverage points. For diagnostic purposes, the residuals \mathbf{v}_j of regression variable \mathbf{x}_j on the remaining controllable variables in a matrix $\mathbf{X}_{(j)}$ which does not contain the column \mathbf{x}_j can be used. Let us use $\mathbf{H}_{(j)}$ to denote the projection matrix which corresponds to the projection into a subspace of columns of matrix $\mathbf{X}_{(j)}$. For diagnosis of influential points from the point of view of multicollinearity, the plot of $v_{ji}^2/(\mathbf{v}_j^T \cdot \mathbf{v}_j)$ against $H_{(j)ii}$ is used, where v_{ij} is the i th component of vector \mathbf{v}_j and $H_{(j)ii}$ is the i th diagonal element of matrix $\mathbf{H}_{(j)}$.

In Fig. 6.52, the points strongly affected by multicollinearity are located in the bottom right-hand corner and the top left-hand corner of the graph. The points located in the top left corner cause multicollinearity only when variable \mathbf{x}_j is included in the model. The points located in the bottom right corner are strongly influential only when variable \mathbf{x}_j is not included in model.

The presence of multicollinearity can be identified on the basis of numerical and

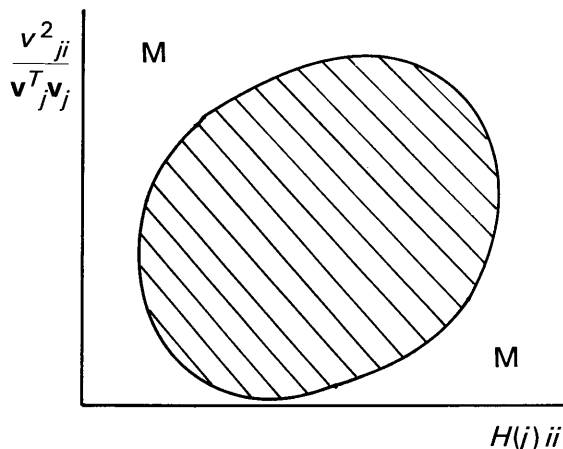


Fig. 6.52—Identification of multicollinearity in data: **M** denotes multicollinearity.

statistical criteria. Instead of the matrix $\mathbf{X}^T\mathbf{X}$, its normalized version \mathbf{R} is used. Matrix \mathbf{R} is formally identical with the correlation matrix of controllable variables.

The following numerical criteria are commonly used.

(a) The determinant of matrix \mathbf{R} is calculated from

$$\det(\mathbf{R}) = \prod_{j=1}^m \lambda_j$$

where the λ_j are eigenvalues of the matrix \mathbf{R} . If $\det(\mathbf{R})$ is small and less than 10^{-3} , multicollinearity is detected.

(b) The conditioning number K is calculated from

$$K = \lambda_{\max}/\lambda_{\min} \quad (6.182)$$

where λ_{\max} and λ_{\min} are the maximal and minimal eigenvalues of a matrix \mathbf{R} . If $K > 10^3$, strong multicollinearity is detected.

The main statistical criterion used is the *VIF factor*, defined as the ratio of the variance of the j th regression coefficient to the same variance for orthogonal variables when \mathbf{R} is the unit matrix. It is given by

$$VIF_j = \tilde{R}_{jj} \quad (6.183)$$

where \tilde{R}_{jj} is the j th diagonal element of matrix \mathbf{R}^{-1} . The *VIF* factors are related to the determination coefficient $\hat{R}_{x_j}^2$ of regression \mathbf{x}_j on $\mathbf{X}_{(j)}$ when \mathbf{x}_j is expressed as a combination of other controllable variables. Then

$$VIF_j = \frac{1}{1 - \hat{R}_{x_j}^2} \quad (6.184)$$

If $VIF_j > 10$, strong multicollinearity is detected.

Problem 6.46. *Testing for multicollinearity in the dependence of the mean activity coefficient on temperature*

The dependence of the logarithm of the mean activity coefficient on the temperature, $\ln \gamma_{\pm} = f(T)$ can be expressed by a polynomial of the third degree. Consider the extent of multicollinearity and use the method of rational rank to decrease the multicollinearity level.

Data: measured for $m_{\text{HCl}} = 0.1$

$T, ^\circ\text{C}$	0	10	20	30	40	50
$\ln \gamma_{\pm}$	0.8067	0.8038	0.8000	0.7964	0.7927	0.7867
	60	70	80	90		
	0.7828	0.775	0.769	0.765		

Solution: For the proposed model, $\ln \gamma_{\pm} = \beta_1 T + \beta_2 T^2 + \beta_3 T^3 + \beta_4$, the regression equation is found to be

$$\ln \gamma_{\pm} = 0.807(\pm 1.06 \times 10^{-3}) - 2.654 \times 10^{-4}(\pm 1.07 \times 10^{-4})T \\ - 3.13 \times 10^{-6}(\pm 2.87 \times 10^{-6})T^2 + 9.44 \times 10^{-9}(\pm 2.09 \times 10^{-8})T^3$$

by the classical least-squares method (estimated standard deviations of parameter estimates are given in brackets). The determination coefficient $\hat{R}^2 = 0.9957$, the quadratic error of prediction $MEP = 3.507 \times 10^{-6}$ and the Akaike criterion $AIC = -132.26$. Table 6.13 lists the eigenvalues and the *VIF* factors. From these numbers, $\det(\mathbf{R}) = 3.97 \times 10^{-4}$ and the conditioning number $K = 1989.73$ are calculated. From *t*-tests at the significance level $\alpha = 0.05$, parameters β_2 and β_3 are statistically insignificant.

Table 6.13. Characteristics detecting multicollinearity

P	Characteristic	$j = 1$	$j = 2$	$j = 3$
10^{-35}	VIF_j	70.42	439.1	184
	λ_j	0.00146	0.0935	2.905
0.05	VIF_j	6.204	0.260	4.373

With the use of the method of rational ranks, the regression equation with precision $P = 0.05$ is estimated in the form

$$\ln \gamma_{\pm} = 0.807(\pm 8.72 \times 10^{-4}) - 3.22 \times 10^{-4}(\pm 3.28 \times 10^{-5})T \\ - 1.476 \times 10^{-6}(\pm 7.18 \times 10^{-8})T^2 - 2.837 \times 10^{-9}(\pm 3.314 \times 10^{-9})T^3$$

with the determination coefficient $\hat{R}^2 = 0.9955$, the mean quadratic error of prediction $MEP = 2.264 \times 10^{-6}$ and Akaike criterion $AIC = -131.7$. Table 6.13 gives the *VIF* factors. The matrix \mathbf{R}^{-1} constructed according to Eq. (6.86), and replacing $j = 1$ with $j = \omega$, removed multicollinearity, and *t*-tests at significance level $\alpha = 0.05$ showed that the parameter β_3 is statistically insignificant.

Figure 6.53 shows the regression model found by the method of least-squares, with the 95% confidence intervals, and Fig. 6.54 shows the model found by the method of rational ranks ($P = 0.05$). From comparison of these figures, it is obvious that elimination of multicollinearity leads to narrower confidence bands.

Conclusion: Significant multicollinearity, as indicated by the *VIF* criterion having a value higher than 10, causes an increase in the estimates of variance, and hence an increase in width of the confidence bands. Elimination of multicollinearity leads to a decrease in goodness-of-fit (a decrease of \hat{R}^2) but to an improvement in the prediction ability of model (the criterion *MEP*), in addition to the decrease in the variance of estimates and narrowing of confidence bands. Elimination of multicollinearity is rather important in calibration in the instrumental methods of analytical chemistry.

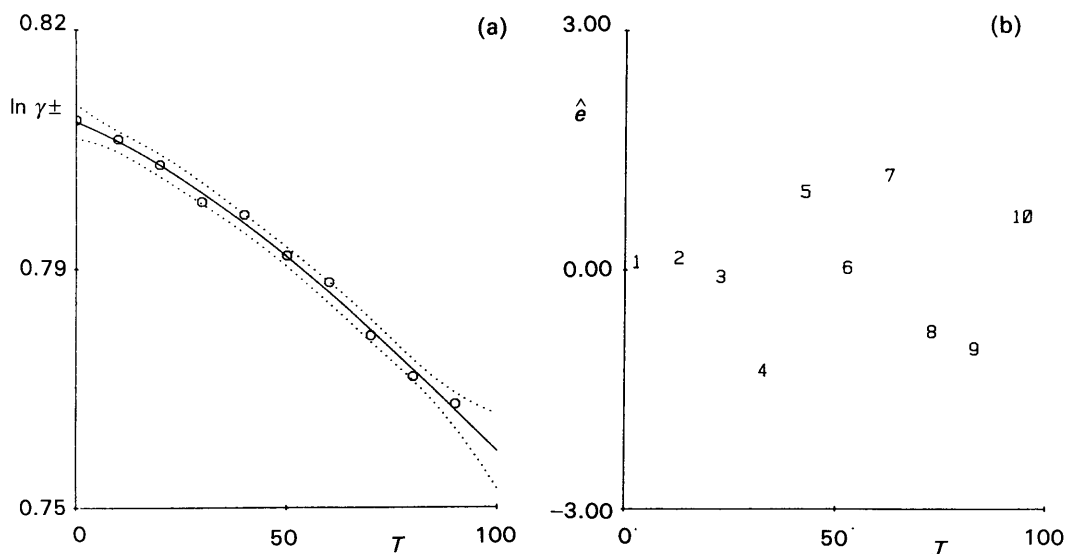


Fig. 6.53—(a) Regression model estimated by the LS method, and (b) graphical examination of residuals.

6.6.4 Variables subject to random errors

In chemometrics problems, both the dependent variable y and controllable variables x are measured quantities. The variances of x are usually significantly smaller than the variance of y , and also smaller than the differences between the locations of individual points. Under such conditions the assumption about the deterministic matrix X may be abandoned. In some cases it is necessary to suppose that instead of variables x_j we measure experimental values t_j given by

$$t_{ij} = x_{ij} + \kappa_{ij} \quad (6.185)$$

where κ_{ij} are errors of measurement of the j th independent variable at the i th point. The result of measurement is the set of n points $\{y_i, t_{ij}; j = 1, \dots, m\}$, $i = 1, \dots, n$. If the x_{ij} are deterministic quantities we speak about *functional models*, but if the x_{ij}

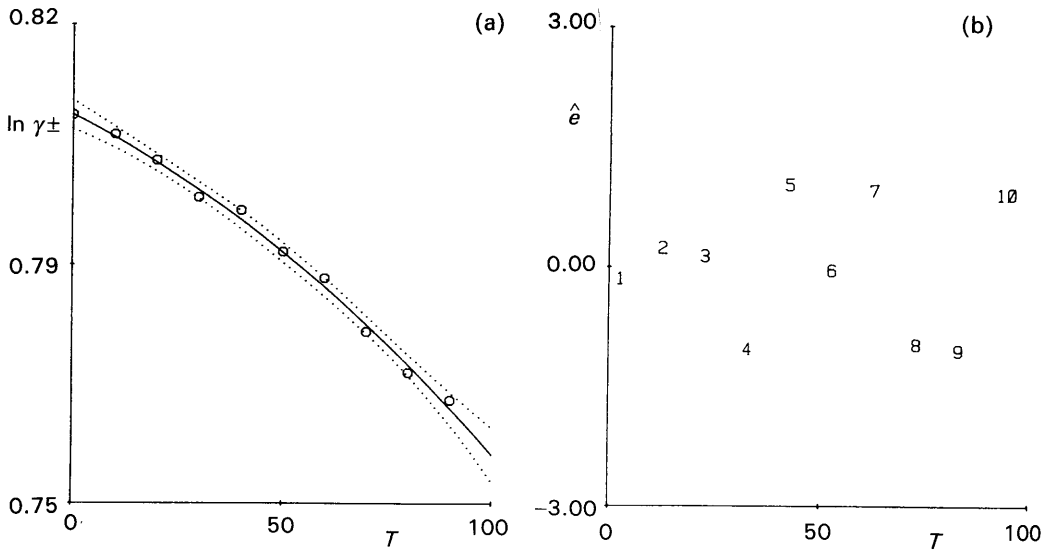


Fig. 6.54—(a) Regression model estimated by the rational rank method and (b) graphical examination of residuals

are random quantities we speak about *structural models*. The errors ε_i and κ_{ij} have the properties:

- Both errors, ε_i and κ_{ij} , have zero means.
- The variance of errors $D(\varepsilon_i^2) = \sigma^2$ and $D(\kappa_{ij}^2) = \tau_j^2$ are at all n points constant (homoscedasticity).
- The errors ε_i and κ_{ij} at different points (i.e. measurements) are uncorrelated, so that $E(\varepsilon_i \varepsilon_j) = 0$, $i \neq j$ and $E(\kappa_{ij} \kappa_{kj}) = 0$, $i \neq k$.
- The errors ε_i and κ_{ij} are mutually uncorrelated so that $E(\varepsilon_i \kappa_{ij}) = 0$, $j = 1, \dots, m$.

If errors ε have a normal distribution $N(0, \sigma^2)$ and errors κ_j also have a normal distribution $N(0, \tau_j^2)$, then according to the maximal likelihood method, the criterion of the *extended least-squares method* (ELS) can be expressed as:

$$U_E(\mathbf{b}, \mathbf{X}) = \sum_{i=1}^n \left[\frac{1}{\sigma^2} \left(y_i - \sum_{j=1}^m b_j x_{ij} \right)^2 + \sum_{j=1}^m \frac{1}{\tau_j^2} (t_{ij} - x_{ij})^2 \right] \quad (6.186)$$

By minimizing the function $U_E(\mathbf{b}, \mathbf{X})$ with respect to \mathbf{b} and to \mathbf{x}_j , we find the extended estimates \mathbf{b}_E and also correct quantities \hat{x}_{ij} of controllable variables x_{ij} . With some simplifying assumptions, Eq. (6.186) can be expressed as

$$U_E(\mathbf{b}) = \frac{\sum_{i=1}^n \left(y_i - \sum_{j=1}^m b_j t_{ij} \right)^2}{\sigma^2 + \sum_{j=1}^{m-1} b_j^2 \tau_j^2} \quad (6.187)$$

where b_m is the intercept term. For a regression straight line, Eq. (6.187) is simplified to:

$$U_E(b_1, b_2) = \frac{\sum_{i=1}^n (y_i - b_1 t_i - b_2)^2}{\sigma^2 + \tau^2 b_1^2} \quad (6.188)$$

Figure 6.55 shows that for $\sigma^2 = \tau^2$, the criterion $U_E(b_1, b_2)$ leads to minimization of the squares of perpendicular distances between the regression function and the experimental points.

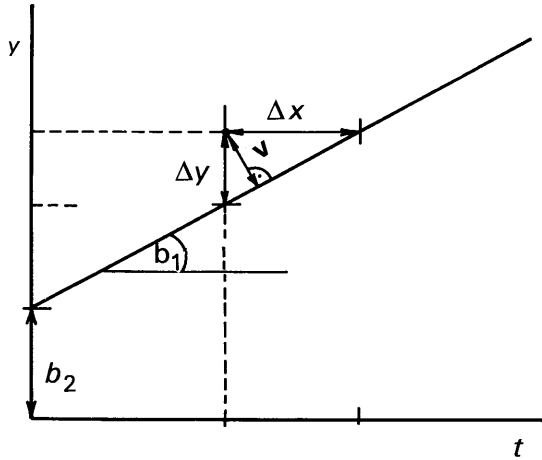


Fig. 6.55—Illustration of the criterion of squares of perpendicular distances $v^2 = \Delta y^2 / (1 + (\Delta y / \Delta x)^2)$ for a straight line $y = b_1 t + b_2$.

The estimate of parameters \mathbf{b} in Eq. (6.187) may be achieved by the iterative procedure of the weighted least-squares method with weight values given by

$$P_i = \left(\sigma^2 + \sum_{j=1}^{m-1} \tilde{b}_{Ej} \tau_j^2 \right)^{-1/2}$$

where \tilde{b}_{Ej} are the parameter values estimated in the previous iteration. However, the procedure requires knowledge of variances σ^2 and τ_j^2 . From the structure of Eqs. (6.187) and (6.188), it may be concluded that a knowledge of the ratios of variances,

$$K_j = \frac{\sigma^2}{\tau_j^2}, \quad j = 1, \dots, m$$

should be useful in application of the iterative procedure.

We restrict ourselves now to the simplest linear case, i.e. the straight line. From Eq. (6.185) we can write

$$y_i = b_1(t_i - \kappa_i) + b_2 + \varepsilon_i = b_1 t_i + b_2 + \varepsilon_i^* \quad (6.189)$$

where $\varepsilon_i^* = \varepsilon_i - b_1 \kappa_i$ represents errors related to a magnitude b_1 . To express the

variance of real x values, the mean quadratic deviation σ_x^2 is used,

$$\sigma_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$$

The ratio

$$K_x = \frac{\sigma_x^2}{\sigma_t^2} = \frac{\sigma_x^2}{\sigma_x^2 + \tau^2}$$

is called the *reliability ratio*. If the true values of x_i and the errors κ_i are not correlated, the mean value $E(b_1)$ estimated by least-squares with neglect of an error ε_i^* structure is given by

$$E(b_1) = \beta_1 K_x \quad (6.190)$$

The corresponding determination coefficient is expressed by

$$R_{yt}^2 = R_{yx}^2 K_x \quad (6.191)$$

If the classical least-squares method is used, measurement errors cause a decrease in slope estimate b_1 and in the correlation coefficient. The magnitude of this decrease depends on the reliability ratio K_x or on the ratio of σ_x^2 to σ_t^2 . If K_x is significantly lower than 1, Eq. (6.188) is used to estimate slope b_{1E} .

To remove the intercept b_2 we introduce the centred variables $(y_i - \bar{y})$ and $(t_i - \bar{t})$. If we know the variance ratio $K = \sigma^2/\tau^2$, Eq. (6.188) may be expressed in the form

$$U_E(b_1) = \frac{\sum_{i=1}^n ((y_i - \bar{y}) - b_1(t_i - \bar{t}))^2}{K + b_1} \quad (6.192)$$

After analytical minimization of $U_E(b_1)$ we obtain

$$b_{1E} = L + \text{sign}(S_{yt}) \times \sqrt{K + L^2} \quad (6.193)$$

where

$$L = \frac{S_y - K \times S_t}{2S_t}$$

and $\text{sign}(S_{yt})$ gives the sign of variable S_{yt} . Symbols S represent the sums of squares

$$S_y = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$S_t = \sum_{i=1}^n (t_i - \bar{t})^2$$

$$S_{yt} = \sum_{i=1}^n (y_i - \bar{y})(t_i - \bar{t})$$

When the slope b_{1E} is known, the intercept b_{2E} of the regression straight line may be calculated from

$$b_{2E} = \bar{y} - b_{1E}\bar{x} \quad (6.194)$$

The influence of the magnitude of K on the set of regression straight lines is evident from Eq. (6.193).

- (a) For $K \rightarrow \infty$, the regression line corresponds to the LS method.
- (b) For $K = 1$, the regression line minimizes the perpendicular distances from experimental points. This is called orthogonal regression.
- (c) For $K \rightarrow 0$ the regression line is an inverse regression i.e. a linear dependence of t on y .

Unsuitable selection of the magnitude of K leads, however, to an increase in variances, so the techniques for simultaneous estimation of parameters and variance ratio are used. Some other procedures of regression analysis, for the case when all variables are subject to random errors, are described by Fuller [36].

We will write an expression for a structural model for which the random variables x , ε and κ have independent normal distributions with variances σ_x^2 , σ^2 and τ^2 , and for which the ratio $K = \sigma^2/\tau^2$ is also known. The variance estimates are then given by

$$\hat{\sigma}_x^2 = \frac{S_{yt}}{K(n-1)} \left[\sqrt{K + L^2} - L \right] \quad (6.195)$$

and

$$\hat{t}^2 = \frac{1}{2K(n-1)} \left[S_y + KS_t - 2S_{yt}\sqrt{K + L^2} \right] \quad (6.196a)$$

and

$$\hat{\sigma}^2 = \hat{t}^2 K \quad (6.196b)$$

The variance of the slope of the regression straight line is given by

$$D(b_{1E}) = \frac{1}{(n-1)\hat{\sigma}_x^4} [\hat{\sigma}_x^2 S_v + \hat{t}^2 S_v - b_{1E}^2 \hat{t}^4] \quad (6.197)$$

where

$$S_v = \frac{n-1}{n-2} (K + b_{1E}^2) \hat{t}^2$$

To test hypotheses about parameter b_{1E} , the test criterion

$$T_E = \frac{|b_{1E} - b_{1E}^*|}{\sqrt{D(b_{1E})}} \quad (6.198)$$

is used. If the null hypothesis $H_0: \beta_1 = b_1^*$ is accepted, this criterion has approximately

the Student t -distribution with $(n - 2)$ degrees of freedom.

The variance of the intercept of the regression straight line is estimated by

$$D(b_{2E}) = \frac{S_V}{n} + \bar{t}^2 D(b_{1E}) \quad (6.199)$$

where

$$\bar{t} = \frac{\sum_{i=1}^n t_i}{n}$$

Problem 6.47. *Validation of a new analytical method when both variables are subject to random errors*

In Problem 6.7 the results of new analytical method (y) are compared with the standard one (x). Estimate both parameters of regression straight line $y = \beta_1 x + \beta_2$ when both variables x and y are loaded by experimental error and the variances of both methods are same, $K = 1$. Test the null hypothesis $H_0: \beta_1 = 1$.

Data: From Problem 6.7

Solution: Sum of squares $S_y = 1.327 \times 10^6$, $S_x = 1.714 \times 10^6$, $S_{yx} = 1.489 \times 10^6$ substituted into Eq. (6.193) give the estimate of slope $b_{1E} = 0.8784$ and into Eq. (6.194) the estimate of intercept $b_{2E} = 11.521$. As the estimates of variances are $\hat{\sigma}_x^2 = 7.367 \times 10^4$ and $\hat{\tau}^2 = 848.53 = \hat{\sigma}^2$, the estimate of variance of the slope from Eq. (6.197) will be $D(b_1) = 9.337 \times 10^{-4}$ and of the intercept from Eq. (6.199) $D(b_2) = 161.53$. When we test the null hypothesis $H_0: \beta_1 = 1$ with Eq. (6.198), we find that the test criterion $T_E = |0.8784 - 1| / \sqrt{9.337 \times 10^{-4}} = 3.9795$ is higher than the quantile $t_{0.975}(22) = 2.074$ and therefore the null hypothesis H_0 is rejected and the slope β_1 differs significantly from 1. Figure 6.56 demonstrates the regression straight line which minimizes perpendicular distances from experimental points.

Conclusion: Correctly recognizing that both variables are subject to random errors does not cause any difficulties in estimation of the parameters of the regression straight line. A set of regression straight lines arranged according to the precision of individual variables (the magnitude K) may be calculated. For $K = 1$ the useful criterion of perpendicular distances from experimental points is obtained.

6.6.5 Other error distributions of the dependent variable

6.6.5.1 The M-estimates method

When the distribution of the errors in the dependent variable y is not normal (violation of condition 7 for the LS method, Section 6.2) the parameter estimates obtained by the LS method are not the best possible estimates. In such a case, instead of the least-squares criterion some other *robust* criterion can be used, that is not so sensitive to violation of the condition about the error distribution, and not sensitive to influential points. The most convenient robust criteria seem to be the group of M-estimates. The M-estimates are maximal likelihood estimates for a given probability density function of errors $p(\epsilon)$. All M-estimates are related to the minimization

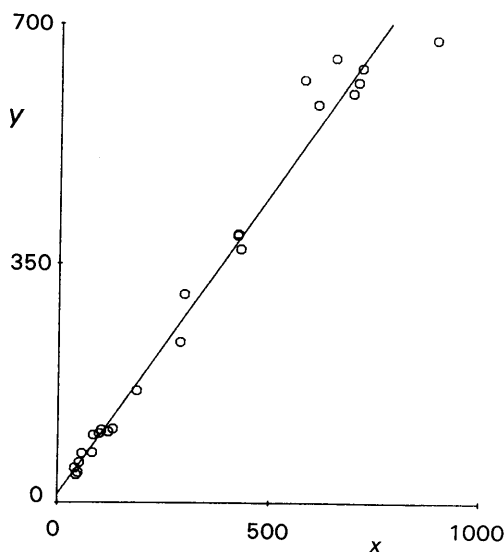


Fig. 6.56—The regression straight line minimizing the perpendicular distances from experimental points.

criterion

$$U_M = \sum_{i=1}^n \rho(e_i/\sigma) = \sum_{i=1}^n \rho((y_i - \mathbf{x}_i \mathbf{b}_M)/\sigma) \quad (6.200)$$

where \mathbf{x}_i is the i th row of matrix \mathbf{X} , σ is the parameter of spread and $\rho(\cdot)$ is a convenient function determined from the probability density $p(\varepsilon)$. By analytical minimization of U_M , (6.200) a set of normal equations is obtained:

$$\sum_{i=1}^n \psi(e_i/\sigma) x_{ij} = 0, \quad j = 1, \dots, m \quad (6.201)$$

where the function

$$\psi(x) = \frac{\delta \rho(x)}{\delta x}$$

represents the derivative of function $\rho(x)$ with respect to x . Then, if $r_i = e_i/\sigma$, Eq. (6.201) may be expressed in a form which corresponds to the weighted least-squares method

$$\sum_{i=1}^n w_i(r) y_i x_{ij} = \sum_{i=1}^n \sum_{k=1}^m w_i(r) x_{ij} x_{ik} b_k \quad j = 1, \dots, m \quad (6.202)$$

where $w_i(r) = \psi(r_i)/r_i$. The parameters are estimated by the iterative method of re-weighted least-squares (IRWLS), by using the following procedure:

- (1) Select $w_i(r) = 1$, $i = 1, \dots, n$ and set $l = 1$.
- (2) Estimate the residuals $r_l = \hat{e}_l/\hat{\sigma}_l$ by the classical least-squares method. In order

to reach convergence, corrected least-squares estimates are used [40].

- (3) Calculate the weights $w_i(r_l)$ from Eq. (6.203), for $l = l + 1$.
- (4) Use the reweighted least-squares to estimate \mathbf{b}_l and the residuals r_l .
- (5) If the estimates \mathbf{b}_l and \mathbf{b}_{l+1} are not close enough, go to step 3, otherwise $b_l = b_M$.

It should be noted that in the j th iteration the weights used have been calculated from residuals \hat{e}_{l-1} in the $(l-1)$ th iteration. By applying this method, the robust estimate of parameter σ can be evaluated. An independent estimate $\hat{\sigma}_l$ from the residuals \hat{e}_{l-1} determined in the previous iteration seems to be most convenient. A useful expression is

$$\hat{\sigma} = \frac{\text{med}(|\hat{e}_l - \text{med}(\hat{e}_l)|)}{0.6745} \quad (6.203)$$

where $\text{med}(\hat{e}_l)$ is the median calculated from all residuals and for sake of simplicity, the indices $(l-1)$ denoting the actual iteration used for residual estimation, are omitted. The constant 0.6745 for large sample size fixes the value $\hat{\sigma}$ to be equal to the residual standard deviation $\hat{\sigma}$ but for a normal error distribution. A simpler option is

$$\hat{\sigma} = 2.1 \text{ med}(|\hat{e}_l|) \quad (6.204a)$$

Hill and Holland [37] recommended the expression

$$\hat{\sigma} = \frac{\text{med}(\text{largest}(n-m)|\hat{e}_l|)}{0.6745} \quad (6.204b)$$

Huber [38] recommends a procedure of simultaneous estimation of \mathbf{b}_j and $\hat{\sigma}_j$ in every iteration. Some variants of IRWLS method are described in a paper by Li [39].

It can be difficult to make the initial guess of the parameters to be estimated. Application of the classical least-squares method can cause difficulties from non-convergence of the estimates. The simple procedure of the corrected least-squares method was suggested by Phillip and Eyring [40]. It starts with estimates \mathbf{b} determined by the classical least-squares method. From residuals \hat{e} the robust parameter of scale is estimated

$$S = \text{med}_i(|\hat{e}_i|) \quad (6.205)$$

and the winsorized residuals are calculated by the rule

$$e_i^w = \begin{cases} -1.5 S & \text{for } \hat{e}_i < -1.5 S \\ \hat{e}_i & \text{for } |\hat{e}_i| \leq 1.5 S \\ 1.5 S & \text{for } \hat{e}_i > 1.5 S \end{cases}$$

The vector of correction $\hat{\mathbf{q}} = (q_1, \dots, q_m)^T$ is calculated as the vector of regression coefficients \mathbf{e}^w on \mathbf{X} from

$$\hat{\mathbf{q}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{e}^w$$

To calculate the quantities r_l and $w_i(r_l)$, the following corrected parameters values are taken as initial values

$$\mathbf{b}^w = \mathbf{b} + \hat{\mathbf{q}}$$

This procedure does not require much computer time, since the matrix $(\mathbf{X}^T \mathbf{X})^{-1}$ is already evaluated.

The statistical analysis of M-estimates is based on fact that estimates \mathbf{b}_M have an asymptotically normal distribution with mean $\boldsymbol{\beta}$ and covariance matrix

$$D(\mathbf{b}_M) = \tau^2 (\mathbf{X}^T \mathbf{X})^{-1} \quad (6.206)$$

where

$$\tau^2 = \frac{E(\psi^2)}{[E(\psi)]^2}$$

Estimate $\hat{\tau}^2$ can be found from the expression

$$\hat{\tau}^2 = K_e \frac{\sum_{i=1}^n \frac{\psi^2(r_i)}{n-m}}{\left[\sum_{i=1}^n \frac{\psi'(r_i)}{n} \right]^2} \quad (6.207)$$

The constant K_e is the correction for finite samples; it may be set equal to one (according to Li [39]) or calculated from an expression suggested by Huber [38].

The advantage of the IRWLS method is the fact that after termination of iterative refinement of parameter estimates the covariance matrix of the LS method is already the estimate $D(\mathbf{b}_M)$.

To examine robustness, functions such as $\rho(r)$ should be selected in order to get the derivative $\psi(r)$ bounded. From Fig. 6.57b it is obvious that for the LS criterion the function $\psi(r)$ is not bounded, because it increases with an increase of r .

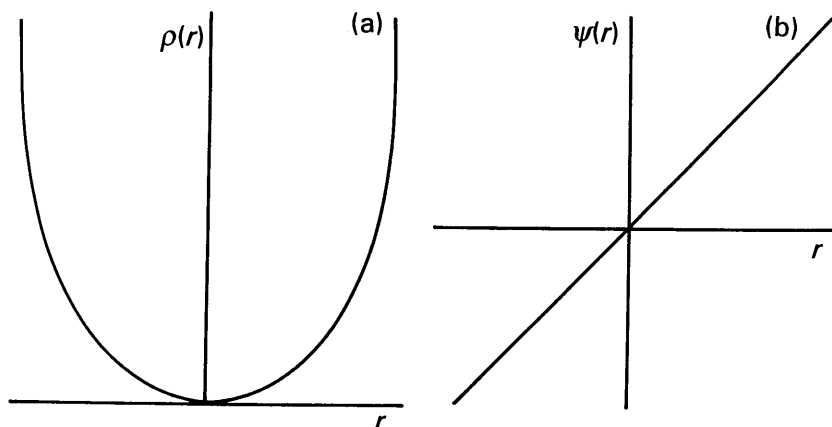


Fig. 6.57—(a) Function $\rho(r)$, and (b) function $\psi(r)$ for the least-squares method.

Of the robust methods, we restrict ourselves here to the L_1 method and the method of combined procedure with “Biweights”, i.e. regression with limited influence. The second method is robust for all types of influential points.

6.6.5.2 The L_1 approximation method

The method of L_1 approximation is also called the method of least absolute residuals. The criterion is in the form

$$L_1(\mathbf{b}) = \sum_{i=1}^n |y_i - \sum_{j=1}^m x_{ij}b_j| \quad (6.208)$$

This is a special case of M-estimates for $\rho(r) = |r|$ and $\psi(r) = \text{sign}(r)$. Both are shown in Fig. 6.58.

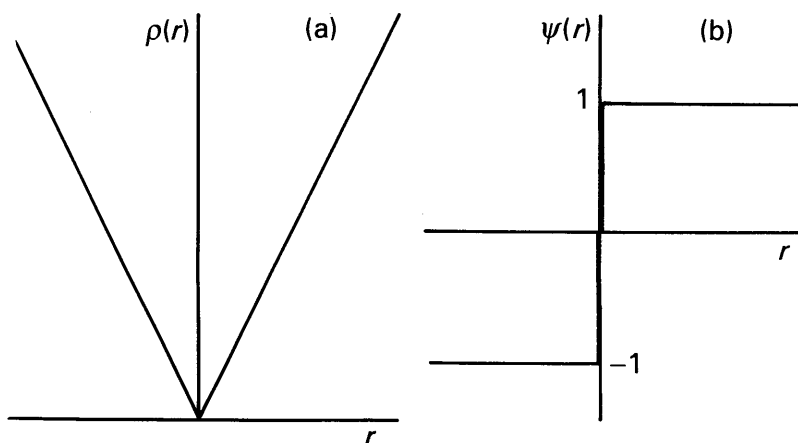


Fig. 6.58—(a) The function $\rho(r)$, and (b) function $\psi(r)$ for the L_1 approximation method.

It can be seen from Fig. 6.58b that the function $\psi(r)$ is bounded for all r by the value ± 1 . This means that the criterion (6.208) is *robust* for all residuals. The estimates \mathbf{b}_L achieved by minimization of the criterion $L_1(\mathbf{b})$ are maximum likelihood estimates when the errors ε have the Laplace distribution. For a symmetric distribution of errors with kurtosis greater than 3, the estimates \mathbf{b}_M are more effective, i.e. they have smaller variances than the estimates \mathbf{b} from the classical LS method. From Eq. (6.208) it arises that the $L_1(\mathbf{b})$ criterion consists of several linear segments. Figure 6.59 shows the dependence of $L_1(\mathbf{b})$ on b for the case of $m = 1$, for a regression straight line passing through the origin.

From Fig. 6.59b it is evident that many different estimates may exist that correspond to a minimum of $L_1(\mathbf{b})$. Minimization of the criterion $L_1(\mathbf{b})$ is a linear programming problem, i.e. to search for the minimum of $[\sum_{i=1}^n (e_i^+ + e_i^-)]$ when

$$\mathbf{X}\mathbf{b}_L + \mathbf{e}^+ + \mathbf{e}^- = \mathbf{y}$$

where $\mathbf{e}^+, \mathbf{e}^- \geq 0$ are vertical deviations from the regression plane $\mathbf{X}\mathbf{b}_L$ (Fig. 6.60).

Estimates for parameters b_L may be obtained by the program IRWLS by using the weights $w_i(r) = 1/|r_i|$ [Eq. (6.202)].

For simple regression models such as the equation of a straight line, we can use the condition that the regression function corresponding to a minimum of the criterion $L_1(\mathbf{b})$ must go through just m experimental points. This fact can be used in writing an algorithm which, for all combinations of m points, determines the parameter

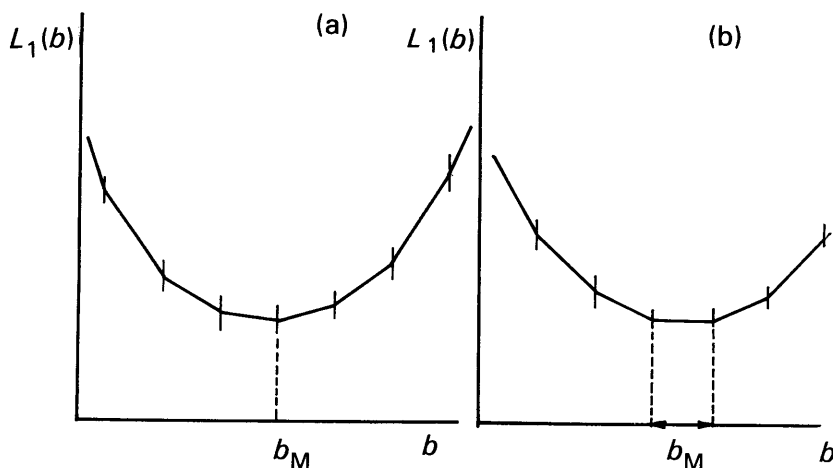


Fig. 6.59—Two possible shapes of the criterion function $L_1(\mathbf{b})$: (a) with an obvious minimum, (b) the minimum covers an interval.

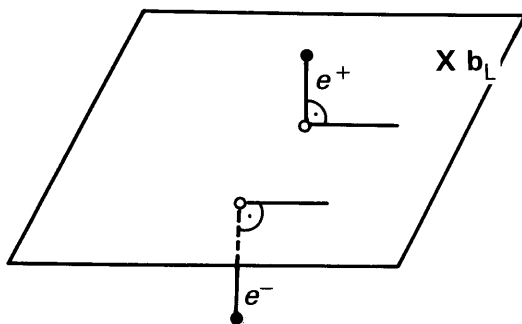


Fig. 6.60—Representation of the vertical deviations e^+ and e^- from the regression plane $X\mathbf{b}_L$.

estimates by solving linear equations for m unknowns, and the values for the estimates \mathbf{b}_M are those for which the criterion $L_1(\mathbf{b})$ has a minimum. For a regression straight line, and a small number of experimental points, this algorithm is rather simple i.e. it formulates a search for the slope and intercept of a straight line going through two points.

Problem 6.48. *Examination of the relationship between the change of surface energy of adsorption and the effective specific surface of a sorbent*

Nikulichev and Kanchenko [41] studied the adsorption of stearic acid from decane solution onto various sorbents including stearates and 12-oxystearates. The effective specific surface S_e of these sorbents and the change of surface energy as a consequence of adsorption, $-\Delta G$, were measured. Compare the L_1 approximation and LS methods and, construct a linear model between the variables $-\Delta G$ and S_e .

Data:

$S_e, \text{m}^2 \cdot \text{kg}^{-1}$	2.6	3.3	4.4	4.2	6.2	6.5
$-\Delta G, \text{kJ} \cdot \text{mol}^{-1}$	17.8	18.6	16.2	17.3	15.8	15.2

Solution: The LS method estimates the regression equation as

$$-\Delta G = -0.7524 S_e + 20.23$$

with the residual sum of squares $RSS = 0.266$ and the mean absolute deviation $A = 0.44$.

By determination of all possible straight lines going through two points and substituting into the $L_1(\mathbf{b})$ criterion, the following model was found:

$$-\Delta G = -0.5556 S_e + 19.24$$

with the residual sum of squares $RSS = 0.352$ and the mean absolute deviation $A = 0.435$. The two straight lines are compared in Fig. 6.61.

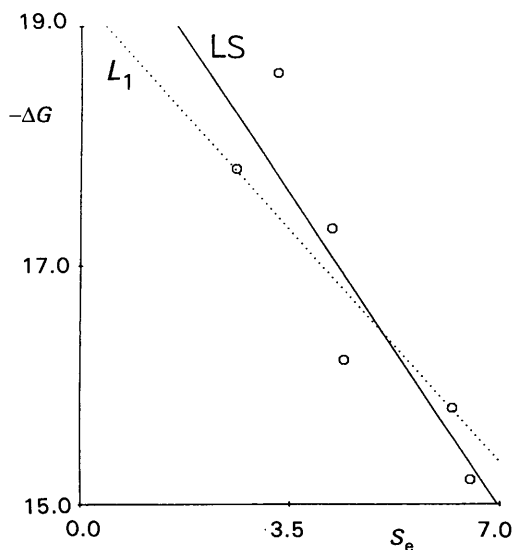


Fig. 6.61—Comparison of the regression straight lines found by the LS method (LS) and the L_1 approximation (L_1).

Conclusion: For a regression straight line, application of the L_1 approximation is simple. This method is robust enough to cope with the outlying point number 2.

Statistical analysis of the results of the L_1 approximation depends on the asymptotic normality of estimates \mathbf{b}_L . The covariance matrix $D(\mathbf{b}_L)$ is calculated from Eq. (6.206), and the variance τ^2 is estimated from

$$\tau = \frac{1}{2p(\bar{\epsilon})}$$

where $p(\tilde{e})$ is the probability density function of errors at the median. It is approximately valid that

$$p(e) = 0.5 (e_{0.75} - e_{0.25})$$

where $e_{0.75}$ is the upper and $e_{0.25}$ is the lower quartile of the residuals (Chapter 2). Statistical analysis is similar to the LS method, but instead of $\hat{\sigma}^2$ the quantity τ^2 is used. Many authors describe the L_1 approximation as a generally robust method, but this method is robust only with reference to outlying points and not to leverages.

Problem 6.49. *Comparison of the robustness of the LS method with the L_1 approximation, in the presence of one influential point*

To illustrate the efficiency of robustness of the two methods, the LS and L_1 approximation, six data points are used. The first data set (A) contains one outlier (y is equal to 10 instead of the correct value, 1) and the second data set (B) contains one leverage (x is equal to 10 instead of the correct value 1). If the regression method is robust enough it should estimate both parameters $\beta_1 = 1$ and $\beta_2 = 0$ in the model $E(y/x) = \beta_1 x + \beta_2$. Estimate b_1 and b_2 by the LS and L_1 approximation methods, and compare the results.
Data:

Data set A	x	1	2	3	4	5	6
	y	10	2	3	4	5	6
Data set B	x	10	2	3	4	5	6
	y	1	2	3	4	5	6

Solution: The estimates of parameters β_1 and β_2 by the LS and L_1 approximation method are listed in Table 6.14, and the regression straight lines are shown in Fig. 6.62.

Table 6.14. Comparison of parameter estimates b_1 and b_2 found by the LS and L_1 -approximation methods for the model $y = 1 \times x + 0$.

	LS method		L_1 approximation	
	b_1	b_2	b_1	b_2
Data set A	−0.2857	6	1	0
Data set B	−1.25	4.125	−0.2857	3.857

The poor robustness of the LS method leads to a change of sign of the slope of the straight line. The L_1 approximation is robust enough towards the outlying point (set A) but not towards the leverage point (outlying in x value) (set B).
Conclusion: The L_1 approximation method is not generally robust enough to cope with all types of influential points. One influential point can be enough for both methods to give a false estimate of slope; the difference may be big enough to result in a change of sign of the slope.

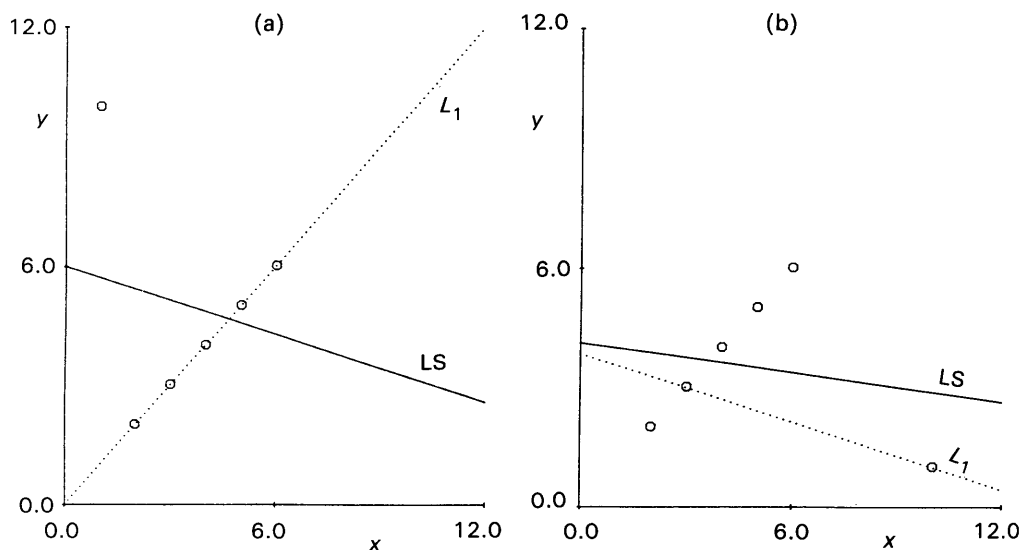


Fig. 6.62—Robustness of the LS method and the L_1 approximation for (a) set A (one outlier) and (b) set B (one leverage and outlier together).

6.6.5.3 Robust estimates with bounded influence

By using a convenient choice of function $\rho(r)$ or $w(r)$, the robust M-estimates may be found by the iterative reweighted least-squares method. Table 6.15 lists the most frequently used types of M-estimates for $\rho(r)$ and $w(r)$ with numerical values of the constant term.

Table 6.15. Functions $\rho(r)$ and $w(r)$ for five selected robust methods

Author of method	$\rho(r)$	$w(r)$	Range	Constant
Andrews	$A^2(1 - \cos(r/A))$	$(A/r) \cdot \sin(r/A)$	$ r \leq A$ $ r > A$	$A = 1.339$
Tukey	$(B^2/2) (1 - (1 - (r/B)^2)^3)$	$(1 - (r/B)^2)^2$	$ r \leq B$ $ r > B$	$B = 4.865$
Huber	$B^2/2$ $r^2/2$ $k r - k^2/2$	1 $k/ r $	$ r \leq k$ $ r > k$	$k = 1.345$
Talwar	$r^2/2$ $T^2/2$	1 0	$ r \leq T$ $ r > T$	$T = 2.795$
Welsch	$(W^2/2) (1 - \exp(-(r/w)^2))$	$\exp(-(r/w)^2)$	—	$W = 2.985$

For analysis of chemometric problems the Tukey “biweight” is recommended. It is suitable for calculation of $r_i = e_i/\sigma$, with $\sigma = S$ from Eq. (6.205). For the normal distribution of errors this estimate is equal to 0.67σ .

The estimates of parameters determined by methods from Table 6.15 are robust only on outliers but not on leverages (outlying in x values). To ensure robustness against all types of influential points, estimates with bounded influence are constructed. In the simplest way the set of equations (6.201) is modified by introducing new

weights $V(\mathbf{x}_i)$ into the expression

$$\sum_{i=1}^n \psi(e_i/\sigma)x_{ij}V(\mathbf{x}_i) = 0$$

The weights $V(\mathbf{x}_i)$ eliminate the influence of leverages and are proportional to the magnitudes H_{ii} of the diagonal elements of the projection matrix \mathbf{H} .

Krasker and Welsch [43] recommend selecting weights by using the expression

$$V(\mathbf{x}_i) = \frac{1 - H_{ii}}{\sqrt{H_{ii}}}$$

Effective procedures for a construction of estimates with bounded influence may be found in the work of Hettmansperger [44]. Introduction of the weights $V(\mathbf{x}_i)$ from Eq. (6.209) into computer programs does not cause any difficulty. In the IRWLS method the weights $w_i(r)$ are replaced by weights $V(\mathbf{x}_i)w_i(r)$.

Problem 6.50. *Examination of the robustness of estimates with bounded influence*

Estimate the parameters of the regression straight line for data from Problem 6.49 by using a combination of Welsch weights $V(\mathbf{x}_i)$ (Table 6.15) in Eq. (6.209).

Data: from Problem 6.49

Solution: Table 6.16 lists estimates of the parameters for the regression straight line for both data sets. Although both sets contain one strongly influential point, the slope estimates are always equal to 1.

Table 6.16. Estimates b_1 and b_2 of parameters $\beta_1 = 1$ and $\beta_2 = 0$ determined by the use of Welsch weights

	b_1	b_2
Set A	1	1.87×10^{-6}
Set B	0.995	0.0196

Conclusion: Estimates with limited influence are robust against all types of influential points.

Some other global robust methods exist. Strong robust methods are methods in which, instead of the sum of squared residuals, the median of squared residuals is sought. These robust methods may be used to locate groups of influential points.

Problem 6.51. *Operation of a plant for the oxidation of ammonia to nitric acid*

The operation of a plant for the oxidation of ammonia to nitric acid was studied [45], and a set of data from 21 days of operation was collected. The dependent variable y represents the percentage of the input ammonia that is lost by escaping as unabsorbed nitric oxides. This is an inverse measure of the yield of nitric acid for the plant. Three independent variables are x_1 the rate of operation, x_2 the temperature of the cooling water in the coils of the absorption tower for the nitric acid, and x_3 the concentration of nitric acid in the absorbing liquid. Investigation of plant

operations indicates that the following sets of runs can be considered as replicates: (1, 2), (4, 5, 6), (7, 8), (11, 12), and (18, 19). While the runs in each set are not exact replicates, the points are sufficiently close to each other in x -space for them to be used as such. Suppose the linear model is $E(y/x) = \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4$, and apply the LS method and the L_1 approximation method to indicate influential points. Data: (*denotes the strongly influential point in output)

Run No.	Stack Loss %	Air Flow x_1	Temperature x_2	[HNO ₃] x_3	Residual \hat{e}_{L1}	Concl.
	y					
1	42	80	27	89	2.53	*
2	37	80	27	88	1.85×10^{-6}	
3	37	75	25	90	2.715	*
4	28	62	24	87	3.814	*
5	18	62	22	87	-0.61	
6	18	62	23	87	-0.89	
7	19	62	24	93	0.50	
8	20	62	24	93	0	
9	15	58	23	87	-0.73	
10	14	58	18	80	-2.8×10^{-6}	
11	14	58	18	89	0.279	
12	13	58	17	88	0.039	
13	11	58	18	82	-1.43	
14	12	58	19	93	-0.887	
15	8	58	18	89	0.601	
16	7	50	18	86	0.008	
17	8	50	19	72	-0.217	
18	8	50	19	79	9.1×10^{-5}	
19	9	50	20	80	0.241	
20	15	56	20	82	0.812	
21	15	70	20	91	-4.722	*

Solution: The classical LS method finds the regression equation

$$\hat{y} = -37.68 + 0.7336x_1 + 1.3883x_2 - 0.2164x_3$$

with determination coefficient $\hat{R}^2 = 0.913$ and residual standard deviation $\hat{\sigma} = 3.243$. The partial regression graphs for the independent variable x_1 and x_2 indicate that points 21 and 4 are outliers. The L-R graph indicates that point 21 is a strongly influential point and points 1, 3 and 4 are less influential.

With the use of the L_1 approximation, the regression equation takes the form

$$\hat{y} = -39.65 + 0.83x_1 + 0.581x_2 - 0.0621x_3$$

with mean absolute deviation $A_p = 2.004$. The last column in the Table shows the

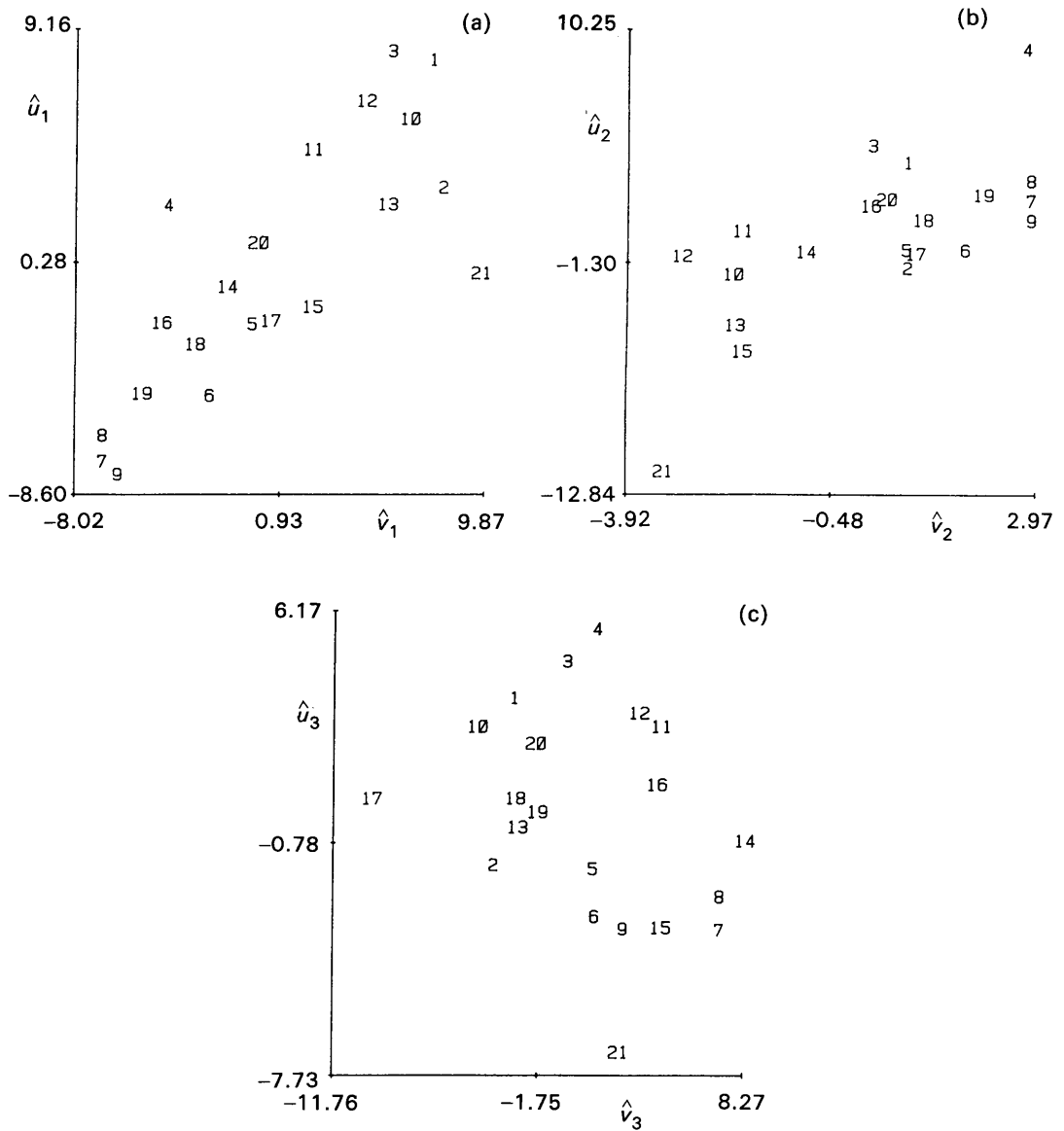


Fig. 6.63—Partial regression graph for independent variables (a) x_1 , (b) x_2 , and (c) x_3 .

residuals $\hat{e}_{L1} = \hat{e}/A_p$, which indicate that points 21 and also 1, 3 and 4 are influential.
Conclusion: The robust methods can be useful for identification of influential points.

Problem 6.52. Examination of the effect of three different factors on the amount of ozone in air

The dependence of the amount of ozone in the air (y) on the intensity of the sun's

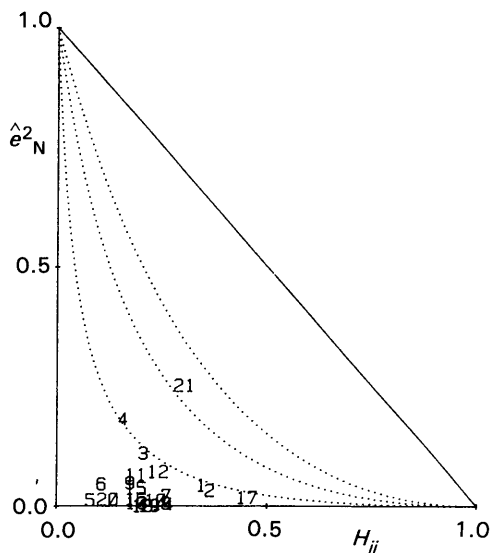


Fig. 6.64—The L-R graph for the diagnostic DF_i .

radiation for the range of wavelengths 400–700 nm (x_1), the mean velocity of wind (x_2) and the highest daytime temperature (x_3) was studied [45]. The linear model $y = \beta_1x_1 + \beta_2x_2 + \beta_3x_3 + \beta_4$ was proposed. Compare the robust estimates of the parameters with the estimates found after strongly influential points are rejected.

Data:

Measurement	y , ppm	x_1	x_2 , miles/hr.	x_3 , °F
1	41	190	74	67
2	36	118	8.0	72
3	12	149	12.6	74
4	18	313	11.5	62
5	23	299	8.6	65
6	19	99	13.8	59
7	8	19	20.1	61
8	16	256	9.7	69
9	11	290	9.2	66
10	14	274	10.9	68
11	18	65	13.2	58
12	14	334	11.5	64
13	34	307	12.0	66
14	6	78	18.4	57
15	30	322	11.5	68
16	11	44	9.7	62

17	1	8	9.7	59
18	11	320	16.6	73
19	4	25	9.7	61
20	32	92	12.0	61
21	23	13	12.0	67
22	45	252	14.9	81
23	115	223	5.7	79
24	37	279	7.4	76

Solution: The classical LS method leads to the regression equation

$$\hat{y} = -79.99 - 0.01868x_1 - 1.996x_2 + 1.963x_3$$

The partial regression graphs and the L–R graph indicate that there is only one influential point, number 23. From these graphs it is also evident that the dependence on the chosen independent variables x_1 , x_2 and x_3 is not strong. When point 23 is omitted, the regression equation becomes

$$\hat{y} = -37.52 + 0.00559x_1 - 0.7488x_2 + 0.9935x_3$$

The L_1 approximation gives the equation

$$\hat{y} = -75.36 + 0.00665x_1 - 0.3391x_2 + 1.527x_3$$

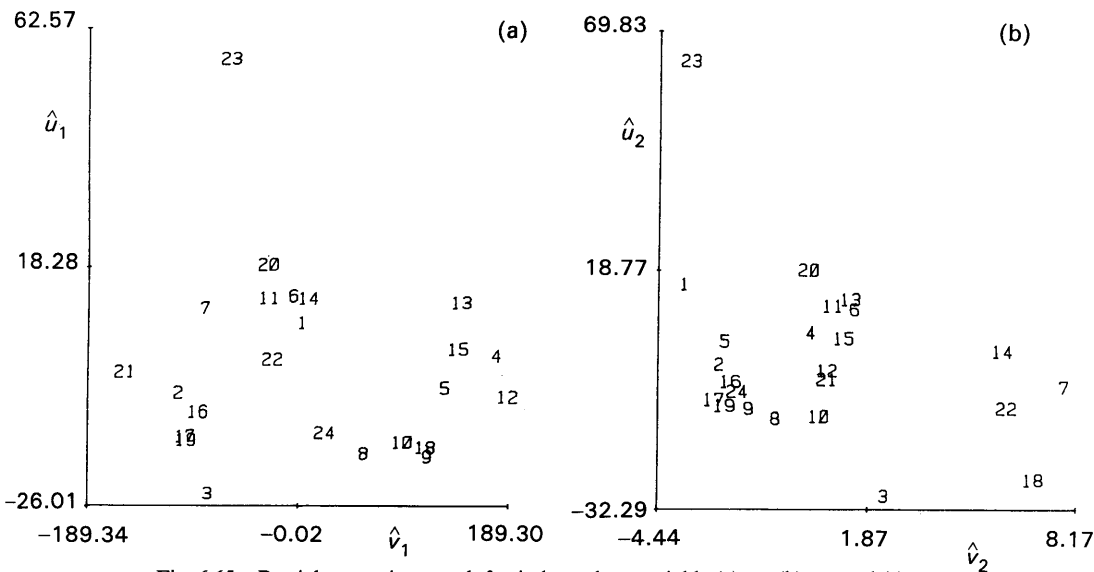


Fig. 6.65—Partial regression graph for independent variable (a) x_1 , (b) x_2 , and (c) x_3 .

It is interesting that when the LS method is used, some predicted points (amount of ozone) have a negative sign: this does not happen with the alternative regression methods. From the physical point of view, predicted values of \hat{y} should always be positive. For example, for the point 7, by the LS method $\hat{y}_7 = -0.7$, by the LS

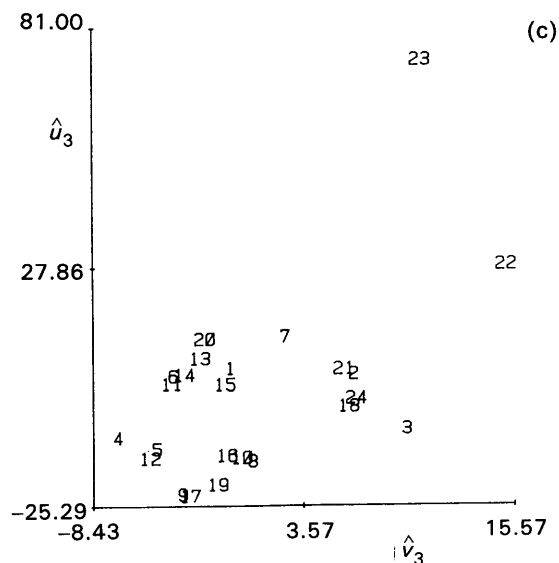
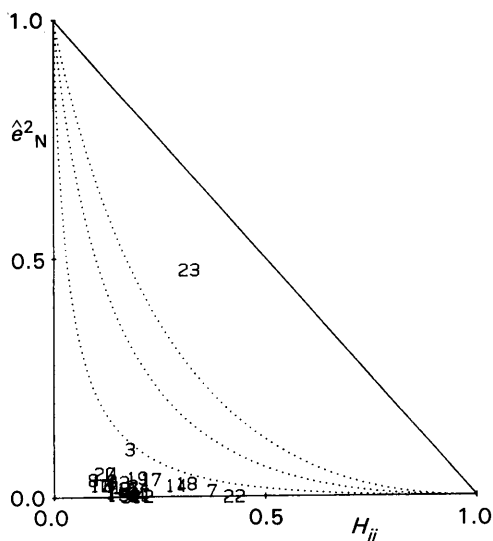


Fig. 6.65—continued

Fig. 6.66—The L-R graph for diagnostic DF_i .

method without point 23, $\hat{y}_7 = 8.14$, and by the L_1 approximation $\hat{y}_7 = 11.13$.

The presence of a single influential point caused the model to be unsuitable for prediction. Omitting the outlier (23) has a more beneficial effect than the use of the L_1 approximation.

Conclusion: Interactive data analysis based on identification of influential points often leads to better and more correct parameter estimates than robust regression.

One of greatest disadvantages of robust methods is a preference for the regression model proposed. If the proposed model is unsuitable, robust methods lead to suppression of the influence of individual points and therefore also to a suppression of the detection of unsuitable proposed models. Therefore, robust methods should be applied only with careful regard to the peculiarities of the model and data.

Sometimes it is falsely believed that the effects of influential points (outliers and leverages) are suppressed in large samples. Let us illustrate the effect of a single influential point (outlier O or leverage L) in a set of 50 points. Samples of 50 points containing one outlier O or one high-leverage point L, were generated, and the influence on the estimate b_1 of the slope of the regression straight line was examined. Data x_i were simulated by generation from a rectangular distribution $R[0, 1]$, and then linearly transformed into the interval (10, 20). The variable y_i was calculated using the relation

$$y_i = \beta_1 \times x_i + \beta_2 + N(0, 1)$$

where $\beta_1 = 1$, $\beta_2 = 1$ and $N(0, 1)$ is a random number with standardized normal distribution. Into this data, for point 37, either an outlier O or a leverage L was introduced. Table 6.17 lists values of the slope estimate b_1 of the regression straight line $y = x + 1$, when point 37 has magnitudes 40, 70, and 80. Regression analysis was performed by classical LS (LS), by the M-estimate with Welsch weights from Table 6.15 (WR) and by the estimate with bounded influence (KWR).

Table 6.17. Estimation of slope b_1 by three different regression methods: LS, WR and KWR for a data set of 50 points with one influential point (O = outlier, L = leverage point). $\beta_1 = 1$.

Value of influential point	LS		WR		KWR	
	L	O	L	O	L	O
40	0.46	1.180	0.926	0.973	0.958	0.962
70	0.17	1.490	0.203	0.969	0.965	0.959
80	0.13	1.602	0.140	0.969	0.965	0.959

Table 6.17 illustrates that one outlier or leverage point in a set of 50 points causes the classical LS method or the M-estimate method to determine a totally false estimate of the slope b_1 .

The method of slope estimate with bounded influence (KWR) is robust and found a true estimate of parameter β_1 .

This example of the influence of a single outlier or high leverage point indicates that without an analysis of influential points in interactive co-operation with the computer, routine data treatment may be totally invalidated by false and meaningless estimates. Just one decimal point falsely written may cause totally erroneous parameter estimates.

6.7 CALIBRATION

Calibration is one of the most important applications in the chemical laboratory for regression analysis. Calibration consists of two steps:

- (1) building a calibration model;
- (2) application of the calibration model.

Building a calibration model is identical with the task of building a regression model. The second step of calibration involves inversion of the first step, i.e. for a measured response y^* the corresponding value x^* and its statistical characteristics are calculated. The main attention in this section is paid to calibration straight lines.

6.7.1 Types of calibration and calibration models

Calibration tasks have been classified according to different criteria by Rossenblatt and Spiegelman [46].

(1) Absolute calibration is the most frequently used procedure in chemical instrumentation. In the construction of a calibration model, the measured quantity η , called the signal (potential, EMF, electric resistance, pH, absorbance, etc.) is related to the quantity ξ which describes a state or a property of the system (composition, concentration, temperature, time, etc.). An example of an absolute calibration is the dependence of the absorbance of a solution (η) on its concentration (ξ).

In a calibration experiment for n samples with known (or precisely measured) values of variable ξ , the corresponding quantities η are measured. Frequently, both variables are monitored instrumentally, and there will be n points $\{x_i, y_i\}$, $i = 1, \dots, n$, where

$$y_i = \eta_i + \varepsilon_i \quad (6.210a)$$

$$x_i = \xi_i + \delta_i \quad (6.210b)$$

where ε_i and δ_i are experimental errors. If the variable ξ_i is measured precisely, or exactly defined standards are used, $\delta_i = 0$, $i = 1, \dots, n$. The quantity η_i is replaced by a calibration model $f(x, \beta)$, and data treatment leads to estimation of parameter β .

In the second phase, there are M repeated measured values of an analytical signal $\{y_j^*\}$, $j = 1, \dots, M$, from which the mean value of property \hat{x}^* with its confidence interval is estimated. An example of a signal that depends on concentration is illustrated in Fig. 6.67, where symbols L_L and L_U denote the lower and upper limits of the confidence interval of concentration. In the rest of this chapter we will consider only absolute calibration.

(2) Comparative calibration is a procedure in which one instrument is calibrated against a second one, and either may be used as the standard. An example is the determination of concentration with the use of absorbance (Lambert–Beer law) as the first method, and potentiometric titration as the second method. Absorbance values are compared with volumes of titrant added. The errors δ_i are not negligible, and to construct the calibration model the regression analysis for the case when both variables are subject to experimental errors must be used.

With reference to the application of the calibration model, the following cases may be distinguished:

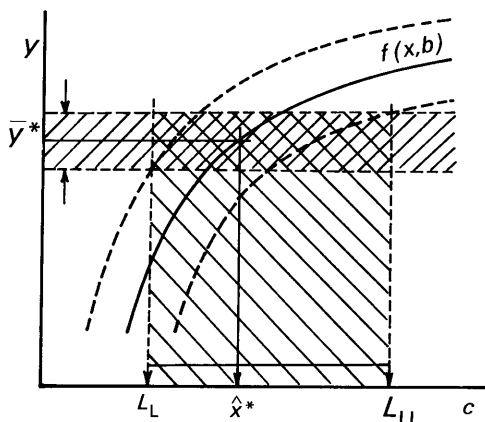


Fig. 6.67—Absolute calibration and a procedure for determination of concentration \hat{x}^* for the mean value of analytical signal \bar{y}^* . L_L and L_U are the lower and upper limits of the confidence interval of concentration.

- single application of the calibration model: the calibration model is constructed from n measured points $\{x_i, y_i\}$, $i = 1, \dots, n$, and then one estimate \hat{x}^* with its confidence interval is calculated from one y^* value;
- multiple application of the calibration model: from the calibration model, several estimates \hat{x}^* are determined from values of the analytical signal.
- single or multiple application in combination with other measurements: the result of the second phase of calibration is used together with other variables and constants for determination of a quantity which is a function of more variables. Here, any bias in the estimates \hat{x}^* which will be included in the final systematic error of the result.

The difficulty of the calibration task depends on the model used. For linear regression models, the confidence bands around the model may be expressed by Eq. (6.45) or for all possible values by Eq. (6.45a). The components of vector \mathbf{x} are functions of a measured property (i.e. usually concentration), and when polynomial models are considered, the individual components correspond to powers of this measured property. To find a value of \hat{x}^* , a root of a polynomial must be found.

For nonlinear regression models the solution is sought in the form

$$\hat{x}^* = f^{-1}(y^*) \quad (6.211)$$

On the base of the Taylor series for this function, the approximate formula for the variance $D(\hat{x}^*)$ may be found in the form [47]

$$D(\hat{x}^*) \approx \left[\frac{\partial f(x, \mathbf{b})}{\partial x} \right]^{-2} \left[\frac{D(y^*)}{M} + D(f(x, \mathbf{b})) \right] \quad (6.212)$$

where $D(y^*)$ is the variance of y^* values, usually equal to σ^2 and $D(f(x, \mathbf{b})) = D(\hat{y})$ is the variance of prediction, estimated from the Taylor series of function $f(x, \mathbf{b})$. For

the linear regression model the variance of prediction is given by

$$D(\hat{y}) = \sigma^2 \left[\frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] = \sigma^2 \left[\frac{1}{n} + \frac{(y^* - \bar{y})^2}{b_1^2 \sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

where b_1 is the estimate of the slope of the regression line. On substituting into Eq. (6.212) we obtain

$$D(\hat{x}^*) \approx \frac{\sigma^2}{b_1^2} \left[\frac{1}{M} + \frac{1}{n} + \frac{(y^* - \bar{y})^2}{b_1^2 \sum_{i=1}^n (x_i - \bar{x})^2} \right] \quad (6.212a)$$

Difficulties are caused by the generally non-symmetric distribution of quantity x^* . Only in the case of a calibration straight line and small residual variance can the distribution of x^* be assumed to be approximately normal [48].

If both y and y^* are random variables with normal distribution, the difference $\Delta = \bar{y}^* - f(x^*, \mathbf{b})$ will also have the normal distribution. The standardized random variable $\Delta/\sqrt{D(\Delta)}$ has the Student distribution with the number of degrees of freedom used for determination of $D(\Delta)$. To find the $100(1 - \alpha)\%$ confidence interval of the quantities \hat{x}^* defined in Eq. (6.211) it is necessary to solve the equation [51]

$$(\bar{y}^* - f(\hat{x}^*, \mathbf{b}))^2 = F_{1-\alpha}(1, r) \times D(\bar{y}^* - f(\hat{x}^*, \mathbf{b}))$$

where $r = n - 2$. The variance $D(\Delta) = D(\bar{y}^* - f(\hat{x}^*, \mathbf{b}))$ may be estimated by the Taylor series expansion of function $f(\hat{x}^*, \mathbf{b})$. It is approximately valid that

$$\begin{aligned} D(\Delta) = D(\bar{y}^*) + \sum_{j=1}^m \left[\frac{\partial f(\hat{x}^*, \mathbf{b})}{\partial b_j} \right]^2 \times D(b_j) \\ + 2 \sum_{i=1}^{m-1} \sum_{j>1}^m \frac{\partial f(\hat{x}^*, \mathbf{b})}{\partial b_i} \times \frac{\partial f(\hat{x}^*, \mathbf{b})}{\partial b_j} \text{cov}(b_i, b_j) \end{aligned} \quad (6.213)$$

where m is the number of regression parameters. The special case represented here is the model of the calibration straight line

$$y = b_1(x - \bar{x}) + \bar{y} \quad (6.213a)$$

for which, after substitution into Eq. (6.213), we obtain

$$D(\Delta) = \hat{\sigma}^2 \left[\frac{1}{M} + \frac{1}{n} + \frac{(\hat{x}^* - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \quad (6.213b)$$

When $D(\Delta)$ is known, the two limits of the confidence interval of \hat{x}^* may be estimated. This involves finding roots of the quadratic equation

$$\bar{y}^* - \bar{y} - b_1(\hat{x}^* - \bar{x}) = F_{1-\alpha}(1, n-2) \times D(\Delta) \quad (6.213c)$$

with respect to the variable \hat{x}^* .

6.7.2 Calibration straight line

A straight line is the usual calibration model in a chemical laboratory. Usually, it is supposed that this model fits all the measured points for a given set of variables x and y . For example, the Lambert–Beer law $A = \epsilon dc$ expresses a linear relationship between absorbance and concentration c where ϵ is the molar absorptivity and d is the path length of the cell.

In some cases, however, the straight-line model is valid only in a limited interval, and above a limiting point $\{x_A, y_A\}$ there is a significant departure from linearity. For example, the Kubelka–Munk relationship between the remission function $(1 - R^2)/2R$ and the concentration c is valid only for low concentrations. The Lambert–Beer law too is valid only up to some limiting concentration, above which curvature occurs.

For statistical data treatment, the model in the form of Eq. (6.213a) may be used, or some other equivalent expression such as

$$y_i = \beta_1 \times x + \beta_2 + \varepsilon_i, \quad i = 1, \dots, n$$

or

$$y_j^* = \beta_1 \times \kappa + \beta_2 + \varepsilon_j^*, \quad j = 1, \dots, M.$$

The task of calibration is to find an estimate of parameter x^* , the primary parameter, and of parameters β_1 and β_2 , the supplementary parameters. The estimation assumes normality of the errors ε_i and ε_j^* . The estimate \hat{x}^* and its confidence interval may be calculated by several procedures.

By substituting into Eq. (6.211) from Eq. (6.213a) we obtain the straight estimate of parameter κ in the form

$$\hat{x}^* = \bar{x} + \frac{(y^* - \bar{y})}{b_1} \quad (6.214)$$

where y^* is the measured signal (or the average \bar{y}^* for $M > 1$ repeated measurements, respectively) and b_1 is the estimate of the slope. This estimate is generally biased and a correction is made by Naszodi's modified estimates [49]

$$\hat{x}_B^* = \bar{x} + (y^* - \bar{y}) \frac{b_1}{b_1^2 + \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \quad (6.215)$$

Krutchhoft [50] proposed the inversion estimate

$$\hat{x}_1^* = \bar{x} + (y^* - \bar{y}) \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (6.216)$$

which refers to the inversion regression model

$$E(x/y) = \alpha_1(y - \bar{y}) + \alpha_2.$$

From analysis of the estimate \hat{x}_1^* it was found that it too is a biased estimate which is not better than the straight estimate \hat{x}^* . Moreover, in the estimation of parameters α_1 and α_2 it is falsely supposed that the y values are measured with negligible errors in comparison with the x -values.

Schwartz [51] proposed the nonlinear estimate given by

$$\hat{x}_N^* = \frac{\sum_{i=1}^n x_i \times \exp\left[\frac{-(y^* - b_2 - b_1 \times x_i)^2}{2\hat{\sigma}^2}\right]}{\sum_{i=1}^n \exp\left[\frac{-(y^* - b_2 - b_1 \times x_i)^2}{2\hat{\sigma}^2}\right]} \quad (6.217)$$

which, however, assumes normality of residuals.

Problem 6.53. *Point estimates of concentration from an AAS calibration line*

The atomic absorbances of solutions of various concentration of lithium were measured. Determine the calibration line and from it then the concentration of lithium for measured absorbance values $A_1 = 0.0002$, $A_2 = 0.5$ and $A_3 = 1.0$.

Data: $n = 16$

C, g of Li in 25 ml	2.5	5.0	7.5	10.0	12.5	15.0	
A	0.063	0.120	0.189	0.251	0.316	0.393	
17.5	20.0	22.5	25.0	27.5	30.0	32.5	35.0
0.442	0.502	0.568	0.639	0.694	0.749	0.821	0.884
37.5	40.0						
0.947	1.010						

Solution: The classical method of least-squares leads to the regression equation

$$A = 0.02525(\pm 0.00011)C + 0.0002(\pm 0.0028)$$

with correlation coefficient $\hat{R} = 0.9999$. Table 6.18 lists the estimates \hat{x}^* , \hat{x}_B^* , \hat{x}_1^* and \hat{x}_N^* for $A = 0.0002$, 0.5 and 1.

Table 6.18. Concentration estimates by various methods

Absorbance	\hat{x}^*	\hat{x}_B^*	\hat{x}_I^*	\hat{x}_N^*
0.0002	0	4.32×10^{-4}	6.05×10^{-3}	2.5
0.5	19.795	19.795	19.795	20
1.0	39.597	39.597	39.592	40

Within experimental error, all estimates except the nonlinear one lead to the same result.

Conclusion: For sufficiently precise data with small spread around the regression straight line, the classical estimate \hat{x}^* is satisfactory.

In the construction of confidence intervals of the estimates \hat{x}^* and \hat{x}_B^* for more scattered data, the simplest is the determination of $D(\hat{x}^*)$ and to use Eq. (6.213) with an assumption of normality. The limits of the 95% confidence interval are calculated by

$$L_L = \hat{x}^* - 1.96\sqrt{D(\hat{x}^*)} \quad (6.218a)$$

$$L_U = \hat{x}^* + 1.96\sqrt{D(\hat{x}^*)} \quad (6.218b)$$

To construct the confidence interval, the ratio

$$\frac{[(b_2 + b_1) \times (\hat{x}^* - y^*)]^2}{\sigma^2 \left[\frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]} \approx F_{1-\alpha}(1, n-2) \quad (6.219)$$

is often used. This ratio exhibits the Fisher–Snedecor distribution with 1 and $(n-2)$ degrees of freedom. The corresponding $100(1-\alpha)\%$ confidence interval of parameter x is calculated from

$$L_{L,U} = \bar{x} + \frac{(y^* - \bar{y}) \pm \hat{\sigma} \sqrt{F_{1-\alpha}(1, n-2)} \left[\frac{1+\lambda}{n} + \frac{(y^* - \bar{y})^2}{b_1^2 \sum_{i=1}^n (x_i - \bar{x})^2} \right]}{b_1 \times (1-\lambda)} \quad (6.220)$$

Parameter λ is given by

$$\lambda = \frac{\hat{\sigma}^2 \times F_{1-\alpha}(1, n-2)}{b_1^2 \sum_{i=1}^n (x_i - \bar{x})^2}$$

is the variation coefficient of slope β_1 . When this ratio is smaller than 0.1 the slope estimate is sufficiently precise for the approximate confidence interval for parameter κ to be used in the form

$$L_{L,U} = x^* \pm t_{1-\alpha/2}(n-2) \times \frac{\hat{\sigma}}{|b_1|} \sqrt{\frac{1}{n} + \frac{(y^* - \bar{y})^2}{b_1^2 \sum_{i=1}^n (x_i - \bar{x})^2}} \quad (6.221)$$

If we want information about the whole regression (calibration) line, we replace in Eq. (6.220) the term $\sqrt{F_{1-\alpha}(1, n-2)}$ by the term $\sqrt{2F_{1-\alpha}(2, n-2)}$.

The Scheffe's confidence interval of one predicted value y^* at \hat{x}^* is calculated by

$$L_{L,U} = y^* \pm \sqrt{2F_{1-\alpha}(2, n-2)} \hat{\sigma} \left[1 + \frac{1}{n} + \frac{(\hat{x}^* - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]^{1/2} \quad (6.222)$$

This confidence interval is larger by the variance σ^2 because the variable y^* is used instead of its mean value $E(y^*)$. By rearrangement of Eq. (6.222) we find the $100(1 - \alpha)\%$ confidence interval of variable κ in the form

$$L_{L,U} = \hat{x}^* \pm \frac{\sigma \sqrt{2F_{1-\alpha}(2, n-2)}}{\lambda_1} \left[\lambda_1 + \frac{\lambda_1}{n} + \frac{(y^* - \bar{y})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]^{1/2} \quad (6.223)$$

where

$$\lambda_1 = \frac{b_1^2 - \sqrt{2F_{1-\alpha}(2, n-2)}}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

When the arithmetic mean \bar{y}^* is used, the term 1 is replaced by term $1/M$ in the brackets of Eq. (6.222). An analogous adjustment can be made to give Eq. (6.223), which corresponds to Eq. (6.213c). The graphical interpretation of the confidence interval of parameter x is shown in Fig. 6.68.

When there are replicate values of y , and \bar{y}^* has been determined, the confidence straight lines U_L and U_U should be calculated. The intersection of straight line U_U with the lower confidence parabola P_L of the calibration straight line leads to point L_U and the intersection of straight line U_L with the upper confidence parabola P_U leads to point L_L .

If the variance of measurement, σ^2 , is known it is easy to define the $100(1 - \alpha)\%$ confidence interval of signal \bar{y}^* in the form

$$U_{L,H} = \bar{y}^* \pm u_{1-\alpha/2} \sigma$$

where $u_{1-\alpha/2}$ is the quantile of the normalized normal distribution. If σ^2 is unknown, the inequality

$$\sigma^2 < \frac{(n-2)\hat{\sigma}^2}{\chi_{\alpha/2}^2(n-2) \times M} \quad (6.224)$$

may be used, where $\chi_{\alpha/2}^2$ is the lower $100 \alpha/2\%$ quantile of the χ^2 -distribution. The

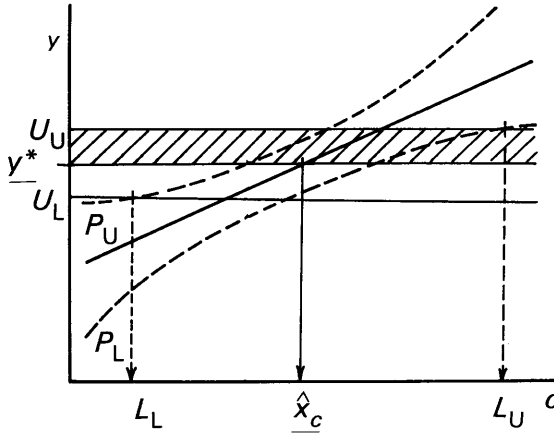


Fig. 6.68—Determination of the confidence interval of parameter x for a calibration straight line. The confidence interval of the signal is indicated by the hatched area.

confidence interval of signal $U_{L,H}$ is then calculated from

$$U_{L,H} = \bar{y}^* \pm u_{1-\alpha/2} \times \frac{\hat{\sigma}}{\sqrt{M}} \sqrt{\frac{n-2}{\chi_{\alpha/2}^2(n-2)}} \quad (6.225)$$

Instead of the quantile $u_{1-\alpha/2}$ in this equation, for $M = 1$ the more convenient quantile of the Student distribution $t_{1-\alpha/2}(n-2)$ is used and the variance σ^2 is replaced by its estimate $\hat{\sigma}^2$.

From Eq. (6.45a) the limiting $100(1 - \alpha)\%$ confidence parabola are given by

$$P_{L,U} = b_1 x + b_2 \pm \sigma \left\{ 2F_{1-\alpha}(2, n-2) \left[\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \right\}^{\frac{1}{2}} \quad (6.226)$$

The limiting value L_U represents the solution of the equation

$$U_U = P_L \quad (6.227a)$$

with respect to variable x . The limiting value L_L is the solution of the equation

$$U_L = P_U \quad (6.227b)$$

Both equations are quadratic with respect to variable x . From Fig. 6.68 it can be seen that in some cases the intersection of a straight line with the parabola does not exist, but in other cases the confidence straight line of the signal may intersect the parabola of calibration line at two points. This indicates that the spread of data is too large when the slope of the calibration straight line is small, and such a calibration straight line is a poor model.

Problem 6.54. Various confidence intervals of concentration from a photometric calibration straight line

For the photometric calibration data from Problem 6.53, calculate the bounds of the 95% confidence interval of concentration, L_L and L_U , for the following absorbance values: $A_1^* = 0.0002$, $A_2^* = 0.5$, $A_3^* = 1$, $\bar{A} = 0.51$ (the mean of two measured values 0.50 and 0.52; $M = 2$) and $\bar{A}_4^* = 0.977$ (the mean of measured values 0.95, 0.98, 1.00; $M = 3$) by using Eqs. (6.218), (6.220), (6.221), (6.223) and (6.227).

Data: from Problem 6.53

Solution: The calculated limits $L_{L,H}$ for the 95% confidence intervals of concentration are listed in Table 6.19.

Table 6.19. The 95% confidence interval of concentration calculated by various expressions

M	A*	Expression for confidence interval				
		(6.218) lower L_L upper L_U	(6.220) L_L L_U	(6.221) L_L L_U	(6.223) L_L L_U	(6.227) L_L L_U
1	0.0002	-0.46 0.46	-0.3018 0.295	-0.299 0.299	-0.646 0.639	-0.398 0.417
1	0.50	19.37 20.22	19.65 19.94	19.65 19.94	19.20 19.94	19.39 19.20
2	(0.50; 0.52)	19.89 20.50	20.05 20.33	20.05 20.33	19.76 20.62	19.90 20.48
1	1.00	39.15 40.05	39.33 39.87	39.33 39.87	38.97 40.23	39.19 40.01
3	(0.95; 0.98; 1.00)	38.37 38.97	38.42 38.93	38.42 38.93	38.26 39.09	38.43 38.91

By using Eq. (6.223) instead of $\lambda_1(1 + 1/n)$ when $M > 1$, the term $\lambda_1(1/M + 1/n)$ is used.

The confidence intervals from Eqs. (6.220) and (6.221) do not reflect a higher precision of determination of \bar{y}^* . The confidence limits (6.227a, b) were evaluated by a simplified expression for the confidence straight line of the signal by

$$U_{L,U} = \bar{y}^* \mp u_{1-\alpha/2} \times \hat{\sigma} / \sqrt{M} \quad (6.228)$$

with $\alpha = 0.05$. From Table 6.19 it is seen that for sufficient precision of data, Eq. (6.220) and its approximation (6.221) lead to the same results. The other confidence intervals are, however, rather different. The approximation (6.218) leads to values $L_{L,U}$ which are close to the values calculated from Eq. (6.227).

Conclusion: For data with a small spread around the regression straight line, the simpler approximation (6.218) should be used. For replicate signal measurements the expressions (6.220) and (6.221) are *not* suitable.

The quality of the confidence interval around the parameter x is improved by

- (1) repeating the signal measurement y^* , i.e. increasing the number of measurements

M . For a sufficient number of replicates, M , the estimate $U_{L,U}$ can be calculated from Eq. (6.228), with σ^2 replaced by the variance σ_y^{2*} and the quantile $u_{1-\alpha/2}$ replaced by the quantile of the Student distribution $t_{1-\alpha/2}$.

- (2) The confidence parabola may be narrowed by elimination of influential points. In polynomial calibration models the confidence bands may be narrowed by the use of biased estimates calculated by the method of the rational ranks.
- (3) decreasing the residual variance $\hat{\sigma}^2$ and so increasing the precision of measurement, or by the use of a correct calibration model.

6.7.3 The precision of calibration

To express the precision of a calibration, limiting values of the concentration for which the measurement signal is still statistically significantly different from the noise are usually defined. To express precision and sensitivity of calibration methods, three levels of signal are identified:

(1) The *critical level* y_C represents the upper limit of the $100(1 - \alpha)\%$ confidence interval of the predicted signal from the calibration model for the concentration equal to zero, i.e. the *blank measurement*. By replacing $\sqrt{2F_{1-\alpha}(2, n-2)}$ by the quantile $t_{1-\alpha/2}(n-2)$ in Eq. (6.222) and setting $x = 0$, we obtain an expression for the critical level y_C in the form

$$y_C = \bar{y} - b_1 \bar{x} + \hat{\sigma} t_{1-\alpha/2}(n-2) \times \sqrt{1 + \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \quad (6.229)$$

Signals above this critical level y_C are significantly different from the noise. The concentration x_C corresponding to this critical level y_C is determined from the calibration model from

$$x_C = \frac{y_C - \bar{y}}{b_1} + \bar{x} \quad (6.229a)$$

(2) The *detection limit* y_D corresponds to the concentration for which the lower $100(1 - \alpha)\%$ confidence interval of signal prediction from the calibration model is equal to y_C . The detection limit y_D and its corresponding concentration x_D are illustrated on Fig. 6.69.

For the linear calibration model we have

$$y_D = y_C + \hat{\sigma} t_{1-\alpha/2}(n-2) \sqrt{1 + \frac{1}{n} + \frac{(x_D - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \quad (6.230)$$

Oppenheimer [52] proposed the following approximation

$$y_D = y_C + \hat{\sigma} t_{1-\alpha/2}(n-2) \times \sqrt{1 + \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \quad (6.231)$$

The corresponding concentration x_D is calculated from

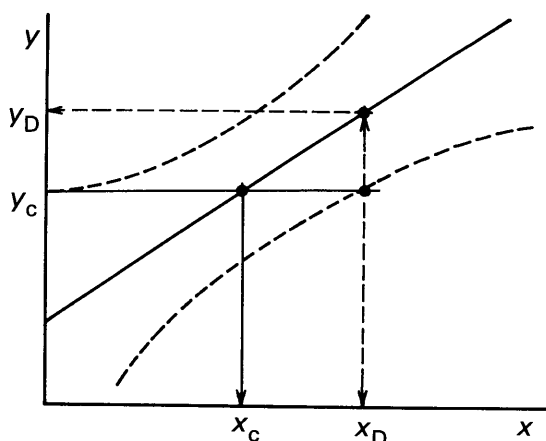


Fig. 6.69—Definition of the critical level y_c , the detection limit y_D and their corresponding concentrations x_c and x_D .

$$x_D = \frac{y_D - \bar{y}}{b_1} + \bar{x} \quad (6.231a)$$

The detection limit gives the lowest true signal level which still permits detection. The quantity x_D gives the minimum concentration which can be distinguished from zero with probability $(1 - \alpha)$.

(3) The *determination limit* y_s is the smallest signal level for which the relative standard deviation of prediction from the calibration model is sufficiently small and equal to the number C , where $C = 0.1$, usually.

If the predicted value at point x_s is given by

$$y(x_s) = \bar{y} + b_1(x_s - \bar{x})$$

and the condition of determination y_s is then equal to

$$\frac{\sqrt{D(y(x_s))}}{\hat{y}(x_s)} = C \quad (6.232)$$

Substitution and rearrangement leads to the expression

$$y_s = \frac{\hat{\sigma}}{C} \sqrt{1 + \frac{1}{n} + \frac{(x_s - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \quad (6.233)$$

In practice, in the chemical laboratory, an approximation is used, as follows.

$$y_s \approx \frac{\hat{\sigma}}{C} \sqrt{1 + \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \quad (6.234)$$

The corresponding concentration x_s is given by

$$x_s = \frac{y_s - \bar{y}}{b_1} + \bar{x} \quad (6.234a)$$

For nonlinear calibration models, Schwartz [53] recommends that the upper L_U and lower L_L limits of the confidence interval of concentration which correspond to different signal levels y^* are determined. Instead of the relative standard deviation of prediction from the calibration model, Schwartz uses the effective relative standard deviation

$$C(x') = \frac{L_U - L_L}{2x' \times t_{1-\alpha/2}(n-2)}$$

(4) The *modified determination limit* y'_s is the value of x' for which $C(x') = C$. This y'_s limit is found graphically by plotting $C(x')$ against x and substituting in the calibration model. Equation (6.235) may be used for linear models as well as nonlinear ones.

All four definitions may be simply used to calculate the detection limit y_D and the determination limit y_s for nonlinear calibration models, and for data for which the variance of measurement is not constant (heteroscedasticity) [52]. Generally, it is valid that

$$y_C \leq y_D \leq y_s$$

Ebel and Kamm [54] have described an alternative procedure of determination of the detection limit y_D , and this is illustrated in Fig. 6.70.

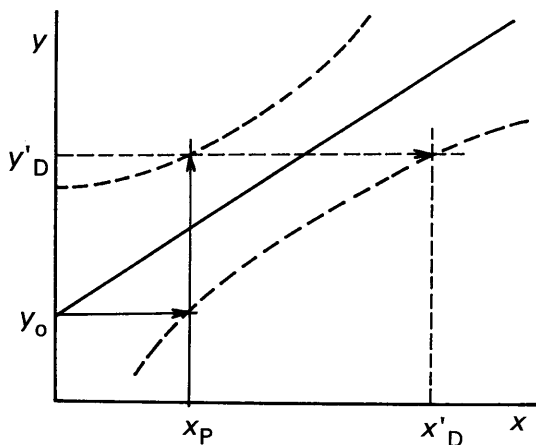


Fig. 6.70—Illustration of the procedure for determination of the detection limit y_D according to Ebel and Kamm [54].

Even for this case for linear calibration models it is easy to determine x_P by the use of L_L from Eqs. (6.220) and (6.221). Substitution of x_P into the expression for L_U leads to

$$L_U = y'_D$$

Problem 6.55. Precision and sensitivity of a photometric calibration model

For the photometric calibration from Problem 6.53 calculate the critical level y_C , the detection level y_D and the determination limit y_s for the relative standard deviation $C = 0.1$.

Data: from Problem 6.53

Solution: The values calculated are $\bar{x} = 21.25$, $\bar{y} = 0.5368$, $\sum_{i=1}^{16} (x_i - 21.25)^2 = 2125$,

$K = 0.0252$, $\hat{\sigma}^2 = 1.722 \times 10^{-6}$ and $t_{0.975}(14) = 2.14$. The limiting levels of absorbance and the corresponding concentrations are determined to be: $y_C = 0.0129$ and $x_C = 0.504$; $y_D = 0.0257$ and $x_D = 1.008$; and $y_s = 0.0593$ and $x_s = 2.339$. To calculate y_D and y_s , the approximate expressions (6.231) and (6.234) were used. Figure 6.71 shows the dependence of $C(x')$ on x' for the interval $0.1 \leq x' \leq 1$. Here, $x'_s \approx 0.6$ and $y'_s = \bar{y} + (0.6 - \bar{x}) \times b_1 = 0.0164$, for $C(x') = 0.1$.

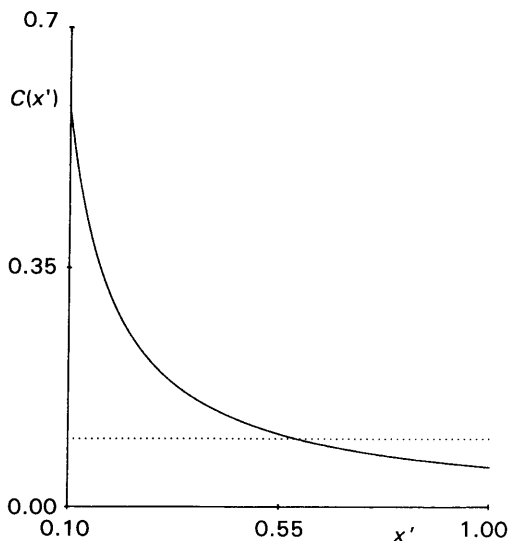


Fig. 6.71—Plot of $C(x')$ vs. x' in a search for the determination limit. Dotted line corresponds to $C = 0.1$.

Conclusion: For linear calibration models, all three limits are easy to calculate.

A modified determination limit x'_s , y'_s may be determined by use of the expression for variance $D(\hat{x})$, Eq. (6.212). Setting $C = 0.1$, the quantity x'_s is the root of the equation

$$0.01x_s'^2 = D(x') \quad (6.236)$$

For linear calibration models with $M = 1$, the variance $D(x')$ is defined by Eq. (6.212a). By a simple rearrangement, Eq. (6.236) can be transformed into a quadratic equation with the following root

$$x'_s = \frac{-\bar{x} + \sqrt{\bar{x}^2 + AD}}{A} \quad (6.237)$$

where

$$A = \frac{b_1^2 \times \sum_{i=1}^n (x_i - \bar{x})^2}{100\sigma^2} - 1$$

and

$$D = \bar{x}^2 + \frac{(n+1) \sum_{i=1}^n (x_i - \bar{x})^2}{n}$$

When x'_s is known, the value y'_s is determined by substitution of the value into the calibration model.

Problem 6.56. *Modified determination limit for the photometric calibration model*

Calculate the modified determination limit for photometric calibration model from Problem 6.53.

Data: from Problem 6.53

Solution: From Eqs. (6.236) and (6.237) the following numerical values are calculated: $A = 7936.5$, $D = 2709.375$ and $x'_s = 0.582$, which corresponds to $y'_s = 0.01596$.

Conclusion: For the linear calibration model the modified determination limit is readily calculated from Eq. (6.237).

6.8 PROCEDURE FOR LINEAR REGRESSION ANALYSIS

The procedure for examination and construction of a linear regression model consists of following steps.

(1) *Proposal of a model*

The procedure should always start from the simplest model, with individual independent controllable variables not raised to powers other than the first, and with no interaction terms of the type $x_j x_k$ included. Only in cases when it is known that the model contains functions of the controllable variables is an exception made.

(2) *Exploratory data analysis in regression*

The scatter of individual variables and all possible pair combinations are examined. The scatter plots of x_j vs. x_k or the index plots x_j vs. j are often used here. In this step of a regression analysis, the significance of individual variables with reference to scatter and the presence of multicollinearity is examined. An approximately linear relationship between variables in scatter plots of x_j vs. x_k indicates multicollinearity.

Also, in this step the influential points causing multicollinearity are detected.

(3) *Parameter estimation*

The parameters of the proposed regression model and the corresponding basic statistical characteristics of this model are determined by the classical least-squares

method (LS). Individual parameters are tested for significance by using the Student t -test, the determination coefficient \hat{R}^2 and the predicted determination coefficient $\hat{\hat{R}}^2$. Other statistical characteristics calculated are the total F -test, the model significance test, the model complexity test, the mean quadratic error of prediction MEP and the Akaike information criterion AIC , to examine the linearity of model.

(4) *Analysis of regression diagnostics*

The statistical analysis of classical residuals leads to estimates of residual variance $\hat{\sigma}^2(\hat{e})$, residual standard deviation $s(\hat{e})$, residual skewness $g_1(\hat{e})$, residual kurtosis $g_2(\hat{e})$, the Pearson χ^2 -test of residual normality and the Jarque–Berra normality test. Different diagnostic graphs are used to examine the regression diagnostics for identification of influential points, and to test the conditions for the least-squares method, namely homoscedasticity, absence of autocorrelation, and normality of error distribution. If influential points are found, it has to be decided whether these points should be eliminated from data. If points are eliminated, the whole data treatment must be repeated. When there are several controllable variables, the significance of each variable and its function is examined by partial-regression graphs and by the partial-residual graph.

(5) *Construction of a more accurate model*

According to the test for fulfilment of the conditions for the least-squares method, and the result of regression diagnostics, a more accurate regression model is constructed as follows.

- (a) When heteroscedasticity is found in the data, the weighted least-squares method (WLS) is used.
- (b) When autocorrelation is found in the data, the generalized least-squares method (GLS) is used.
- (c) When some restrictions apply to the parameters, the conditioned least-squares method (CLS) is used.
- (d) When multicollinearity is found in the data, the method of rational ranks (MRV) is used.
- (e) When all variables are subject to random errors, the extended least-squares method (ELS) is used.
- (f) When the data have an error distribution other than normal, or the data contain outliers or high leverage points, some robust methods are used.

(6) *Evaluation of the quality of the model proposed*

With the use of classical tests, regression diagnostics and some supplementary information about the “model + data + method”, the quality of the proposed linear regression model is evaluated.

(7) *Analysis of calibration models*

For a calibration model proposed for the given signal value y^* , the quality of the independent variable x^* together with its confidence interval is estimated. Before

application of the calibration model, the detection limit and the determination limit should be estimated. These limits determine the allowable lower limit of the calibration model.

(8) *Statistical hypothesis testing*

In some cases, to compare several straight lines, statistical hypothesis testing is performed.

6.9 ADDITIONAL PROBLEMS

Problem 6.57. *The effect of influential points on the detection limit and determination limit in a photometric calibration model*

The relationship between absorbance A and the concentration of nitrate c in solution is described by the Lambert–Beer law,

$$A = \varepsilon \times d \times c + a$$

where d is the cuvette length in cm (here $d = 1$ cm), ε is the molar absorptivity and a is the absorbance of the blank. Estimate (i) the parameters of the Lambert–Beer calibration model, (ii) the detection limit and determination limit for nitrate.

Data: $n = 16$

c , mg of NO_3^- in 25 ml	0.005	0.0161	0.0165	0.0213	0.0275
A	0.110	0.272	0.224	0.274	0.338
0.0324	0.0382	0.0453	0.0523	0.0575	0.0632
0.389	0.449	0.522	0.595	0.649	0.708
0.0803	0.0862	0.0918	0.0982		
0.885	0.946	1.005	1.067		

Solution: The classical least-squares method gives the regression model

$$A = 10.20(\pm 0.11)c + 0.06478(\pm 0.00635),$$

where the standard deviations of the parameter estimates are given in brackets. Both parameter estimates are statistically significant at significance level $\alpha = 0.05$. The determination coefficient $\hat{R}^2 = 0.9984$, the mean quadratic error of prediction $MEP = 0.0001904$ and residual standard deviation $\hat{\sigma} = 0.01257$. Heteroscedasticity is identified in the data by the criterion S_f , and the sign test indicates a trend in residuals. Figure 6.72a shows the regression model with the 95% confidence intervals and Fig. 6.72b is a plot of the classical residuals *vs.* c .

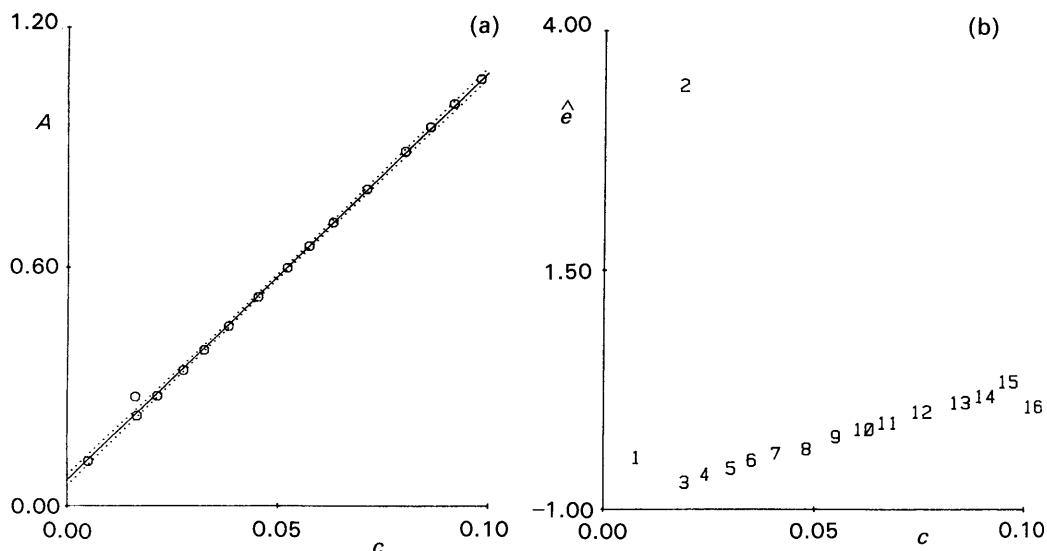


Fig. 6.72—(a) Regression model with the 95% confidence interval, and (b) graph of classical residuals \hat{e} vs. concentration c .

The graphs indicate that point 2 is an outlier and causes the heteroscedasticity in data and the trend in residuals. When point 2 is compared with 3 and 4, it is obvious that 2 is a mistake. The rankit graph of normalized residuals (Fig. 6.73a) leads to the same conclusion. The rankit graph for recursive residuals is also interesting (Fig. 6.73b) as a consequence of one outlier in a base from which the first estimates of the parameters are calculated.

From Fig. 6.74 it is evident that, apart from point 2, which is strongly masking the influence of other points, there are other influential points such as 3 and 16, and to some extent 4 and 1.

When point 2 is omitted, the classical least-squares method gives the residual regression model

$$A = 10.33(\pm 0.013) \times c + 0.05497(\pm 0.00078)$$

with determination coefficient $\hat{R}^2 = 0.9999$, and $MEP = 2.92 \times 10^{-6}$. The residual standard deviation $\hat{\sigma} = 0.00143$ demonstrates a significant improvement in the statistical regression characteristics, too. In residuals there is no evidence of trends nor heteroscedasticity. Omitting point 2 caused just small changes in the numerical estimates of the parameters. Figure 6.75a shows the regression model and Fig. 6.75b the residual plot without point 2.

Despite the excellent degree of fit of the regression straight line to the experimental points (Fig. 6.75a), the residual plot (Fig. 6.75b) indicates the presence of some other influential points, i.e. points 1 and 15 (previously 16) and to some extent 14. Valuable information about influential points is found from the L-R graph for D_i and McCulloh–Meeter graph, in Fig. 6.76.

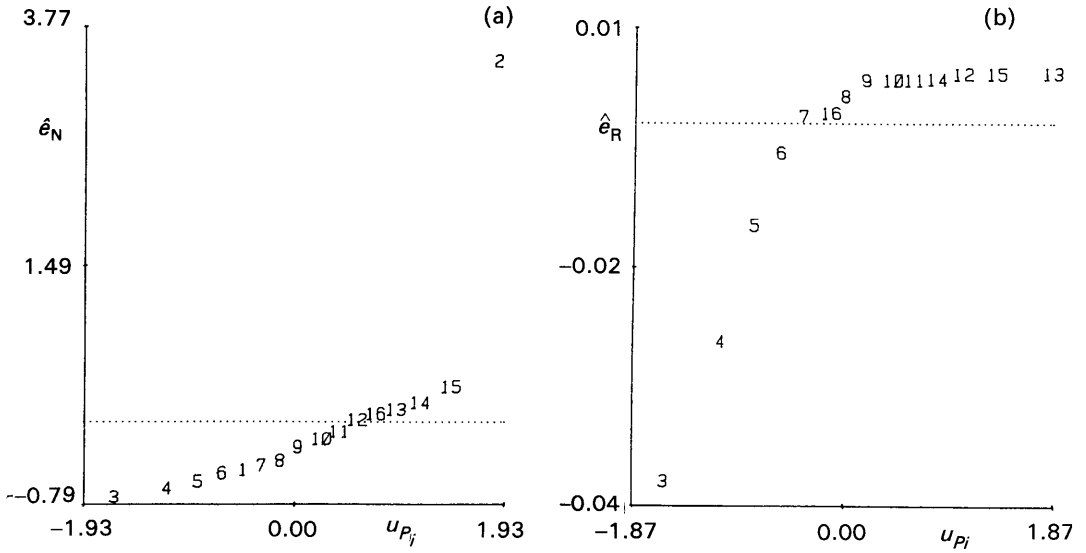


Fig. 6.73—Rankit graph for (a) normalized residuals, and (b) recursive residuals.

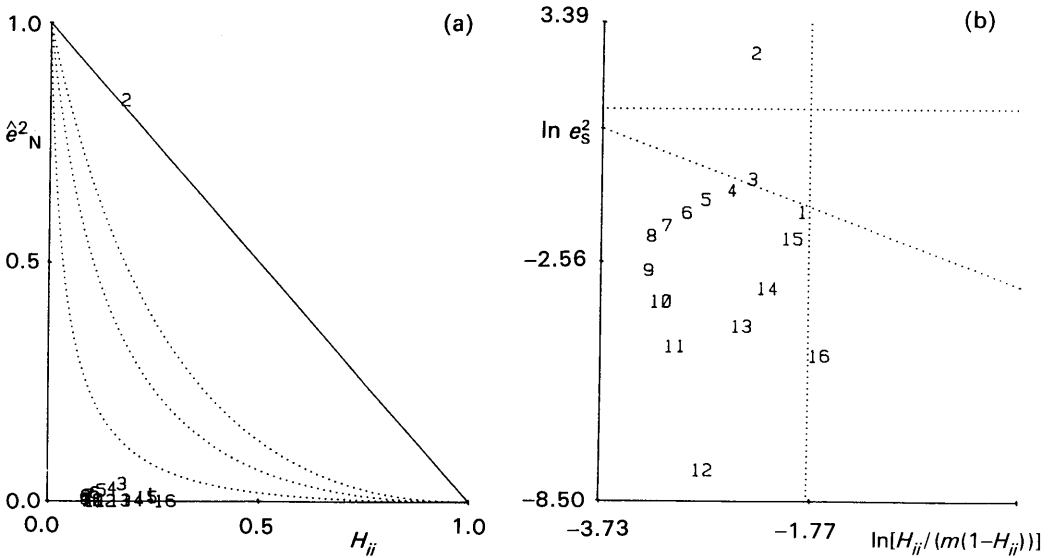


Fig. 6.74—(a) The L-R graph for D_i and (b) McCulloh-Meeter graph.

The Jack-knife residuals show that point 1 is strongly influential, having $\hat{e}_{j1} = 4.072$ and, also point 15 is suspect, with $\hat{e}_{j15} = -1.976$.

Since calibration requires the highest precision, points 1 and 15 were omitted; i.e. from the original data set, points 1, 2 and 16 were discarded. The regression model by classical least-squares has the form

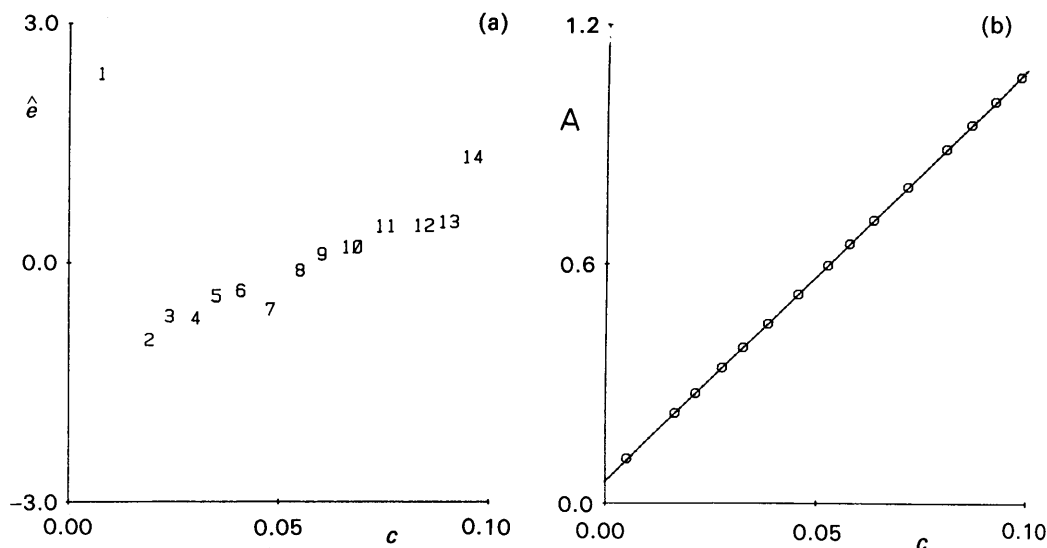
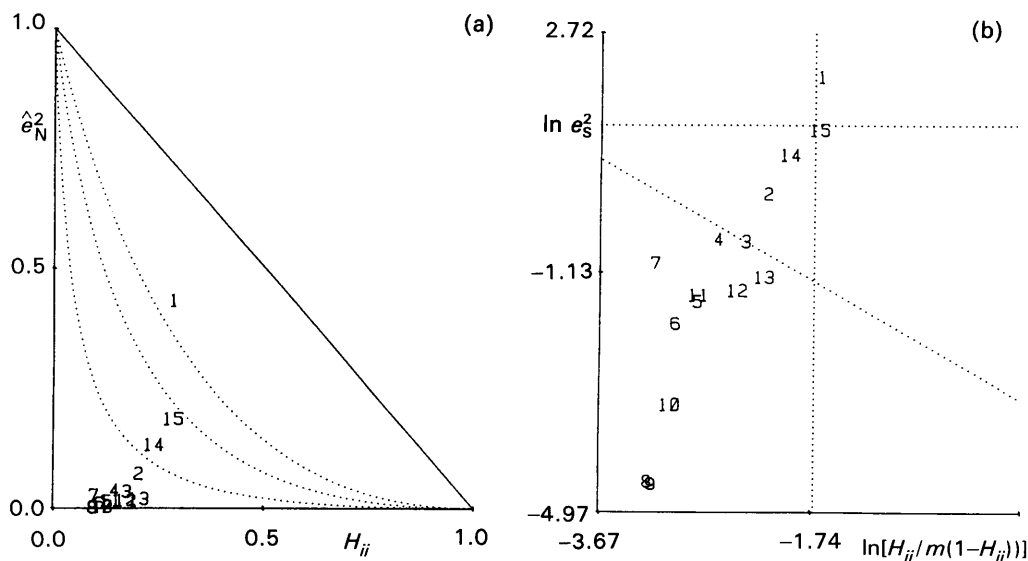


Fig. 6.75—(a) Regression model without point 2, and (b) the residual plot.

Fig. 6.76—(a) The L-R graph for D_i , and (b) McCulloh-Meeter graph for the data without point 2.

$$A = 10.364(\pm 0.0033)c + 0.053(\pm 1.91 \times 10^{-4})$$

with standard deviation $\hat{\sigma} = 0.00029$.

The estimates for all three data sets of the detection limits (A_D, c_D) and the determination limits (A_s, c_s) are listed in Table 6.20.

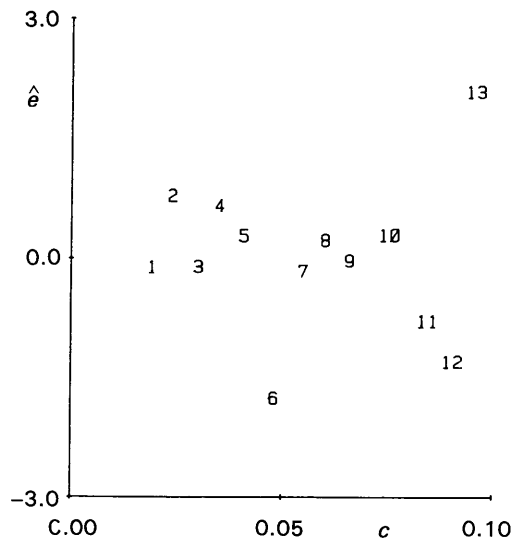


Fig. 6.77—Plot of residuals vs. variable c .

Table 6.20. The effect of influential points on detection and determination limits

Set	Size	Points omitted	Detection limit		Determination limit	
			A_D	c_D	A_s	c_s
1	16	none	0.125	5.93×10^{-3}	0.1410	7.458×10^{-3}
2	15	2	0.062	6.83×10^{-4}	0.0163	-3.74×10^{-3}
3	12	1,2,16	0.055	1.46×10^{-4}	0.0034	-4.79×10^{-3}

Problem 6.58. *Influence of instrument precision on the detection limit*

Three instruments A, B and C, with different precisions, were used to measure the signal S for 20 levels of concentration c . For the proposed model of the calibration straight line

$$E(S/c) = \beta_1 \times c + \beta_2$$

calculate the detection limit and determination limit. For signal levels $S^* = 6$, with one signal measurement, ($M = 1$) estimate the concentration and the corresponding 95% confidence limit.

Data: $n = 20$

$c, \text{ mol dm}^{-3}$	Signal S from the instrument		
	A, %	B, %	C, %
0.01	4.823	5.108	5.101
0.02	5.197	5.207	5.199
0.03	5.937	5.311	5.300

0.04	5.424	5.399	5.400
0.05	5.255	5.497	5.500
0.06	5.702	5.606	5.599
0.07	5.790	5.700	5.699
0.08	5.962	5.769	5.800
0.09	5.734	5.889	5.899
0.10	5.786	6.008	5.999
0.11	6.117	6.099	6.101
0.12	6.555	6.203	6.200
0.13	6.570	6.289	6.299
0.14	6.815	6.399	6.400
0.15	6.187	6.494	6.500
0.16	6.552	6.606	6.598
0.17	6.947	6.697	6.702
0.18	7.090	6.807	6.801
0.19	7.159	6.896	6.900
0.20	7.291	6.994	7.001

Solution: Since the main task of this problem is to examine the influence of the instrument precision on the estimates b_1 and b_2 , the regression diagnostics were not used. Table 6.21 lists the parameter estimates b_1 and b_2 , correlation coefficient \hat{R} and the residual standard deviation $\hat{\sigma}$ estimated by the classical least-squares method.

Table 6.21. Parameter estimates b_1 and b_2 of the calibration model from data measured by three instruments with different precisions

Instrument	Precision	b_1	b_2	\hat{R}	$\hat{\sigma}$
A	Fair	0.1120	4.968	0.937	0.058
B	Good	0.0997	5.002	0.9999	0.022
C	High	0.1000	5.000	1.0000	0.0002

Table 6.22 shows the calculated detection limit and the determination limit for all three instruments. The values of the limits are affected by the relatively large value of the intercept of the calibration straight line, which shows that the signal measurement always gives a large blank value.

Table 6.22. Limit values for three instruments

Instrument	Precision	S_D	c_D	S_s	c_s
A	Fair	6.143	10.48	2.792	-1.942
B	Good	5.046	0.445	0.105	-4.91
C	High	5.004	0.0448	0.0106	-49.86

Table 6.23 lists the point estimates of concentration $\hat{c} = \hat{x}$ for the signal levels $S^* = 6$, with $M = 1$.

Table 6.24 lists the 95% confidence interval of the estimated concentration \hat{c} at the signal level $S^* = 6$.

Conclusion: The spread of points around the calibration straight line is related to the

Table 6.23. Concentration estimates for signal level $S^* = 6$

Instrument	Precision	\hat{x}^*	\hat{x}_B^*	\hat{x}_T^*	\hat{x}_N^*
A	Fair	9.209	9.219	9.366	9.209
B	Good	10.009	10.009	10.010	10
C	High	10.001	10.001	10.001	10

Table 6.24. The 95% confidence interval of concentration $L_L \leq \hat{c} \leq L_U$ for the signal level $S^* = 6$ by five methods.

Instrument	Precision	(6.218)	(6.220)	(6.221)	(6.223)	(6.227)
A	Fair L_L	4.662	7.71	7.83	2.77	4.68
	L_U	13.76	10.56	10.59	15.5	13.9
B	Good L_L	9.817	9.95	9.95	9.746	9.82
	L_U	10.2	10.07	10.07	10.27	10.2
C	High L_L	9.98	9.995	9.995	9.974	9.98
	L_U	10.02	10.007	10.007	10.03	10.02

precision of the instrument. It has a significant affect on the detection and determination limits, and also on the confidence interval for the concentration. In evaluating calibration experiments, attention should be paid to the model quality and to the data quality.

Problem 6.59. *Determination of the degree of polynomial in the approximation of analytical data*

There are two sets of analytical data, the spectrum of molar absorptivities and the titration data. Both sets of data should be approximated by a polynomial, and the degree of polynomial should be examined with regard to prediction ability.

Data: (S) the spectrum of molar absorption coefficients as a function of wavelength; $n = 15$

λ , nm	460	470	480	490	500	510	520	530	540	550	560	570	580	590	600
ϵ , mol ⁻¹ cm ⁻¹	3.0	3.4	4.3	5.0	6.0	6.8	8.1	9.2	10.7	11.6	12.9	13.6	14.6	15.3	15.5

(T) the titration curve, $n = 13$

v , ml	0.12	0.56	0.83	1.36	1.48	1.73	2.20
y , mV	3.85	9.42	12.90	17.36	19.31	22.73	32.89
	2.57	2.83	3.01	3.32	3.62	3.90	
	44.51	53.01	62.09	81.00	102.11	124.00	

Solution: (1) The spectral data:

Because the x values are in the range 460–600, possibly resulting in some numerical difficulties with polynomials of higher degree in the precision of building the

covariance matrix, the linear transformation $\lambda^* = (\lambda - 460)/100$ was applied, to give data in the interval $\langle 0, 14 \rangle$.

Figure 6.78a shows the straight line regression model and Fig. 6.78b a plot of the classical residuals \hat{e} vs. λ^* . Obviously, the nonlinear pattern in the residuals shows that the approximation polynomial must be of degree greater than 2. Table 6.25 presents the statistics MEP , AIC and $\hat{\sigma}$ for increasing degree of polynomial, together with the conclusion from the sign test of residuals. The best seems to be the polynomial of 4th degree. The polynomial of the 3rd degree differs only slightly, and all its parameters are significantly different from zero. Therefore, on the base of Student t -tests the polynomial of the 3rd degree was chosen.

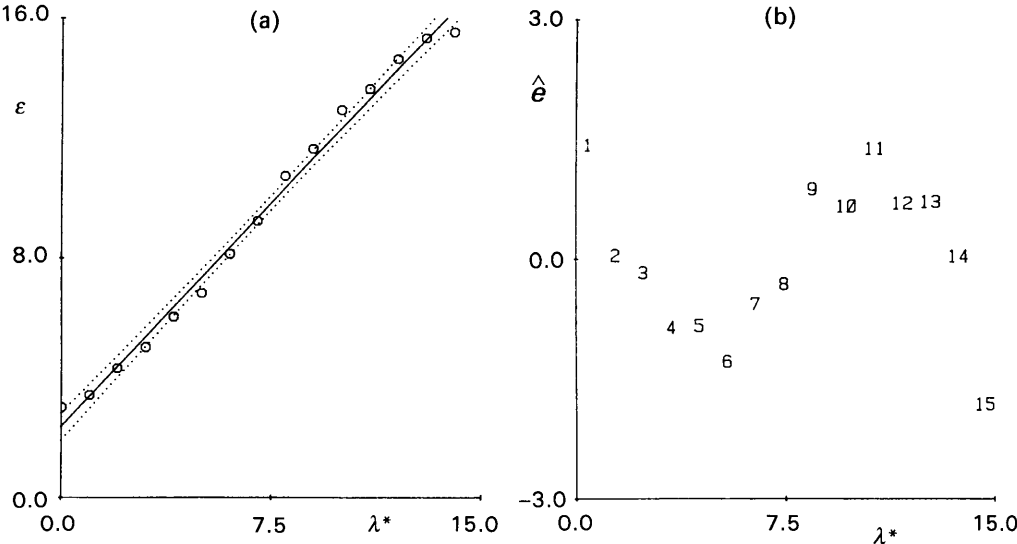


Fig. 6.78—(a) The linear regression model with the 95% confidence interval, and (b) the residual plot vs. λ^* .

Table 6.25. A search for the most convenient degree of polynomial, on the basis of four statistical characteristics MEP , ADC , $\hat{\sigma}$ and the sign test.

Polynomial degree	MEP	AIC	$\hat{\sigma}$	Trend in residuals	Conclusion
1	0.2280	14.26	0.427	yes	Rejected
2	0.3502	16.26	0.445	yes	Rejected
3	0.0284	-18.1	0.138	no	Accepted
4	0.0277	-18.89	0.132	no	Accepted
5	0.0560	-17.14	0.138	yes	Rejected
6	0.2570	-15.19	0.146	yes	Rejected

Figure 6.79a shows the degree of fit of the polynomial model of the 3rd degree with the 95% confidence interval, and Fig. 6.79b the residuals plot. The polynomial model is

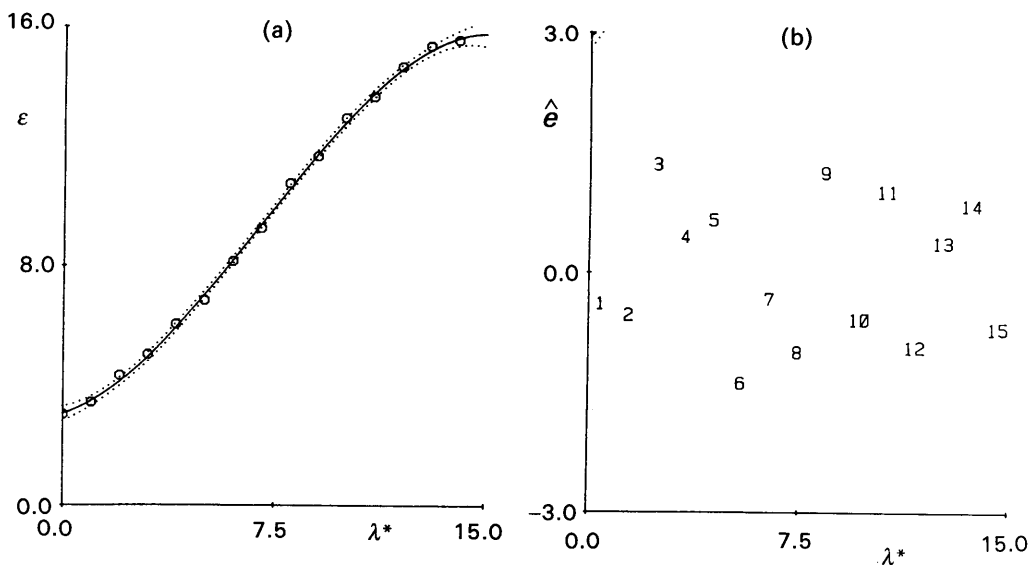


Fig. 6.79—(a) The 3rd degree polynomial, and (b) the residuals plot.

$$\varepsilon = 3.054(\pm 0.113) + 0.297(\pm 0.0727)\lambda^* + 0.1288(\pm 0.0123)\lambda^{*2} - 6.146(\pm 0.578) \times 10^{-3}\lambda^{*3}.$$

(2) The titration data:

The titration curve data covers a small range of x values and therefore does not need any initial data transformation. Table 6.26 lists the statistical characteristics. All characteristics are statistically significant for the polynomial of the 3rd degree.

Table 6.26. The search for the most convenient degree of polynomial by analysis of titration data, on the basis of the statistical characteristics MEP , ADC , $\hat{\sigma}$ and the sign test

Polynomial degree	MEP	AIC	$\hat{\sigma}$	Trend in residuals	Conclusion
1	264.1	101.3	14.04	yes	Rejected
2	40.97	71.63	4.357	yes	Rejected
3	1.07	25.85	0.731	no	Accepted
4	1.04	26.36	0.732	no	Accepted
5	2.497	27.17	0.747	no	Rejected

Figure 6.80a shows the curve for the polynomial model of the 3rd degree, with the 95% confidence interval, and Fig. 6.80b shows the residuals plot. The regression model is

$$y = 2.145(\pm 0.814) + 16.51(\pm 1.71)v - 7.715(\pm 0.977)v^2 + 2.959(\pm 0.159)v^3$$

Conclusion: In the search for the best degree of polynomial, several statistical characteristics of regression quality should be considered.

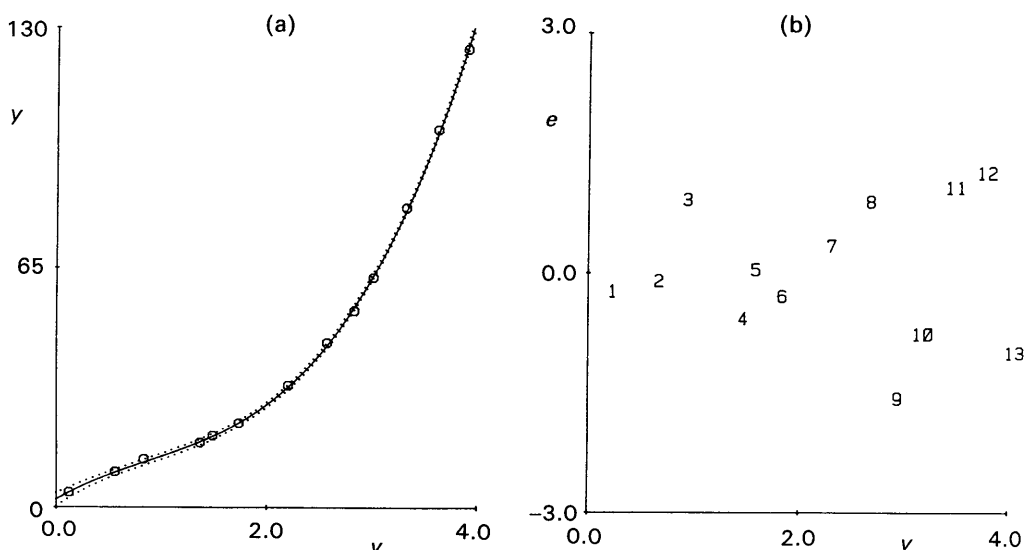


Fig. 6.80—(a) The 3rd degree polynomial, and (b) the residuals plot.

Problem 6.60. *Influence of statistical weights on the search for the degree of an approximated polynomial*

Experimental data were obtained by repeated measurements of the dependent variable y . For each measured point, the standard deviations $s(y)$ were also calculated. Try to find the best approximate polynomial with regard to its prediction ability.

Data: $n = 13$

x	0.12	0.56	0.83	1.36	1.48	1.73	2.20	2.57	2.83
y	3.85	9.42	12.90	17.36	19.31	22.73	32.89	44.51	62.09
$s(y)$	0.09	0.15	0.42	0.42	0.23	0.27	0.36	0.83	0.52
	3.01	3.32	3.62	3.90					
	81.00	102.11	124.00	124.00					
	0.61	0.93	0.86	0.71					

Solution: Table 6.27 compares the results of a search for the polynomial degree by two methods, the classical least-squares (LS) and weighted least-squares (WLS) with weights $1/s(y)$. By the WLS method, from the point of view of the mean quadratic error of prediction MEP , it was found that the best degree is the third degree, whereas by the classical LS method it is the 4th degree.

Conclusion: The precision of each point measurement may affect the search of the most convenient polynomial degree. Often differences among several polynomial degrees are quite small. In this case, the polynomial with the lowest degree is chosen,

Table 6.27. Choice of best degree of polynomial on the basis of selected statistical characteristics, by the LS and WLS methods

Polynomial degree	Method	MEP	AIC	$\hat{\sigma}$	Trend in residuals	Conclusion
2	LS	40.82	71.57	4.346	yes	Rejected
	WLS	1168	69.46	4.007	yes	Rejected
3	LS	1.04	25.29	0.715	no	Accepted
	WLS	0.626	16.2	0.504	no	Accepted
4	LS	1.341	25.71	0.714	no	Rejected
	WLS	0.826	17.55	0.521	no	Accepted
5	LS	4.469	18.71	0.539	no	Rejected
	WLS	—	—	—	no	Rejected

or the polynomial for which most characteristics are statistically significant. Another way is to use specialized techniques such as a stepwise regression.

Problem 6.61. Calibration model of the polarographic determination of clotiazepine Ebel and Brockmeyer [55] made a polarographic determination of the drug clotiazepine, which is derived from benzodiazepine. For calibration the height of the polarographic peak was measured for samples of exactly determined concentration. Find a suitable calibration model (the authors proposed a linear one).

Data: Measurements were repeated three times.

$c, \mu\text{g/ml}$	Polarographic current $I, \mu\text{A}$		
	1st run	2nd run	3rd run
0.76	5.81	5.63	6.46
1.48	9.96	9.96	10.47
2.16	14.50	14.72	15.09
2.82	19.24	18.84	18.82
3.44	23.12	23.47	23.19
4.04	27.59	28.32	28.33
4.62	31.50	32.57	31.67
5.17	35.38	36.34	35.40
5.69	40.19	38.99	39.20
6.20	43.53	42.95	43.41

Solution: In the first step, the regression model was analysed by the classical least-squares method, and the regression equation was found.

$$I = 6.923(\pm 0.0545) \times c - 0.03085(\pm 0.2197)$$

The intercept of this regression straight line is not statistically significant at significance level $\alpha = 0.05$. Some characteristics $\hat{R}^2 = 0.9983$, $\hat{\sigma} = 0.5179$, $MEP = 0.2946$ and $AIC = 61.09$ do not point to any deviations from linearity, but the sign test shows a trend in residuals. No influential points are indicated.

Figure 6.81a shows the calibration straight line and Fig. 6.81b the residual plot. The nonrandom pattern of residuals shows the need to introduce a nonlinear term, i.e. x^2 .

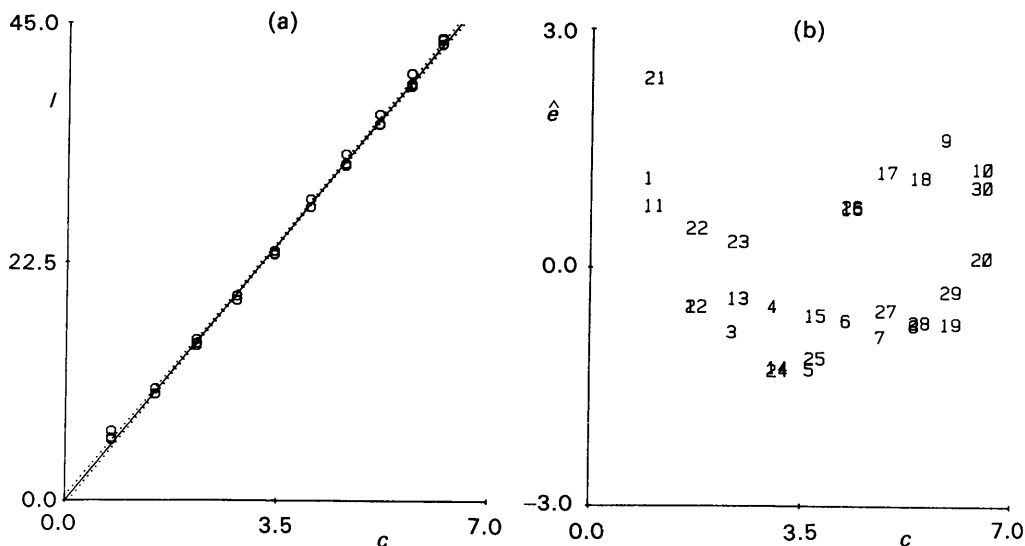


Fig. 6.81—(a) The calibration straight line of peak height vs. concentration of clotazepine, and (b) the residuals plot.

The second step was to do a regression analysis for a quadratic model by the classical least-squares method, with the result

$$I = 1.017(\pm 0.325) + 6.132(\pm 0.209)c + 0.1126(\pm 0.0291)c^2.$$

All coefficients of polynomial are significantly different from zero at $\alpha = 0.05$. Characteristics $\hat{R}^2 = 0.9989$, $\hat{\sigma} = 0.423$, $MEP = 0.1989$ and $AIC = 49.86$ show that the quadratic model fits the data better than the linear one. The residuals do not exhibit any trend.

Figure 6.82a shows the quadratic model and Fig. 6.82b the residuals plot.

To test the significance of the difference between the two models, the Fisher-Snedecor F -test is used. The test criterion

$$F = \frac{(7.511 - 4.883) \times 27}{4.883 \times 1} = 14.53$$

gives a higher value than the quantile $F_{0.95}(1, 27) = 4.21$, so the quadratic model is more suitable than the linear one.

Conclusion: The calibration equation is better expressed by a quadratic model than a linear one, because the quadratic model has significantly lower value of the residual quadratic error. The estimated instrumental error is also lower for this model.

Problem 6.62. Investigation of the influence of three factors on the percentage of conversion of *n*-heptane

The influence of the reaction temperature (x_1), the ratio of hydrogen to *n*-heptane (x_2) and the reaction contact time (x_3) on the yield (%) of acetylene (y) from

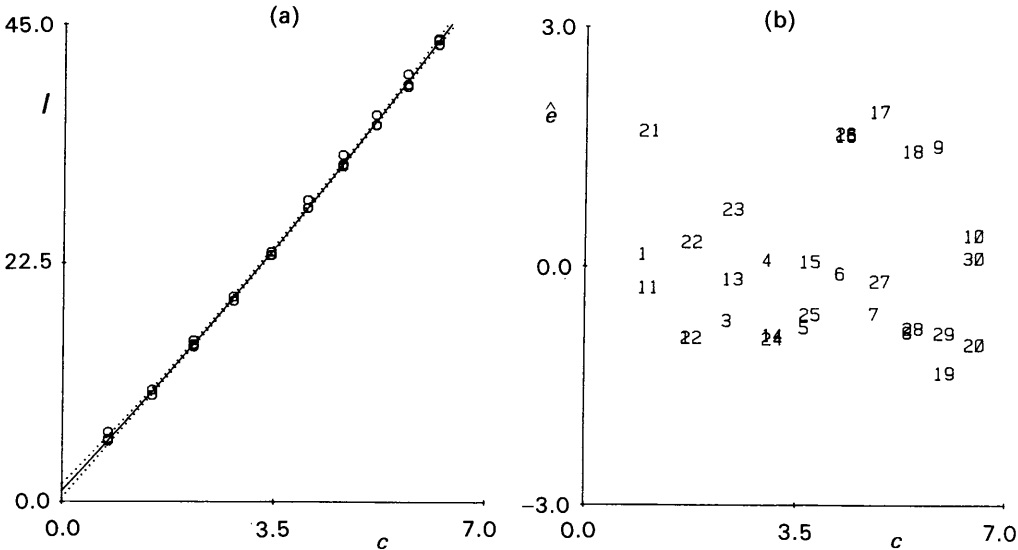


Fig. 6.82—(a) The quadratic regression model, and (b) the residuals plot.

n-heptane has been studied [56]. Determine the influence of the three factors x_1 , x_2 and x_3 on y , assuming the linear regression model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

Test also the validity of the conditions for the least-squares method.

Data: $n = 16$

i	$x_1, ^\circ\text{C}$	x_2	x_3, sec	$y, \%$
1	1300	7.5	0.012	49.0
2	1300	9.0	0.012	50.2
3	1300	11.0	0.0115	50.2
4	1300	13.5	0.013	48.5
5	1300	17.0	0.0135	47.5
6	1300	23.0	0.012	44.5
7	1200	5.3	0.040	28.0
8	1200	7.5	0.038	31.5
9	1200	11.0	0.032	34.5
10	1200	13.5	0.026	35.0
11	1200	17.0	0.034	38.0
12	1200	23.0	0.041	38.5
13	1100	5.3	0.084	15.0
14	1100	7.5	0.098	17.0
15	1100	11.0	0.092	20.5
16	1100	17.0	0.086	29.5

Solution: There is significant correlation between y and x_1 , and y and x_3 . There is also correlation between x_1 and x_3 which shows up as a multicollinearity. Also, the individual VIF factors ($VIF_1 = 12.2$, $VIF_2 = 1.06$ and $VIF_3 = 12.32$) indicate multicollinearity (see Table 6.28).

Table 6.28. Preliminary statistical analysis of data

Variable	Mean	Standard deviation	Partial regression coefficient
y	36.11	11.90	1.000
x_1	1213.0	80.62	0.945
x_2	12.44	5.662	0.370
x_3	0.04031	0.03164	-0.914

Correlation between independent variables:

x_1 vs. x_2 : 0.2236

x_1 vs. x_3 : -0.9582

x_2 vs. x_3 : -0.2402

Table 6.29 lists the parameter estimates obtained by the classical least-squares method and the statistical tests of these parameters. It is discovered that parameters β_2 and β_3 are nearly equal to zero.

Table 6.29. The estimates of the four parameters of the proposed regression model, and their statistical analysis at $\alpha = 0.05$

i	Parameter β_i	Estimate b_i	Standard deviation $s(b_i)$	t_{exp}	Conclusion $\beta_i = 0$
0	β_0	-121.3	55.44	-2.188	no
1	β_1	0.1269	0.04218	3.007	no
2	β_2	0.3482	0.1770	1.967	yes
3	β_3	19.02	107.5	-0.1762	yes

The basic statistical characteristics of the proposed linear model are listed in Table 6.30.

Table 6.30. Basic statistical characteristics of the linear regression and results of statistical tests

Determination coefficient, R^2	0.9198
Standard deviation of prediction, $s(y)$	3.767
Total F -test, F_{exp}	45.88
Criterion M	0.7718
Multicollinearity test	Worse model modification requested
Significance of model	Model is significant at the level $\alpha = 0.05$
Heteroscedasticity test	Insignificant heteroscedasticity
Complexity of model	2.775×10^6
Quadratic error of prediction	21.02
Autocorrelation coefficient	0.4445
Sign test	No trend in residuals

Examination of the linear model by statistical testing leads to the conclusion that the proposed model is as a whole statistically significant. The conversion is affected

mostly by temperature (x_1), then by the ratio of n-heptane to acetylene (x_2) and the smallest influence is that of the reaction time (x_3). Although the results are a little distorted by multicollinearity, the distortion is not very significant.

Table 6.31 shows the predicted values \hat{y} , the standard deviation of prediction $s(\hat{y})$, and the relative residual $\hat{e}_R(\hat{y})$. Table 6.32 gives an analysis of residuals.

Table 6.31. Comparison of predicted values \hat{y} with experimental values y

i	y	\hat{y}	$s(\hat{y})$	$\hat{e}_R(\hat{y}), \%$
1	49.00	46.02	1.892	6.075
2	50.20	46.55	1.702	7.280
3	50.20	47.25	1.550	6.433
4	48.50	48.09	1.537	0.839
5	47.50	49.30	1.679	-3.794
6	44.50	51.42	2.245	-15.550
7	28.00	32.04	1.684	-14.426
8	31.50	32.84	1.519	-4.265
9	34.50	34.18	1.729	0.939
10	35.00	35.16	2.255	-0.459
11	38.00	36.23	1.686	4.666
12	38.50	38.18	2.137	0.824
13	15.00	18.52	1.905	-23.447
14	17.00	19.02	2.372	-11.863
15	20.50	20.35	1.941	0.735
16	29.50	22.55	2.047	23.5551

The mean of absolute residuals : 2.475

The mean of relative residuals, % : 7.822

The residual sum of squares : 170.3

Table 6.32. The various residuals

i	\hat{e}	\hat{e}_S	\hat{e}_N	\hat{e}_J
1	2.977	0.790	0.906	0.899
2	3.654	0.970	1.087	1.096
3	3.249	0.862	0.946	0.942
4	0.407	0.108	0.118	0.113
5	-1.802	-0.478	-0.535	-0.518
6	-6.920	-1.837	-2.287	-2.916
7	-4.039	-1.072	-1.199	-1.223
8	-1.343	-0.357	-0.390	-0.375
9	0.324	0.086	0.097	0.093
10	-0.161	-0.043	-0.053	-0.051
11	1.773	0.471	0.526	0.510
12	-0.317	0.084	0.102	0.098
13	-3.517	-0.934	-1.082	-1.091
14	-2.017	-0.535	-0.698	-0.673
15	0.151	0.040	0.047	0.045
16	6.948	1.844	2.197	2.721

From these tables it is evident that points 6 and 16 are outliers, and this should be investigated further.

The partial regression graphs in Fig. 6.83 show that excluding both outliers 6 and 16 does not affect the parameter estimates, so the slope of remaining points in graphs will not change much. A more detailed analysis of individual groups of points

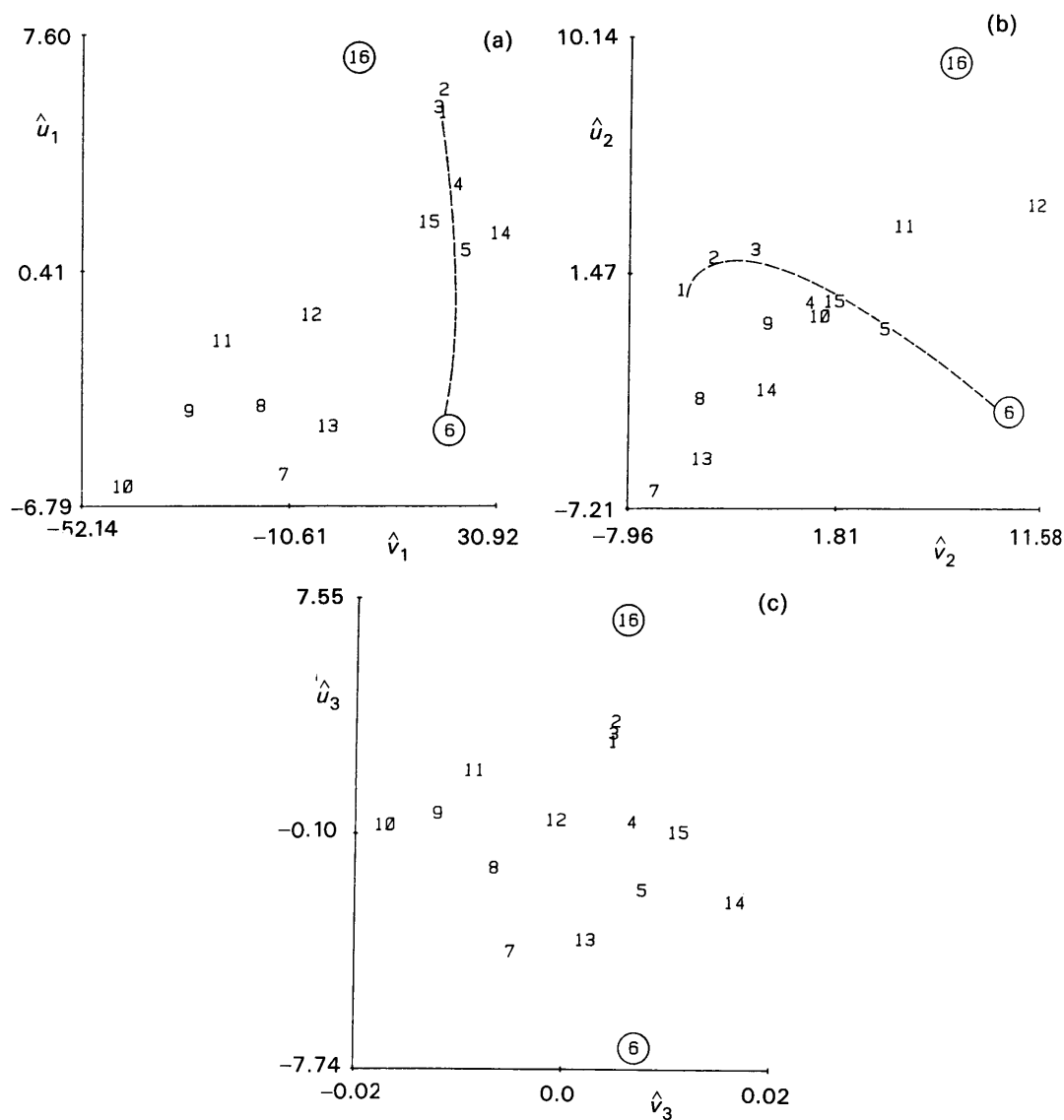


Fig. 6.83—The partial regression graphs for (a) variable x_1 , (b) variable x_2 , and (c) variable x_3 .

(according to different temperatures) shows that points 1 to 6 (temperature 1300°C) have a different trend (dotted) from the points of the other two groups.

Omitting this group (temperature 1300°C) increases the significance of factor x_2 . *Conclusion:* Despite the statistical significance of the linear model, the data should be experimentally investigated again and the number of data points increased. With a small number of points it is rather difficult to examine the regression model fully.

It may be concluded that

- (a) the temperature 1300°C has a different effect on the reaction from temperatures 1200°C and 1100°C;
- (b) besides temperature, the yield is affected also by the ratio of hydrogen to n-heptane;
- (c) for the time range studied, the contact time has no influence on the yield.

Problem 6.63. *Investigation of the dependence between phosphorus content in maize and in soil*

The content of the inorganic phosphorus (x_1) and organic phosphorus (x_2) in a soil affecting the content of phosphorus (y) in the maize was studied [45]. Examine the influence of factors x_1 and x_2 on variable y .

Data: $n = 18$

i	x_1	x_2	y
1	0.4	53.0	64.0
2	0.4	23.0	60.0
3	3.1	19.0	71.0
4	0.6	34.0	61.0
5	4.7	24.0	54.0
6	1.7	65.0	77.0
7	9.4	44.0	81.0
8	10.1	31.0	93.0
9	11.6	29.0	93.0
10	12.6	58.0	51.0
11	10.9	37.0	76.0
12	23.1	46.0	96.0
13	23.1	50.0	77.0
14	21.6	44.0	93.0
15	23.1	56.0	95.0
16	1.9	36.0	54.0
17	26.8	58.0	168.0
18	29.9	51.0	99.0

Solution: In the first step of regression analysis the linear model

$$E(y/x) = \beta_3 + \beta_1 x_1 + \beta_2 x_2$$

was proposed. The method of least squares gives

$$\hat{y} = 56.25(\pm 16.31) + 1.79(\pm 0.557)x_1 + 0.0867(\pm 0.415)x_2.$$

At the significance level $\alpha = 0.05$, parameter β_2 is statistically insignificant. The determination coefficient is $\hat{R}^2 = 0.482$ and the residual standard deviation is $\hat{\sigma} = 20.68$. The proposed linear model is, as a whole, statistically significant since

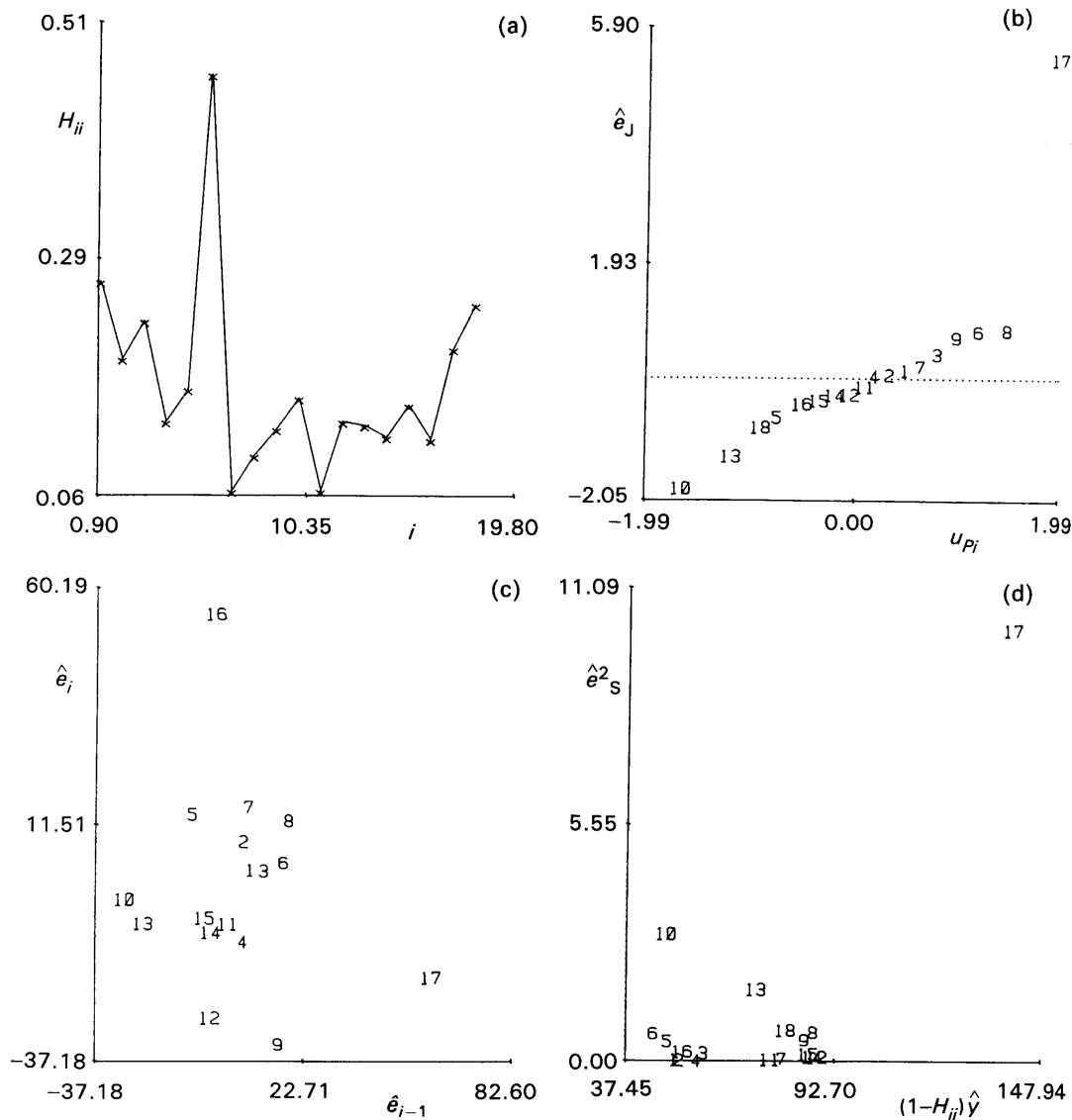


Fig. 6.84—Graphical examination for influential points (a) the index graph for H_{ii} ; (b) the rankit graph for \hat{e}_i ; (c) the autocorrelation graph; (d) the heteroscedasticity graph.

$F_{\text{exp}} = 6.988$ is larger than corresponding quantile of the F -distribution. The test criterion of the Jarque–Berra test $L_J = 11.52$ proves the strong non-normality of residuals. One strong outlier, point 17, is found (see Fig. 6.84) for which the standardized Jack-knife residual is $\hat{e}_{j,17} = 5.36$.

The graphical analysis of residuals shows that the outlier 17 causes the heteroscedasticity of data. Figure 6.85 shows the partial regression graphs for variables x_1 and

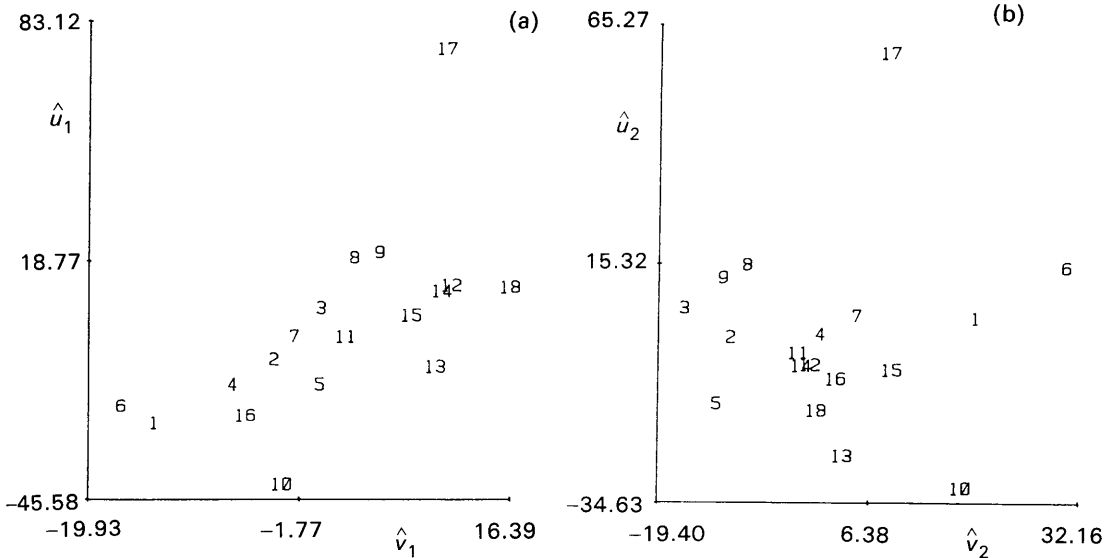


Fig. 6.85—The partial regression graph for (a) variable x_1 , and (b) variable x_2 .

x_2 . Excluding point 17 decreases the significance of factor x_2 , as a random pattern of points is now formed. The significance of the factor x_1 remains unchanged.

In the second step of regression analysis, point 17 and factor x_2 were omitted. The calculated regression model is

$$y = 62.57(\pm 4.452) + 1.229(\pm 0.306)x_1.$$

Both parameters are statistically significant, the determination coefficient is increased now to $\hat{R}^2 = 0.519$ and the residual standard deviation is decreased to $\hat{\sigma} = 11.92$. The normality test $L_J = 10.10$ still demonstrates the non-normality of residuals. There is no heteroscedasticity or trend in residuals. The regression diagnostics discover one influential point, number 10, with $\hat{e}_{J,10} = -2.84$.

Figure 6.86 shows the linear regression model with its 95% confidence interval and the classical residual plot.

Conclusion: The content of phosphorus in maize stem is affected only by the content of inorganic phosphorus in soil. Excluding a strongly influential point may cause a decrease in the significance of some variables in a model.

REFERENCES

- [1] N. R. Draper and H. Smith, *Applied Regression Analysis*, 2nd Ed., Wiley, New York, 1981.
- [2] G. A. F. Seber, *Linear Regression Analysis*, Wiley, New York, 1977.
- [3] I. Guttman, *Linear Models — An Introduction*, Wiley, New York, 1982.
- [4] S. R. Searle, *Linear Models*, Wiley, New York, 1971.
- [5] F. J. Anscombe, *Am. Statist.*, 1973, **27**, 17.
- [6] J. Utts, *Commun. Statist.*, 1982, **11**, 2801.
- [7] W. Kramer and H. Sonnberger, *The Linear Regression Model under Tests*, Physica Verlag, Heidelberg, 1986.

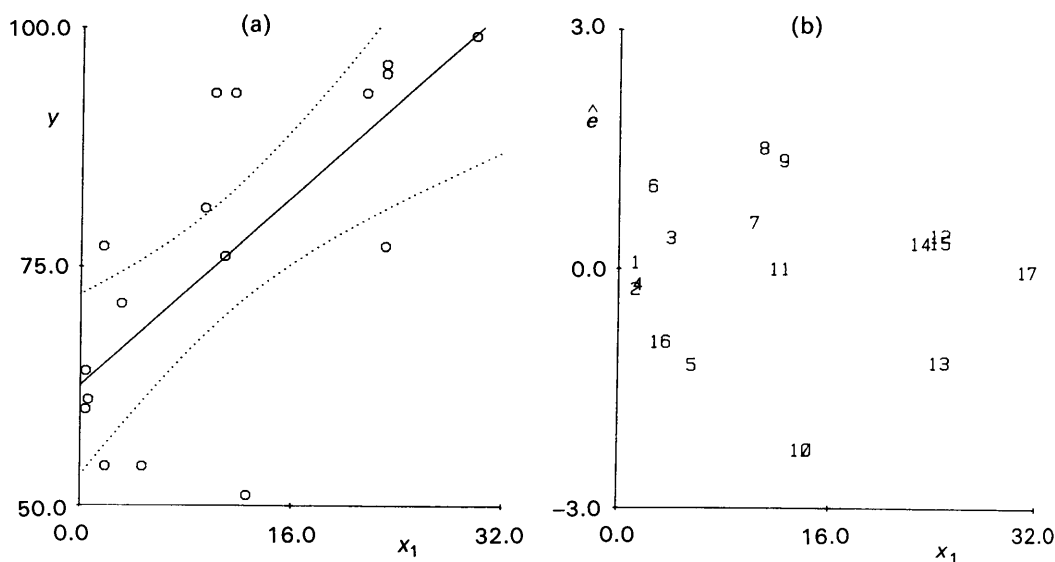


Fig. 6.86—(a) The regression straight line, and (b) the residual plot of \hat{e} on x_1 .

- [8] J. R. Scott, *Appl. Statist.*, 1975, **24**, 42.
- [9] J. Cassela, *Am. Statist.*, 1983, **37**, 147.
- [10] R. Suich, G. C. Derringer, *Technometrics*, 1977, **19**, 213.
- [11] A. N. Kornilov and L. B. Smenina, *Zh. Fiz. Khim.*, 1970, **44**, 1932.
- [12] J. Militky, *Proc. Conf. European Simulation Conference 87*, Prague, September, 1987.
- [13] G. R. Phillip, J. M. Harris and E. M. Eyring, *Anal. Chem.*, 1982, **54**, 2053.
- [14] J. W. Neil and D. E. Johnson, *Commun. Stat.*, 1984, **13**, 485.
- [15] J. R. Green and D. Margerison, *Statistical Treatment of Experimental Data*, Elsevier, Amsterdam, 1978.
- [16] J. C. Nash, *Compact Numerical Algorithms for Computer*, A. Hilger, Bristol, 1979.
- [17] A. M. Antila and M. L. Sikvonen, *Z. Anal. Chem.*, 1987, **327**, 799.
- [18] Ch. Lawson and R. Hanson, *Solving Least-Squares Problems*, Englewood Cliffs, New Jersey, 1974.
- [19] A. G. Dahlquist and A. Bjorck, *Numerical Methods*, Prentice-Hall, Englewood Cliffs, 1974.
- [20] D. M. Marquardt, *Technometrics*, 1970, **12**, 591.
- [21] D. A. Belsey, E. Kuh and R. E. Welsch, *Regression Diagnostics*, Wiley, New York, 1980.
- [22] A. C. Atkinson, *Plots, Transformations and Regression*, Clarendon Press, Oxford, 1985.
- [23] S. Weisberg, *Technometrics*, 1983, **25**, 219.
- [24] R. D. Cook and S. Weisberg, *Residuals and Influence in Regression*, Chapman and Hall, New York, 1982.
- [25] B. Joiner, *Am. Statist.*, 1981, **35**, 227.
- [26] S. Chatterjee and A. S. Hadi, *Statist. Sci.*, 1986, **1**, 379.
- [27] M. A. O'Gorman and R. M. Myers, *Commun. Statist.*, 1987, **16**, 771.
- [28] J. B. Gray, *Proc. Statist. Comput. Sect.*, p. 159, ASA, Washington 1983.
- [29] C. L. Mallows, *Technometrics*, 1986, **28**, 313.
- [30] R. D. Cook and S. Weisberg, *Biometrika*, 1983, **70**, 1.
- [31] N. Querry, *Technometrics*, 1964, **6**, 225.
- [32] C. M. Jarque and A. K. Bera, *Int. Stat. Rev.*, 1987, **55**, 163.
- [33] G. G. Judge and M. E. Bock, *Statistical Implications of Pre-test and Stein Rule Estimators in Econometrics*, North Holland, Amsterdam, 1978.
- [34] J. J. Leary and E. B. Messick, *Anal. Chem.*, 1985, **57**, 956.
- [35] S. D. Horn, R. A. Horn and D. B. Duncan, *J. Am. Statist. Assoc.*, 1975, **70**, 380.
- [36] W. A. Fuller, *Measurement Error Models*, Wiley, New York, 1987.
- [37] R. W. Hill and D. W. Holland, *J. Am. Statist. Assoc.*, 1977, **72**, 828.
- [38] D. J. Huber, *Robust Statistics*, Wiley, New York, 1981.

- [39] Li G., in Hoaglin D. C. *et al.*, eds., *Exploring Data Tables, Trends and Shapes*, Wiley, New York, 1985.
- [40] G. R. Phillip and E. M. Eyring, *Anal. Chem.*, 1983, **55**, 1134.
- [41] Yu. G. Nikulichev, Yu. A. Kanchenko, A. E. Mysak and I. F. Kobilinskaya, *Kolloidn. Zh.*, 1988, **50**, 473.
- [42] J. Mala and I. Slama, *Chem. Pap.*, 1988, **42**, 319.
- [43] W. S. Krasker and R. E. Welsch, *J. Am. Statist. Assoc.*, 1982, **77**, 595.
- [44] T. P. Hettmansperger, *Aust. J. Statist.*, 1987, **29**, 1.
- [45] P. J. Rousseuw and A. M. Leroy, *Robust Regression and Outliers Detection*, Wiley, New York, 1987.
- [46] J. R. Rosenblatt and C. H. Spiegelman, *Technometrics*, 1981, **23**, 329.
- [47] S. Ebel and U. Becht, *Z. Anal. Chem.*, 1987, **327**, 157.
- [48] L. M. Schwartz, *Anal. Chem.*, 1976, **48**, 2287.
- [49] L. J. Naszodi, *Technometrics*, 1978, **20**, 201.
- [50] R. G. Krutchkoff, *Technometrics*, 1967, **9**, 425.
- [51] L. M. Schwartz, *Anal. Chem.*, 1977, **49**, 2062.
- [52] L. Oppenheimer, T. P. Capizzi, R. M. Weppelman and H. Mehta, *Anal. Chem.*, 1983, **55**, 638.
- [53] L. M. Schwartz, *Anal. Chem.*, 1983, **55**, 1424.
- [54] S. Ebel and U. Kamm, *Z. Anal. Chem.*, 1984, **318**, 293.
- [55] S. Ebel and R. Brockmeyer, *Z. Anal. Chem.*, 1970, **326**, 770.
- [56] D. Himmelblau, *Process Analysis by Statistical Methods*, Wiley, New York, 1969.

7

Correlation

Chapter 5 considers the characteristics and procedures of multivariate data analysis, and Chapter 6 describes the construction of linear regression models. In this chapter we describe relationships expressing dependencies among the components ξ_1, \dots, ξ_m of an m -dimensional vector ξ by using regression. The difference from construction of linear regression models in Chapter 6 is that here the data form a random sample from the m -dimensional distribution of random vector ξ . There is no consideration about which component ξ_j of the random vector ξ is the response (in the linear model, the dependent variable) and which components of vector ξ are controllable (in the linear model, the independent variable).

The random sample $\{x_i\}$, $i = 1, \dots, n$ of size n represents an $(n \times m)$ array of data

$$\begin{bmatrix} x_{11} & \dots & x_{12} & \dots & x_{1m} \\ x_{21} & \dots & x_{22} & \dots & x_{2m} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{n2} & \dots & x_{nm} \end{bmatrix}$$

where the number of rows n (i.e., the number of m -dimensional "points" x_i) is larger than the number of columns m (i.e., the number of "variables" or components of vector x). The characteristic fact is that all components of the data vector are measured and not controllable by the experimenter.

In the regression models of Chapter 6, some independent variables such as temperature, concentration, etc. are also measured (and therefore random) variables, but the experimenter could adjust and control their magnitude.

Although in chemometric practice, correlation problems do not often require detailed analysis, we find that problems such as (a) comparison of various analytical methods on different samples or (b) searching for relationships among various properties or characteristics of compounds, are more problems of correlation than regression.

7.1 CORRELATION MODELS

As for univariate random variables, the components of random vectors can be characterized by use of means and variances. A *measure of intensity* of the dependence between components ξ_i and ξ_j , $i \neq j$ is given by the second central mixed moment $\text{cov}(\xi_i, \xi_j)$, denoted as the covariance. The *standardized covariance* or *correlation coefficient* $\rho(\xi_i, \xi_j)$, is more useful.

Covariance and correlation coefficients and methods for their estimation are described in Chapter 5. Here we use the covariance matrix **C** with elements formed by individual covariances or the correlation matrix **R** with elements formed by individual correlation coefficients. The covariance matrix **C** has the variances on the diagonal while the correlation matrix has ones.

A random vector is characterized by the vector of mean values $\mathbf{u} = (u_1, \dots, u_m)^T$ where $u_j = E(\xi_j)$, and by the covariance matrix **C**. This information is generally not sufficient. Analogously to the mean values $E(\xi_j)$, conditional means or conditional variances can also be defined. We will define these characteristics for a case of two random quantities ξ_1 and ξ_2 , and then for the general random vector ξ .

7.1.1 Correlation models for two random variables

For two random variables ξ_1 and ξ_2 the conditional means are given by

$$E(\xi_1/x_2) = \int_{-\infty}^{\infty} x_1 f(x_1/x_2) dx_1 \quad (7.1a)$$

$$E(\xi_2/x_1) = \int_{-\infty}^{\infty} x_2 f(x_2/x_1) dx_2 \quad (7.1b)$$

where $f(x_2/x_1)$ and $f(x_1/x_2)$ are the conditional probability densities (cf. Chapter 5).

From Eq. (7.1a,b) it is evident that the conditional mean value $E(\xi_1/x_2)$ is in fact a mean value of the random variable ξ_1 , with condition that the random variable ξ_2 lies in the infinitely small interval around the value x_2 . The conditional mean value $E(\xi_2/x_1)$ is defined similarly. Because they are conditioned by a random variable, the conditional mean values are *random variables* which may be characterized by the means and variances. The means of the conditional mean values do not provide any new information because

$$E(E(\xi_1/x_2)) = E(\xi_1)$$

and

$$E(E(\xi_2/x_1)) = E(\xi_2)$$

By introducing the variances of the conditional mean values ($D(E(\xi_1/x_2))$ and $D(E(\xi_2/x_1))$) the total variances $D(\xi_1)$ and $D(\xi_2)$ may be decomposed into the components

$$D(\xi_1) = E(D(\xi_1/x_2)) + D(E(\xi_1/x_2)) \quad (7.2a)$$

$$D(\xi_2) = E(D(\xi_2/x_1)) + D(E(\xi_2/x_1)) \quad (7.2b)$$

The first terms on the right hand sides represent the mean values of the conditional variances which may be defined in a similar way to the conditional mean values (Eq. 7.1), with the use of conditional probability densities.

The conditional mean values have the same properties as the unconditional. For the conditional mean value $E(\xi_2/x_1)$ it is also true that

- (1) For any x_1 of random variable ξ_1 , the values $E(\xi_2/x_1)$ exist if $E(\xi_2) < \infty$.
- (2) If a random variable ξ_1 does not depend on the random variable ξ_2 , the conditional mean value is independent of the condition and $E(\xi_2/x_1) = E(\xi_2)$.
- (3) If $\xi_2 = g(\xi_1)$ where $g(\cdot)$ is a function notation, then $E(\xi_2/x_1) = g(x_1)$.
- (4) The conditional mean value is a not symmetric function of the arguments, so that $E(\xi_2/x_1) \neq E(\xi_1/x_2)$.

Property (3) shows that the conditional mean value is a function of quantity x_1 of condition ξ_1 and therefore it is denoted as the regression of variable ξ_2 on the variable ξ_1 .

Generally, two types of regression are distinguished [1].

- (1) The *theoretical regression* is a conditional mean value derived from knowledge of a conditional probability density $f(x_2/x_1)$ or the knowledge of a joint probability density $f(x_1, x_2)$ and both marginal densities $f(x_1)$ and $f(x_2)$. It is valid that for all elliptic conditional distributions, including the normal one, the theoretical regression is a linear one [2]. For some conditional distributions, however, the theoretical regression may be nonlinear [3].
- (2) The *empirical regression* is any conveniently selected function $g(x_1)$ which approximates the behaviour of the conditional mean value $E(\xi_2/x_1)$. To find the function $g(\cdot)$ and the parameters estimates, the methods in Chapter 6 may be used.

We will now deal with the theoretical regression when all the components of vector ξ have the normal distribution, and also the joint distribution of vector ξ is normal.

Problem 7.1. *Deriving a theoretical regression for the normalized normally distributed random variables*

Let us assume that the random variables ξ_1 and ξ_2 have the normalized normal distribution $N(0, 1)$ with zero mean and variance equal to one. The joint distribution of the variables is also normal with the probability density function

$$f(x_1, x_2) = \frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left[-\frac{x_1^2 - 2\rho \times x_1x_2 + x_2^2}{2(1-\rho^2)}\right]$$

where $\rho = \rho(\xi_1, \xi_2)$ is the correlation coefficient between the random variables ξ_1 and ξ_2 . Derive the theoretical regression $E(\xi_2/x_1)$.

Solution: In the first step, the conditional probability density $f(x_2/x_1)$ should be calculated.

$$f(x_2/x_1) = \frac{f(x_1, x_2)}{f(x_1)}$$

On substitution and rearrangement, we get

$$f(x_2/x_1) = \frac{1}{\sqrt{2\pi(1-\rho^2)}} \times \exp\left[-\frac{(x_2 - \rho x_1)^2}{2(1-\rho^2)}\right] \quad (7.3)$$

By substituting into Eq. (7.1) and analytical differentiation we get

$$E(\xi_2/x_1) = \int_{-\infty}^{\infty} \frac{x_2}{\sqrt{2\pi(1-\rho^2)}} \exp\left[-\frac{(x_2 - \rho x_1)^2}{2(1-\rho^2)}\right] dx_2 = \rho x_1 \quad (7.4)$$

Conclusion: For this case of random variables, the theoretical regression is linear with zero intercept, and the slope corresponds to the correlation coefficient ρ .

The conditional variances are the characteristics of variability of conditional distributions. The conditional variance $D(\xi_2/x_1)$ expresses the variability of the random variable ξ_2 around the theoretical regression $E(\xi_2/x_1)$, on condition that ξ_1 has a realization x_1 , where function x_1 is called the scedastic function. If $D(\xi_2/x_1)$ is a constant independent of ξ_1 (or x_1) it is a homoscedastic function. The homoscedastic and heteroscedastic functions are illustrated in Fig. 7.1.

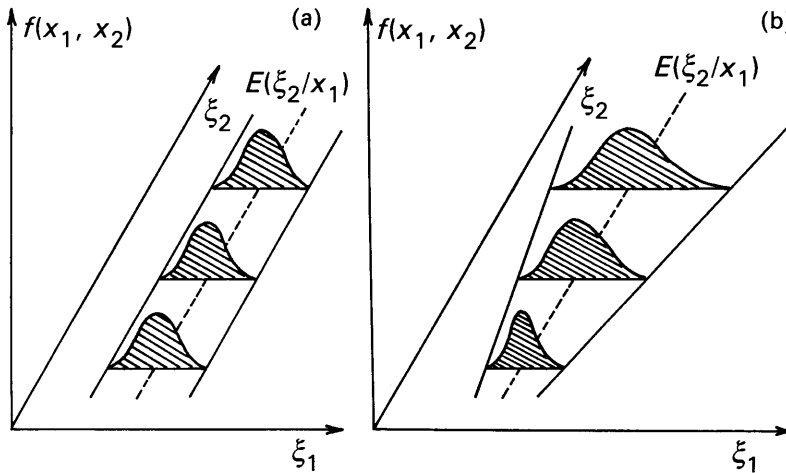


Fig. 7.1—(a) Homoscedastic, and (b) heteroscedastic relationship between two linearly dependent random variables.

For independent random variables the following expressions are valid

$$E(\xi_2/x_1) = E(\xi_2)$$

and

$$D(\xi_2/x_1) = D(\xi_2).$$

The theoretical regression $E(\xi_2/x_1)$ and $E(\xi_1/x_2)$ then represents two mutually perpendicular straight lines, parallel with the co-ordinates when the scedastic functions are constant.

For dependent random quantities, either the conditional mean value (Fig. 7.1) or conditional variance, or both, is/are non-constant.

Problem 7.2. Determination of a scedastic function

Determine the scedastic function $D(\xi_2/x_1)$ for the variables ξ_1 and ξ_2 defined in Problem 7.1.

Solution: From the definition of variance we can write

$$D(\xi_2/x_1) = \int_{-\infty}^{\infty} (x_2 - E(\xi_2/x_1))^2 f(x_2/x_1) dx_2$$

After substitution from Eq. (7.3) and analytical differentiation, we find

$$D(\xi_2/x_1) = (1 - \rho^2) \quad (7.5)$$

Conclusion: For normalized normal quantities with the normal distribution, when their joint distribution is also normal the scedastic function is constant.

Conditional variances are also random variables (dependent on condition) which may be characterized by the mean values and variances. In the linear regression of ξ_2 on ξ_1 the mean conditional variance is

$$E(D(\xi_2/x_1)) = D(\xi_2)[1 - \rho^2] \quad (7.6)$$

and the variance of conditional mean value is

$$D(E(\xi_2/x_1)) = D(\xi_2)\rho^2 \quad (7.7)$$

For homoscedastic functions, the mean conditional variance is equal to the conditional variance which does not depend on the conditional values. Conditional variances have all the properties of unconditional variances.

The mean conditional variances generally characterize a stochastic dependence between random variables which can be nonlinear.

If $E(D(\xi_2/x_1)) = D(\xi_2)$, ξ_1 and ξ_2 are independent.

If $E(D(\xi_2/x_1)) < D(\xi_2)$, there is a stochastic relationship between the variables.

From Eq. (7.2b) and the definition of regression it follows that the variance of the conditional mean is that part of the total variance concerned with "the theoretical regression" caused by the influence of variable ξ_1 on the variability of variable ξ_2 . The mean value of the conditional variance expresses the influence of all other (not considered) variables which cause variability in output variable ξ_2 .

For a measure of regression quality, we use the ratio R_R^2

$$R_R^2 = \frac{D(E(\xi_2/x_1))}{D(\xi_2)}$$

which determines the part of the variability of random variable ξ_2 explained by theoretical regression. From Eq. (7.7) it follows that R_R^2 for linear regression is equal to the square of correlation coefficient or to the determination coefficient.

Let us mention the theoretical regression and conditional variances for two random variables, ξ_1 with distribution $N(\mu_1, \sigma_1^2)$ and ξ_2 with distribution $N(\mu_2, \sigma_2^2)$, when

their joint distribution is also normal. On the basis of Problem 7.1, the theoretical regression $E(\xi_2/x_1)$ may be expressed in the form

$$E(\xi_2/x_1) = \mu_2 + \rho \frac{\sigma_2}{\sigma_1}(x_1 - \mu_1) \quad (7.9)$$

This is a straight line with slope $b_1 = \rho\sigma_2/\sigma_1$ and intercept $b_2 = \mu_2 - b_1\mu_1$, which passes through the centre of gravity of co-ordinates $[\mu_1, \mu_2]$.

Similarly, on the basis of the results of Problem 7.2, the conditional variance may be determined as

$$D(\xi_2/x_1) = \sigma_2^2(1 - \rho^2) \quad (7.10)$$

From the definition it follows that

$$D(\xi_2/x_1) = E[\xi_2 - E(\xi_2/x_1)]^2 \quad (7.11)$$

and the conditional variance is equal to the mean value of the square of deviations of random quantity ξ_2 .

In a similar way, the linear expression for the regression $E(\xi_1/x_2)$ may be found

$$E(\xi_1/x_2) = \mu_1 + \rho \frac{\sigma_1}{\sigma_2}(x_2 - \mu_2) \quad (7.12)$$

The conditional variance

$$D(\xi_1/x_2) = \sigma_1^2(1 - \rho^2) \quad (7.13)$$

It is obvious that both theoretical regressions go through the same point. The product of their slopes is equal to the square of the correlation coefficient. If the correlation coefficient $\rho = 1$, both slopes of theoretical regressions will be equal to one and both regressions will be identical. If $\rho = 0$, both slopes will be equal to zero and the regressions will be parallel with the axis of the co-ordinate system and will have an angle of 90° . The angle φ between the theoretical regressions gives a measure of the linear dependence between the random quantities ξ_1 and ξ_2 . For this angle

$$\tan \varphi = \frac{\sigma_1 \sigma_2 (1 - \rho^2)}{\rho(\sigma_1^2 + \sigma_2^2)} \quad (7.14)$$

Figure 7.2 shows the relationship between the two theoretical regressions.

In cases when the correlation coefficient is not equal to zero or one, there exist two different regressions. Often it is possible to determine which variable is a response and which is the controllable one, and according to that, to select a suitable type of regression. When it is not possible to determine the type of variables, to determine the linear relationship between ξ_1 and ξ_2 we can use:

- (1) the *principal axis* which corresponds to the minimum of the squares of the perpendicular distances of points from the regression straight line.
- (2) the *reduced principal axis* of the ellipses of constant densities, corresponding to the minimum of the products of deviations in the two variables. The slope of

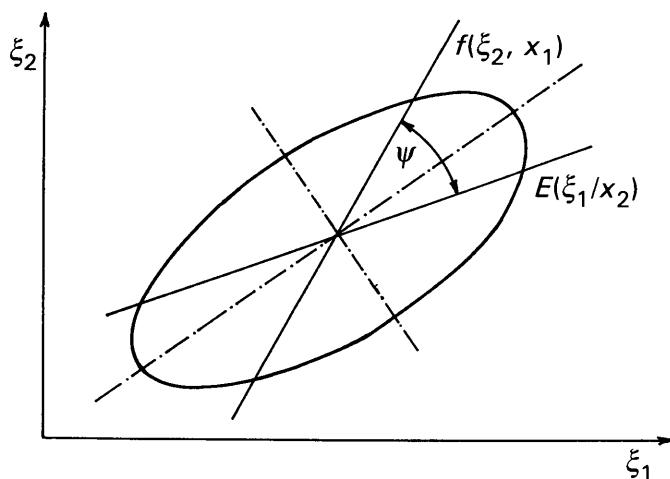


Fig. 7.2—The relationship between theoretical regressions.

the corresponding reduced principal axis is directly equal to the variance ratio $d = \sigma_2/\sigma_1$ and the regression goes through the centre of gravity (μ_1, μ_2) .

These expressions may be used for practical purposes. The means μ_1 and μ_2 may be replaced by the arithmetic averages \bar{x}_1 and \bar{x}_2 , the variances σ_1^2 and σ_2^2 by the sample variance estimates s_1^2 and s_2^2 , and the correlation coefficient ρ by the sample estimate of the correlation coefficient:

$$R = \sqrt{\frac{\sum_{i=1}^n (x_{1i} - \bar{x}_1)(x_{2i} - \bar{x}_2)}{\sum_{i=1}^n (x_{1i} - \bar{x}_1)^2 \times \sum_{i=1}^n (x_{2i} - \bar{x}_2)^2}} \quad (7.15)$$

The slope b_1 and the intercept b_2 of the regression $E(\xi_2/x_1)$ correspond to estimates found by the least-squares method, and $D(\xi_2/x_1)$ corresponds to the residual sum of squares.

Problem 7.3a. *Influence of solvent type on the degree of polymerization of cotton*

For 17 differently degraded samples of cotton the relative viscosity was determined in (a) a solution of the ethylenediamine complex of copper (CUEN), and (b) an alkaline solution of copper tetra-ammine hydroxide (CUOXAN). From viscosity values, the degree of polymerization values \overline{DP} were calculated. Examine the relationship between the degree of polymerization \overline{DP}_1 in solution CUEN and \overline{DP}_2 in solution CUOXAN.

Data: $n = 17$

$\overline{DP}_1 \times 1000$	5.913	1.837	5.732	3.792	5.823	2.837
$\overline{DP}_2 \times 1000$	2.341	1.740	2.863	1.648	2.608	1.391

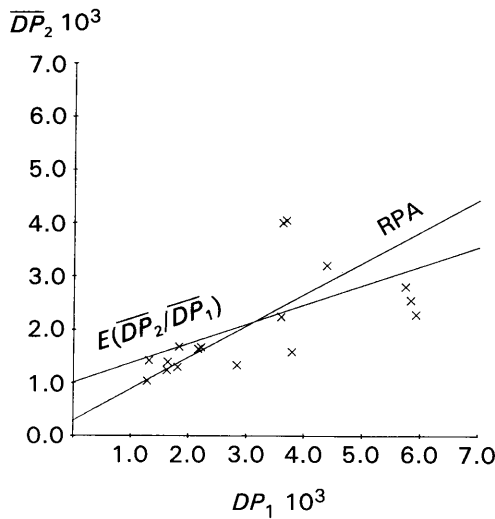


Fig. 7.3—The straight line corresponding to the reduced principal axis (RPA) and to the regression $E(\xi_2/x_1)$.

$$\mu = \begin{bmatrix} \mu_1 \\ \mu^* \end{bmatrix} \quad C = \begin{bmatrix} \sigma_1^2 & c_1^T \\ c_1 & C^* \end{bmatrix}$$

where $c_1^T = [\text{cov}(\xi_1, \xi_2), \dots, \text{cov}(\xi_1, \xi_m)]$ is the vector containing the covariance between the response variable ξ_1 and the explanatory variables ξ^* . The symbol C^* represents the covariance matrix of the explanatory variables.

As for the case of two variables, the conditioned probability density $f(x_1/x^*)$ may be determined, and this has a normal distribution, so the conditional mean value is

$$E(\xi_1/x^*) = \mu_1 + c_1^T C^{*-1} (x^* - \mu^*) \quad (7.16)$$

Let us introduce the vector $\mathbf{a} = (a_1, \dots, a_{m-1})^T$, and the expression

$$\mathbf{a} = C^{*-1} c_1 \quad (7.17)$$

Then Eq. (7.16) represents the linear function of variables x^* in the form

$$E(\xi_1/x^*) = \mu_1 + \mathbf{a}^T (x^* - \mu^*) = \mu_1 + \sum_{i=1}^{m-1} a_i (x_{i+1} - \mu_{i+1}) \quad (7.18)$$

If the joint distribution of the random vector is normal and the distribution of all ξ_j components is also normal, the resulting theoretical regression is linear.

The vector of regression coefficients \mathbf{a} is estimated here as a solution of the set of $(m-1)$ linear equations

$$\mathbf{a} C^* = c_1 \quad (7.19)$$

where C^* and c_1 contain individual covariances. The corresponding conditional

variance is given by the relation

$$D(\xi_1/\mathbf{x}^*) = \sigma_1^2 - \mathbf{c}_1^T \mathbf{C}^{*-1} \mathbf{c}_1 = \sigma_1^2 - \sum_{i=1}^m c_i^1 C_{ij}^{22} c_j^1 \quad (7.20)$$

where c_i^1 , c_j^1 are elements of vector \mathbf{c}_1 , and C_{ij}^{22} are elements of matrix \mathbf{C}^{*-1} . If all components of the random vector $\boldsymbol{\xi}^*$ are mutually independent, matrix \mathbf{C} is diagonal with variances σ_j^2 on the diagonal. For individual regression coefficients, then

$$a_{j-1} = \frac{\text{cov}(\xi_1, \xi_j)}{\sigma_j^2} = \rho(\xi_1, \xi_j) \frac{\sigma_1}{\sigma_j}, \quad j = 2, \dots, m \quad (7.21)$$

Similarly, the expression for the conditional variance will simplify to

$$D(\xi_1/\mathbf{x}^*) = \sigma_1^2 - \sum_{j=2}^m \rho^2(\xi_1, \xi_j) \sigma_1^2 = \sigma_1^2 [1 - R_{1(2, \dots, m)}^2] \quad (7.22)$$

where $R_{1(2, \dots, m)}$ is the multiple correlation coefficient between ξ_1 and the vector $\boldsymbol{\xi}^*$. For this correlation coefficient

$$R_{1(2, \dots, m)} = \sqrt{\frac{D(E(\xi_1/\mathbf{x}^*))}{\sigma_1^2}} = \sqrt{1 - \frac{\det(\mathbf{R})}{\det(\mathbf{R}_{11})}} \quad (7.23)$$

where \mathbf{R}_{ij} is the matrix formed by leaving out the i th row and the j th column of the correlation matrix \mathbf{R} .

Basic properties of the multiple correlation coefficient are

- (1) $0 \leq R_{1(2, \dots, m)} \leq 1$;
- (2) if $R_{1(2, \dots, m)} = 1$, the random quantity ξ_1 is exactly a linear combination of quantities ξ_2, \dots, ξ_m ;
- (3) if $R_{1(2, \dots, m)} = 0$, all pairwise correlation coefficients $\rho(\xi_1, \xi_j) = 0$, $j = 2, \dots, m$;
- (4) for the case of a single explanatory variable, the multiple correlation coefficient is identical with the absolute value of the paired correlation coefficient, $R_{1(2)} = |\rho(\xi_1, \xi_2)|$;
- (5) as the number of explanatory variables increases, the multiple correlation coefficient never decreases:

$$R_{1(2)}^2 \leq R_{1(2,3)}^2 \leq R_{1(2,3,4)}^2 \leq \dots \leq R_{1(2, \dots, m)}^2$$

Problem 7.4. Multiple correlation coefficient for two explanatory variables

Estimate the multiple correlation coefficient $R_{1(2,3)}$ between a variable ξ_1 and two variables ξ_2, ξ_3 .

Solution: For correlation matrices \mathbf{R} and \mathbf{R}_{11} we can write

$$\mathbf{R} = \begin{bmatrix} 1 & R_{12} & R_{13} \\ R_{12} & 1 & R_{23} \\ R_{13} & R_{23} & 1 \end{bmatrix} \quad \text{and} \quad \mathbf{R}_{11} = \begin{bmatrix} 1 & R_{23} \\ R_{23} & 1 \end{bmatrix}$$

In these expressions, the symmetry $R_{ij} = R_{ji}$ of paired correlation coefficients is used. After substitution into Eq. (7.23) and some rearrangement we get

$$R_{1(2,3)} = \sqrt{\frac{R_{12}^2 + R_{13}^2 - 2R_{12} \times R_{13} \times R_{23}}{1 - R_{23}^2}} \quad (7.24)$$

Equation (7.24) shows that paired correlation coefficients can not reach any value in the range $-1 \leq R_{ij} \leq 1$, but they are mutually bounded by the condition $R_{1(2,3)} \leq 1$.

If $R_{23} = 0$, the explanatory variables are mutually uncorrelated, and

$$R_{1(2,3)}^2 = R_{12}^2 = R_{13}^2 \quad (7.25a)$$

Conclusion: The multiple correlation coefficient may be estimated as the function of paired correlation coefficients. When the explanatory variables ξ_2, \dots, ξ_m are mutually uncorrelated, the square of the multiple correlation coefficient is equal to the sum of squares of paired correlation coefficients.

In some cases the centred random variables or normalized random variables are used. For centred random variables

$$\xi_{cj} = \xi_j - \mu_j, \quad j = 1, \dots, m$$

and for normalized random variables

$$\xi_{Nj} = \frac{\xi_j - \mu_j}{\sigma_j}, \quad j = 1, \dots, m$$

The regression defined by Eq. (7.18) may be expressed with the use of centred random variables in the form

$$E(\xi_1/\mathbf{x}_c^*) = \mathbf{c}_1^T \mathbf{C}^{*-1} \mathbf{x}_c^* = \sum_{i=1}^{m-1} a_i x_{c,i+1} \quad (7.25b)$$

It can be seen that centring does not change the estimates of the regression coefficients, but the intercept term is equal to zero.

With normalized random variables, the regression $E(\xi_1/\mathbf{x}^*)$ takes the form

$$E(\xi_{N1}/\mathbf{x}_N^*) = \mathbf{R}^T \mathbf{R}^{*-1} \mathbf{x}_N^* = \frac{\mathbf{a}^T \mathbf{D} \mathbf{x}_c^*}{\sigma_1} = \mathbf{b}^T \mathbf{x}_N^* \quad (7.26)$$

where \mathbf{R} is a vector of size $((m-1) \times 1)$ containing paired correlation coefficients $\rho(\xi_1, \xi_j)$, $j = 2, \dots, m$, and \mathbf{R}^* is the correlation matrix of the vector of explanatory variables of size $(m-1) \times (m-1)$, \mathbf{D} denotes a diagonal transformation matrix with elements σ_j , $j = 2, \dots, m$, on the main diagonal. The coefficients $\mathbf{b}_j = \mathbf{R}^{*-1} \mathbf{R}$ are called the normalized regression coefficients. From Eq. (7.26) it follows that a relationship exists between non-normalized (a_j) and normalized (b_j) regression coefficients

$$a_{j-1} = b_{j-1} \frac{\sigma_1}{\sigma_j}, \quad j = 2, \dots, m \quad (7.27)$$

The normalization changes the magnitude of the regression coefficients. The advantage

of normalized regression coefficients is the fact that they concern directly the paired correlation coefficients and are easier to interpret.

Problem 7.5. Regression for two explanatory variables

Estimate coefficients **a** and **b** of the linear regression of one response variable ξ_1 and two explanatory variables ξ_2, ξ_3 .

Solution: The vector **R** is (R_{12}, R_{13}) . Matrix **R*** is identical to the matrix **R**₁₁ from Problem 7.4. For the matrix **R***⁻¹

$$\mathbf{R}^{*-1} = (1 - R_{23}^2)^{-1} \begin{bmatrix} 1 & -R_{23} \\ -R_{23} & 1 \end{bmatrix}$$

The normalized regression coefficients then are

$$\begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = \frac{1}{1 - R_{23}^2} \begin{bmatrix} 1 & -R_{23} \\ -R_{23} & 1 \end{bmatrix} \begin{bmatrix} R_{12} \\ R_{13} \end{bmatrix}$$

leading to the expressions

$$b_1 = \frac{R_{12} - R_{23} \times R_{13}}{1 - R_{23}^2} \quad (7.28a)$$

$$b_2 = \frac{R_{13} - R_{23} \times R_{12}}{1 - R_{23}^2} \quad (7.28b)$$

Non-normalized regression coefficients are expressed by

$$a_1 = b_1 \frac{\sigma_1}{\sigma_2}$$

$$a_2 = b_2 \frac{\sigma_1}{\sigma_3}$$

If the explanatory variables are mutually uncorrelated, $R_{23} = 0$ and the normalized regression coefficients correspond to the paired correlation coefficients

$$b_1 = R_{12} \quad (7.29a)$$

and

$$b_2 = R_{13} \quad (7.29b)$$

Conclusion: It is obvious that the regression coefficients are functions of the only paired correlation coefficients. When the explanatory variables are mutually uncorrelated, the normalized regression coefficients are directly equal to the paired correlation coefficients between the *j*th response and the explanatory random variable.

The regression coefficients **a** and **b** are such that the correlation between random variables ξ_1 and $\hat{\xi}_1 = (\xi_1/\mathbf{x}^*)$ is maximal. Random variable $\hat{\xi}_1$ is the linear combination of components $\xi_2, \xi_3, \dots, \xi_m$ of random vector $\boldsymbol{\xi}$ for realization \mathbf{x}^* . It is also valid that

$$\rho(\xi_1, \hat{\xi}_1) = R_{1(2,3,\dots,m)} \quad (7.30)$$

In regression analysis $\hat{\xi}_1$ is called the *prediction*. Equation (7.30) shows that the multiple correlation coefficient is, in fact, the paired correlation coefficient between the vector ξ_1 and vector $\hat{\xi}_1$.

The random variable $\varepsilon = \xi_1 - \hat{\xi}_1$ is called the residual in point \mathbf{x}^* . The residuals are uncorrelated with individual explanatory variables because

$$\text{cov}(\varepsilon, \xi_j) = E(\varepsilon \xi_j) = E(\xi_1 - \hat{\xi}_1 \xi_j) = 0$$

The covariance of residuals with a controllable variable ξ_1 is equal to the conditioned variance

$$\text{cov}(\varepsilon, \xi_1) = D(\xi_1/\mathbf{x}^*) \quad (7.31)$$

which is also the residual variance $D(\varepsilon)$. As for the case of one explanatory variable (Section 7.1.1) the estimates of the regression parameters and other random variables may be obtained on the basis of the sample means $\bar{x}_1, \dots, \bar{x}_m$, and the sample covariance, or correlation matrix, respectively. It is also valid that these estimates are identical with the estimates obtained by the least-squares method.

Problem 7.6. *The effect of inorganic and organic nitrogen in soil on the nitrogen content of corn*

The effect of the concentration of inorganic nitrogen (x_2) and organic nitrogen (x_3) in the soil on the content of nitrogen in corn has been studied [4]. Estimate the linear regression model for the regression of x_1 on x_2 and x_3 , and calculate the multiple correlation coefficient.

Data:

x_1	64	60	71	61	54	77	81	93	93	51
x_2	0.4	0.4	3.1	0.6	4.7	1.7	9.4	10.1	11.6	12.6
x_3	53	23	19	34	24	65	44	31	29	58
	76	96	77	93	95	54	168	99		
	10.6	23.1	23.1	21.6	23.1	1.9	26.8	29.9		
	37	46	50	44	56	36	59	51		

Solution: From the expressions for the mean, variance and pair correlation coefficient we estimate

$$\begin{aligned} \bar{x}_1 &= 81.28 & \bar{x}_2 &= 11.94 & \bar{x}_3 &= 42.11 \\ s_1^2 &= 728.8 & s_2^2 &= 103.6 & s_3^2 &= 185.6 \end{aligned}$$

and

$$\mathbf{R} = \begin{bmatrix} 1 & 0.6934 & 0.3545 \\ 0.6934 & 1 & 0.4616 \\ 0.3545 & 0.4616 & 1 \end{bmatrix}$$

On substitution into Eq. (7.24) we obtain

$$R_{1(2,3)} = \sqrt{\frac{0.6934^2 + 0.3545^2 - 2 \times 0.6934 \times 0.3545 \times 0.4616}{1 - 0.4616^2}} = 0.6945$$

From Eq. (7.28a) we obtain

$$b_1 = \frac{0.6934 - 0.4616 \times 0.3545}{1 - 0.4616^2} = 0.6732$$

and from Eq. (7.28b)

$$b_2 = 0.04375$$

For non-normalized regression coefficients

$$a_1 = \sqrt{\frac{728.8}{103.6}} \times 0.6732 = 1.7855$$

$$a_2 = 0.08669$$

For the intercept term, from Eq. (7.18), we have

$$a_0 = \bar{x}_1 - a_1 \bar{x}_2 - a_2 \bar{x}_3 = 56.31$$

The linear regression model has the form

$$x_1 = 56.31 + 1.7855x_2 + 0.08669x_3$$

The multiple correlation coefficients are the same, and the coefficients a_0 , a_1 , a_2 correspond to the estimates by the least-squares method.

Conclusion: The coefficients of linear regression models and the multiple correlation coefficient may be calculated directly from the definitions.

From a practical point of view, it is convenient to use computer programs for linear regression. However, these expressions show that regression models can be directly derived from the random vector, and moreover they often aid the interpretation of the statistical characteristics.

In some cases it is useful to examine a relationship between two components ξ_1 and ξ_j of a random vector, when the other components of vector ξ are excluded. To express the intensity of this dependence, the *partial correlation coefficients* of various orders are used. The simplest are the partial correlation coefficients of zero order, which correspond to the paired correlation coefficients.

The partial correlation coefficient of the first order $R_{1,3(2)}$ corresponds to the paired correlation coefficient between the residuals

$$\varepsilon_2 = \xi_1 - E(\xi_1/x_2)$$

and the residuals

$$\kappa_2 = \xi_3 - E(\xi_3/x_2)$$

Then

$$R_{1,3(2)} = \frac{R_{13} - R_{12} \times R_{23}}{\sqrt{(1 - R_{12}^2)(1 - R_{23}^2)}} \quad (7.32)$$

Similarly, other partial correlation coefficients $R_{1i(j)}$ of the first order can be defined from the paired correlation coefficients between residuals

$$\varepsilon_j = \xi_1 - E(\xi_1/x_j)$$

and residuals

$$\kappa_j = \xi_i - E(\xi_i/x_j) \quad (7.33)$$

Then, it can be shown that

$$R_{1,i(j)} = \frac{R_{1i} - R_{1j}R_{ij}}{\sqrt{(1 - R_{1i}^2)(1 - R_{ij}^2)}}$$

The partial correlation coefficients of the second order $R_{1,i(j,k)}$ are the same as the paired correlation coefficients of residuals

$$\varepsilon_{j,k} = \xi_1 - E(\xi_1/(x_j, x_k))$$

and residuals

$$\kappa_{j,k} = \xi_i - E(\xi_i/(x_j, x_k))$$

To estimate these, an equation analogous to Eq. (7.33) may be used, where instead of the paired correlation coefficients the partial correlation coefficients of the first order are used

$$R_{1,i(j,k)} = \frac{R_{1i(j)} - R_{1j(k)} \times R_{ij(k)}}{\sqrt{(1 - R_{1j(k)}^2)(1 - R_{ij(k)}^2)}} \quad (7.34)$$

The partial correlation coefficient of the $(m - 1)$ th order $R_{1,i(2,3,\dots,m)}$ corresponds to the paired correlation coefficient between residuals

$$\varepsilon_{2,\dots,m} = \xi_1 - E(\xi_1/\mathbf{x}^*)$$

and residuals

$$\kappa_{2,\dots,m} = \xi_i - E(\xi_i/\mathbf{x}^*)$$

where the vector \mathbf{x}^* contains the components $x_2, x_3, \dots, x_{i-1}, x_{i+1}, \dots, x_m$.

Generally, the partial correlation coefficients of the higher orders are estimated according to a recursive formula

$$R_{1,j(2,3,\dots,j-1)} = \frac{A - B \times C}{\sqrt{(1 - B^2)(1 - C^2)}} \quad (7.35)$$

where

$$A = R_{1,j(2,3,\dots,j-2)}$$

$$B = R_{1,j-1(2,3,\dots,j-2)}$$

and

$$C = R_{j,j-1(2,3,\dots,j-2)}$$

When individual partial correlation coefficients of all orders are known, the multiple correlation coefficient can be estimated from

$$R_{1(2,\dots,m)}^2 = 1 - (1 - R_{1,2}^2)(1 - R_{1,3(2)}^2)(1 - R_{1,4(2,3)}^2) \dots (1 - R_{1,m(2,3,\dots,(m-1))}^2)$$

All these expressions can be evaluated on a pocket calculator.

In the computer estimation of the partial correlation coefficients, matrix notation is convenient.

$$R_{1,i(2,3,\dots,m)} = \frac{(-1)^i \times \det(\mathbf{R}_{1,i})}{\sqrt{\det(\mathbf{R}_{1,1}) \times \det(\mathbf{R}_{i,i})}} \quad (7.37)$$

where \mathbf{R} is the correlation matrix corresponding to the vector ξ and $\mathbf{R}_{i,j}$ is the matrix formed by leaving out the i th row and the j th column of matrix \mathbf{R} .

Problem 7.7. *Partial correlation coefficients of the first order*

For the random variables ξ_1, ξ_2, ξ_3 , estimate with the use of Eq. (7.37) the partial correlation coefficients $R_{1,2(3)}$ and $R_{1,3(2)}$.

Solution: For calculation of the correlation coefficients the matrices $\mathbf{R}, \mathbf{R}_{11}, \mathbf{R}_{12}, \mathbf{R}_{22}, \mathbf{R}_{13}$ and \mathbf{R}_{33} are necessary. We determined

$$\begin{aligned} \mathbf{R} &= \begin{bmatrix} 1 & R_{12} & R_{13} \\ R_{12} & 1 & R_{23} \\ R_{13} & R_{23} & 1 \end{bmatrix} & \mathbf{R}_{11} &= \begin{bmatrix} 1 & R_{23} \\ R_{23} & 1 \end{bmatrix} \\ \mathbf{R}_{12} &= \begin{bmatrix} R_{12} & R_{23} \\ R_{13} & 1 \end{bmatrix} & \mathbf{R}_{33} &= \begin{bmatrix} 1 & R_{12} \\ R_{12} & 1 \end{bmatrix} \\ \mathbf{R}_{13} &= \begin{bmatrix} R_{12} & 1 \\ R_{13} & R_{23} \end{bmatrix} & \mathbf{R}_{22} &= \begin{bmatrix} 1 & R_{13} \\ R_{13} & 1 \end{bmatrix} \end{aligned}$$

On substitution into Eq. (7.37) we find

$$R_{1,2(3)} = \frac{R_{12} - R_{23} \times R_{13}}{\sqrt{(1 - R_{23}^2)(1 - R_{13}^2)}} \quad (7.38)$$

or

$$R_{1,3(2)} = \frac{R_{13} - R_{12} \times R_{23}}{\sqrt{(1 - R_{23}^2)(1 - R_{12}^2)}} \quad (7.39)$$

Conclusion: The partial correlation coefficients may be estimated directly from Eq. (7.37).

It is interesting that by use of the partial correlation coefficients, the normalized regression coefficients may also be found. The intensity of the mutual relationship

between components of a random vector may be better estimated by the partial correlation coefficients than by the paired correlation coefficients.

Problem 7.8. *Partial correlation between the nitrogen content in corn and in soil*

For the data from Problem 7.6, calculate the partial correlation coefficients between the nitrogen content in corn and (a) the content of inorganic nitrogen in soil $R_{1,2(3)}$, and (b) the content of organic nitrogen in soil, $R_{1,3(2)}$.

Data: from Problem 7.6

Solution: By direct substitution into Eq. (7.38), we find

$$R_{1,2(3)} = \frac{0.6934 - 0.4616 \times 0.3545}{\sqrt{(1 - 0.4616^2)(1 - 0.3545^2)}} = 0.6386$$

and from Eq. (7.39)

$$R_{1,3(2)} = 0.05325$$

Conclusion: The nearly zero value of the partial correlation coefficient $R_{1,3(2)}$ shows that the influence of organic nitrogen in soil on the nitrogen content in corn is negligible. The relatively high value of the paired correlation coefficient $R_{13} = 0.3545$ is strongly affected by the correlation $R_{23} = 0.462$ between the organic and inorganic nitrogen in soil. Detailed data analysis shows that point 17 is an outlier, so the analysis should be repeated with that point omitted.

7.2 CORRELATION COEFFICIENTS

Correlation coefficients serve as basic measures for expressing “closeness” of the linear stochastic dependence between components of the random vector ξ . In the literature [5] many other characteristics are used which can also cover nonlinear stochastic dependences. We restrict ourselves here to a description of the distributions of the sample correlation coefficients, and some selected tests.

7.2.1 Paired correlation coefficient

The paired correlation coefficient $\rho(\xi_i, \xi_j) = R_{ij}$ is a measure of the linear stochastic dependence between the random variables ξ_i and ξ_j . For sample size n , the same correlation coefficient may be estimated from Eq. (7.15). For simplicity, we denote the paired correlation coefficient by the letter ρ and the sample paired correlation coefficient by R .

At first we restrict discussion to a case when the joint distribution of quantities ξ_1 and ξ_2 is normal and $\rho = 0$. Then the probability density of random quantity R is symmetrical around zero and has the shape [6] expressed by

$$f(R) = \frac{2^{n-3}}{\pi} \times \frac{\left[\Gamma\left[\frac{n-1}{2}\right] \right]^2}{\Gamma[n]} \times (1 - R^2)^{(n-4)/2} \quad (7.40)$$

where $\Gamma(\cdot)$ is the Gamma function. For $n = 5, 9$ and 51 , the curves $f(R)$ are illustrated in Fig. 7.4.

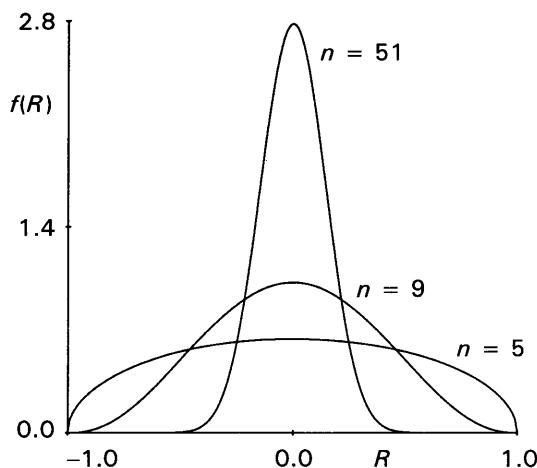


Fig. 7.4—The probability density of the sample correlation coefficient for $\rho = 0$ and for sample size $n = 5, 9$ and 51 .

In construction of significance tests the following test criterion is used

$$t = \frac{R\sqrt{n-2}}{\sqrt{1-R^2}} \quad (7.41)$$

which for $\rho = 0$ has the Student distribution with $(n-2)$ degrees of freedom. This may be used for testing the independence between a pair of random variables. If their distribution is two-dimensionally normal, lack of correlation is identical to independence.

The null hypothesis $H_0: \rho = 0$ is tested *vs.* various alternatives. If the criterion $|t|$ from Eq. (7.41) is larger than the corresponding quantile of the Student distribution, the null hypothesis is rejected and the random variables are not correlated. This test is strongly non-robust and is valid only in the case of two-dimensional normality of ξ_1 and ξ_2 . To speed up the convergence of $f(R)$ to the normal distribution, various transformations are used. The simple Ruben transformation has the form

$$R(R) = \frac{R\sqrt{n-2.5}}{\sqrt{1-0.5R^2}} \quad (7.42)$$

The random variable $R(R)$ has, even for small sample sizes, the normalized normal distribution $N(0, 1)$.

Problem 7.9. *Significance of the degree of polymerization of cotton in two solutions*
Determine the significance of the correlation coefficient between the degrees of polymerization of cotton determined in solution CUEN and CUOXAN (Problem 7.3) when the sample correlation coefficient $R = 0.6142$ was estimated from 17 data values.

Data: from Problem 7.3

Solution: We select the significance level $\alpha = 0.05$. Since there must be a positive linear relationship between the degrees of polymerization in the two solutions, we know that $\rho \geq 0$. We examine the null hypothesis $H_0: \rho = 0$ vs. $H_A: \rho > 0$.

- (a) By substituting into Eq. (7.41), we calculate the test criterion $t = 3.104$. This value is higher than the quantile of the Student distribution $t_{0.95}(15) = 1.753$, so the inequality $\rho > 0$ is accepted.
- (b) By substituting into Eq. (7.42) we calculate the test criterion of the Ruben transformation $R(R) = 2.596$. This value is higher than the quantile of the normalized normal distribution $u_{0.95} = 1.645$, so the null hypothesis H_0 is rejected.

Conclusion: Both tests prove that the population correlation coefficient is, with 95% probability, positive, and therefore correlation exists between the two degrees of polymerization of the two solutions.

A common case is a simultaneous distribution of two random variables that are two-dimensionally normal with $\rho \neq 0$. For $n > 3$ the probability density function of the sample correlation coefficient may be expressed in the form [6].

$$f(R/\rho) = \frac{2^{n-3}}{\pi(n-3)!} (1 - \rho^2)^{(n-1)/2} \times (1 - R^2)^{(n-4)/2} \\ \times \sum_{j=0}^{\infty} \frac{(2\rho R)^j}{j!} \left[\Gamma \left[\frac{n+j-1}{2} \right] \right]^2 \quad \text{for } -1 < R < 1, \quad (7.43)$$

$$f(R/\rho) = 0 \quad \text{elsewhere}$$

The probability density function $f(R/\rho)$ is rather asymmetrical for small sample sizes (see Fig. 7.5).

For sufficiently large sample sizes ($n > 500$), the distribution of $f(R/\rho)$ can be approximated by a normal distribution with mean $E(R) = \rho$ and variance $D(R) = (1 - \rho^2)^2/(n - 1)$.

If random variables ξ_1 and ξ_2 have a two-dimensional *elliptic* distribution with correlation coefficient ρ and kurtosis g_2 , the random variable

$$u_n = \frac{|R - \rho| \times \sqrt{n-1}}{1 - \rho^2} \quad (7.44)$$

has an asymptotically normal distribution with zero mean value and variance equal to $(1 + g_2)$.

Problem 7.10. *The confidence interval of the correlation coefficient*

For 600 random samples, the content of iron was determined by two analytical methods with correlation coefficient $R = 0.85$. Estimate the 95% confidence interval of the correlation coefficient ρ . Examine the null hypothesis $H_0: \rho = 0.9$ against the alternative $H_A: \rho \neq 0.9$.

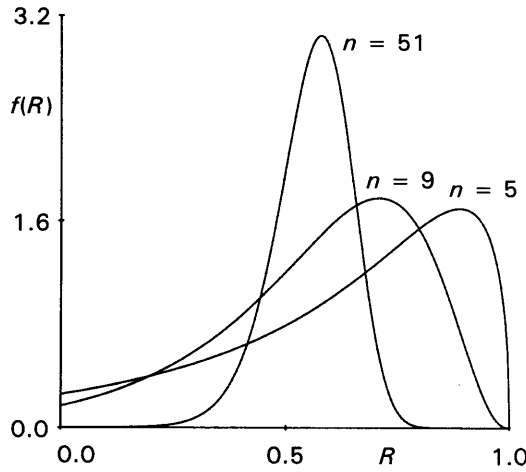


Fig. 7.5—The probability density function of the sample correlation coefficient (for $\rho = 0.6$) for sample sizes $n = 5, 9$ and 51 .

Solution: We make use of the asymptotic normality of the distribution of the correlation coefficient. It is valid that

$$R - u_{1-\alpha/2} \frac{(1 - R^2)}{\sqrt{n-1}} \leq \rho \leq R + u_{1-\alpha/2} \frac{(1 - R^2)}{\sqrt{n-1}}$$

where $u_{1-\alpha/2}$ is the quantile of the normalized normal distribution. On substituting, we get

$$0.828 \leq \rho \leq 0.872$$

To test the H_0 hypothesis, we use the test criterion u_n [Eq. (7.44)]

$$u_n = \frac{|0.85 - 0.9| \sqrt{600 - 1}}{1 - 0.9^2} = 6.44$$

which is higher than the quantile $u_{0.975} = 1.96$, so we reject the null hypothesis H_0 : $\rho = 0.9$ at the significance level $\alpha = 0.05$. The correlation between the two analytical methods for iron determination differs significantly from the value 0.9.

Conclusion: For large sample sizes the asymptomatic normality can be utilized to permit a test of paired correlation coefficients.

If the sample correlation coefficient is estimated from the sample size n , some measures of location, spread and distribution shape may be approximated [6]. The mean may be expressed in the form

$$E(R) = \rho + (1 - \rho^2) \left[-\frac{\rho}{2(n-1)} + \frac{\rho - 9\rho^3}{8(n-1)^2} + \frac{\rho + 42\rho^3 - 75\rho^5}{16(n-1)^3} + \dots \right] \quad (7.45)$$

the variance

$$D(R) = \frac{1}{n}(1 - \rho^2)^2 \left[1 + \frac{11\rho^2}{2(n-1)} + \frac{-24\rho^2 + 75\rho^4}{2(n-1)^2} + \dots \right] \quad (7.46)$$

the skewness

$$g_1(R) = \frac{6\rho}{n-1} \left[1 + \frac{-30 + 70\rho^2}{12(n-1)} + \dots \right] \quad (7.47)$$

and the kurtosis

$$g_2(R) \approx \frac{6(12\rho^2 - 1)}{n-1} \quad (7.48)$$

The bias $E(R) - \rho$ is, as an initial approximation, equal to

$$E(R) - \rho = \frac{-\rho(1 - \rho^2)}{2(n-1)}$$

and the estimate R calculated by Eq. (7.15) is rather underestimated for $\rho > 0$. For very small sample sizes ($n < 15$), the corrected correlation coefficient is used for practical calculations. This is given by

$$R^* = R \left[1 + \frac{1 - R^2}{2(n-4)} \right] \quad (7.49)$$

The square of the correlation coefficient is strongly *overestimated* in cases when the sample is not random. For larger intervals of sample values and more uniform scaling, the R^2 value is more overestimated.

To improve the statistical properties of a distribution of the sample correlation coefficient, many transformations which speed up convergence to normality are used. The best known is the Fisher transformation [7] which takes the form

$$\begin{aligned} Z(R) &= \arctan(R) \\ &= 0.5 \ln \left[\frac{1+R}{1-R} \right] \end{aligned} \quad (7.50)$$

This transformation stabilizes the variance. For $n > 50$, the distribution of quantity $Z(R)$ is approximately normal, with mean value $E(Z)$ and variance $D(Z)$ calculated from

$$E(Z) = Z(\rho) \quad (7.51a)$$

$$D(Z) = \frac{1}{n-3} \quad (7.51b)$$

More exact estimates of the mean, variance, skewness and kurtosis are given by

$$E(Z) = Z(\rho) + \frac{0.5\rho}{n-1} + \frac{\rho(5 + \rho^2)}{8(n-1)^2} + \frac{\rho(11 + 2\rho^2 + 3\rho^4)}{16(n-1)^3} + \dots \quad (7.52)$$

$$D(Z) = \frac{1}{n-1} \left[1 + \frac{0.5(4-\rho^2)}{n-1} + \dots \right] \approx \frac{1}{n-3} - \frac{\rho^2}{2(n-3)^2} \quad (7.53)$$

$$g_1(Z) = \frac{\rho^3}{(n-1)^{2/3}} + \dots \quad (7.54)$$

$$g_2(Z) = \frac{2}{n+1} + \dots \quad (7.55)$$

For small sample sizes, the Sammiunddin transformation [8] is recommended

$$S = \frac{(R - \rho)\sqrt{n-2}}{\sqrt{(1-R^2)(1-\rho^2)}} \quad (7.56)$$

The quantity S has approximately the Student distribution with $n-2$ degrees of freedom for $\rho \neq 0$.

Kraemer [9] replaces the correlation coefficient ρ in Eq. (7.56) with the median $\tilde{\rho}$ of a distribution of the sample correlation coefficient, for which

$$\tilde{\rho} \approx \rho + (1-\rho^2)\rho \left[\frac{2}{n-1} + \frac{(-7\rho^2+15)}{24(n-1)^2} + \dots \right] \quad (7.57)$$

There are many other transformations [6] which are recommended for small or large sample sizes.

Problem 7.11. *Examination of the correlation coefficient between the degrees of polymerization of cotton in two solutions*

Suppose that if the correlation coefficient between the degrees of polymerization of cotton in CUEN and CUOXAN solutions is not smaller than $\rho_0 = 0.85$, a significant linear association between the results exists. Examine the null hypothesis $H_0: \rho = 0.85$ against alternative $H_A: \rho < 0.85$, with the use of various transformations.

Solution: (a) Fisher transformation (7.50) leads to the test criterion

$$u_F = |Z(R) - Z(\rho_0)| \times \sqrt{n-3}$$

with approximate distribution $N(0, 1)$. From Eq. (7.50) we get

$$Z(R) = 0.5 \ln \left[\frac{1 + 0.6142}{1 - 0.6142} \right] = 0.7156$$

and

$$Z(\rho_0) = Z(0.85) = 1.256$$

Then

$$u_F = |0.7156 - 1.256| \sqrt{17-3} = 2.021$$

is higher than the quantile $u_{0.95} = 1.64$ and therefore the null hypothesis $H_0: \rho = 0.85$ cannot be accepted.

(b) The Sammiunddin transformation (7.56) leads to the test criterion

$$S = \frac{|0.6142 - 0.85|\sqrt{17 - 2}}{\sqrt{(1 - 0.614^2)(1 - 0.85^2)}} = 2.197$$

is higher than the quantile $t_{0.95}(15) = 1.725$, so the null hypothesis $H_0: \rho = 0.85$ is rejected.

(c) Kraemer modification (7.57) leads to the median

$$\tilde{\rho} \approx 0.85 + (1 - 0.85^2)0.85 \left[\frac{2}{16} + \dots \right] = 0.879$$

The test criterion

$$S = \frac{|0.6142 - 0.879|\sqrt{17 - 2}}{\sqrt{(1 - 0.879^2)(1 - 0.6142^2)}} = 2.726$$

is higher than the quantile $t_{0.95}(15) = 1.725$ and therefore the null hypothesis $H_0: \rho = 0.85$ is rejected.

Conclusion: All three tests used show that the correlation coefficient ρ is significantly lower than the value 0.85 and therefore the stochastic dependence between the degrees of polymerization is not very strong.

When the joint distribution of random variables is not normal and the sample contains strong outliers, the normalized transformation is not valid and the correlation coefficient is not suitable for expressing a stochastic association. We can then use various robust estimates of correlation coefficients, which apply robust estimates of parameters of location, spread and covariance. Some techniques have been described [10].

The correlation coefficients should be interpreted very carefully. As a general rule, a significant paired correlation is not the proof of a causal dependence. Sometimes false correlations are formed when either ξ_1 or ξ_2 strongly correlate with some unconsidered random variable ξ_3 , and a high value of $\rho(\xi_1, \xi_2)$ is the consequence of high values of $\rho(\xi_1, \xi_3)$ and $\rho(\xi_2, \xi_3)$. In the interpretation of correlation coefficients, the partial correlation coefficients should also be considered.

7.2.2 Partial correlation coefficients

For calculation of partial correlation coefficients either the recursive formulae [Eq. (7.35)] or the matrix method [Eq. (7.37)] can be used.

For statistical testing and building of the confidence interval, we use a rule that a distribution of the partial correlation coefficient of the order $(m - 1)$ is identical to the distribution of the paired correlation coefficients for sample size $(n - m + 1)$. Thus, techniques described in Section 7.2.1, with modified sample size, may also be used.

Problem 7.12. *Significance of the dependence between the organic nitrogen in soil and the content of nitrogen in corn*

For the data from Problem 7.8, examine the significance of the correlation coefficient $R_{1,3(2)}$ as an expression of the association between organically bound nitrogen in soil and content of nitrogen in corn.

Data: from Problem 7.8

Solution: To examine the significance of the null hypothesis $H_0: R_{1,3(2)} = 0$ against $H_A: R_{1,3(2)} \neq 0$ we use the above relationship. Because the partial correlation coefficient is of the first order, we have $m - 1 = 1$ and the reduced sample size is $n - 1$. From Eq. (7.41) the test criterion

$$t_P = \frac{\hat{R}_{1,3(2)}\sqrt{n-3}}{\sqrt{1-R_{1,3(2)}^2}} = 0.227$$

is smaller than the quantile $t_{0.975}(15) = 2.13$, so the null hypothesis $H_0: R_{1,3(2)} = 0$ is accepted. In calculation of t_P , the partial correlation coefficient $R_{1,3(2)} = 0.05858$ for $n = 18$ was used.

Conclusion: On the basis of the test of the partial correlation coefficient, it is concluded that there is no significant correlation between the organic nitrogen content in soil and in corn.

Partial correlation coefficients can be used to elucidate some false correlations. Consider a case when the paired correlation coefficient between ξ_1 and ξ_2 is $R_{12} = H$ where $H \rightarrow 1$. Suppose there is a random variable ξ_3 which strongly correlates with ξ_1 and ξ_2 , so that $R_{13} = H^2$ and $R_{23} = H$. Then the multiple correlation coefficient

$$R_{1(2,3)} = H$$

may be estimated, and the partial correlation coefficients are equal to

$$R_{1,3(2)} = 0$$

$$R_{1,2(3)} = \frac{H}{\sqrt{1+H^2}}$$

Despite the high value of the paired correlation coefficient ($R_{13} = H^2$ is close to one) the quantity ξ_3 does not contribute to the explanation of the variability of ξ_1 and ξ_2 . It is a typical *parasite variable*. When $R_{23} = R_{13}/R_{12}$, the variable ξ_3 is a parasite.

This situation can arise from the neglect of a significant variable such as, for example, time or temperature. For example, at various time values during a degradation process, the mechanical or optical properties of the materials will be different. If time is ignored, a significant "false" correlation among these properties appears. When time is included as a variable, the optical properties do not contribute to explaining the variability of the mechanical properties. A high value of the paired correlation coefficient is not always a guarantee of a significant association between variables.

Similarly, there are cases when a low value of the paired correlation coefficient leads to high partial correlation coefficients and a high multiple correlation coefficient.

When $R_{13} = 0$ and $R_{12} = \varepsilon$ ($\varepsilon \rightarrow 0$), but variables ξ_2 and ξ_3 are strongly correlated i.e.

$$R_{23} = \sqrt{(1 - \varepsilon^2)}$$

then $R_{12(3)} = 1$, $R_{13(2)} = -1$ and $R_{1(2,3)} = 1$. From this it follows that a zero paired correlation coefficient does not mean automatically that a given random variable is insignificant or parasite and may be excluded.

Moreover it is not valid that strongly correlated random variables are always redundant. These examples demonstrate that often no conclusions can be made from the paired correlation coefficients. It is because other variables are not considered that a "false" correlation may be concluded.

7.2.3 Multiple correlation coefficient

The multiple correlation coefficient, denoted as $R_{1(2,\dots,m)}$, is a measure of the overall linear stochastic association of one random variable ξ_1 with the best linear combination of the other components ξ_2, \dots, ξ_m , of the random vector ξ . The sample correlation coefficient $R_{1(2,\dots,m)}$ may be readily calculated from Eq. (7.23) by replacing the correlation matrix \mathbf{R} by the sample correlation matrix $\hat{\mathbf{R}}$. For the sake of simplicity, we refer to the multiple correlation coefficient of the population as R_m and its sample estimate as \hat{R}_m .

Let us suppose that the vector ξ has an m -multidimensional normal distribution and that all its components have a normal distribution.

Case $R_m = 0$:

The probability density of random variable \hat{R}_m^2 is given by

$$f(\hat{R}_m^2) = K_m \times (\hat{R}_m^2)^{(m-3)/2} \times (1 - \hat{R}_m^2)^{(n-m-2)/2} \quad (7.58)$$

where K_m is a constant dependent on m and n . The distribution defined by Eq. (7.58) is a beta-distribution $\text{Be}[(m-1)/2, (n-m)/2]$. Then the random variable given by

$$F_R = \frac{(n-m) \times \hat{R}_m^2}{(m-1)(1 - \hat{R}_m^2)} \quad (7.59)$$

has the F -distribution with $(m-1)$ and $(n-m)$ degrees of freedom.

For large sample sizes, the distribution of

$$C_R = (n-1) \times \hat{R}_m^2 \quad (7.60)$$

is χ^2 with $(m-1)$ degrees of freedom.

For the mean value of the sample squared multiple correlation coefficient, we have

$$E(\hat{R}_m^2) = \frac{m-1}{n-1} \quad (7.61)$$

and for the variances

$$D(\hat{R}_m^2) = \frac{2(n-m+2)(m-1)}{(n^2-1)(n-1)} \approx \frac{2(m-1)}{n^2} \quad (7.62)$$

Equation (7.61) shows that for a small number of measurements and a large number of explanatory variables, the quantity \hat{R}_m^2 will be significantly different from zero even in cases when the population multiple correlation coefficient R_m^2 is equal to zero. For example, if $n = 12$ and $m = 9$, then $E(\hat{R}_m^2) = 0.727$ even though $R_m^2 = 0$. This negative effect may be removed by decreasing m or increasing n . A sample size higher than $n_{\min} = (1 + 100 \times m)$ ensures that $\hat{R}^2 \approx 0.1$ for uncorrelated random variables.

In the case of a multi-variable normal distribution, the null hypothesis $H_0: R_m = 0$ against the alternative $H_A: R_m \neq 0$, with use of criterion F_R , is suitable as a test for independence.

Problem 7.13. *Significance of the relationship between the nitrogen content in soil and in corn*

In Problem 7.6 the multiple correlation coefficient expressing the relationship between the nitrogen content in corn and a linear combination of organically bound nitrogen and inorganically bound nitrogen in soil is equal to $\hat{R}_{1(2,3)} = 0.6945$. Examine the null hypothesis $H_0: R_{1(2,3)} = 0$.

Solution: According to Eq. (7.59), the test criterion

$$F_R = \frac{(18 - 3)0.6945^2}{2 \times (1 - 0.6945^2)} = 6.988$$

is higher than the quantile of the Fisher-Snedecor distribution $F_{0.95}(2, 15) = 3.682$, and therefore the null hypothesis $H_0: R_{1(2,3)} = 0$ is rejected at significance level $\alpha = 0.05$.

Conclusion: The content of nitrogen in soil significantly affects the content of nitrogen in corn. Inorganically bound nitrogen contributes predominantly.

Case $R_m > 0$:

To calculate the sample multiple correlation coefficient, \hat{R}_m , the complicated exact expression or a convenient approximation may be used. Gurland [6] has proposed a relatively precise approximation

$$\frac{\hat{R}_m^2}{1 - \hat{R}_m^2} \approx \frac{\frac{(n-1)R_m^2}{1 - R_m^2} + (m-1)}{n-m} \times F_{r,n-m} \quad (7.63)$$

where the quantity $F_{r,n-m}$ has the F -distribution with r and $(n-m)$ degrees of freedom. Then

$$r = [K(n-1) + (m-1)]/Z \quad (7.63a)$$

where

$$Z = \frac{(n-1)K(K+2) + (m-1)}{(n-1)K + (m-1)} \quad (7.63b)$$

and

$$K = \frac{R_m^2}{1 - R_m^2} \quad (7.63c)$$

For large sample sizes, the square of the multiple correlation coefficient reaches approximately a normal distribution with the mean value

$$E(\hat{R}_m^2) = R_m^2$$

and variance

$$D(\hat{R}_m^2) = \frac{4R_m^2(1 - R_m^2)^2}{n - 1}$$

The random variable

$$u_R = \frac{\sqrt{n - 1}(\hat{R}_m^2 - R_m^2)}{2R_m(1 - R_m^2)} \quad (7.64)$$

has the normalized normal distribution. Also, the Fisher and other transformations for speeding up convergence to normality can be used.

For the mean value of the squared multiple correlation coefficient

$$E(\hat{R}_m^2) = R_m^2 + \frac{m - 1}{n - 1}(1 - R_m^2) - \frac{2(n - m)}{n^2 - 1}R_m^2(1 - R_m^2) + \dots \quad (7.65)$$

The variance is given by

$$D(\hat{R}_m^2) = \frac{4R_m^2(1 - R_m^2)^2 \times (n - m)^2}{(n^2 - 1)(n + 3)} \approx \frac{4R_m^2(1 - R_m^2)^2}{n} \quad (7.66)$$

For smaller sample sizes, the estimate \hat{R}_m^2 is overestimated. The corrected multiple correlation coefficient is expressed by

$$\hat{R}_m^{*2} = \hat{R}_m^2 - \frac{m - 3}{n - m}(1 - \hat{R}_m^2) - \left[\frac{2(n - 3)}{(n - m)^2}(1 - \hat{R}_m^2) + \dots \right] \quad (7.67)$$

It can be seen that $\hat{R}_m^{*2} < \hat{R}_m^2$. For small values of \hat{R}_m^2 , the corrected \hat{R}_m^{*2} can even be negative and therefore it should be restricted to the interval $<0, 1>$.

7.2.4 Rank correlation

In some cases, the classical paired correlation coefficient can be replaced by the rank correlation coefficient, which is not very sensitive to the presence of outliers.

The rank of the i th element of a sample is equal to the index of the order statistic. Let us write the sample ranks for variable ξ_1 as x_{1si} and sample ranks regarding to the variable ξ_2 as x_{2si} . The Spearman rank correlation coefficient is then expressed by

$$\tilde{\rho}_s = 1 - \frac{6}{n(n^2 - 1)} \sum_{i=1}^n (x_{1si} - x_{2si})^2 \quad (7.68)$$

For $\rho_s = 0$, the distribution of $\hat{\rho}_s$ is symmetric with mean value $E(\hat{\rho}_s) = 0$ and variance $D(\hat{\rho}_s) = 1/(n - 1)$.

For $n > 10$, the quantity

$$t_s = \frac{|\hat{\rho}_s| \sqrt{n-2}}{\sqrt{1-\hat{\rho}_s^2}} \quad (7.69)$$

has the Student distribution, asymptotically, with $n-2$ degrees of freedom, if the theoretical coefficient $\rho_s = 0$.

Problem 7.14. *Correlation between the effective specific surface and the change of surface energy of adsorption*

For six different stearates, the effective surface ξ_1 and the change of surface energy of adsorption ξ_2 were evaluated. Estimate the Spearman correlation coefficient and examine its significance.

Data:

x_1	2.6	3.3	4.4	4.2	6.2	6.5
x_2	17.8	18.6	16.2	17.3	15.8	15.2

Solution: Table 7.1 lists the ranks of x_{1si} and x_{2si} .

Table 7.1. The order of quantities x_1 and x_2

x_{1s}	1	2	3	4	5	6
x_{2s}	5	6	3	4	2	1

From Eq. (7.68), we can calculate

$$\hat{\rho}_s = 1 - \frac{6}{6(6^2 - 1)} (4^2 + 4^2 + 1 + 1 + 3^2 + 5^2) = -0.943$$

Then substitution into Eq. (7.69) leads to

$$t_s = \frac{0.943 \sqrt{4}}{\sqrt{1 - 0.943^2}} = 5.66$$

Because the quantile $t_{0.975}(4) = 2.776$ is lower than 5.66, the null hypothesis $H_0: \rho_s = 0$ is rejected.

Conclusion: The nonparametric test used showed significant negative correlation between the effective specific surface and the change of surface energy of adsorption. For small sample sizes, the conclusion is of little consequence.

In practical problems, often several elements of a sample have the same rank. In this case, these elements have the same mean rank as if they had different values, and the Spearman correlation coefficient is then estimated from

$$\hat{\rho}_s = \frac{\frac{n(n^2-1)}{6} - \sum_{i=1}^n (x_{1si} - x_{2si})^2 - a - b}{\left\{ \left[\frac{n(n^2-1)}{6} - 2a \right] \left[\frac{n(n^2-1)}{6} - 2b \right] \right\}^{1/2}} \quad (7.70)$$

where a and b are correcting coefficients for rank, expressed by

$$a = \frac{1}{12} \sum_j (a_j^3 - a_j) \quad (7.70a)$$

$$b = \sum_k (b_k^3 - b_k) \quad (7.70b)$$

where j is the number of clusters of the same rank for x_1 and a_j is the number of values of the same rank in the j th cluster. The definitions for k and b_k are similar.

The rank correlation coefficient ρ_s lies in the interval $-1 \leq \rho_s \leq 1$. If the sample comes from a two-dimensional normal distribution and $n \geq 30$, then

$$R_{12} = \rho(\xi_1, \xi_2) = 2 \sin\left(\frac{\pi}{6} \times \rho_s\right) \quad (7.71)$$

When rank correlation coefficients are used, it should be remembered that transforming data from x_{1i} and x_{2i} into x_{1si} and x_{2si} always causes loss of information. Robustness and a decrease in sensitivity to deviations from normality are the compensation.

7.3 PROCEDURE FOR CORRELATION ANALYSIS

The procedure of correlation analysis assumes some mutual relationships (bounds) among the components of the random vector. Besides a pair correlation coefficient, a partial correlation coefficient should also be computed, to enable deeper analysis of mutual bounds. Interpretation should be made carefully, especially when the sample size is not large.

REFERENCES

- [1] V. V. Gubarev, *Algoritmy statisticeskich izmerenij*, Energoatomizdat, Moskva, 1986.
- [2] I. Nimo-Smith, *Biometrika*, 1979, **66**, 390.
- [3] Ch. J. Kovalski, *Am. Statist.*, 1973, **27**, 103.
- [4] P. Prescott, *Technometrics*, 1975, **17**, 129.
- [5] G. J. Mirskij, *Charakteristiki stochasticeskoj vzaimnosvzaji i ich izmerenije*, Energoizdat, Moskva, 1982.
- [6] M. Siotani, T. Hyakawa and Y. Fujikoshi, *Modern Multivariate Statistical Analysis*, American Science Press, 1985.
- [7] R. A. Fisher, *Metron*, 1921, **1**, 1.
- [8] M. S. Srivastava, *Commun. Statist.*, 1983, **A12**, 125.
- [9] H. Ch. Kraemer, *J. Am. Statist. Assoc.*, 1973, **68**, 1004.
- [10] S. J. Devlin, R. Gnanadesikan and J. R. Kettenring, *Biometrika*, 1975, **62**, 531.

8

Nonlinear regression models

Nonlinear models are often used in the chemical laboratory. There are three main ways in which nonlinear models are utilized in chemometrics.

- (1) Construction of *calibration models* when the measured variable y is a nonlinear function of the independent (adjustable, controllable) variable x .
- (2) Construction of *chemical models* describing the stoichiometry, concentration and equilibrium constants of all the products of chemical reactions at equilibrium, or the kinetics of chemical reactions.
- (3) Construction of *empirical models* based on a study of the nonlinear dependence between the dependent variable y and independent explanatory variables x .

According to the actual type of task, an approach to building the regression model $f(\mathbf{x}, \boldsymbol{\beta})$ is chosen. The regression model $f(\mathbf{x}, \boldsymbol{\beta})$ is a function of a vector of controllable independent variables \mathbf{x} and of a vector of unknown parameters $\boldsymbol{\beta}$ of dimension $(m \times 1)$, $\boldsymbol{\beta} = \{\beta_1, \dots, \beta_m\}^T$. Nonlinear regression considers the set of points $\{y_i, \mathbf{x}_i^T\}$, $i = 1, \dots, n$, where y represents the response (dependent) variable.

The dimension of vector \mathbf{x}_i does not affect the dimension of vector $\boldsymbol{\beta}$. The regression problem is formulated with regard to a regression triplet:

- (1) the data set,
- (2) a proposed model, and
- (3) a regression criterion.

The regression problem consists of a search for the best model $f(\mathbf{x}, \boldsymbol{\beta})$ on a basis of the data set $\{y_i, \mathbf{x}_i\}$, $i = 1, \dots, n$, such that the model sufficiently fulfils the given regression criterion.

In chemometrics, the model $f(\mathbf{x}, \boldsymbol{\beta})$ is usually known, so the regression problem consists of searching for the best estimates of unknown parameters $\boldsymbol{\beta}$. In contrast to linear regression models, the parameters $\boldsymbol{\beta}$ play a very important role in nonlinear models. In linear regression models, the regression parameters have no physical meaning but are just numerical coefficients; the parameters in a nonlinear model can

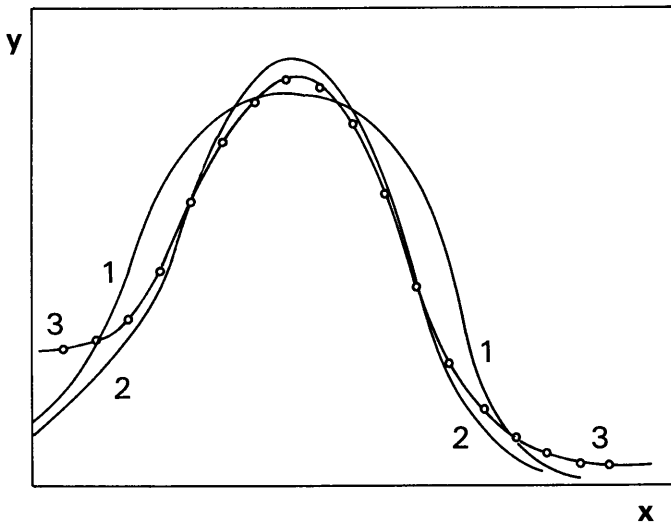


Fig. 8.1—Regression model, $y = f(x; \beta)$ at three stages of a search for unknown parameters β : (1) for an initial guess of parameters $\mathbf{b}^{(0)}$, (2) in the k -th iteration of parameters refinement $\mathbf{b}^{(k)}$, (3) for the best estimates \mathbf{b} .

have a specific physical meaning. Finding the numerical values is often the main purpose of the regression analysis. Examples are equilibrium constants (dissociation constants, stability constants, solubility products) of reactions, rate constants in kinetic models, or unknown concentrations in titration curves. In the interpretation of estimates of model parameters, it must be remembered that they are *random variables* which have variance, and which are often strongly correlated.

Problem 8.1. *Formulation of the parameters and variables of a regression model*

The dependence of a rate constant for a chemical reaction, k , on temperature T is described by the Arrhenius equation

$$k = k_0 \exp(-E/RT) \quad (8.1)$$

where k_0 is the activation entropy of the chemical reaction, E is the activation energy, and R is the universal gas constant. Formulate the model parameters and examine their correlation.

Solution: The rate constant k (response variable, y) was measured at various temperatures, T (explanatory variable, x). The unknown parameters in the regression model [Eq. (8.1)] are $\beta_1 = k_0$ and $\beta_2 = -E/R$. If the additive errors model is used

$$y_i = f(x_i, \beta) + \varepsilon_i = \beta_1 + \exp(\beta_2/x_i) + \varepsilon_i \quad (8.2)$$

Here, b_1 and b_2 are estimates of β_1 and β_2 and are determined from experimental data $\{k_i, T_i\}$, $i = 1, \dots, n$, based on a regression criterion. When errors ε_i in Eq. (8.2) are independent (values of y_i are mutually independent) random variables of the same distribution, and have constant variance $\sigma^2(y)$. The regression criterion corresponding

to the least-squares method may be used. That is

$$U(\boldsymbol{\beta}) = \sum_{i=1}^n (y_i - f(x_i, \boldsymbol{\beta}))^2 \quad (8.3)$$

The estimates \mathbf{b} then minimize the criterion $U(\boldsymbol{\beta})$.

It can be shown here that a strong correlation exists between estimates $\ln b_1$ and b_2 . This correlation may be expressed by the paired correlation coefficient

$$r = \sqrt{c^2 + c} \times \ln(1 + 1/c) \quad (8.4)$$

where $c = T_1/(T_n - T_1)$ is related to minimum (T_1) and maximum (T_n) temperatures. When $T_1 = 300\text{K}$ and $T_2 = 360\text{K}$, then $r = 0.9986$ and $\ln b_1$ and b_2 are nearly linearly dependent. This means that the ratio $(\ln b_1)/b_2$ is constant and hence, individual parameters cannot be estimated independently. When parameters are correlated, unfortunately a change in the first parameter is often compensated for by a change in the second one. There may be several different pairs of parameter estimates $(\ln b_1, b_2)$ which give nearly same values of the least squares criterion $U(\mathbf{b})$. The parameter estimates achieved by various regression programs may differ by some orders of magnitude, but nevertheless apparently a “best” fit to the experimental data, and low values of $U(\mathbf{b})$ are reached.

Conclusion: Even a simple nonlinear model may lead to difficulties in the accuracy of the parameter estimates and also in their interpretation.

Often, attempts are made to apply nonlinear regression models in situations which are totally inappropriate. Models are often applied outside the range of their validity, and generally, it is supposed that they can substitute for missing data. In chemical kinetics, for example, attempts may be made to estimate parameters, from data far from equilibrium. The calculated parameters then differ significantly from the true equilibrium parameters, and should be interpreted as model parameters only.

The result of nonlinear regression depends on the quality of the regression triplet, i.e. (1) the data, (2) the model, and (3) the regression criterion. Correct formulation leads to parameter estimates which have meaning not only formally but also physically.

In this chapter, we solve problems for which a regression model is known. For solving calibration problems or searching for empirical models, the regression model is appropriate. Some procedures are mentioned in Chapter 9.

8.1 FORMULATION OF A NONLINEAR REGRESSION MODEL

A *linear* regression model is a model which is formed by a linear combination of model parameters. This means that linear regression models can, with reference to the model functions, be nonlinear. For example, the model $f(x, \boldsymbol{\beta}) = \beta_1 + \beta_2 \times \sin x$ is sinusoidal, but with regards to parameters it is a linear model. For linear regression models, the following condition is valid

$$g_j = \frac{\delta f(x, \boldsymbol{\beta})}{\delta \beta_j} = \text{constant}, \quad j = 1, \dots, m \quad (8.5)$$

If for any parameter, β_j , the partial derivative is not a constant, we say that the regression model is nonlinear. Nonlinear regression models may be divided into the following groups:

(1) *Non-separable models*, when condition (8.5) is not valid for any parameter. For example, in the model

$$f(x, \beta) = \exp(\beta_1 x) + \exp(\beta_2 x).$$

(2) *Separable models*, when condition (8.5) is valid for one model parameter. For example, the model

$$f(x, \beta) = \beta_1 + \beta_2 \exp(\beta_3 x)$$

is nonlinear only with regards to the parameter β_3 .

(3) *Intrinsically linear models* are nonlinear, but by using a correct transformation they can be transformed into linear regression models. For example, the model

$$f(x, \beta) = \beta^2 x$$

is nonlinear in parameter β , but the shape of the model is a straight line. With the use of the reparameterization

$$\gamma = \beta^2$$

the nonlinear model is transformed into a linear one.

Reparameterization means transformation of parameters β into parameters γ which are related to the original ones by a function

$$\gamma = g(\beta) \tag{8.6}$$

By reparameterization, many numerical and statistical difficulties of regression may be avoided or removed and non-separable models transformed into separable models. The model of the Arrhenius equation (8.2) is separable, i.e. linear with regard to β_1 , and by the reparameterization, $f(x, \gamma) = \exp(\gamma_1 + \gamma_2/x)$ is transformed into a non-separable model, where $\gamma_1 = \ln \beta_1$ and $\gamma_2 = \beta_2$. Each regression model may be reparameterized in many ways, one of which is described in Section 8.5.

In chemometrics, we often distinguish models that are linearly transformable, which can, by use of an appropriate transformation, be transformed into linear regression models. For example, the Arrhenius regression model (8.2), may be transformed into the form (if random errors ε are neglected)

$$\ln y = \gamma_1 + \gamma_2 z$$

where $\gamma_1 = \ln \beta_1$, $\gamma_2 = \beta_2$ and $z = 1/x$. The resulting model is a linear model with respect to parameters γ . For finite errors ε , however, this transformation is not correct, and causes heteroscedasticity. When the measured rate constants k_i have constant variance $\sigma^2(k_i)$, then the quantities $\ln k_i$ have non-constant variance $\sigma^2(\ln k) = \sigma^2(k_i)/(k_i)^2$, i.e. constant relative error. The linear transformation is useful for simplification of the search for parameters, but it leads to biased estimates and is therefore used to find a guess for initial estimates of unknown parameters (Section

8.5). The derivatives g_j in Eq. (8.5) are *sensitivity measures* of parameter β_j in model $f(\mathbf{x}, \boldsymbol{\beta})$.

From the sensitivity measures of individual parameters, a preliminary analysis of nonlinear regression models can be made, classifying their quality and identifying any redundancy caused by an excessive number of parameters. A model should not contain excessive parameters and its parameters may be unambiguously estimated if the sensitivity measures, g_j , for given data are found to be linearly independent. This means that it is not possible to determine non-zero coefficients $v_j, j = 1, \dots, m$, such that the Eq. (8.7) is fulfilled.

$$\sum_{j=1}^m g_j v_j = 0 \quad (8.7)$$

However, if at least one non-zero coefficient, $v_j \neq 0$, exists for which Eq. (8.7) is fulfilled, the regression model is redundant and should be simplified by *excluding* some parameters. If Eq. (8.7) is valid, all parameters may not be individually estimable.

Problem 8.2. *Examination of redundancy of a regression model*

Test for redundant parameters in the regression model $f(x, \boldsymbol{\beta}) = \beta_1 \exp(\beta_2 + \beta_3 x)$. Apply the sensitivity measures, g_j .

Solution: We first compute sensitivity measures,

$$g_1 = \exp(\beta_2 + \beta_3 x)$$

$$g_2 = \beta_1 \exp(\beta_2 + \beta_3 x)$$

and

$$g_3 = \beta_1 x \exp(\beta_2 + \beta_3 x).$$

On substituting into Eq. (8.7), we get

$$(v_1 + v_2 \beta_1 + v_3 \beta_1 x) \exp(\beta_2 + \beta_3 x) = 0$$

For $v_1 = -\beta_1$, $v_2 = 1$ and $v_3 = 0$, Eq. (8.7) is fulfilled, so the model contains redundant parameters.

To confirm the redundant parameters, reparameterization of the model may be used i.e.

$$f(x, \boldsymbol{\gamma}) = \exp(\gamma_1 + \gamma_2 x)$$

or

$$f(x, \boldsymbol{\delta}) = \delta_1 \exp(\delta_2 x)$$

where $\gamma_1 = \ln \beta_1 + \beta_2$, $\gamma_2 = \beta_3$ or $\delta_1 = \beta_1 \exp(\beta_2)$ and $\delta_2 = \beta_3$.

Conclusion: Parameters β_1 and β_2 cannot be estimated separately. Only their functions γ_1, γ_2 or δ_1, δ_2 may be estimated.

Examination for redundancy in regression models should be part of the investigation of any regression model. Some models exhibit redundancy for only some combinations of parameters $\boldsymbol{\beta}$. The model cannot be simplified without knowledge of preliminary estimates of some parameters.

Problem 8.3. *Influence of the magnitude of parameters on redundancy in a regression model*

Examine the redundancy in the regression model

$$f(x, \beta) = \exp(\beta_1 x) + \exp(\beta_2 x)$$

with regard to the magnitude of its parameters β_1 and β_2 .

Solution: The sensitivity measures

$$g_1 = x \exp(\beta_1 x)$$

and

$$g_2 = x \exp(\beta_2 x)$$

are substituted into Eq. (8.7), resulting in the expression

$$x(v_1 \exp(\beta_1 x) + v_2 \exp(\beta_2 x)) = 0$$

There are no values for v_1 and v_2 that satisfy the equation, unless $\beta_1 = \beta_2$, in which case $v_1 = 1$ and $v_2 = -1$ fulfils the condition.

Conclusion: When in a search for the best estimates of parameters β_1 and β_2 , estimate b_1 is nearly equal to b_2 , the model is ill-conditioned, and if $b_1 = b_2$, the model is redundant.

There are models which exhibit local redundancy for selected points or values of the independent variable x .

Problem 8.4. *Examination of local redundancy of parameters of the Arrhenius equation*

Find the conditions under which the model of the Arrhenius equation (8.2) is redundant.

Solution: The sensitivity measures substituted into Eq. (8.7) lead to

$$(v_1 + v_2 \beta_1/x) \exp(\beta_2/x) = 0$$

For $v_1 = -\beta_1/x$ and $v_2 = 1$ this equation is satisfied. Local redundancy occurs when the value β_1 is of the same magnitude as some of the experimental quantities x_i, y_i , $i = 1, \dots, n$.

Conclusion: Redundant parameters in the Arrhenius equation occur when $\beta_1 \approx x_i$.

Redundancy always leads to singularity of matrix $\mathbf{J}^T \mathbf{J}$ (cf. Section 8.4). This means that algorithms for the inversion of this matrix by classical procedures will fail (Section 8.5). The local redundancy of parameters may be avoided by using pseudoinversion of matrix $\mathbf{J}^T \mathbf{J}$.

Ill-conditioned nonlinear models cause problems when Eq. (8.7) is fulfilled only approximately. It is analogous to multicollinearity in linear regression models. Although parameter estimates may be found when $\mathbf{J}^T \mathbf{J}$ is ill-conditioned, some numerical difficulties appear during its inversion. If we know the approximate magnitude of the parameter estimates $\mathbf{b}^{(0)}$, we may construct the matrix $\mathbf{L} = n^{-1}(\mathbf{J}^T \mathbf{J})$ with elements

$$L_{jk} = \frac{1}{n} \sum_{i=1}^n \frac{\delta f(x_{i\cdot} \mathbf{b})}{\delta b_j} \frac{\delta f(x_{i\cdot} \mathbf{b})}{\delta b_k} \bigg|_{\mathbf{b} = \mathbf{b}^{(o)}} \quad (8.8)$$

Matrix \mathbf{L} corresponds to the matrix $(1/n) \mathbf{X}^T \mathbf{X}$ for linear regression models. To estimate the ill-conditioning, matrix \mathbf{L} is transformed into the standardized form \mathbf{L}^* with elements

$$L_{ij}^* = \frac{L_{ij}}{\sqrt{L_{ii} L_{jj}}} \quad (8.9)$$

The conditioning of matrix \mathbf{L}^* gives a guide to the conditioning of parameters $\mathbf{b}^{(o)}$ in a given model for a given experimental data set.

A simple measure of ill-conditioning is the determinant of matrix \mathbf{L}^* , $\det(\mathbf{L}^*)$. When the determinant is less than 0.01, i.e. $\det(\mathbf{L}^*) < 0.01$, the nonlinear model is ill-conditioned and hence has to be simplified [1].

In many regression programs, the inversion of matrix $(\mathbf{J}^T \mathbf{J})$ involves its eigenvalues, $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$. (An indication of redundancy is the zero values of some eigenvalues.) For a measure of ill-conditioning the ratio $\lambda_p = \lambda_1 / \lambda_m$ may be used. If $\lambda_p > 900$ the regression model is ill-conditioned [2].

Problem 8.5. *Examination of ill-conditioning of the Arrhenius equation for a chemical reaction in the solid phase*

Examine the conditioning of a model of the Arrhenius equation (8.2) for a simulated data set [3] of a chemical reaction in the solid phase. Guessed values of initial estimate are $\beta_1^{(0)} = 10^7 \text{ min}^{-1}$ and $\beta_2^{(0)} = -15047$.

Data:

$k, \text{ min}^{-1}$	0.0112	0.0120	0.0325	0.0535
$T, \text{ K}$	730	750	770	790

Solution: From Eq. (8.9), the elements of matrix \mathbf{L} are estimated by

$$L_{11} = \frac{1}{4} \sum_{i=1}^4 \exp(2\beta_2^{(0)}/T_i) = 1.105458 \times 10^{-17}$$

$$L_{12} = L_{21} = \frac{\beta_1^{(0)}}{4} \sum_{i=1}^4 \frac{1}{T_i} \exp \frac{2\beta_2^{(0)}}{T_i} = 1.417623 \times 10^{-13}$$

$$L_{22} = \frac{\beta_1^{(0)2}}{4} \sum_{i=1}^4 \frac{1}{T_i^2} \exp \frac{2\beta_2^{(0)}}{T_i} = 1.818692 \times 10^{-9}$$

By using a standardization procedure and Eq. (8.9), we get

$$\mathbf{L}^* = \begin{bmatrix} 1 & 0.9979 \\ 0.9979 & 1 \end{bmatrix}$$

The determinant of this matrix, $\det(\mathbf{L}^*) = 4.1182 \times 10^{-4}$, shows significant ill-conditioning.

Conclusion: Ill-conditioning is caused by the small range of experimental temperatures.

8.2 MODELS OF MEASUREMENT ERRORS

Suppose that the experimental data $\{\mathbf{x}_i^T, y_i\}$, $i = 1, \dots, n$, and the regression model are known. The response variable y is the variable measured and subject to various kinds of errors. Common errors include measurement errors ε_M , errors of model formulation ε_T , errors of adjusting the independent controllable variable \mathbf{x} , ε_x and the random errors of the experiment ε_N . The total error, ε , of the dependent variable y , is the sum of the individual errors. It is assumed that the total error of measurement, for all values of y_i , $i = 1, \dots, n$, has a mean value equal to zero, i.e. $E(\varepsilon_i) = 0$. When $E(\varepsilon_i) = \text{constant}$, the intercept term in model is missing and when $E(\varepsilon_i) \neq \text{constant}$, the model is falsely proposed. The general regression model can be expressed in the form

$$y_i = Z_i(\mathbf{x}_i, \varepsilon_i, \boldsymbol{\beta}), \quad i = 1, \dots, n \quad (8.10)$$

where the function Z_i depends on the type of errors and on the form of the regression function.

When the data represent the results of experimental measurement, the *additive model* of measurement errors is usually assumed:

$$Z_i = f(\mathbf{x}_i, \boldsymbol{\beta}) + \varepsilon_i \quad (8.11)$$

In many experiments, there are some restrictions on the measured variable y_i , $i = 1, \dots, n$. For example, y_i may take only positive values, with non-constant variance, $\sigma^2(y_i)$, but with constant relative error, $\sigma^2(y)/y$. Such conditions are valid in the *multiplicative model* of measurement errors

$$Z_i = f(\mathbf{x}_i, \boldsymbol{\beta}) \exp(\varepsilon_i) \quad (8.12)$$

In chemical practice, the *combined model* of measurement errors

$$Z_i = f(\mathbf{x}_i, \boldsymbol{\beta}) \exp(v_i) + \varepsilon_i \quad (8.13)$$

is also used. The errors v_i and ε_i in Eq. (8.13) are assumed to be independent.

In a chemical laboratory the measurement is usually made on just one experimental system. For example, in the investigation of the equilibria of reaction products, the voltage of the glass electrode cell (or absorbance) is monitored during a titration after each addition of a volume of titrant. Cumulative errors can appear in such an experimental procedure. Instrumental measurements are often subject to a constant relative error

$$v(y) = \frac{\sigma(y)}{f(\mathbf{x}, \boldsymbol{\beta})}$$

so that the variance of measured variable y is proportional to the square of the value

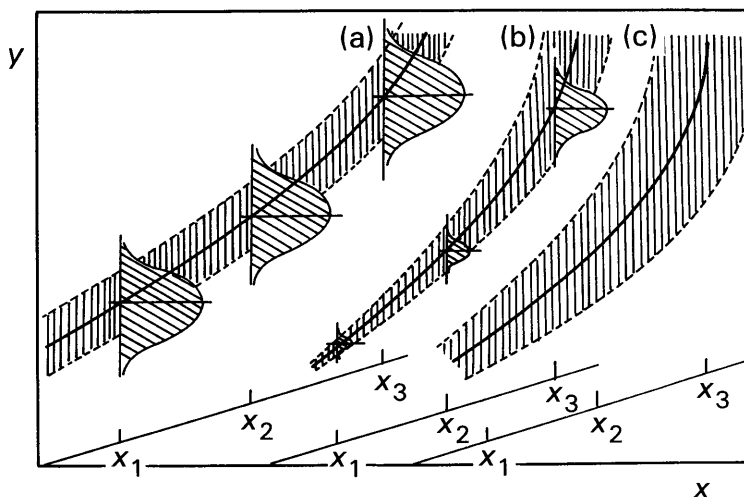


Fig. 8.2—Three models of measurement errors: (a) the additive model, (b) the multiplicative model, and (c) the combined model.

of function $f(x, \beta)$,

$$\sigma^2(y) \approx f^2(x, \beta)$$

The total error ε_i is then expressed by

$$\varepsilon_i = \sum_{j=1}^i u_j + v_i \quad (8.14)$$

where v_i represents the measurement error and u_j is the process error. Process errors are caused by fluctuations in experimental conditions such as temperature, pressure, purity of reagents, etc., and they are cumulative. The total error ε_i expressed by Eq. (8.14) is therefore additive.

In order to find a proper criterion for regression and to make a statistical analysis, the distribution of the random quantities y_i must be determined. This distribution is closely related to the distribution of errors ε_i given by the probability density function $p_\varepsilon(\varepsilon)$. This function depends on distribution parameters such as the variance σ^2 , etc.

In chemometric problems, the error distribution is assumed to be unimodal and symmetrical, with the maximum at $E(\varepsilon) = 0$. It is often assumed, but falsely, that the measurement errors ε_i are mutually independent. The point probability density function $p_\varepsilon(\varepsilon)$ is then given by the product of the *marginal* densities $p_\varepsilon(\varepsilon_i)$

$$p_\varepsilon(\varepsilon) = \prod_{i=1}^n p_\varepsilon(\varepsilon_i) \quad (8.15)$$

Several distributions, including the normal, rectangular, Laplace and trapezoidal ones may be expressed by the probability density function

$$p_\varepsilon(\varepsilon_i) = Q_N \exp(-|\varepsilon_i|^p/\alpha) \quad (8.16)$$

where Q_N is the normalizing constant and α is a parameter proportional to the variance. If $p = 1$, the resulting distribution is Laplace. When $p = 2$, the distribution is normal and when $p \rightarrow \infty$, rectangular. The disadvantage of describing distribution $p_\epsilon(\epsilon_i)$ by Eq. (8.16) is that for $p < 2$, in the neighbourhood of origin, the distribution is not locally quadratic. Therefore, alternative probability density functions, such as the generalized Student distribution, are used [4].

In some cases, the errors are not independent but are characterized by a covariance matrix of errors C_ϵ . When the errors ϵ come from a symmetric and unimodal distribution with the mean $E(\epsilon) = 0$, the probability density function is chosen from a class of elliptic distributions

$$p_\epsilon(\epsilon) = Q_N \sqrt{\det(\mathbf{B})} h \sqrt{\epsilon^T \mathbf{B} \epsilon} \quad (8.17)$$

where Q_N is the normalizing coefficient, \mathbf{B} is the covariance matrix of errors C_ϵ and $h(\cdot)$ is a positive function defined on the interval $\langle 0, \infty \rangle$ with finite moments up to $(n + 1)$. The most widely used distribution is the multivariate normal distribution, $N(0, C_\epsilon)$, which for $h(x) = \exp(-0.5x^2)$ gives the probability density function

$$p_\epsilon(\epsilon) = (2\pi)^{-n/2} (\det C_\epsilon)^{-1/2} \exp(-0.5 \epsilon^T C_\epsilon^{-1} \epsilon) \quad (8.18)$$

It is also possible to use the multivariate Laplace, the Student or other distributions [4].

The form of the covariance matrix of errors C_ϵ depends on the type of error dependence. A simple example is the case of heteroscedasticity, when errors are mutually independent but have non-constant variance $E(\epsilon_i^2) = \sigma_i^2$. The matrix C_ϵ is then diagonal with the elements σ_i^2 on a diagonal and the probability density function (8.17) is transformed into Eq. (8.15). For other types of autocorrelation, the matrix C_ϵ is not diagonal and their *off-diagonal* elements C_{ij} correspond to the covariance between ϵ_i and ϵ_j , $C_{ij} = E(\epsilon_i \epsilon_j)$.

Problem 8.6. *Covariance matrix of errors for a combination of measurement errors and process errors*

Derive the covariance matrix of errors for a case when the errors ϵ_i result from errors of measurement v_i and process errors u_j according to Eq. (8.14). Make the following assumptions:

- the process errors u_j and errors of measurement v_i are mutually independent, $E(u_j v_i) = 0$;
- the process errors u_j are independent, $E(u_j u_i) = 0$ for $j \neq i$ and have constant variance, $E(u_i^2) = \sigma^2$;
- the measurement errors v_i are independent, $E(v_i v_j) = 0$ for $j \neq i$ and have non-constant variance, $E(v_i^2) = \sigma_0^2 f^2(\mathbf{x}, \boldsymbol{\beta})$.

Solution: Equation (8.14) is rewritten as

$$\epsilon_i = \epsilon_{i-1} + u_i + v_i = \epsilon_{i-1} + w_i \quad (8.19a)$$

where w_i is the total error on changing from the $(i - 1)$ th state to the i th state. This error has zero mean $E(w_i) = 0$ but non-constant variance, $\tau_i^2 = \sigma^2 + \sigma_0^2 f(x_i, \boldsymbol{\beta})$. From

Eq. (8.19a), the individual errors ε_i may be written in the form

$$\varepsilon_1 = w_1,$$

$$\varepsilon_2 = w_1 + w_2$$

...

$$\varepsilon_n = w_1 + w_2 + \dots + w_n$$

and in matrix notation

$$\boldsymbol{\varepsilon} = \mathbf{A}\mathbf{w} \quad (8.19b)$$

where \mathbf{A} is the lower triangular matrix of ones, on and under the main diagonal

$$\mathbf{A} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 1 & 1 & 0 & \dots & 0 \\ 1 & 1 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 1 & \dots & 1 \end{bmatrix}$$

The covariance matrix of errors \mathbf{C}_ε is given by

$$\mathbf{C}_\varepsilon = E(\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T) = E(\mathbf{A} \mathbf{w} \mathbf{w}^T \mathbf{A}^T) = \mathbf{A}E(\mathbf{w}\mathbf{w}^T)\mathbf{A}^T = \mathbf{A}\mathbf{V}\mathbf{A}^T \quad (8.20)$$

where $E(\mathbf{w}\mathbf{w}^T) = \mathbf{V}$ is the covariance matrix of errors \mathbf{w} .

With the given assumptions, the errors w_i are independent so that \mathbf{V} is the diagonal matrix with elements on the diagonal $V_{ii} = \tau_i^2$. Substitution into Eq. (8.20) results in

$$\mathbf{C}_\varepsilon = \begin{bmatrix} \tau_1^2 & \tau_1^2 & \dots & \tau_1^2 \\ \tau_1^2 & \tau_1^2 + \tau_2^2 & \dots & \tau_1^2 + \tau_2^2 \\ \vdots & \vdots & \ddots & \vdots \\ \tau_1^2 & \tau_1^2 + \tau_2^2 & \dots & \sum_{i=1}^n \tau_i^2 \end{bmatrix} \quad (8.21)$$

with the general element of this covariance matrix $C_{ij} = \sum_{j=1}^k \tau_j^2$ where $k = \min(i, j)$.

Conclusion: Knowledge of error composition is important in covariance matrix building.

If we know the point probability density function of the measurement errors $p_\varepsilon(\boldsymbol{\varepsilon})$ or the marginal densities $p_\varepsilon(\varepsilon_i)$ we can determine the probability density $p(\mathbf{y})$ or $p(y_i)$ from the expression for the probability density for a function of a random variable.

In the case of independent random errors $\boldsymbol{\varepsilon}$

$$p(y_i) = p_\varepsilon[Z_i^{-1}(\mathbf{x}_i, y_i, \boldsymbol{\varepsilon})] \cdot \left| \frac{\delta Z_i^{-1}(\cdot)}{\delta y_i} \right| \quad (8.22)$$

where $Z_i^{-1}(\cdot)$ denotes the inverse of the function $Z(\cdot)$. For the additive model of measurements [Eq. (8.11)] the following function may be written

$$Z_i^{-1}(\cdot) = y_i - f(\mathbf{x}_i, \boldsymbol{\beta})$$

with the derivative

$$\left| \frac{\delta Z_i^{-1}(\cdot)}{\delta y_i} \right| = 1$$

Substitution into Eq. (8.22) gives

$$p(y_i) = p_e(y_i - f(x_i, \beta)) \quad (8.23)$$

Hence, it may be concluded that the additive model does not cause any deformations of the distribution of the measured quantities with regard to the error distribution.

In the case of the multiplicative model of measurements [Eq. (8.12)], the equation obtained is

$$Z_i^{-1}(\cdot) = \ln y_i - \ln f(x_i, \beta)$$

with the derivative

$$\left| \frac{\delta Z_i^{-1}(\cdot)}{\delta y_i} \right| = \frac{1}{y_i}$$

where only positive values of the measured variable y are allowed. Substitution into Eq. (8.22) gives

$$p(y_i) = \frac{1}{y_i} p_e(\ln y_i - \ln f(x_i, \beta)) \quad (8.24)$$

The probability density obtained does not correspond to the probability density of the errors $p_e(\cdot)$.

Problem 8.7. *Distribution of the variable y for combined errors ε*

Determine the distribution of the vector of measured variables \mathbf{y} , for a case of combined errors (8.14) assuming that the errors ε have multivariate normal distribution $N(0, \mathbf{C}_\varepsilon)$ and the additive model of measurements is valid.

Solution: According to Eq. (8.23), the probability density function $p_3(\varepsilon)$ is defined for a general covariance matrix of errors \mathbf{C}_ε in Eq. (8.18). It is necessary to evaluate $\det(\mathbf{C}_\varepsilon)$ and $\mathbf{C}_\varepsilon^{-1}$ for \mathbf{C}_ε defined by Eq. (8.21). If $\mathbf{C}_\varepsilon = \mathbf{A}\mathbf{V}\mathbf{A}^T$, then its inverse is

$$\mathbf{C}_\varepsilon^{-1} = (\mathbf{A}^{-1})^T \mathbf{V}^{-1} (\mathbf{A}^{-1}) \quad (8.25)$$

From Eq. (8.20), the matrix \mathbf{A}^{-1} is

$$\mathbf{A}^{-1} = \begin{bmatrix} 1 & 0 & \cdot & \cdot & 0 & 0 \\ -1 & 1 & \cdot & \cdot & 0 & 0 \\ 0 & -1 & \cdot & \cdot & 0 & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \cdot & \cdot & -1 & 1 \end{bmatrix}$$

The bidiagonal matrix \mathbf{A}^{-1} is a matrix with one diagonal and one underdiagonal

band. Since the matrix \mathbf{V} is diagonal, the matrix \mathbf{V}^{-1} is also diagonal with elements τ_i^{-2} on a diagonal. Substituting into Eq. (8.25) leads to

$$\mathbf{C}_e^{-1} = \begin{bmatrix} (\tau_1^{-2} + \tau_2^{-2}) & -\tau_2^{-2} & . & . & 0 & 0 \\ -\tau_2^{-2} & (\tau_2^{-2} + \tau_3^{-2}) & . & . & . & . \\ 0 & -\tau_3^{-2} & . & . & . & . \\ . & . & . & . & \tau_{n-1}^{-2} & . \\ . & . & . & . & (\tau_{n-1}^{-2} - \tau_n^{-2}) & \tau_n^{-2} \\ 0 & 0 & . & . & \tau_n^{-2} & \tau_n^{-2} \end{bmatrix} \quad (8.26)$$

This matrix is a tridiagonal matrix. Its determinant $\det(\mathbf{C}_e)$ is calculated from Eq. (8.21)

$$\det(\mathbf{C}_e) = \det(\mathbf{A}) \det(\mathbf{V}) \det(\mathbf{A}^T) = \det(\mathbf{V}) = \prod_{i=1}^n \tau_i^2$$

Conclusion: The joint probability density function of a vector \mathbf{y} is, according to Eqs. (8.18) and (8.23), given by the expression

$$p(\mathbf{y}) = (2\pi)^{-n/2} \left[\prod_{i=1}^n \tau_i^2 \right]^{-1/2} \exp[-0.5(\mathbf{y} - \mathbf{f})^T (\mathbf{A}^{-1})^T \mathbf{V}^{-1} \mathbf{A}^{-1} (\mathbf{y} - \mathbf{f})] \quad (8.27)$$

where the vector \mathbf{f} contains the elements $f(x_i, \beta)$, $i = 1, \dots, n$. The variable y also has a multivariate normal distribution with the same covariance matrix of errors as \mathbf{C}_e .

From a survey of error models, it follows from experimental conditions and assumptions about various types of errors, that the distribution of the measured variable y can be derived. In the measurements made in a chemical laboratory, most of the observed errors have the normal distribution, and follow the additive model of errors. Any differences are characterized by the covariance matrix \mathbf{C}_e , which may contain only diagonal elements, or off-diagonal elements in addition.

8.3 FORMULATION OF THE REGRESSION CRITERION

For the vector of measured values $\mathbf{y} = \{y_1, \dots, y_n\}^T$, the joint probability density function is denoted by the likelihood function $L(\boldsymbol{\theta})$. This function depends on the vector of parameters, $\boldsymbol{\theta}$, which contains the model parameters, $\boldsymbol{\beta}$, and distribution parameters, $\boldsymbol{\sigma}$. The maximum likelihood estimates of parameters, $\hat{\boldsymbol{\theta}}$, are determined by maximization of the logarithm of the function

$$\ln L(\hat{\boldsymbol{\theta}}) = \ln p(\mathbf{y}) = \sum_{i=1}^n \ln p(y_i) \quad (8.28)$$

The maximum likelihood estimates $\hat{\boldsymbol{\theta}}$ have an asymptotic variance equal to the inverse of the expected *Fisher information matrix*

$$D(\hat{\boldsymbol{\theta}}) = \mathbf{I}^{-1}(\hat{\boldsymbol{\theta}}) \quad (8.29)$$

The elements of matrix $\mathbf{I}(\boldsymbol{\theta})$ are given by

$$II_{ij} = -E \left[\frac{\delta^2 \ln L(\theta)}{\delta \theta_i \delta \theta_j} \right] \quad (8.30)$$

For practical purposes, the Fisher information matrix $\mathbf{I}(\theta)$ is replaced by the *estimated information matrix* $\hat{\mathbf{I}}(\theta)$ with elements

$$\hat{I}_{ij} = - \left[\frac{\delta^2 \ln L(\theta)}{\delta \theta_i \delta \theta_j} \right]_{\theta = \hat{\theta}} \quad (8.31)$$

The estimated information matrix can be used to construct confidence intervals more conveniently. For maximum likelihood estimates, some important properties may be derived:

- (1) The estimates $\hat{\theta}$ are asymptotically ($n \rightarrow \infty$) unbiased. Therefore the bias

$$\mathbf{h} = \theta - E(\hat{\theta}) = 0 \quad (8.32)$$

is the zero vector. For a finite sample size n , the estimates $\hat{\theta}$ are biased and the magnitude \mathbf{h} depends on the degree of non-linearity of the regression model.

(2) The estimates $\hat{\theta}$ are asymptotically efficient and the variance estimates are minimal of all unbiased estimates. The covariance matrix $D(\hat{\theta})$ lies on the lower limit of the Cramer–Rao inequality [5]. For finite samples, this property is generally not fulfilled.

(3) The random vector $\sqrt{n}(\hat{\theta} - \theta)$ has, asymptotically, the normal distribution $N(0, \mathbf{I}^{-1})$ with zero mean and variance equal to the inverse of the Fisher information matrix. When the error distribution is approximately normal, the normality of estimates is valid for finite samples.

For sufficiently large sample sizes, many interesting properties of the estimates $\hat{\theta}$ may be used. For finite sample sizes, some difficulties arise from the bias estimates $\hat{\theta}$. If the probability density function $p(\mathbf{y})$ is known, the maximum likelihood estimates or a criterion for their determination (the *regression criterion*) may be found.

Problem 8.8. Regression criterion for additive errors

Derive the regression criterion for the case when measurements errors are independent, with zero mean, constant variance, and the normal distribution $N(0, \sigma^2 \mathbf{E})$; and with the assumption that the additive model of measurement errors (8.11) is valid.

Solution: Let $f_i = f(x_i; \beta)$ and $\theta^T = (\beta^T, \sigma^2)$. If the distribution of measured variable y_i is normal, $N(f_i, \sigma^2)$, then

$$p(y_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \times \exp \left[-\frac{(y_i - f_i)^2}{2\sigma^2} \right] \quad (8.33)$$

The logarithm of the likelihood function, $\ln L(\theta)$, has the form

$$\ln L(\theta) = \sum_{i=1}^n \ln p(y_i) = \frac{-n}{2} \ln(2\pi\sigma^2) - \frac{2}{(2\sigma)^2} \times U(\beta) \quad (8.34)$$

where $U(\beta)$ is the *least-squares criterion* or the residual sum of squares of deviations

defined as

$$U(\boldsymbol{\beta}) = \sum_{i=1}^n (y_i - f_i)^2 \quad (8.35)$$

Analytical maximization of $\ln L(\boldsymbol{\theta})$ according to σ^2 leads to

$$\frac{\delta \ln L(\boldsymbol{\theta})}{\delta \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} U(\boldsymbol{\beta}) = 0 \quad (8.36a)$$

and therefore

$$\hat{\sigma}^2 = \frac{U(\boldsymbol{\beta})}{n} \quad (8.36b)$$

The estimate $\hat{\sigma}^2$ is biased for a small number of measurements. The unbiased form is

$$\hat{\sigma}^2 = \frac{U(\boldsymbol{\beta})}{n - m} \quad (8.37)$$

On substituting from Eq. (8.36b) into Eq. (8.34), the concentrated likelihood function $\ln L(\boldsymbol{\beta})$ is formulated as

$$\ln L(\boldsymbol{\beta}) = -\frac{n}{2}(1 + \ln(2\pi)) - 0.5 \ln U(\boldsymbol{\beta}) \quad (8.38)$$

The maximum of $\ln L(\boldsymbol{\beta})$ corresponds to the minimum of the regression criterion $U(\boldsymbol{\beta})$ which is, in fact, a condition for the least-squares method (LS). That is, the method of maximum likelihood is identical to the least-squares method.

On the basis of Eq. (8.29), the covariance matrix of estimates $D(\hat{\boldsymbol{\theta}})$ is given by

$$D(\hat{\boldsymbol{\theta}}) = \begin{bmatrix} \sigma^2(\mathbf{J}^T \mathbf{J})^{-1} & 0 \\ 0 & 2\sigma^4/n \end{bmatrix} \quad (8.39)$$

where \mathbf{J} is the $(n \times m)$ Jacobi matrix of the first derivatives of the model, with elements

$$J_{ij} = \frac{\delta f(x_i, \boldsymbol{\beta})}{\delta \beta_j} \quad (8.40)$$

From Eq. (8.39), it follows that the estimates $\hat{\sigma}^2$ and \mathbf{b} are independent and the parameter covariance matrix is $D(\mathbf{b}) = \sigma^2(\mathbf{J}^T \mathbf{J})^{-1}$.

Conclusion: The maximum likelihood method enables either the formulation of a regression criterion or the determination of the covariance matrix of estimates. It may be concluded that with the use of the properties of the maximum likelihood method, we can simplify the construction of the confidence intervals and carry out statistical hypothesis testing.

Maximization of the likelihood function leads to the problem of *nonlinear optimization*. When the covariance matrix of errors \mathbf{C}_ε is known, we can for the additive model of measurement errors and normal error distribution, $\varepsilon \approx N(0, \mathbf{C}_\varepsilon)$,

find the maximum likelihood estimates \mathbf{b} of parameters $\boldsymbol{\beta}$ by minimizing the criterion of the generalized least-squares

$$U(\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{f})^T \mathbf{C}_\varepsilon^{-1} (\mathbf{y} - \mathbf{f}) = \text{Tr}[\mathbf{C}_\varepsilon^{-1} \hat{\mathbf{e}} \hat{\mathbf{e}}^T] \quad (8.41a)$$

where $\hat{\mathbf{e}} = \mathbf{y} - \mathbf{f}$ is the deviation and the symbol $\text{Tr}(\mathbf{A})$ denotes the trace of matrix \mathbf{A} . If the matrix \mathbf{C}_ε is diagonal, the situation is much simpler. The least-squares criterion [Eq. (8.41a)] transforms into the relationship

$$U(\boldsymbol{\beta}) = \sum_{i=1}^n w_i^2 (y_i - f_i)^2 = \sum_{i=1}^n w_i^2 (y_i - f(x_i, \boldsymbol{\beta}))^2 \quad (8.41b)$$

where $w_i^2 = 1/C_{ii}$ is the weight equal to the reciprocal value of the elements of covariance matrix. If the variables $y_i^* = w_i y_i$ and $f_i^* = w_i f(x_i, \boldsymbol{\beta})$ are introduced, Eq. (8.41b) takes the form of the classical least-squares method

$$U(\boldsymbol{\beta}) = \sum_{i=1}^n (y_i^* - f_i^*)^2 \quad (8.41c)$$

When the weights are known, a weighted least-squares problem can be converted into a classical least-squares problem with modified variables. This procedure can also be used for the unknown matrix \mathbf{C}_ε when its elements are estimated separately, for example, because of heteroscedasticity.

For an unknown matrix \mathbf{C}_ε , the technique of consecutive maximization is used. First, the estimate \mathbf{C}_ε is computed and substituted into the likelihood function. The resulting concentrated likelihood function contains only the parameters $\boldsymbol{\beta}$. Bard [6] derived the following derivatives

$$\frac{\delta \ln \det(\mathbf{C}_\varepsilon)}{\delta \mathbf{C}_\varepsilon} = (\mathbf{C}_\varepsilon)^{-1} \quad (8.42a)$$

and

$$\frac{\delta \text{Tr}[\mathbf{C}_\varepsilon^{-1} \times \hat{\mathbf{e}} \hat{\mathbf{e}}^T]}{\delta \mathbf{C}_\varepsilon} = -(\mathbf{C}_\varepsilon^T)^{-1} \hat{\mathbf{e}} \hat{\mathbf{e}}^T (\mathbf{C}_\varepsilon^T)^{-1} \quad (8.42b)$$

The derivative of the likelihood function is

$$\ln L = -\frac{n}{2} \ln 2\pi - \frac{1}{2} \ln \det(\mathbf{C}_\varepsilon) - \frac{1}{2} \hat{\mathbf{e}}^T \mathbf{C}_\varepsilon^{-1} \hat{\mathbf{e}} \quad (8.43a)$$

and since

$$\mathbf{C}_\varepsilon = \hat{\mathbf{e}} \hat{\mathbf{e}}^T \quad (8.43b)$$

then

$$\text{Tr}[\hat{\mathbf{e}} \hat{\mathbf{e}}^T]^{-1} = \text{Tr}(\mathbf{E}) = n \quad (8.43c)$$

and the concentrated likelihood function has the form

$$\ln L(\boldsymbol{\beta}) = -\frac{n}{2}[\ln 2\pi + 1] - \frac{1}{2} \ln \det(\hat{\mathbf{e}} \hat{\mathbf{e}}^T) \quad (8.44)$$

Maximization of the function $\ln L(\boldsymbol{\beta})$ is the same as minimization of the criterion $U_D(\boldsymbol{\beta})$ where

$$U_D(\boldsymbol{\beta}) = \det(\mathbf{e} \mathbf{e}^T) = \det[(\mathbf{y} - \mathbf{f})(\mathbf{y} - \mathbf{f})^T] \quad (8.45a)$$

When matrix \mathbf{C}_e is diagonal, matrix $\hat{\mathbf{C}}_e$ is also diagonal and Eq. (8.45a) converts into the form

$$U_D(\boldsymbol{\beta}) = \sum_{i=1}^n [y_i - f(x_i, \boldsymbol{\beta})]^2 \quad (8.45b)$$

Hence, in the case of heteroscedasticity, application of the classical least-squares method leads to unbiased estimates \mathbf{b} but the estimates of the covariance matrix $\hat{\mathbf{C}}_e$ are biased.

Problem 8.9. *Regression criterion for combined errors*

Derive a regression criterion for a case of combined errors [Eq. (8.14)] assuming that errors $\boldsymbol{\varepsilon}$ have the multivariate normal distribution and that

- (1) the measurement errors v_i are negligible in comparison to the process errors u_j ;
- (2) the process errors u_j are negligible with respect to the errors of measurement v_i .

Solution: The joint probability density function $p(\mathbf{y})$ is expressed by Eq. (8.27). The logarithm of the likelihood function may be expressed as

$$\ln L(\boldsymbol{\theta}) = \frac{n}{2} \ln(2\pi) - \frac{1}{2} \sum_{i=1}^n \ln \tau_i^2 - \frac{1}{2} (\mathbf{y} - \mathbf{f})^T (\mathbf{A}^{-1})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{f}) \quad (8.46)$$

The last term in this equation can be expressed as

$$\begin{aligned} U_c(\boldsymbol{\beta}) &= \sum_{i=1}^n \tau_i^{-2} [(y_i - y_{i-1}) - (f(x_i, \boldsymbol{\beta}) - f(x_{i-1}, \boldsymbol{\beta}))]^2 \\ &= \sum_{i=1}^n \tau_i^{-2} [L_i - K_i(\boldsymbol{\beta})]^2 \end{aligned} \quad (8.47)$$

where $y_0 = 0$ and $f(x_0, \boldsymbol{\beta}) = 0$. Moreover, $L_i = y_i - y_{i-1}$ and $K_i(\boldsymbol{\beta}) = f(x_i, \boldsymbol{\beta}) \times f(x_{i-1}, \boldsymbol{\beta})$. Equation (8.47) corresponds to the weighted least-squares method for first differences.

The maximum of the function $\ln L(\boldsymbol{\theta})$ generally corresponds to a minimum of

$$Q = \sum_{i=1}^n \ln[\sigma^2 + \sigma_o^2 \times f(x_i, \boldsymbol{\beta})] + U_c(\boldsymbol{\beta}) \quad (8.48)$$

Maximization of Q in terms of σ^2 , σ_o^2 and $\boldsymbol{\beta}$ may be achieved by general minimization methods:

- (1) Small measurement errors. For $\sigma^2 \gg \sigma_o^2 \times f(x_i, \boldsymbol{\beta})$, $\tau_i^2 = \sigma^2 = \text{constant}$ and we have a *model of pure process errors*. On substituting into the likelihood function (8.46)

and differentiating with respect to σ^2 , we get

$$\frac{\delta \ln L(\boldsymbol{\theta})}{\delta \sigma^2} = -\frac{0.5n}{\sigma^2} + \frac{0.5}{\sigma^4} \sum_{i=1}^n [L_i - K_i(\boldsymbol{\beta})]^2 = 0 \quad (8.49)$$

On rearrangement, we find

$$\hat{\sigma}^2 = \frac{U_p(\boldsymbol{\beta})}{n} \quad (8.50)$$

Here

$$U_p(\boldsymbol{\beta}) = \sum_{i=1}^n [L_i - K_i(\boldsymbol{\beta})]^2 \quad (8.51)$$

is the regression criterion for the minimum of the sum of the squared first differences. If estimate $\hat{\sigma}^2$ [Eq. (8.50)] is substituted into the likelihood function (8.46), the concentrated likelihood function is obtained

$$\ln L(\boldsymbol{\beta}) = -\frac{n}{2}[1 + \ln 2\pi] - \frac{1}{2} \ln U_p(\boldsymbol{\beta}) \quad (8.52)$$

The maximum of the function $\ln L(\boldsymbol{\beta})$ corresponds to a minimum of $U_p(\boldsymbol{\beta})$. The function $U_p(\boldsymbol{\beta})$ [Eq. (8.51)] may be minimized by many nonlinear regression programs, after a simple rearrangement of y_i and $f(x_i, \boldsymbol{\beta})$ into variables L_i and $K_i(\boldsymbol{\beta})$.

When the process errors (fluctuations of the system) are small in comparison to errors of measurement, we have $\sigma^2 \ll \sigma_o^2 \times f(x_i, \boldsymbol{\beta})$, and the variance $\tau^2 \approx \sigma_o^2 \times f^2(x_i, \boldsymbol{\beta})$. On substituting into the likelihood function [Eq. (8.46)] and differentiating with respect to σ_o^2 , we get

$$\frac{\delta \ln L(\boldsymbol{\theta})}{\delta \sigma_o^2} = -\frac{n}{2\sigma_o^2} + \frac{0.5}{\sigma_o^4} \sum_{i=1}^n f^{-2}(x_i, \boldsymbol{\beta})(L_i - K_i(\boldsymbol{\beta}))^2 = 0 \quad (8.53)$$

On rearrangement, we get

$$\hat{\sigma}_o^2 = \frac{1}{n} \sum_{i=1}^n f^{-2}(x_i, \boldsymbol{\beta})(L_i - K_i(\boldsymbol{\beta}))^2 \quad (8.53b)$$

and on substituting $\hat{\sigma}_o^2$ into the likelihood function [Eq. (8.46)], the concentrated likelihood function becomes

$$\ln L(\boldsymbol{\beta}) = -\frac{n}{2}[1 + \ln n + \ln 2\pi] - \sum_{i=1}^n \ln f(x_i, \boldsymbol{\beta}) - \frac{1}{2} \ln U_w(\boldsymbol{\beta}) \quad (8.54)$$

where $U_w(\boldsymbol{\beta})$ is the criterion of the weighted least-squares method for first differences

$$U_w(\boldsymbol{\beta}) = \sum_{i=1}^n \left[\frac{L_i - K_i(\boldsymbol{\beta})}{f(x_i, \boldsymbol{\beta})} \right]^2 \quad (8.55)$$

The maximum of the function $\ln L(\boldsymbol{\beta})$ [Eq. (8.54)] corresponds to the minimum of function $U_c(\boldsymbol{\beta})$

$$U_c(\boldsymbol{\beta}) = U_w(\boldsymbol{\beta}) \prod_{i=1}^n f^2(x_i, \boldsymbol{\beta}) \quad (8.56)$$

The criterion $U_c(\boldsymbol{\beta})$ may be minimized either by general algorithms or by the iterative method of weighted least-squares.

It is obvious that this situation is more complicated than the case of small measurement errors. When the accumulated process errors are negligible then $\varepsilon_i = v_i$ in Eq. (8.44) and the problem of heteroscedasticity has to be solved. The corresponding regression criterion will have the form (8.56) but the function $U_w(\boldsymbol{\beta})$ is expressed in terms of variables y_i , $f(x_i, \boldsymbol{\beta})$ instead of L_i , $K_i(\boldsymbol{\beta})$.

Conclusion: When the errors are complicated, a suitable regression criterion may be derived by using the maximum likelihood method.

8.4 GEOMETRY OF NONLINEAR REGRESSION

Although some chemometric problems lead to criteria different from the classical least-squares method (LS), the LS method is still the most commonly used method in chemical practice. In Problem 8.8, it was shown that the LS method is really a special case of the maximum likelihood method for an additive model of measurement errors and the normal distribution of independent errors ε with zero mean and constant variance. For the purpose of geometric interpretation, the least-squares criterion $U(\boldsymbol{\beta})$ in Eq. (8.35) is rewritten in vector notation as

$$U(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{f}\|^2 \quad (8.57)$$

where $\mathbf{y} = (y_1, \dots, y_n)^T$, $\mathbf{f} = (f(x_1, \boldsymbol{\beta}), \dots, f(x_n, \boldsymbol{\beta}))^T$ and the symbol $\|\mathbf{x}\| = \sqrt{\mathbf{x}^T \mathbf{x}}$ means the Euclidean norm.

Examination of the shape of the criterion function $U(\boldsymbol{\beta})$ in the space of the estimators helps to explain why the search for the function minimum is so difficult. In this $(m+1)$ -dimensional space, values of criterion $U(\boldsymbol{\beta})$ are plotted against the parameters β_1, \dots, β_m .

For linear regression models, the criterion function $U(\boldsymbol{\beta})$ is an elliptic hyperparaboloid with its centre at $[\mathbf{b}, U(\mathbf{b})]$, the place where $U(\mathbf{b})$ reaches a minimum (the "pit point"). For linear models, the criterion function $U(\boldsymbol{\beta})$ has a quadratic form of type $\boldsymbol{\beta}^T (\mathbf{X}^T \mathbf{X}) \boldsymbol{\beta}$, and the matrix $\mathbf{X}^T \mathbf{X}$ is positive-definite.

In some cases, the parameter space is used for interpretation of $U(\boldsymbol{\beta})$. Parameter space is an m -dimensional space with the components of vector $\boldsymbol{\beta}$ on the axes. The value $U(\boldsymbol{\beta})$ is a perpendicular projection of an $(m+1)$ -dimensional object into this m -dimensional space. For the two estimated parameters β_1 and β_2 where $m=2$, the criterion function $U(\boldsymbol{\beta})$ for a linear model is drawn in the $(m+1)$ -dimensional (i.e. 3-dimensional) space, in Fig. 8.3.

Rather complicated shapes can occur in nonlinear models, as a result of the nonlinear function $f(\mathbf{x}, \boldsymbol{\beta})$; there may be a number of extremes and saddle points. Figure 8.4 is an illustration of a criterion function $U(\boldsymbol{\beta})$ with two minima and one saddle point.

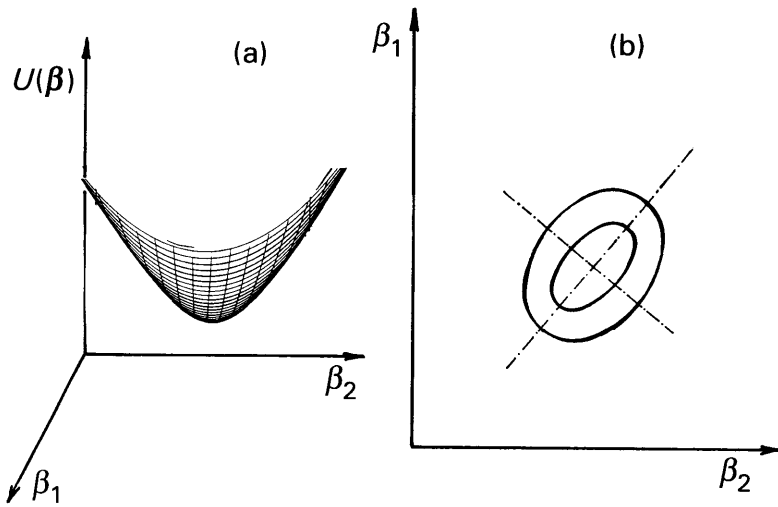


Fig. 8.3—Interpretation of the criterion function $U(\beta)$ for a linear model ($m = 2$): (a) the elliptic hyperparaboloid in the $(m + 1)$ -dimensional space of estimators, (b) the concentric ellipses as contours in the m -dimensional parameter space.

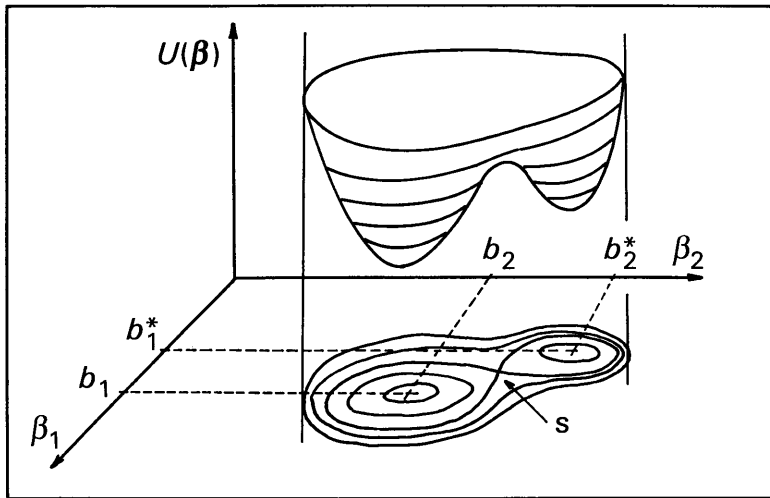


Fig. 8.4—A criterion function $U(\beta)$ with two local minima β and β^* and one saddle point S .

Quantitative information on the local behaviour of the criterion function $U(\beta)$ in the vicinity of any point β_j may be obtained from a Taylor series expansion up to quadratic terms:

$$U(\beta) = U(\beta_j) + \Delta\beta_j^T \mathbf{g}_j + \frac{1}{2} \Delta\beta_j^T \mathbf{H}_j \Delta\beta_j \quad (8.58)$$

where $\Delta\beta_j = \beta - \beta_j$ and \mathbf{g}_j is the gradient vector of a criterion function containing the components

$$\mathbf{g}_k = \frac{\delta U(\beta)}{\delta \beta_k}, \quad k = 1, \dots, m \quad (8.59a)$$

The matrix \mathbf{H}_j of dimension $(m \times m)$, is the symmetric Hessian matrix defined by the second derivative of the criterion function $U(\beta)$ with components

$$H_{lk} = \frac{\delta^2 U(\beta)}{\delta \beta_l \delta \beta_k}, \quad l, k = 1, \dots, m \quad (8.60)$$

Equation (8.59a) is valid for any criterion function $U(\beta)$. In the least-squares method the gradient of the criterion function $U(\beta)$ from Eq. (8.57) has the form

$$\mathbf{g}_j = -2\mathbf{J}^T \mathbf{e} \quad (8.59b)$$

where \mathbf{e} is the difference vector with elements

$$e_i = y_i - f(x_i, \beta), \quad i = 1, \dots, n.$$

The matrix \mathbf{J} of dimension $(n \times m)$ is called the Jacobian matrix with elements corresponding to the first derivative of the regression model in terms of the individual parameters at given points. These elements have the form

$$J_{ik} = \frac{\delta f(x_i, \beta)}{\delta \beta_k}, \quad \begin{matrix} i = 1, \dots, n \\ k = 1, \dots, m \end{matrix} \quad (8.61)$$

With the least-squares method, a similar relationship involving the Hessian matrix may be derived:

$$\mathbf{H}_j = 2[\mathbf{J}^T \mathbf{J} + \mathbf{B}] \quad (8.62)$$

where \mathbf{B} is a matrix containing the second derivatives of the regression function with elements

$$B_{kj} = \sum_{i=1}^n e_i \frac{\delta^2 f(x_i, \beta)}{\delta \beta_k \delta \beta_j}, \quad k, j = 1, \dots, m \quad (8.63)$$

In the vicinity of local minima \mathbf{b} , the gradient \mathbf{g} is approximately equal to zero. This means that

- (1) the error vector $\hat{\mathbf{e}}$ is perpendicular to the columns of a matrix \mathbf{J} in m -dimensional space;
- (2) the criterion function $U(\beta)$ is proportional to the quadratic form $\Delta\beta_i^T \mathbf{H}_i \Delta\beta_i$.

The type of local extreme is distinguished by a matrix \mathbf{H} . When the matrix \mathbf{H} is

- (a) positive-definite, the extreme is a minimum and $U(\beta)$ approximates to an elliptic hyperparaboloid;
- (b) negative-definite, the extreme is a maximum;
- (c) indefinite, no extreme is present.

Definiteness of matrices is examined by the Sylvester conditions. For practical

calculation, it is necessary that the positive-definite matrix is regular, has rank m and has eigenvalues that are all positive. It is useful to compare the Taylor series expansion of criterion function $U(\boldsymbol{\beta})$ by Eq. (8.58) with the criterion function $U(\boldsymbol{\beta})$ into which the Taylor series expansion of a model function is substituted. With the use of Taylor series expansion, the function $f(x, \boldsymbol{\beta})$ in the vicinity of the point $\boldsymbol{\beta}_j$ may be approximated by

$$f(x_i, \boldsymbol{\beta}) = f(x_i, \boldsymbol{\beta}_j) + \mathbf{J}_i(\boldsymbol{\beta} - \boldsymbol{\beta}_j) + \frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\beta}_j)^T \mathbf{G}_i(\boldsymbol{\beta} - \boldsymbol{\beta}_j) \quad (8.64)$$

where \mathbf{G}_i is the matrix of second derivatives of a model function $f(x_i, \boldsymbol{\beta})$ with elements

$$G_{jk} = \frac{\delta^2 f(x_i, \boldsymbol{\beta})}{\delta \beta_j \delta \beta_k}, \quad j, k = 1, \dots, m \quad (8.65)$$

and \mathbf{J}_i is the i th row of the matrix \mathbf{J} . Generally, a vector \mathbf{f} may be approximated by the Taylor series expansion into quadratic terms

$$\mathbf{f} \approx \mathbf{f}(\boldsymbol{\beta}_j) + \mathbf{J}\Delta\boldsymbol{\beta}_j + \frac{1}{2}\Delta\boldsymbol{\beta}_j^T \vec{\mathbf{G}}_i \Delta\boldsymbol{\beta}_j \quad (8.66)$$

where $\vec{\mathbf{G}}_i$ is an $(n \times m \times n)$ -dimensional array with layers formed by the matrices \mathbf{G}_i .

Usually, a linearization of the function \mathbf{f} is used:

$$\mathbf{f} \approx \mathbf{f}(\boldsymbol{\beta}_j) + \mathbf{J}\Delta\boldsymbol{\beta}_j \quad (8.66a)$$

Substituting Eq. (8.66a) into (8.57) we get the criterion for the "linearly-transformed" least-squares method

$$U_L(\boldsymbol{\beta}) = \mathbf{e}^T \mathbf{e} - 2\Delta\boldsymbol{\beta} \mathbf{J}^T \mathbf{e} + \Delta\boldsymbol{\beta}^T (\mathbf{J}^T \mathbf{J}) \Delta\boldsymbol{\beta} \quad (8.67)$$

The first term of this equation is equal to $U(\boldsymbol{\beta}_j)$

$$\mathbf{e}^T \mathbf{e} = U(\boldsymbol{\beta}_j) \quad (8.67a)$$

and the second one to $\Delta\boldsymbol{\beta}_j^T \mathbf{g}$,

$$-2\Delta\boldsymbol{\beta} \mathbf{J}^T \mathbf{e} = \Delta\boldsymbol{\beta}_j^T \mathbf{g} \quad (8.67b)$$

Equation (8.67) differs from (8.58) only in the third term containing matrix $2\mathbf{J}^T \mathbf{J}$ instead of matrix \mathbf{H} . It follows from Eqs. (8.62) and (8.63) that for small error values e_i , the matrix \mathbf{B} may be neglected, making Eqs. (8.67) and (8.58) identical. This means that the linearization of the regression model corresponds to the Taylor series expansion of the criterion function $U(\boldsymbol{\beta})$ into quadratic terms, assuming that matrix \mathbf{B} is negligible. From Eq. (8.67), it also follows that for nonlinear regression models the matrix $\mathbf{J}^T \mathbf{J}$ corresponds to $\mathbf{X}^T \mathbf{X}$ in linear models. If the linearization (8.67) is sufficiently precise, the statistical analysis may be performed in a similar way to that used for linear regression models.

(A) *T. geometry of linear least-squares*

For the interpretation of the geometry of linear regression, n -dimensional sample space is used. The vector of observations $\mathbf{y} = (y_1, \dots, y_n)^T$ defines a line $\overline{\mathbf{OY}}$ from

the origin O to the point Y with co-ordinates (y_1, \dots, y_n) . The X matrix has m column vectors \mathbf{x}_i , $i = 1, \dots, m$, each containing n elements. The elements of the j th column define the co-ordinates $(x_{j1}, x_{j2}, \dots, x_{jn})$ of a point \mathbf{X}_j in the sample space, and the j th column vector of matrix X defines the vector $\overrightarrow{OX_j}$ in sample space. The m vectors $\overrightarrow{OX_1}, \overrightarrow{OX_2}, \dots, \overrightarrow{OX_m}$ define a subspace of m dimensions called the estimation space which is contained within the sample space. Any point in this subspace can be represented by the termination of a vector which is a linear combination of the vectors defining the space — that is, a linear combination of the columns of X, such as $\mathbf{X}\boldsymbol{\beta}$ where $\boldsymbol{\beta} = (\beta_1, \dots, \beta_m)$ is an $m \times 1$ vector. Suppose the vector $\mathbf{X}\boldsymbol{\beta}$ defines the point T. Then the squared distance YT^2 is given by

$$(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = U(\boldsymbol{\beta})$$

as defined in Chapter 6. Thus the sum of squares $U(\boldsymbol{\beta})$ represents, in the sample space, the squared distance of Y from a general point T in the estimation space. Minimization of $U(\boldsymbol{\beta})$ with respect to $\boldsymbol{\beta}$ implies finding that value of $\boldsymbol{\beta}$, say \mathbf{b} , which provides a point P (defined by the vector $\hat{\mathbf{y}} = \mathbf{X}\mathbf{b}$) in the estimation space closest to the point Y. Then, geometrically, P must be the foot of the perpendicular from Y to the estimation space, that is, the foot of a line passing through Y and the orthogonal to all the columns of matrix X. In terms of the vectors from the origin, we can write

$$\begin{aligned}\mathbf{y} &= \hat{\mathbf{y}} + (\mathbf{y} - \hat{\mathbf{y}}) \\ &= \mathbf{y} + \hat{\mathbf{e}}\end{aligned}$$

where \mathbf{e} is the vector of residuals. The vector \mathbf{y} is thus divided into two orthogonal components:

- (1) $\hat{\mathbf{y}}$, which lies entirely in the estimation space, and
- (2) $\hat{\mathbf{e}}$, the vector of residuals, which lies in the *residual space*. The residual space is defined as the $(n \times m)$ -dimensional subspace, which is the remainder of the full n -dimensional space, after the m -dimensional estimation space has been defined. The estimation and residual spaces are thus orthogonal.

If T is a general point in the estimation space and YP is orthogonal to the space, then

$$YT^2 = YP^2 + PT^2$$

or

$$U(\boldsymbol{\beta}) = U(\mathbf{b}) + PT^2$$

Thus, the contours for which $U(\boldsymbol{\beta}) = \text{constant}$ must be such that

$$PT^2 = U(\boldsymbol{\beta}) - U(\mathbf{b}) = \text{constant}.$$

In the sample space, then, the contours defined by $U(\boldsymbol{\beta}) = \text{constant}$ consist of all points T such that $PT^2 = \text{constant}$; that is, points in the estimation space with the form $\mathbf{X}\boldsymbol{\beta}$ which lie on an m -dimensional sphere centered at the point P defined by $\mathbf{X}\mathbf{b}$. The radius of this sphere is $\sqrt{U(\boldsymbol{\beta}) - U(\mathbf{b})}$.

In order to illustrate the geometry of the linear least-squares, we will look at a sample space with $n = 3$ (Fig. 8.5a) and $m = 2$. That is, there are three components (y_1, y_2, y_3) of the vector y , two parameters β_1 and β_2 of the parametric vector β and a three by two matrix X of the form

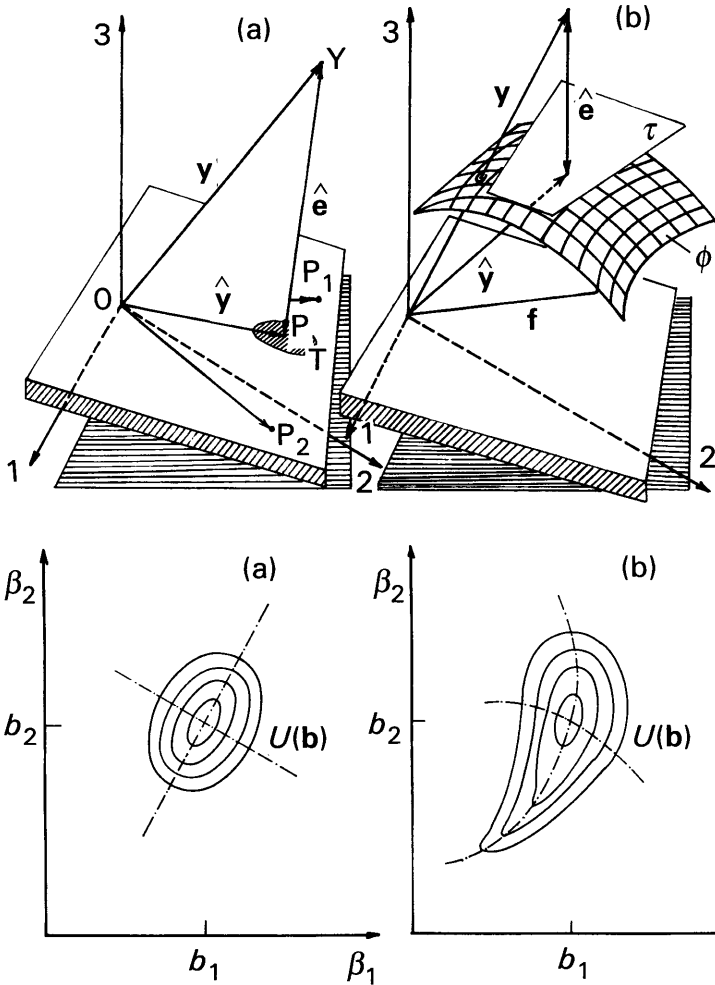


Fig. 8.5—Geometrical representation of (a) linear least-squares, and (b) nonlinear least-squares. Upper diagram shows the sample space and the lower diagram illustrates the parameter space.

$$X = \begin{bmatrix} x_{11} & x_{21} \\ x_{12} & x_{22} \\ x_{13} & x_{23} \end{bmatrix}$$

The columns of X define two points P_1 and P_2 with co-ordinates (x_{11}, x_{12}, x_{13}) and (x_{21}, x_{22}, x_{23}) , respectively, and the vectors \vec{OP}_1 and \vec{OP}_2 define a plane which

represents the 2-dimensional estimation space in which vector $\hat{\mathbf{y}} = \mathbf{X}\mathbf{b}$ must lie. The point \mathbf{Y} lies above this plane and the perpendicular \mathbf{YP} from \mathbf{Y} to the plane $\mathbf{OP}_1\mathbf{P}_2$ meets the plane at \mathbf{P} . Thus \mathbf{YP} is the shortest distance from \mathbf{Y} to any point in the estimation space, \mathbf{P} , and is defined by $\hat{\mathbf{y}} = \mathbf{X}\mathbf{b}$ and $U(\mathbf{b}) = \mathbf{YP}^2$. Then, from Pythagoras's theorem:

$$\mathbf{OY}^2 = \mathbf{OP}^2 + \mathbf{YP}^2$$

If $\mathbf{y} = \hat{\mathbf{y}} + \mathbf{e}$, then the geometrical vector equation is $\overrightarrow{\mathbf{OY}} = \overrightarrow{\mathbf{OP}} + \overrightarrow{\mathbf{YP}}$. We recall that, in general, contours with constant $U(\boldsymbol{\beta})$ are represented by m -dimensional spheres in the estimation space. Here, the contours must be circles on the plane $\mathbf{OP}_1\mathbf{P}_2$. It is evident that if \mathbf{T} is a general point $\mathbf{X}\mathbf{b}$ on the plane, $U(\boldsymbol{\beta}) = \text{constant}$ means that $\mathbf{YT}^2 = \text{constant}$, so that $\mathbf{PT}^2 = \mathbf{YT}^2 - \mathbf{YP}^2 = \text{constant}$. Hence, we obtain circles about \mathbf{P} as shown in Fig. 8.5.

The parameter space is an m -dimensional space in which a set of parameter values $(\beta_1, \dots, \beta_m)$ defines a point. The minimum value of $U(\boldsymbol{\beta})$ is attained at the point $\mathbf{b} = (b_1, \dots, b_m)$. We recall that

$$U(\boldsymbol{\beta}) - U(\mathbf{b}) = (\boldsymbol{\beta} - \mathbf{b})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \mathbf{b})$$

All values of $\boldsymbol{\beta}$ which satisfy $U(\boldsymbol{\beta}) = \text{constant} = K$ are given by

$$(\boldsymbol{\beta} - \mathbf{b})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \mathbf{b}) = K - U(\mathbf{b})$$

It can be shown that this is the equation of a closed ellipsoidal contour surrounding the point \mathbf{b} . When $K_1 > K_2$, the contour $U(\boldsymbol{\beta}) = K_1$ completely encloses the contour $U(\boldsymbol{\beta}) = K_2$ and \mathbf{b} lies in the centre of these nested m -dimensional "eggs". A $100(1 - \alpha)\%$ confidence region for the true (but unknown) value of $\boldsymbol{\beta}$ is enclosed by the contour given by

$$\frac{[U(\boldsymbol{\beta}) - U(\mathbf{b})]/m}{U(\mathbf{b})/(n - m)} F_{1-\alpha}(m, n - m)$$

only if errors are normally distributed, i.e. \mathbf{e} comes from $N(0, \sigma^2)$. The equation can be rearranged as

$$U(\boldsymbol{\beta}) = U(\mathbf{b}) \left[1 + \frac{m}{n - m} F_{1-\alpha}(m, n - m) \right]$$

where the expression on the right-hand side is a constant value that defines the contour. The outer contour shown in Fig. 8.5a is labelled as the $100(1 - \alpha)\%$ confidence contour, defined by the above equation. In the 2-dimensional space (β_1, β_2) , the contours are concentric ellipses about the point (b_1, b_2) . Note that contours of this type are obtained, irrespective of the value of n (the number of observations), since the dimension of the parameter space depends on m alone.

(B) The geometry of nonlinear least-squares

When the model is nonlinear there is no \mathbf{X} matrix as in the linear model. Although there is still an estimation space, it is not one that is defined by a set of vectors and

it may be very complex. This estimation space is called the *solution locus* and it consists of all points with co-ordinates of the form $\{f(x_1, \beta), f(x_2, \beta), \dots, f(x_n, \beta)\}$. Since the sum of squares $U(\beta)$ still represents the square of the distance from the point of the solution locus, minimization of $U(\beta)$ still corresponds geometrically to finding the point P of the solution locus which is nearest to Y.

Figure 8.5b shows the sample space for an example involving $n = 3$ observations, y_1, y_2 and y_3 , taken at x_1, x_2 and x_3 respectively, and two parameters β_1 and β_2 . The curved lines $f(\beta_j)$, $j = 1, 2$, also called the *estimation space curves*, indicate the co-ordinate system of parameters on the estimation space or solutions locus. It consists of all points of the form $\{f(x_1, \beta_1, \beta_2), f(x_2, \beta_1, \beta_2), f(x_3, \beta_1, \beta_2)\}$, as β_1 and β_2 vary with x_1, x_2 and x_3 fixed. Generally, this co-ordinate system is formed by all possible combinations of parameters values β in vector \mathbf{f} , where \mathbf{f} denotes $f(\beta_j)$ in which all parameters β_k , for $k \neq j$, are constant. The co-ordinate system of the estimation space curves forms the *estimation surface* Φ (the hatched part in Fig. 8.5b) of all possible solutions. From Fig. 8.5b, it is obvious that the termination points of all vectors \mathbf{f} lie on this estimation surface. When a solution lies on the estimation surface, a vector of parameter estimates \mathbf{b} exists for which $\mathbf{y} = f(\mathbf{b})$ i.e. $U(\beta) = 0$ and the regression model goes through all experimental points. The tangent plane τ in a location \mathbf{b} denoted by the unhatched region, is expressed by Eq. (8.66a) when β_j is replaced by b_j . From the geometry shown in Fig. 8.5b, the following conclusions are drawn:

(1) The minimum $U(\beta)$ corresponds to the minimum distance between a vector \mathbf{y} and the estimation surface.

(2) When the tangent plane τ sufficiently approximates the estimation surface Φ in the vicinity of the point \mathbf{b} , the vector $\mathbf{f} = [f(x_1, \mathbf{b}), \dots, f(x_n, \mathbf{b})]$, called the *prediction*, is a perpendicular projection of a vector \mathbf{y} onto the tangent plane. The corresponding projection matrix has the form

$$\mathbf{P} = \mathbf{J}(\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \quad (8.68)$$

(3) The residual vector $\hat{\mathbf{e}}$ with components $\hat{e}_i = y_i - f(x_i, \mathbf{b})$ is perpendicular to the tangent plane τ . Therefore a condition for the existence of a minimum of $U(\beta)$ is the validity of the expression $\mathbf{J}^T \hat{\mathbf{e}} = 0$. It is then important to determine how precisely the tangent plane approximates the estimation surface.

In the linear model, the contours of constant $U(\beta)$ in parameter space consist of concentric ellipses. When the model is nonlinear, the contours are sometimes banana-shaped, often elongated. Sometimes the contours stretch to infinity and do not close, or they may have multiple loops surrounding a number of stationary values. When several stationary values exist they may have different levels or provide alternative minima for $U(\beta)$.

It is convenient to draw both the solution locus and the parameter space, simultaneously (Fig. 8.6). From m -dimensional parameter space a projection is made onto the estimation surface.

Estimation of the parameters \mathbf{b} requires a matrix inversion:

$$\mathbf{b} = \mathbf{f}^{-1}(\mathbf{y}) \quad (8.69)$$

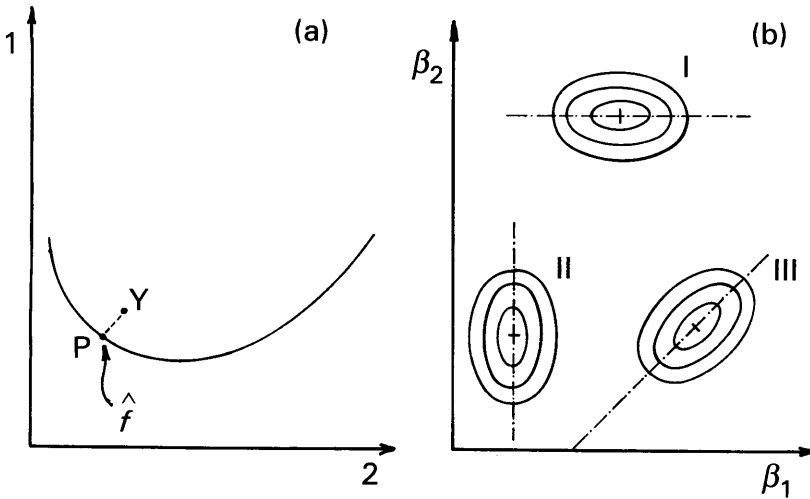


Fig. 8.6—(a) The solution locus; (b) the parameter space for $m = 2$ and $n = 3$. (I) β_2 well determined but not β_1 ; (II) β_1 well determined but not β_2 ; (III) relationship between β_1 and β_2 .

A condition of unambiguity is that each y value in sample space always corresponds to just one point in parameter space. When the estimation surface is nonlinear, the arbitrary straight line in the parameter space through $\beta^{(o)}$ given by

$$\beta = \beta^{(o)} + l\mathbf{h} \quad (8.70)$$

where $\mathbf{h} = (h_1, \dots, h_m)^T$ is any nonzero vector, generates a curve or “lifted line” on the solution locus given by

$$\mathbf{f}_h = \mathbf{f}(\beta^{(o)} + l\mathbf{h}) \quad (8.71)$$

In both Eqs. (8.70) and (8.71), \mathbf{h} represents the direction vector and l is the parameter of the straight line, and

$$\Delta\beta = l\mathbf{h}$$

The tangent \mathbf{f}'_h to this curve \mathbf{f}_h at $\beta^{(o)}$, is found from Eq. (8.66a), to be

$$\mathbf{f}'_h = \mathbf{J}\mathbf{h} \quad (8.72)$$

The set of all such linear combinations is referred to as the tangent plane τ at $\beta^{(o)}$. To express the curvature of the estimation surface, the vector of second partial derivatives, known as the acceleration of the lifted line \mathbf{f}''_h , may be shown to be

$$\mathbf{f}''_h = \mathbf{h}^T \vec{\mathbf{G}} \mathbf{h} \quad (8.73)$$

Each element of \mathbf{f}''_h has the form $\mathbf{h}^T \mathbf{G}_i \mathbf{h}$ where \mathbf{G}_i is the i th plane of the $\vec{\mathbf{G}}$ array.

The acceleration vector \mathbf{f}''_h comprises three components: the first component \mathbf{f}''_h determines the change in direction of the instantaneous velocity vector \mathbf{f}'_h normal to the tangent plane. The second and the third components, which can be added together

to give $\mathbf{f}_h''^P$, determine the change in direction of \mathbf{f}_h' parallel to the tangent plane and the change in speed of the moving point respectively. From the projection matrix \mathbf{P} defined by Eq. (8.68) it follows that

$$\mathbf{f}_h''^N = \mathbf{P}\mathbf{f}_h'' \quad (8.74a)$$

$$\mathbf{f}_h''^P = (\mathbf{E} - \mathbf{P})\mathbf{f}_h'' \quad (8.74b)$$

The acceleration components may be converted into curvatures, namely the *intrinsic curvature*

$$K_h^N = \frac{\|\mathbf{f}_h''^N\|}{\|\mathbf{f}_h'\|^2} \quad (8.75a)$$

and the *parameter-effects curvature*

$$K_h^P = \frac{\|\mathbf{f}_h''^P\|}{\|\mathbf{f}_h'\|^2} \quad (8.75b)$$

Only the latter depends on the particular parameterization chosen.

Interpreted geometrically, K_h^N represents the reciprocal of the radius of a circle which approximates the estimation surface in the direction \mathbf{h} . This curvature depends on the actual type of regression model and on the data used. It is not affected by reparameterization.

The curvature K^P corresponds to the nonparallelity of curves formed by the projection of uniformly spaced points on parallel straight lines from the parameter space into non-uniformly spaced points on the estimation surface. This curvature may be removed by reparameterization.

For characterization of nonlinear behaviour of regression models, we look for a value of vector \mathbf{h} such that the value of K_h^N and K_h^P have maximum values. Thus, both curvatures may be converted into response-invariant standardized relative curvatures Γ^N and Γ^P respectively. Multiplication by the standard radius $\varrho = \hat{\sigma}\sqrt{m}$, where $\hat{\sigma}$ is the square root of the estimated residual variance $\hat{\sigma}^2$, results in the *maximum intrinsic curvature*

$$\Gamma^N = \hat{\sigma}\sqrt{m} \max(K_h^N) \quad (8.75c)$$

and the *maximum parameter-effects curvature*

$$\Gamma^P = \sigma\sqrt{m} \max(K_h^P) \quad (8.75d)$$

8.5 NUMERICAL PROCEDURE FOR PARAMETER ESTIMATION

If a regression model $f(\mathbf{x}, \boldsymbol{\beta})$ is nonlinear in at least one model parameter β_r , substitution into the criterion function [Eq. (8.57)] leads to a task of *nonlinear minimization*. The application of maximum likelihood (Section 8.3) leads to the task of *nonlinear maximization*. The application of any regression criterion leads to the problem of finding an extreme, where the regression parameters $\boldsymbol{\beta}$ are “variables”.

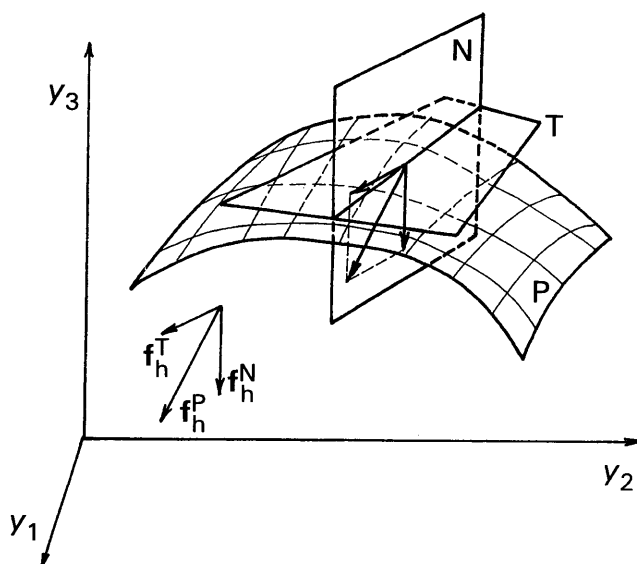


Fig. 8.7—Geometrical illustration of curvature; decomposition of the acceleration vector f_h'' into the components $f_h^{P''}$ and $f_h^{N''}$.

This task can be solved by the application of general optimization methods to search for a free extreme (if no restrictions are placed on the regression parameters) or a search for a constrained extreme if the regression parameters are subject to certain restrictions.

Owing to the great variability of regression models, regression criteria and data, ideal algorithms that can achieve convergence to a global extreme sufficiently fast cannot be found. Most algorithms for many numerical methods often fail, i.e. they converge very slowly or diverge. The more complicated procedures for complicated problems are rather slow and require a large amount of computer memory.

Any program for solving nonlinear regression problems should contain procedures for

- (1) searching for extremes in a given direction (one-dimensional optimization);
- (2) inversion of matrices;
- (3) numerical differentiation (in derivative methods);
- (4) methods of overcoming local areas of divergence.

Programs that use the same algorithms may differ in practical applicability. Comparison of individual programs requires special problems [8] which allows testing under approximately similar conditions. Some typical programs for nonlinear regression have been compared [9].

Even with recommended algorithms, the correct result may not be reached. Kuesters and Mize [10] and Wolfe [11] have proposed schemes for solving partial problems such as numerical differentiation, a search in a given direction, etc. Some practical recommendations for the construction of optimization programs have been

suggested by Gill, Murray and Wright [12], and Schmidt [13] has proposed a program for minimization of the least-squares criterion.

For nonlinear regression, it has been recommended that a package of various regression algorithms is applied either individually or in combination [14].

Nonlinear regression algorithms may be classified into the following principal groups:

- (1) Derivative-free optimization methods;
- (2) Derivative methods for the least-squares method (LS);
- (3) General derivative methods;
- (4) Algorithms for special cases.

The selection of a particular group depends on many factors. Generally, when a criterion function cannot be differentiated, the derivative-free methods should be used. The derivative methods use a special form of the least-squares criterion (LS) which are based on a quadratic approximation of the regression criterion. The general derivative methods enable solution of the task of maximization of likelihood function, for any regression model.

In this chapter we concentrate on the procedures of the first two groups. The general derivative methods are the most commonly used today [10,12]. The algorithms for special cases are determined either by other regression procedures, such as the robust or L_p approximation methods, or by sums of exponential, etc. From our experience, the first two groups of regression programs can solve most types of chemometrics problems.

8.5.1 Non-derivative optimization procedures

Non-derivative optimization procedures allow a search for extremes in the general criterion $U(\beta)$ in terms of parameters β . They use information about the form of function $U(\beta)$ obtained by "mapping" a parameter space. An extreme is selected according to various heuristic rules. These procedures are, in practice, quite popular for their simplicity, but a disadvantage is that most of the derivative-free methods converge very slowly especially in the vicinity of an optimum and have complicated forms of $U(\beta)$. For example, in a skewed banana shape, the algorithms terminate before reaching a minimum. These algorithms are more appropriate for functions $U(\beta)$ that cannot be differentiated.

From the various possible strategies and procedures, we select the following methods:

- (a) methods of direct search;
- (b) simplex methods;
- (c) random search techniques;
- (d) special methods for least-squares.

Among the derivative-free methods are some that involve the numerical calculation of the derivatives. Here, such methods are classed as derivative, since information about the criterion function requires use of its derivatives.

From the statistical point-of-view, the main disadvantage of the derivative-free

methods is the fact that after termination of the minimization process, no information is available about the covariance matrix of estimates, determined from the matrix $(\mathbf{J}^T \mathbf{J})^{-1}$.

In the next section, we concentrate on methods for searching for a minimum of $U(\boldsymbol{\beta})$.

8.5.1.1 Direct search methods

This group includes many heuristic† procedures. One of the simplest is the Hooke–Jeeves algorithm [16] which can be described by the following two-phase procedures:

- (a) a step increment in the direction of the individual co-ordinate axes from the estimate $\boldsymbol{\beta}^{(i)}$ in the i th iteration, then a search for an improved estimate $\boldsymbol{\beta}_p^{(i)}$ for which $U(\boldsymbol{\beta}_p^{(i)}) < U(\boldsymbol{\beta}^{(i)})$;
- (b) a step search for a local minimum $\boldsymbol{\beta}^{(i+1)}$ in the direction determined by points $\boldsymbol{\beta}^{(i)}$ and $\boldsymbol{\beta}_p^{(i)}$. Such a search, for $m = 2$, is illustrated in Fig. 8.8.

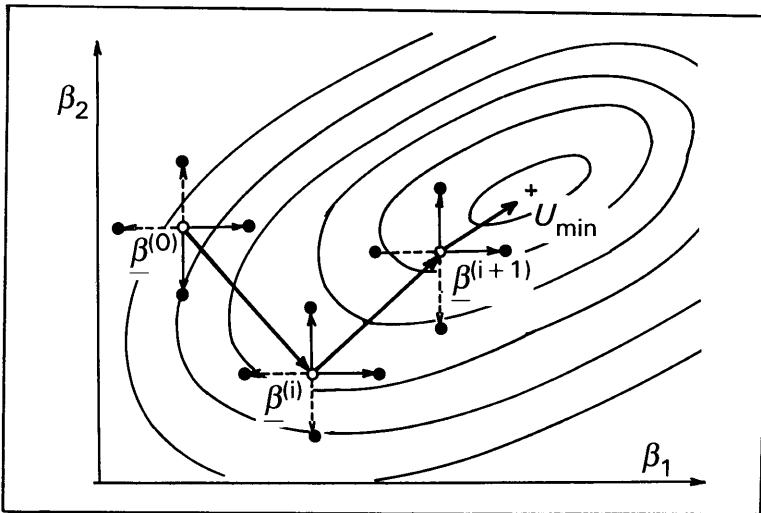


Fig. 8.8—The procedure for co-ordinate search (dashed lines denote unsuccessful directions).

Individual variants of the Hooke–Jeeves method differ in their strategy for shortening the step-increments in the first phase. The increments should not only shorten but also elongate in the given direction when necessary. This is particularly useful in applications such as inverse parabolic interpolation for points $\boldsymbol{\beta}^{(i)}$, $\boldsymbol{\beta}_p^{(i)} + \boldsymbol{\beta}^{(i)} + 2(\boldsymbol{\beta}_p^{(i)} - \boldsymbol{\beta}^{(i)})$ (Section 8.5.2).

The Rosenbrock method [15], instead of involving a step-by-step search for a local minimum in direction \mathbf{s} defined by points $\boldsymbol{\beta}^{(i)}$ and $\boldsymbol{\beta}_p^{(i)}$, rotates the co-ordinate system so that one axis is identical to the direction $\mathbf{s} = \boldsymbol{\beta}^{(i+1)} - \boldsymbol{\beta}^{(i)}$. In the next cycle, a co-ordinate search in the new axis is performed. The advantage of this procedure is that it works well even for cases when the function $U(\boldsymbol{\beta})$ has a narrow skewed (banana)

†heuristic means trial-and-error

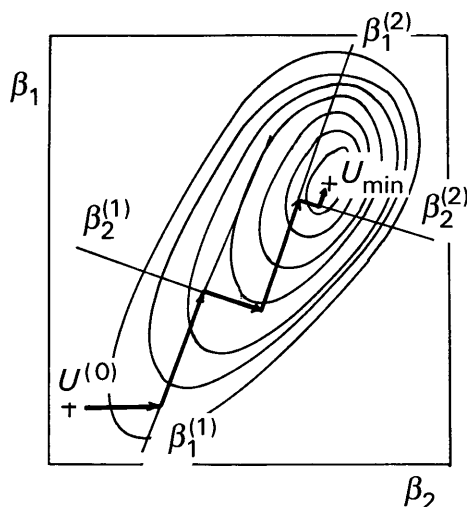


Fig. 8.9—The procedure for Rosenbrock's minimization.

shape. The procedure, when $m = 2$, is illustrated in Fig. 8.9.

The most effective methods in this group are procedures of conjugated directions when, instead of a rotation of the co-ordinate system, the new direction to the local extreme is conjugate to the original one (Powell method [16]).

8.5.1.2 Simplex methods

The simplex methods are the most commonly used optimization methods in analytical chemistry, process engineering and applied statistics. The original non-adaptive simplex method, proposed by Spendley *et al.* [17] is only rarely used now. The first useful modification, by Nelder and Mead [18], led to a simple and widely applicable algorithm. Further modifications have speeded up convergence and removed some limitations of the Nelder and Mead algorithm.

The method starts from an initial guess of parameters $\beta^{(0)}$ from which the simplex is formed. The simplex is a polyhedron, having $(m + 1)$ vertices constructed in m -dimensional parameter space. For $m = 2$, the simplex is a triangle and for $m = 3$, a tetrahedron (Fig. 8.10).

The process of minimization by a simplex method involves three steps.

- (1) Construction of the initial simplex.
- (2) Iterative search for a minimum.
- (3) Identification of a search termination.

In the minimization procedure in each cycle, steps (2) and (3) are repeated. Step (1) affects the speed of convergence.

(1) Construction of the initial simplex

The co-ordinates of the vertices of a simplex create rows of the matrix \mathbf{V} of dimension $[(m + 1) \times m]$. A regular simplex is defined by the magnitude of its edge t and a

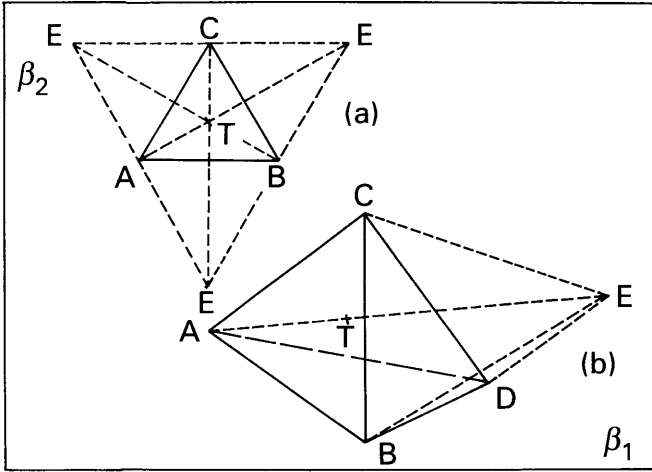


Fig. 8.10—The simplex for (a) $m = 2$, two parameters, and (b) $m = 3$, three parameters. The simplex ABC may be reflected into three possible positions CBE, ABE and ACE.

suitable location in parameter space. For a regular simplex, starting at the origin of co-ordinates and with edges of equal lengths, matrix V has the form

$$V = \begin{bmatrix} 0 & 0 & 0 & . & . & . & 0 \\ c & a & a & . & . & . & a \\ a & c & a & . & . & . & . \\ . & . & . & . & . & . & . \\ . & . & . & . & . & . & . \\ a & a & a & . & . & . & c \end{bmatrix} \quad (8.76)$$

where

$$c = \frac{t(m-1) + \sqrt{(m+1)}}{m\sqrt{2}} \quad (8.77)$$

$$a = \frac{t(\sqrt{m+1} - 1)}{m\sqrt{2}} \quad (8.78)$$

When an initial guess $\beta^{(o)}$ is known, the initial simplex is constructed such that the first row of the matrix V contains this initial guess as the co-ordinates of the component. The j th row ($j = 2, \dots, m+1$), is given by

$$V_{ji} = \beta_i^{(o)} + \frac{0.5\beta_i^{(o)}}{m\sqrt{2}}(m-1 + \sqrt{m+1}), \quad i \neq j \quad (8.79)$$

$$V_{ji} = \beta_i^{(o)} + \frac{0.5\beta_j^{(o)}}{m\sqrt{2}}(\sqrt{m+1} - 1), \quad i = j \quad (8.80)$$

Sometimes, it is more convenient to have the initial simplex constructed such that

the initial guess $\beta^{(0)}$ is in its centre of gravity or is suitably turned into a direction of the steepest descent.

(2) *Iterative search for a minimum*

The principle idea of simplex methods is very simple: the direction to a minimum is on the connecting line between the vertex V_U of a maximum value $U(V_U)$ and its reflection.

The procedure calculates the criterion $U(\beta_j)$ for all vertices of the simplex $V_j = \beta_j$ and enables the vertex V_U to be reflected through the centre of gravity of the other vertices to form a new simplex. To speed up the search for an optimum, five main operations are used: reflection, expansion, contraction, reduction and transfer. These operations (for $m = 2$) are illustrated in Fig. 8.11.

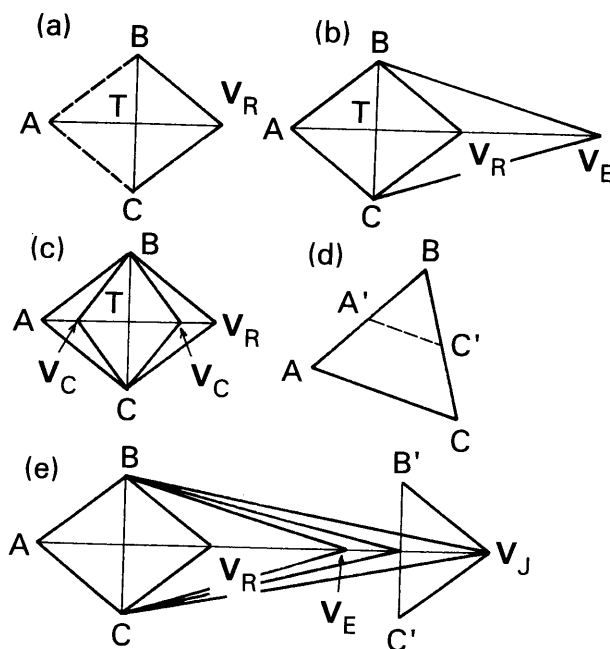


Fig. 8.11—The simplex operations: (a) reflection, (b) expansion, (c) contraction, (d) reduction, and (e) transfer.

The numerical expression of these operations assumes that after $(k - 1)$ steps, the simplex has vertices $\beta_i = V_i$, $i = 0, \dots, m$ with corresponding values of criterion function $U(V_i) = U_i$. Let us say that $d(V_i, V_j)$ is the distance between vertices V_i and V_j , $U_L = \min U_i$ for a minimum of the regression criterion at the vertex V_L , and $U_U = \max U_i$ for a maximum of the regression criterion at the vertex V_U . The minimum and maximum are selected from all the vertices of the simplex. The centre of gravity, of all vertices except the U th (V_U) is calculated from

$$\mathbf{P} = \frac{\sum_{i \neq U}^m \mathbf{V}_i}{m} \quad (8.81)$$

The vertex \mathbf{V}_R , being a reflection of the vertex \mathbf{V}_U , may be expressed in the form

$$\mathbf{V}_R = \mathbf{P} + \alpha(\mathbf{P} - \mathbf{V}_U) \quad (8.82)$$

Whenever possible, the regularity of the simplex should be retained. When, for at least one vertex $\mathbf{V}_i (i \neq L, i \neq U)$, $U_i > U_R > U_L$, the vertex \mathbf{V}_U is replaced by the vertex \mathbf{V}_R and the k th cycle is finished.

If, however, $U_L > U_R$, an expansion is made, to obtain point \mathbf{V}_E :

$$\mathbf{V}_E = \mathbf{P} + \gamma(\mathbf{V}_R - \mathbf{P}) \quad (8.83)$$

In practice, it has been found that the most suitable value of γ is 2.9 or in general, any value for which $\gamma > 2$. If $U_E \leq U_R$, then \mathbf{V}_U or \mathbf{V}_R may be replaced by \mathbf{V}_E . Sometimes additional expansion to the point \mathbf{V}_J is performed, that is

$$\mathbf{V}_J = (J + 1)\mathbf{V}_E - J\mathbf{P} \quad (8.84)$$

where $J = 2, 3, 4, \dots$, until $U_{J+1} > U_J$. Then \mathbf{V}_U is replaced by \mathbf{V}_J and a transfer of another vertex \mathbf{V}_i into a new position \mathbf{V}_i^N is performed in order to keep the original form of the simplex,

$$\mathbf{V}_i^N = \mathbf{V}_i + J(\mathbf{V}_E - \mathbf{P}), \quad i \neq U \quad (8.85)$$

until the k th cycle is once again completed.

When $U_R > U_i$ for all $i \neq U$ after a reflection, a contraction is performed by using the point \mathbf{V}_C given by

$$\mathbf{V}_C = \mathbf{P} + \beta(\mathbf{O} - \mathbf{P}) \quad (8.86)$$

Equation (8.83), $\mathbf{O} = \mathbf{V}_U$ if $U_U < U_R$ or $\mathbf{O} = \mathbf{V}_R$ if $U_U > U_R$. The shortening parameter $\beta < 1$ is usually set equal to $\beta \approx 0.55$. If $U_C < U_U$, \mathbf{V}_U is replaced by \mathbf{V}_C , completing the k th iteration.

If, despite the contraction, $U_C \geq U_U$, the simplex is shortened around \mathbf{V}_L with the smallest criterion value. Shortening involves replacement of vertices $\mathbf{V}_i (i \neq L)$ with the new vertices \mathbf{V}_i^Z such that

$$\mathbf{V}_i^Z = \mathbf{V}_L + \lambda(\mathbf{V}_i - \mathbf{V}_L) \quad (8.87)$$

For this, $\lambda = 0.5$ is usually chosen. This shortening procedure completes the k th cycle. It is useful, however, to test whether the shortening process is not too excessive. Later, enlargement of the simplex magnitude ($\lambda > 1$) can be used to overcome the local unevennesses of the surface of criterion function.

(3) Identification of a termination criterion

Nelder and Mead [20] recommend that at the end of each cycle, an examination of the magnitude of the decrease of the criterion function and of the relative changes of the simplex vertices be made by use of

$$|U_U - U_L| < \varepsilon_1 \quad (8.88)$$

and

$$\frac{1}{m+1} \sum_{i=0}^m d[\mathbf{V}_i^{(k)}, \mathbf{V}_i^{(k-1)}]^2 < \varepsilon_2 \quad (8.89)$$

The constant ε_1 should have the value 10^{-4} and ε_2 , 10^{-8} . The term $\mathbf{V}_i^{(k)}$ denotes the i th vertex of the simplex in the k th iteration.

Sometimes [18] instead of Eq. (8.88), the criterion

$$\sum_{i=0}^m \frac{(U_i - U_L)^2}{m+1} \leq \varepsilon_3 \quad (8.90)$$

is used.

Routh, Schwartz and Denton [19] have proposed a procedure called the *supermodified simplex* (SMS), which is an improved algorithm for an optimum search in the directions \mathbf{V}_U , and \mathbf{V}_R , using the inversion parabolic interpolation. The procedure starts with three known points (\mathbf{V}_U, U_U) , (\mathbf{P}, U_P) , (\mathbf{V}_R, U_R) ; a parabola is fitted through these points and, by an analytical differentiation, its minimum is found. If the parabola is concave and has a maximum, the simplex is contracted or reduced. The magnitude of the coefficient, α , in Eq. (8.82) is selected according to the slope of the parabola.

Ryan, Barr and Fodd [20] devised a generalized simplex method in which, in an attempt to retain the simplicity of the method, the vertex \mathbf{V}_U is reflected in the direction of the approximate gradient of the criterion function $U(\boldsymbol{\beta})$. This procedure is called the *weighted centroid method* (WCM), and uses the weighted centre of gravity \mathbf{P}_W defined by

$$\mathbf{P}_W = \frac{\sum_{i=0}^m (U_i - U_U) \mathbf{V}_i}{\sum_{i=0}^m (U_i - U_U)} \quad (8.91)$$

This equation assumes that the gradient of the function U is close to the direction of the line joining \mathbf{V}_L and \mathbf{V}_U . The next procedure is identical to those described above but, instead of \mathbf{P} , the weighted centre of gravity, \mathbf{P}_W , is used.

The application of Eq. (8.91) in some cases leads to degeneration of the simplex, i.e. an object of lower dimension than the original simplex is formed, because the angle between the directions $(\mathbf{V}_L, \mathbf{V}_U)$ and $(\mathbf{V}_U, \mathbf{P}_W)$ is very small.

The measure of the deviation of \mathbf{P}_W from \mathbf{P} is given by

$$V_o = \frac{\|\mathbf{P}_W - \mathbf{P}\|}{\|\mathbf{V}_L - \mathbf{P}\|} \quad (8.92)$$

where $\|\mathbf{x}\| = \sqrt{\sum x_i^2}$ is the norm of vector \mathbf{x} . When \mathbf{P}_W is close to \mathbf{V}_L , then $V_o \rightarrow 1$. The reduced centre of gravity \mathbf{P}_C has the form

$$\mathbf{P}_C = (1 - \kappa)\mathbf{P} + \kappa\mathbf{P}_W \quad (8.93)$$

and should be used when the parameter κ takes the value 0.3.

Another modification of the simplex method is based on the reflection of more vertices in every iteration. Evans [24] separates simplex vertices in the k th iteration into two groups. The first group contains vertices that have large values of the criterion function, U_i , and these are reflected. The second group contains vertices with relatively small U values, and these remain unchanged in the k th iteration. Rules for simultaneous changes of more simplex vertices have been proposed by Volker *et al.* [22].

One modification [21] of the simplex method involves a random search for the initial guess of parameters $\beta^{(0)}$.

The simplex method is useful when the initial guesses of parameters are rather far from the true values. Its disadvantage is a slow convergence in the vicinity of the minimum. Therefore, simplex methods are often combined with derivative methods for locating extremes. For example, in program FUNMIN [23] the simplex method is combined with the Marquardt algorithm.

Spendley [24] proposed a simple procedure which combines the simplex with the Gauss–Newton method (Section 8.5.2). The procedure starts from a linearized function $f(x_i, \beta)$, defined by Eq. (8.64), for $\beta = \mathbf{V}_j$ or $\beta = \mathbf{V}_L$. It can be shown that

$$e_i(\mathbf{V}_j) - e_i(\mathbf{V}_L) \approx \mathbf{J}_i^T (\mathbf{V}_j - \mathbf{V}_L) \quad (8.94)$$

where \mathbf{J}_i^T is the i th row of the Jacobian \mathbf{J} and the symbol $e_i(\mathbf{V}_j) = y_i - f(x_i, \mathbf{V}_j)$ denotes the i th residual for the estimate \mathbf{V}_j . Similarly $e_i(\mathbf{V}_L)$ denotes the i th residual for the estimate, \mathbf{V}_L . In matrix notation, Eq. (8.94) can be written as

$$\mathbf{T} \approx \mathbf{J}\mathbf{A} \quad (8.95)$$

where \mathbf{T} is the $(n \times m)$ matrix with elements

$$\mathbf{T}_{ij} = e_i(\mathbf{V}_j) - e_i(\mathbf{V}_L), \quad i = 1, \dots, n, \quad j = 1, \dots, m \quad (8.96)$$

and \mathbf{A} is the $(m \times m)$ matrix with elements

$$\mathbf{A}_{jk} = V_{jk} - V_{Lk}, \quad j = 1, \dots, m, (j \neq L) \quad k = 1, \dots, m \quad (8.97)$$

If the simplex vertices in the k th iteration are known, the matrix \mathbf{J} may be estimated from Eq. (8.95). Let us assume that the criterion for the least-squares method (LS) is valid. The increment vector, \mathbf{L} , of the Gauss–Newton method (Section 8.5.2) for this criterion may be calculated from the approximate expression

$$\mathbf{L} = \mathbf{A} - \mathbf{D}^{-1}\mathbf{w} \quad (8.98)$$

where elements of the matrix \mathbf{D} are given by

$$\mathbf{D}_{jk} = \sum_{i=1}^n [e_i(\mathbf{V}_j) - e_i(\mathbf{V}_i)][e_i(\mathbf{V}_k) - e_i(\mathbf{V}_L)], \quad j, k = 1, \dots, m (j \neq L) \quad (8.99)$$

and those of vector \mathbf{w} by

$$w_j = \sum_{i=1}^n [e_i(\mathbf{V}_j) - e_i(\mathbf{V}_i)]e_i(\mathbf{V}_L), \quad j \neq L \quad (8.100)$$

In the k th iteration of a given simplex optimization, the procedure determines \mathbf{L} , from Eq. (8.98) and calculates the criterion value, $U(\mathbf{V}_L + \mathbf{L})$, from Eq. (8.57). This value then determines whether the procedure continues according to the original simplex method or replaces vertex \mathbf{V}_U of the maximum value $U(\mathbf{V}_U)$ by a vertex $(\mathbf{V}_L + \mathbf{L})$, and then uses an approximate Gauss–Newton method. Programming the combined procedures is complex because it involves an inversion of the matrix \mathbf{D} ; on the other hand, it has the advantage that after the termination of the minimization process and the determination of an optimum, the estimate of the matrix $(\mathbf{J}^T \mathbf{J})^{-1}$ may be calculated from

$$(\mathbf{J}^T \mathbf{J})^{-1} \approx \mathbf{A} \mathbf{D}^{-1} \mathbf{A}^T \quad (8.101)$$

In the nonlinear least-squares method, statistical analysis may then be applied as for derivative methods.

Problem 8.10. Search for a minimum by the simplex method

Estimate the minimum of the function

$$f(\boldsymbol{\beta}) = 100(\beta_1^2 - \beta_2)^2 + (1 - \beta_1)^2$$

using the modified simplex method with initial guesses $\beta_1^{(0)} = -1.2$ and $\beta_2^{(0)} = 1$.

Solution: After 179 iterations, a minimum was reached with parameter estimates $b_1 = 0.9934$ and $b_2 = 0.987$ and $f(\mathbf{b}) = 4.24 \times 10^{-5}$. The values are $b_1 = b_2 = 1$.

Conclusion: In comparison with the derivative methods, the modified simplex method is rather slow. However, the values found for the estimates are quite close to the true values.

8.5.1.3 Random optimization

Random methods are convenient for complicated polymodal criterion functions when other algorithms fail. They are very simple and do not require knowledge of good initial guesses of parameters $\boldsymbol{\beta}$. We restrict ourselves to a few selected procedures of adaptive random optimization which are simple and sufficiently effective.

CRS Algorithm (controlled random search) was originated by Price [25]. It uses a combination of controlled random search and the simplex method. It starts from the matrix $\boldsymbol{\beta}$ of dimension $(Z \times m)$, where $Z = 10(m + 1)$. The individual rows of this matrix comprise randomly selected points $\boldsymbol{\beta}_j$ in the parameter space. The minimization procedure consists of three steps:

(1) The estimation of points $\boldsymbol{\beta}_L = \min U(\boldsymbol{\beta}_j)$ and $\boldsymbol{\beta}_U = \max U(\boldsymbol{\beta}_j)$ where the minimum and maximum are searched for over all the rows of the matrix $\boldsymbol{\beta}$.

(2) A random sample of $(m + 1)$ rows $\mathbf{V}_1, \dots, \mathbf{V}_{m+1}$ is taken from matrix $\boldsymbol{\beta}$ [except the row corresponding to $(\boldsymbol{\beta}_L)$] and \mathbf{V}_1 is set equal to $\boldsymbol{\beta}_L$.

(3) Vectors $\mathbf{V}_i, i = 1, \dots, m + 1$ form the vertices of the simplex. A new test point, \mathbf{V}_R , is searched for as a reflection of vertex \mathbf{V}_{m+1} over the centre of gravity \mathbf{P} of the other points. Similarly, as in the classical simplex method, we have

$$\mathbf{V}_R = 2\mathbf{P} - \mathbf{V}_{m+1} = 2 \sum_{i=1}^m \frac{\mathbf{V}_i}{m} - \mathbf{V}_{m+1} \quad (8.102)$$

If $U(\mathbf{V}_R) \geq U(\mathbf{V}_U)$, the procedure goes to the 2nd step, otherwise \mathbf{V}_U is replaced by \mathbf{V}_R . When a minimum is not reached, the procedure returns to step (1).

GSA Algorithm (Generalized simulation annealing) [26] is based on a simulation of the behaviour of physical systems with many degrees of freedom and which lead to a state of minimum energy. An example of such a system is the annealing of metals at different rates of cooling. The procedure starts with the construction of a random point, \mathbf{V} , on the surface of an m -dimensional sphere of radius Δr , followed by a calculation of the criterion function $U(\mathbf{V})$.

The point \mathbf{V} is selected with a certain probability which depends on the difference between $U(\mathbf{V})$ and the current lowest value of criterion function $U(\mathbf{V}_o)$. If $\Delta U = U(\mathbf{V}) - U(\mathbf{V}_o) \leq 0$, point \mathbf{V} has the probability of acceptance of $P_p = 1$. Otherwise, the probability of acceptance is determined from

$$P_p = \exp(-\kappa \Delta U) \quad (8.103)$$

where κ is a positive constant affecting the convergence rate. The acceptance of \mathbf{V} is decided from a pseudorandom number R from a rectangular distribution $R[0, 1]$. If $R < P_p$, point \mathbf{V} is accepted as the centre of other m -dimensional spheres.

The minimization procedure requires a preliminary value for the global minimum U_{\min} . A method for selection has been described in [26]. The function, $U(\boldsymbol{\beta}) - U_{\min}$ is minimized when it reaches a global minimum zero value.

The algorithm GSA consists of the following steps: (the values $\boldsymbol{\beta}^{(i)}$ and $U(\boldsymbol{\beta}^{(i)}) = U_i$ are known)

(1) A search in a random direction \mathbf{S} , with the use of m random numbers N_i from the normalized normal distribution $N[0, 1]$. The components of the vector \mathbf{S} are defined by the normalized expression

$$S_j = \frac{N_j}{\sqrt{\sum_{i=1}^m N_i^2}}, \quad j = 1, \dots, m, \quad (8.104)$$

(2) Determination of a random point \mathbf{V} on the surface of the hypersphere, with the centre at the point $\boldsymbol{\beta}^{(i)}$ and radius Δr , (a value of $\Delta r = 0.15$ is usually selected) according to

$$\mathbf{V} = \Delta r \mathbf{S} + \boldsymbol{\beta}^{(i)} \quad (8.105)$$

(3) If $\Delta U = U(\mathbf{V}) - U_i \leq 0$, $\boldsymbol{\beta}^{(i+1)}$ take the value of \mathbf{V} , and the i th step is complete.

(4) If $\Delta U > 0$, the probability P_p is estimated [Eq. (8.103)] where usually $\kappa \approx 3.5$. Then a random number, R , is generated from the rectangular distribution $R[0, 1]$. If $R < P_p$, point \mathbf{V} , is accepted, i.e. $\boldsymbol{\beta}^{(i+1)} = \mathbf{V}$, and the i th step is completed. If $R \geq P_p$, then point \mathbf{V} is rejected and the procedure returns to step (1) i.e. to generation of a new random vector.

This minimization procedure can usually identify a round cluster of accepted points \mathbf{V} , around the global minimum.

Adaptive Random Search Technique Algorithm [27] is one of the simplest available.

Random vectors \mathbf{V} are generated around the point $\beta^{(i)}$, with values that depend on the earlier rate of convergence. The vectors \mathbf{V} correspond to a special probability distribution. The procedure for search in the i th iteration may be expressed by following steps:

- (1) Locate point \mathbf{V} from

$$\mathbf{V} = \beta^{(i)} + \frac{D(2R - 1)^K}{K} \quad (8.106)$$

where D is the magnitude of the increment (usually $D \approx 0.5$) and R is a pseudorandom number from a rectangular distribution $R[0, 1]$. The parameter K determines the type of probability distribution of the second term in Eq. (8.106). When $K = 1$, the distribution is rectangular. The higher the value of K , the more peaked the distribution and the smaller the variance.

(2) Calculate $U(\mathbf{V})$. If $U(\mathbf{V}) < U(\beta^{(i)})$, the substitution $\beta^{(i+1)} = \mathbf{V}$ is performed, which completes the i th iteration. If $U(\mathbf{V}) \geq U(\beta^{(i)})$, the value of K is adapted, then the procedure returns to step (1). The following strategy is recommended for adaptation parameter K :

- (a) Start with $K = 1$;
- (b) After five successful steps K is increased to 3;
- (c) After another 15 successful steps, K is increased to 5;
- (d) After a further 10 successful steps K is increased to 7;
- (e) When convergence becomes slow, K is halved until $K = 1$ is reached.

Adaptive Random Search Algorithm. Pronzato *et al.* [28] have proposed an improved version of the last algorithm, applying the principle of adaptation. This also begins with the generation of random test points

$$\mathbf{V} = \beta^{(i)} + \mathbf{Z} \quad (8.107)$$

Random vectors \mathbf{Z} , are generated from an m -dimensional normal distribution with zero mean and a diagonal covariance matrix \mathbf{C} .

Each step consists of two parts. In the first part, optimum values of variance C_j , in a covariance matrix \mathbf{C} , are searched for. In the second part, using around 100 optimum values, \mathbf{V} are generated such that they have minimal values $U(\mathbf{V})$ less than $U(\beta^{(i)})$. The search is terminated when five consecutive iterations produce the same estimates of the variance $\hat{\sigma}_j$.

Adaptive search for the optimum variances is based on set upper β_U and lower β_L limits determining the acceptable range of parameter estimates. Five vectors of variances $\sigma_1^2, \dots, \delta_5^2$ are estimated from the equations

$$\begin{aligned} \sigma_1 &= \beta_U - \beta_L \\ \sigma_j &= 0.1^{(j-1)} \sigma_1, \quad j = 2, 3, 4, 5 \end{aligned} \quad (8.108)$$

For each vector σ_j , the point \mathbf{V} from Eq. (8.107) is calculated $f_j = 100/j$ times, for $\mathbf{Z} = \sigma_j \mathbf{N}$, where \mathbf{N} is the vector of independent random quantities of the normalized normal distribution. If $U(\mathbf{V}) < U(\beta^{(i)})$, the replacements $\beta^{(i)} = \mathbf{V}$ and $C_j = C_j + 1$ are

made. The value of $\hat{\sigma}_j$ for which \mathbf{C} has a maximum is taken into stage 2 as the maximum.

In cases where the initial guess of parameters $\beta^{(0)}$ is used, the vector σ_1 is calculated from

$$\sigma_1 = \mathbf{I} 0.3 + 2|\beta^{(0)}| \quad (8.109)$$

where \mathbf{I} is unit vector.

Generally, these algorithms are simple and lead to global extremes. However, they involve lengthy computations, and hence they are rather time-consuming.

8.5.1.4 Special procedures for the least-squares method (LS)

Some non-derivative procedures are based on the assumption that the criterion function $U(\beta)$ may be approximated in the vicinity of a given point by an m -dimensional hyperparaboloid. This is in accordance with the definition of the criterion function $U(\beta)$ of the least-squares method by Eq. (8.58). Two algorithms commonly used in chemistry are, LETAGROP and DUD.

LETAGROP algorithm. The principle of “pit-mapping” (in Swedish leta-grop) originated by Sillén and Ingri [29, 30] is the approximation of the criterion function $U(\beta)$ in the vicinity of $\beta^{(i)}$ in the i th iteration by an m -dimensional elliptic hyperparaboloid. The coefficients of this hyperparaboloid, which in derivative methods are expressed by components of matrices \mathbf{J} and \mathbf{H} [Eq. (8.58)], are calculated from $(m+1)(m+2)/2$ points $\{\beta_{(j)}, U(\beta_{(j)})\}$. Substitution of these points (called “shots” in LETAGROP) into the equation for an m -dimensional hyperparaboloid leads to a set of $(m+1)(m+2)/2$ linear equations for the estimation of their coefficients. If these coefficients are found from analytical differentiating, a minimum of the “approximate” paraboloid may be calculated and hence the vector $\beta^{(i+1)}$ is established.

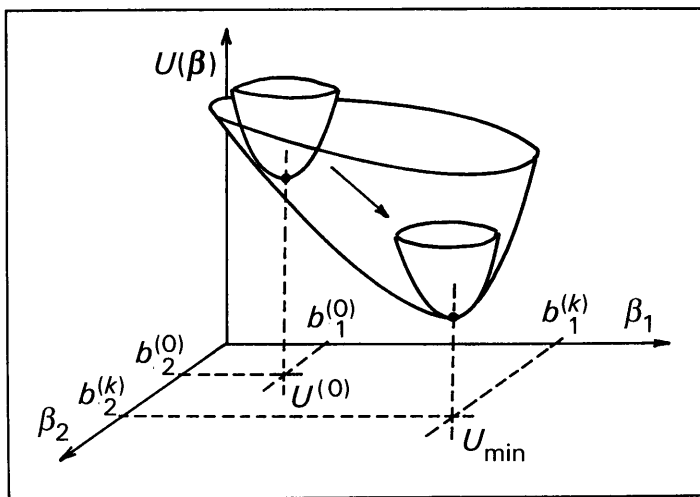


Fig. 8.12—Geometric illustration of the search for a minimum in $U(\beta)$ by the LETAGROP method.

Problems 8.11. *Equations for the minimum of a parabola*

Derive the equations for the minimum of a parabola approximating the criterion function $U(\beta)$ in the vicinity of the point $\beta^{(i)}$, if a regression function $f(x, \beta)$ contains only one parameter.

Solution: The first step involves linearization of the function $f(x, \beta)$ in the vicinity of point $\beta^{(i)}$. On substituting into Eq. (8.66a) we get

$$f(x, \beta) \approx f(x, \beta^{(i)}) + \frac{\delta f(x, \beta)}{\delta \beta} (\beta - \beta^{(i)})$$

If $\Delta = \beta - \beta^{(i)}$, then substitution of this approximation into the least-squares criterion function gives,

$$U(\beta) \approx \sum_{j=1}^n \left[y_j - f(x_j, \beta^{(i)}) - \Delta \frac{\delta f(x, \beta)}{\delta \beta} \right]^2 \approx K_0 + K_1 \Delta + K_2 \Delta^2$$

where

$$K_0 = \sum_{j=1}^n [y_j - f(x_j, \beta^{(i)})]^2$$

$$K_1 = -2 \sum_{j=1}^n \frac{\delta f(x_j, \beta)}{\delta \beta} [y_j - f(x_j, \beta^{(i)})]$$

and

$$K_2 = \sum_{j=1}^n \left[\frac{\delta f(x_j, \beta)}{\delta \beta} \right]^2$$

For estimation of the coefficients K_0 , K_1 and K_2 of an elliptic hyperparaboloid, knowledge of three values β_1 , β_2 and β_3 and their corresponding values $U(\beta_1)$, $U(\beta_2)$ and $U(\beta_3)$ are necessary. The desired coefficients are a solution of the set of three linear equations:

$$U(\beta_1) = K_0 + K_1(\beta_1 - \beta^{(i)}) + K_2(\beta_1 - \beta^{(i)})^2$$

$$U(\beta_2) = K_0 + K_1(\beta_2 - \beta^{(i)}) + K_2(\beta_2 - \beta^{(i)})^2$$

$$U(\beta_3) = K_0 + K_1(\beta_3 - \beta^{(i)}) + K_2(\beta_3 - \beta^{(i)})^2$$

The minimum of the "approximate" parabola is given by

$$\Delta_{\min} = \beta^{(i+1)} - \beta^{(i)} = -\frac{K_1}{2K_2}$$

Conclusion: With knowledge of suitable points β_j , $j = 1, \dots, (m+1)(m+2)/2$, the coefficients and the minimum of the parabola approximating the criterion function, $U(\beta)$ are easily found.

LETAGROP can also be modified for cases when $U(\beta)$ has the shape of a skewed narrow valley, such as when a strong correlation exists between parameters. In this case, the co-ordinate axes are rotated so that the axes of one co-ordinate lies along

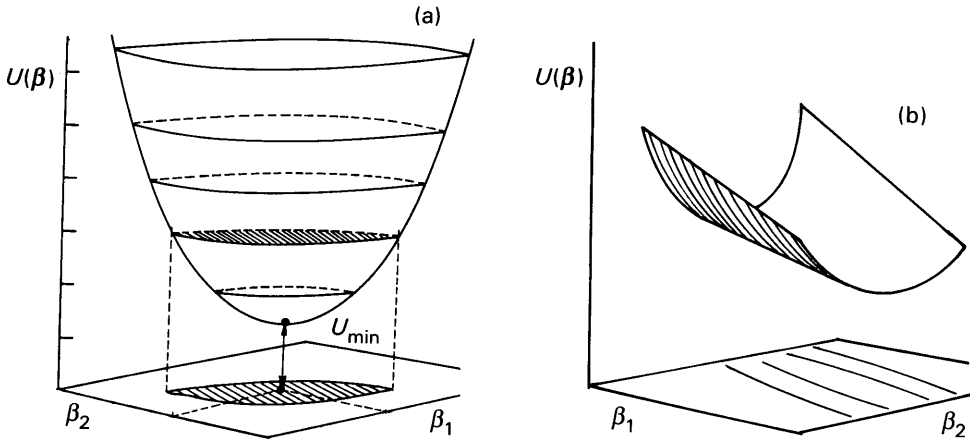


Fig. 8.13—Graphical illustration of a search for the minimum of an elliptic hyperparaboloid in the cases of (a) well-conditioned parameters without any correlation between β_1 and β_2 , and (b) ill-conditioned parameters with a strong correlation between β_1 and β_2 .

the axis of the valley.

LETAGROP is effective when there are a small number of parameters to be estimated. With increasing m , the number of evaluations of the criterion function $U(\beta)$ increases and the search for coefficients of the approximate paraboloid becomes more laborious. A major disadvantage of non-derivative methods is the absence of the matrix $(\mathbf{J}^T \mathbf{J})^{-1}$, which is useful for the statistical analysis of the parameter estimates. For the calculation of $(\mathbf{J}^T \mathbf{J})^{-1}$, the same procedure as in the Nelder-Mead simplex algorithm can be used. The Peckham method [31] uses a similar principle.

Algorithm DUD (Doesn't use derivatives). Procedure DUD is a nonderivative analogy of the Gauss-Newton algorithm [32]. Where linearization of the regression model is performed in the Gauss-Newton method, the function $f(x, \beta)$ is replaced by an affine function $l_i(\alpha)$ which is consistent with it at $(m+1)$ preceding points $\beta^{(k)}$, $k = i - m, \dots, i$, i.e. the results of previous iterations. Geometrically, this can be represented by drawing the function $f(x, \beta)$ in the second hyperplane instead of the estimation space.

To simplify the notation, let us renumber results of the previous iterations $\beta^{(k)}$, $k = i - m, \dots, i$, as $\beta^{(j)}$, $j = 1, \dots, m+1$ where $j = i - m - 1 + k$. The affine function $l_i(\alpha)$ for the i th point $\{x_i, y_i\}$ may be written in the form

$$l_i(\alpha) = f(x_i, \beta^{(m+1)}) + \sum_{j=1}^m \alpha_j [f(x_i, \beta^{(j)}) - f(x_i, \beta^{(m+1)})] \quad i = 1, \dots, n \quad (8.110)$$

After substitution of $l_i(\alpha)$ into the criterion function, the unknown coefficients α_j , $j = 1, \dots, m$ may be found in the form

$$\alpha = (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T (\mathbf{y} - \mathbf{f}) \quad (8.111)$$

where \mathbf{F} is the matrix of dimension $(n \times m)$ with elements

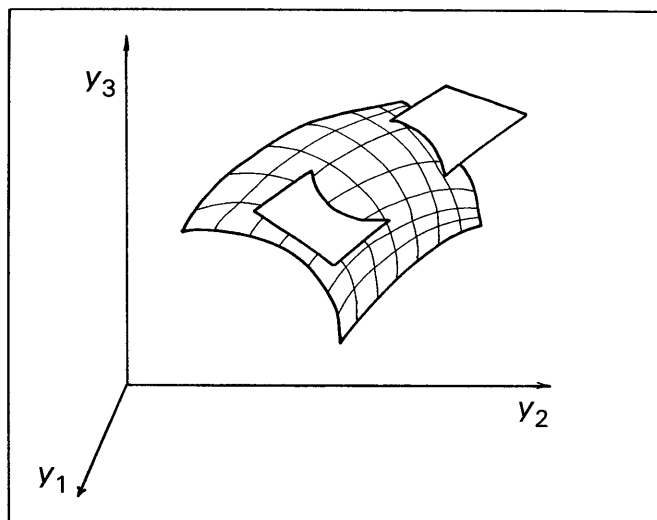


Fig. 8.14—Replacement of the estimation space by the secant hyperplane for $n = 3$ and $m = 2$.

$$F_{ij} = f(x_i, \beta^{(j)}) - f(x_i, \beta^{(m+1)})$$

and the vector \mathbf{f} has elements $f(x_i, \beta^{(m+1)})$. If vector α is known, a better evaluation β_N of $\beta^{(i+1)}$ may be calculated from

$$\beta_N = \beta^{(i+1)} + \mathbf{L}\alpha \quad (8.112)$$

where the j th column of the matrix \mathbf{L} has the form

$$\mathbf{L}_j = \beta^{(i)} - \beta^{(m+1)}, \quad j = 1, \dots, m$$

If $U(\beta_N) < U(\beta^{(m+1)})$, the i th iteration is completed and the substitution $\beta_N \rightarrow \beta^{(m+1)}$ is made. In the reverse case, the distance between β_N and $\beta^{(m+1)}$ is shortened by

$$\beta_N = p_z \beta_N + (1 - p_z) \beta^{(m+1)}$$

where $p_z = 1$ for $z = 1$ and $p_z = -(-1/2)^z$ for $z = 2, 3, 4$ and 5 . The stepwise regression can be used [32] for selecting a suitable vector α .

8.5.2 Derivative procedures for the least-squares method

Algorithms of this group are very commonly used, not only because the least-squares method is a frequent regression criterion but also to provide information necessary for subsequent statistical analysis of the regression results.

These algorithms are useful for all model functions which are twice differentiable. They have the disadvantage that the local convergence depends on the choice of the initial guess $\beta^{(0)}$. All algorithms of this group are of iterative nature. In the i th iteration the procedure starts from the estimates $\beta^{(i)}$ to which a suitable increment vector Δ_i is added by

$$\boldsymbol{\beta}^{(i+1)} = \boldsymbol{\beta}^{(i)} + \Delta_i$$

Generally, the procedure which searches for a minimum of $U(\boldsymbol{\beta})$ consists of four steps:

- (1) Determination of an initial guess $\boldsymbol{\beta}^{(0)}$.
- (2) A search for a convenient directional vector \mathbf{V}_i .
- (3) Determination of scalar α_i satisfying the condition, $\Delta_i = \alpha_i \mathbf{V}_i$.
- (4) Examination of the minimum obtained.

The vector Δ_i is usually considered to be acceptable if

$$U(\boldsymbol{\beta}^{(i)} + \Delta_i) < U(\boldsymbol{\beta}^{(i)}) \quad (8.113)$$

Some algorithms also allow equality of $U(\boldsymbol{\beta}^{(i+1)})$ and $U(\boldsymbol{\beta}^{(i)})$, or even a small increase. Individual algorithms differ in the realization of steps (2) and (3). Let us discuss each of the four steps:

(1) *Determination of the initial guess*

For many algorithms, this step is decisive for success of the minimization procedure. With a good initial guess, $\boldsymbol{\beta}^{(0)}$, even simple unsophisticated methods usually converge. With a poor initial guess, either a minimum can not be found at all or the minimum obtained is a local one. When the regression model can be linearly transformed, the initial guess $\boldsymbol{\beta}^{(0)}$ may be found by the linear least-squares method. In some cases, the initial guess may be obtained from physical or geometrical characteristics.

A transformation into stable parameters expressing geometrically defined characteristics of a regression model may be made, with, for example, function values or their derivatives at selected points, etc. With a personal computer, the path of the function $f(x, \mathbf{b}^{(0)})$ with given data may readily be tested and therefore the quality of the initial guess examined.

(2) *Determination of a directional vector*

The derivative of a criterion function $U(\boldsymbol{\beta})$ at a point $\Delta = \boldsymbol{\beta} + \alpha \mathbf{V}$, has the form

$$\frac{\delta U(\boldsymbol{\beta})}{\delta \alpha} = \left[\frac{\delta U(\boldsymbol{\beta})}{\delta \boldsymbol{\beta}} \right]^T \frac{\delta \boldsymbol{\beta}}{\delta \alpha} \quad (8.114)$$

For $\alpha \rightarrow 0$ we get, from Eq. (8.114), the *directional derivative*,

$$S_D = \left. \frac{\delta U(\boldsymbol{\beta})}{\delta \alpha} \right|_{\alpha \rightarrow 0} = \mathbf{g}^T \mathbf{V} \quad (8.115)$$

where \mathbf{g} is the gradient for which Eqs. (8.59a) and (8.59b) are valid.

From Fig. 8.15, it is evident that the gradient vector is the vector perpendicular to the tangent vector \mathbf{t} . All directional derivatives S_D of the hatched area are acceptable

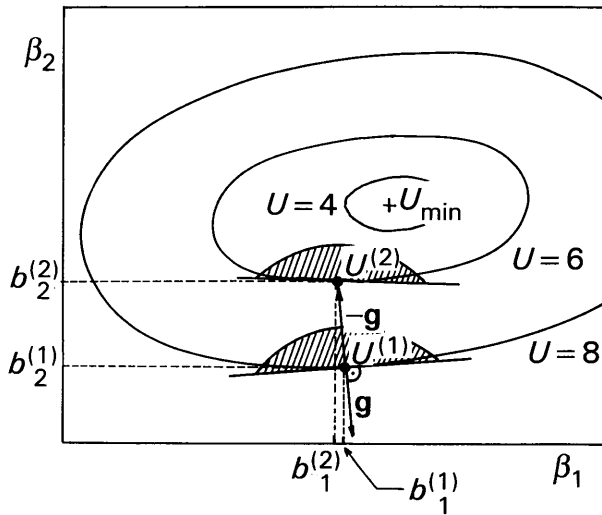


Fig. 8.15—Geometrical illustration of the gradient of a criterion function $U(\beta)$ for $m = 2$. Hatched areas denote admissible directional vectors.

because they do not cause an increase in $U(\beta)$.

The steepest decrease of the criterion function is in the direction $-\mathbf{g}$. The condition of acceptability of a given directional vector \mathbf{V} requires the directional derivative not to be positive. Any direction for which $\mathbf{g}^T \mathbf{V} > 0$ is unsuitable.

Moreover if the directional vector \mathbf{V} is admissible a positive definite matrix \mathbf{R} exists such that

$$\mathbf{V} = -\mathbf{R}\mathbf{g}$$

The directional derivative S_D is then

$$S_D = -\mathbf{g}^T \mathbf{R} \mathbf{g} \quad (8.116)$$

For a positive definite matrix \mathbf{R} , the quadratic forms are positive so that S_D in Eq. (8.115) is negative.

(3) Calculation of an optimum increment

In searching for an optimal increment $\alpha \mathbf{V}$ in the direction \mathbf{V} , an approximation of $U(\beta)$ in this direction by the Taylor series up to the second order can be used. This leads to the form

$$U(\beta + \alpha \mathbf{V}) \approx U(\beta) + \alpha \mathbf{g}^T \mathbf{V} + \frac{\alpha^2}{2} \mathbf{V}^T \mathbf{H} \mathbf{V} \quad (8.117)$$

where \mathbf{H} is the Hessian matrix defined by Eq. (8.62). Equation (8.117) is approximately quadratic with respect to α , so that an optimal α may be estimated by setting the first derivative with respect to α , $U(\beta + \alpha \mathbf{V})$, equal to zero. Hence, we get

$$\alpha^* = -\frac{\delta U(\mathbf{b})}{\delta \alpha} \bigg/ \frac{\delta^2 U(\mathbf{b})}{\delta \alpha^2} = -\mathbf{g}^T \mathbf{V} [\mathbf{V}^T \mathbf{H} \mathbf{V}]^{-1} \quad (8.118)$$

After substitution from Eq. (116) we obtain the *Raleigh coefficient*

$$\alpha^* = \mathbf{g}^T \mathbf{R} \mathbf{g} [\mathbf{g}^T \mathbf{R}^T \mathbf{H} \mathbf{R} \mathbf{g}]^{-1} \quad (8.119)$$

The Raleigh coefficient, α^* , is restricted to a region in which an approximation of type (8.117) can be used.

Another possibility in the search for an optimal α_i value in the direction \mathbf{V}_i is the one-dimensional minimization of the function $U(\boldsymbol{\beta} + \alpha_i \mathbf{V}_i)$.

(4) Termination of the iteration process

The natural criterion for an optimum \mathbf{b} is a zero value of the gradient \mathbf{g} of the criterion function. Many methods terminate the iterative process searching for a minimum when the norm of the gradient

$$\|\mathbf{g}\|^2 = \sum_{j=1}^m g_j^2$$

is sufficiently small. It is possible to select a critical value e.g. 10^{-4} , at which the point $\mathbf{b}^{(i)}$ is considered to be an extreme \mathbf{b} . Often, an iteration is terminated when the changes in the parameter estimates are very small. None of these criteria lead to termination at a true minimum.

Minimization may terminate less heuristically if the residual vector $\hat{\mathbf{e}}$ is approximately perpendicular to the columns of the matrix \mathbf{J} . From Fig. 8.5, it follows that $\mathbf{J}^T \hat{\mathbf{e}} = \mathbf{0}$. The angle, α_j , between the residual vector $\hat{\mathbf{e}}$ and the j th column \mathbf{J}_j of matrix \mathbf{J} , is given by the following expression:

$$\cos \alpha_j = \hat{\mathbf{e}}^T \mathbf{J}_j [\mathbf{J}_j^T \mathbf{J}_j \hat{\mathbf{e}}^T \hat{\mathbf{e}}]^{-1/2} \quad (8.120)$$

When the maximum value of $\cos \alpha_j$ is sufficiently small (e.g. smaller than 10^{-9}) it is assumed that a minimum of $U(\boldsymbol{\beta})$ is found. Many other termination criteria have been proposed [33].

We concentrate our attention on the following derivative algorithms for the least-squares method:

- (a) Gauss–Newton methods;
- (b) Marquardt methods;
- (c) the dog-leg method.

There is a wide spectrum of different improvements and modifications to these methods, but we restrict ourselves here to some simple and efficient techniques.

8.5.2.1 Gauss–Newton methods

To determine a convenient directional vector \mathbf{V}_i , the quadratic approximation of a criterion function $U(\boldsymbol{\beta})$ from Eq. (8.58) may be used, and this also corresponds to Eq. (8.117) for $\alpha = 1$. From

$$\frac{\delta U(\boldsymbol{\beta} + \mathbf{V})}{\delta \mathbf{V}} = 0$$

the optimum direction sector $\mathbf{V}_i = \mathbf{N}_i$ can be computed. The result takes the form

$$\mathbf{N}_i = -\mathbf{H}^{-1}\mathbf{g} = (\mathbf{J}^T\mathbf{J} + \mathbf{B})^{-1}\mathbf{J}^T\hat{\mathbf{e}} \quad (8.121)$$

On substituting into Eq. (8.119), we estimate that $\alpha^* = 1$ and \mathbf{N}_i is directly an increment vector Δ_i . This method is called the *Newton method*. When the criterion $U(\boldsymbol{\beta})$ is a quadratic function, (an elliptic paraboloid), the minimum \mathbf{b} will be found in one step. However, for other forms of criterion function $U(\boldsymbol{\beta})$ and estimates $\boldsymbol{\beta}^{(0)}$ far from \mathbf{b} , this method does not converge sufficiently fast. Moreover it requires knowledge of the matrix of second derivatives \mathbf{G}_i for determination of the matrix \mathbf{B} in Eq. (8.63). Neglecting matrix \mathbf{B} , which is equivalent to the linearization of the regression model, is theoretically acceptable for a case when the residual vector $\hat{\mathbf{e}}$ is small. The corresponding direction vector \mathbf{L}_i has the form

$$\mathbf{L}_i = (\mathbf{J}^T\mathbf{J})^{-1}\mathbf{J}^T\hat{\mathbf{e}} \quad (8.122)$$

Methods applying the directional vector \mathbf{L}_i are called *Gauss–Newton methods*. The methods are simple, and are the most frequently used procedure of nonlinear regression. Substituting $\mathbf{H} \approx (\mathbf{J}^T\mathbf{J})$ into Eq. (8.119) leads to $\alpha^* = 1$. From a practical point, the Gauss–Newton method works well if some of following conditions are fulfilled:

- (1) The residuals $\hat{e}_i = y_i - f(x_i, \boldsymbol{\beta})$ are small.
- (2) The model function $f(x, \boldsymbol{\beta})$ is nearly linear, i.e. the Hessian matrix, \mathbf{H} , has a small norm and its elements are nearly zero.
- (3) The residuals \hat{e}_i have alternating signs so that \mathbf{B} is an approximate zero matrix. This condition is valid in the vicinity of the optimum \mathbf{b} .

It is possible to use other methods to extend the region of convergence of this very simple method. The principle ones include:

- (a) Inversion of the matrix $\mathbf{J}^T\mathbf{J}$ and solution of the set of linear equations

$$(\mathbf{J}^T\mathbf{J})\mathbf{L} = \mathbf{J}^T\hat{\mathbf{e}} \quad (8.123)$$

- (b) Improving the matrix $(\mathbf{J}^T\mathbf{J})$ in order to be close to the Hessian matrix \mathbf{H} .
- (c) Choice of the optimal length of the step α .

We shall describe some successful methods which, on combination, lead to more effective modification of the original Gauss–Newton method.

- (a) *Inversion of the matrix $(\mathbf{J}^T\mathbf{J})$*

When the matrix $(\mathbf{J}^T\mathbf{J})$ is well conditioned, the columns of the matrix \mathbf{J} are linearly independent. Then, for a solution to the set of linear equations (8.123), various procedures may be applied. One of the simplest techniques, with minimal requirement on computer memory, is *Choleski decomposition*. In many practical problems (involving exponential and other nonlinear models), matrix \mathbf{J} has some nearly

collinear columns and therefore the matrix $(\mathbf{J}^T \mathbf{J})$ is ill-conditioned. The length of vector \mathbf{L} estimated from Eq. (8.123) is usually too large and its components have “inconvenient” signs. When some columns of a matrix \mathbf{J} are linearly dependent, the matrix $(\mathbf{J}^T \mathbf{J})$ is singular. These problems may be eliminated by pseudoinversion of matrix $\mathbf{J}^T \mathbf{J}$ or generalized inversion of matrix \mathbf{J} by the algorithm SVD.

Jennrich and Sampson [35] solved the set of linear equations in Eq. (8.123) by stepwise regression. Components of a vector \mathbf{L} are found and these significantly decrease the function $U(\boldsymbol{\beta})$.

In our programs, the matrix $\mathbf{J}^T \mathbf{J}$ is decomposed into eigenvalues and eigenvectors. To invert matrix $(\mathbf{J}^T \mathbf{J})$, we use the technique of rational ranks described in Chapter 6.

One of the most advanced procedures for searching for a suitable vector \mathbf{L} was devised by Schmidt [36]. Instead of the matrix \mathbf{J} , Schmidt constructs matrix \mathbf{M} containing only those columns of a matrix \mathbf{J} for which it is valid that

- they are not linearly dependent,
- they cause the largest decrease in $U(\boldsymbol{\beta})$
- at least one of them is not orthogonal to vector $\hat{\mathbf{e}}$.

This procedure protects the task against difficulties arising from ill-conditioning.

(b) *Improvement of the Hessian matrix*

This group includes methods of variable metric also known as the quasi-Newton methods. Here, each step treats the matrix $(\mathbf{J}^T \mathbf{J})$ so that it approximates the Hessian matrix \mathbf{H} . These methods are suitable for cases with large residuals, that is when $U(\mathbf{b}) \gg 0$. The main idea is simple, and comes from the fact that the Hessian matrix is the derivative of the gradient with respect to a parameter vector, and this derivative is approximated by the difference

$$H_{i+1} \approx \frac{\mathbf{g}_{i+1} - \mathbf{g}_i}{\boldsymbol{\beta}^{(i+1)} - \boldsymbol{\beta}^{(i)}} = \Delta \mathbf{g}_i \mathbf{S}_i^{-1} = \mathbf{B}_{i+1} \quad (8.124)$$

where \mathbf{S}_i is the increment vector and index i denotes the i th iteration. The matrix \mathbf{B}_{i+1} is an approximation of the matrix \mathbf{H}_{i+1} calculated only from information about gradients and values of the vector $\boldsymbol{\beta}^{(i)}$. The course of the procedure is that, instead of \mathbf{B}_{i+1} , the increment $\Delta \mathbf{B}_i = \mathbf{B}_{i+1} - \mathbf{B}_i$ is calculated in the individual iterations. In many cases $\Delta \mathbf{B}_i$ is calculated from [40]

$$\Delta \mathbf{B}_i = \frac{(\Delta \mathbf{g}_i - \mathbf{B}_i \mathbf{S}_i) \mathbf{C}_i^T + \mathbf{C}_i (\Delta \mathbf{g}_i - \mathbf{B}_i \mathbf{S}_i)^T}{\mathbf{C}_i^T \mathbf{S}_i} - \frac{\mathbf{S}_i^T (\Delta \mathbf{g}_i - \mathbf{B}_i \mathbf{S}_i) \mathbf{C}_i \mathbf{C}_i^T}{(\mathbf{C}_i^T \mathbf{S}_i)^2} \quad (8.125)$$

where $\Delta \mathbf{g}_i = \mathbf{g}_{i+1} - \mathbf{g}_i$ and the vector \mathbf{C}_i allows the choice of various strategies of improving the Hessian matrix. From a theoretical point of view, the best option is $\mathbf{C}_i = \mathbf{g}_i$. The process of improvement starts with the zero matrix \mathbf{B} .

Instead of approximating all components of the Hessian matrix \mathbf{H} , it is possible to improve only the part \mathbf{B} containing second derivatives. The matrix \mathbf{B}_{i+1} has the form

$$\mathbf{B}_{i+1} = \mathbf{J}_{i+1}^T \mathbf{J}_{i+1} + \mathbf{K}_{i+1} \quad (8.126)$$

The matrix \mathbf{K}_{i+1} is symmetrical and corresponds to the condition

$$\mathbf{K}_{i+1} \mathbf{S}_i = \mathbf{V}_i \quad (8.127)$$

The choice of \mathbf{V}_i causes variable final results. The matrix \mathbf{K}_{i+1} is again improved in each iteration according to the expression

$$\mathbf{K}_{i+1} = \mathbf{K}_i + \Delta \mathbf{K}_i \quad (8.128)$$

For determination of the matrix $\Delta \mathbf{K}_i$, Eq. (8.125) may be used, but with \mathbf{B}_i , replaced by \mathbf{K}_i and $\Delta \mathbf{g}_i$ replaced by \mathbf{V}_i . The vector \mathbf{V}_i may be computed from the *Broyden-Dennis formula* [37] where

$$\mathbf{V}_i^{(\text{BD})} = \Delta \mathbf{g}_i - \mathbf{J}_{i+1}^T \mathbf{J}_{i+1} \mathbf{S}_i \quad (8.129)$$

or the *Betts formula*

$$\mathbf{V}_i^{(\text{B})} = \Delta \mathbf{g}_i - \mathbf{J}_i^T \mathbf{J}_i \mathbf{S}_i \quad (8.130)$$

It is possible to use a linear combination of these two formulae or to use more complicated procedures of adaptive improvement of the matrix \mathbf{H} , described in detail in the literature [40].

The adaptive improvement of the Hessian matrix, by applying Eq. (8.126), does not automatically result in positive-definiteness of matrix \mathbf{B}_{i+1} . Therefore, this technique should be combined with procedures of pseudoinversion.

A mixed strategy is sometimes used when, according to parameter α , a direction \mathbf{V}_s is selected between the linearization \mathbf{L} and approximately Newton direction \mathbf{N} .

Gill and Murray [38] propose calculation of an approximation of the Hessian matrix by the difference formulae and then application of the SVD procedure for determination of significant components of the gradient. Many authors [39] recommend the method of variable metrics as a standard part of a library of programs for the minimization of the criterion function $U(\boldsymbol{\beta})$. It is best to restart the calculation of matrix \mathbf{B}_{i+1} when matrix \mathbf{H} , becomes unsuitable because of cumulated errors.

(c) Selection of step length

Many variants of the Gauss-Newton method use, for a selection of the optimal step α^* , the quadratic approximation $U(\mathbf{b} + \alpha \mathbf{L})$ in the direction \mathbf{L} . With values $U(\boldsymbol{\beta}^{(i)}) = U_i$ and $U(\boldsymbol{\beta}^{(i)} + \mathbf{L}_i) = U_{i+1}$ and the direction derivative S_D

$$\mathbf{g}_i^T \mathbf{L}_i = -2\mathbf{e}^T (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \mathbf{e}^T = S_D$$

the optimal step length of α^* is estimated from

$$\alpha_i^* = \frac{-S_D}{2(U_{i+1} - S_D - U_i)} \quad (8.131)$$

If $|S_D|$ is small, the value α_i^* is also small. Therefore $\alpha_i^* = \max(0.25, \alpha_i^*)$ may be

selected. There are various heuristic strategies for selection of convenient α values, and these can speed up an iterative search for a minimum.

8.5.2.2 Marquardt-type methods

The obvious selection of a directional vector \mathbf{V}_i is the direction of steepest descent, $-\mathbf{g}$. From Eq. (8.119), an optimum coefficient α^* is given by

$$\alpha^* = \mathbf{g}^T \mathbf{g} [\mathbf{g}^T \mathbf{H} \mathbf{g}]^{-1} \approx \mathbf{g}^T \mathbf{g} [\mathbf{g}^T (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{g}]^{-1} \quad (8.132)$$

The increment vector $\Delta_i = -\alpha^* \mathbf{g}$ corresponds to the gradient method.

The gradient method converges slowly in the vicinity of an optimum. On the other hand, in cases when $\boldsymbol{\beta}^{(i)}$ is far from \mathbf{b} , the direction leading to a minimum can be found. It is effective to use a combination of directions of the Newton method \mathbf{N}_i or the direction of linearization \mathbf{L}_i , together with the direction $-\mathbf{g}$, to construct a more robust algorithm. These procedures are called *hybrid procedures*. The best known example is the Marquardt method, which calculates the directional vector $\mathbf{V}_i(\lambda)$ from

$$\mathbf{V}_i(\lambda) = (\mathbf{J}^T \mathbf{J} + \lambda \mathbf{D}_i^T \mathbf{D}_i)^{-1} \mathbf{J}^T \hat{\mathbf{e}} \quad (8.133)$$

where λ is the parameter and \mathbf{D}_i is the diagonal matrix which eliminates the influence of various magnitudes of the components of the matrix \mathbf{J} . Usually the diagonal elements D_{ii} are equal to diagonal elements of matrix $(\mathbf{J}^T \mathbf{J})$. According to the magnitude of λ , the vector $\mathbf{V}_i(\lambda)$ has following properties:

- The length $\|\mathbf{V}_i(\lambda)\|$ is a decreasing function of λ . For $\lambda \rightarrow \infty$, $\|\mathbf{V}_i(\lambda)\| \rightarrow 0$. The parameter λ operates similarly to the parameter α , i.e. it enables a change in the length of the increment vector.
- The cosine of the angle between the vector $\mathbf{V}_i(\lambda)$ and the negative gradient $-\mathbf{g}_i$ increases as a function of λ . As $\lambda \rightarrow \infty$ it approaches a value of one. It then follows that for large λ values, the directional vector $\mathbf{V}_i(\lambda)$ approaches the directional vector of the gradient method.
- The cosine of the angle between the direction of linearization \mathbf{L}_i and the vector $\mathbf{V}_i(\lambda)$ is a decreasing function of λ . When $\lambda = 0$, it reaches value 1 and $\mathbf{V}_i(\lambda)$ is identical with the direction of the Gauss–Newton method.

The curve $\mathbf{V}_i(\lambda)$ in parameter space begins at the point $\boldsymbol{\beta}^{(i)} + \mathbf{L}_i$ and ends at the point $\boldsymbol{\beta}^{(i)}$, where it has direction $-\mathbf{g}$. However, the space curve does not lie in the plane of vectors \mathbf{L}_i and $-\mathbf{g}$. Appropriate selection of parameter λ ensures:

- positive definiteness of the matrix $\mathbf{R} = (\mathbf{J}^T \mathbf{J} + \lambda \mathbf{D}_i^T \mathbf{D}_i)$, which ensures that its inverse can be found;
- a shortening step $\mathbf{V}_i(\lambda)$ moving from the direction of linearization \mathbf{L}_i ;
- the possibility of choosing between the direction \mathbf{L}_i and approximate direction $-\mathbf{g}$;
- restriction of the magnitude of the incremental vector \mathbf{V}_i to the “admissible” region in the vicinity of $\boldsymbol{\beta}^{(i)}$.

The necessity of repeated matrix inversion for each λ is a disadvantage of this

procedure, as it is rather time-consuming. Moreover, it may occur that a large λ results in a very small magnitude of V_i . Therefore, the use of the maximum value of the magnitude of λ is limited. Individual modifications of the Marquardt method differ, especially in the strategy of the adaptive setting of λ . The original algorithm begins with $\lambda_0 = 0.01$. After each successful step, the calculation $\lambda_{i+1} = \lambda_i/10$ is performed, and after an unsuccessful step, $\lambda_{i+1} = 10\lambda_i$.

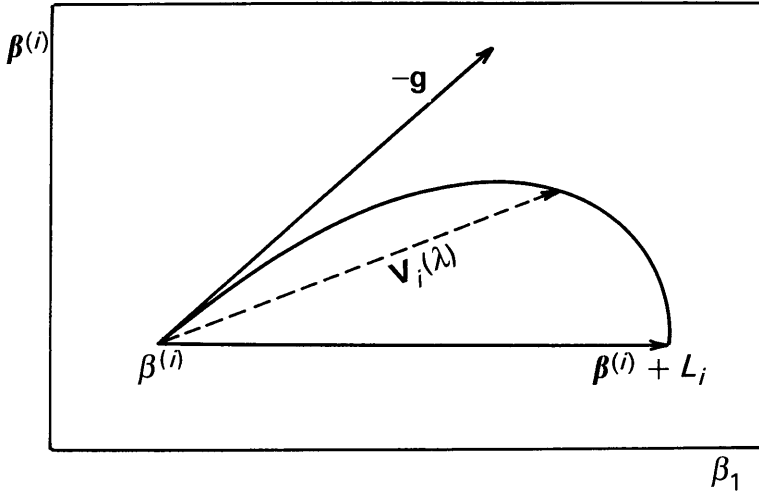


Fig. 8.16—Schematic path of function $V_i(\lambda)$ for two parameters ($m = 2$).

In the Nash algorithm [40], with minimal need for memory for solving a set of linear equations (i.e. determination of $V_i(\lambda)$), the Choleski decomposition is used. The adaptive adjustment of λ starts from $\lambda = 10^{-4}$ and after a successful step $\lambda_{i+1} = \max(10^{-6}, \lambda_i/10)$ is chosen. If a step is unsuccessful, $\lambda_{i+1} = \max(10^6, 4\lambda_i)$ is selected.

In another procedure [41], λ_i is selected according to the magnitude of the maximal diagonal element of matrix $(J^T J)$. Let us denote this element, J_{\max} . In each iteration, the procedure begins with $\lambda_1 = C_{\min} J_{\max}$. If there is no decrease in $U(\beta)$, this coefficient is increased according to

$$\lambda_{i+1} = \lambda_i (C_{\max}/C_{\min})^{1/4}$$

This process of increasing λ is performed until a decrease in $U(\beta)$ occurs or until λ is equal to J_{\max} . Then another search method is used, where $C_{\max} = 1$ and $C_{\min} = 0.01$ are chosen.

A very good and effective method was proposed by More [42]; this forms part of the program NL2SOL.

Some authors recommend choosing the optimal value of λ_{opt} as the one which leads to a maximal decrease in $U(\beta)$. This local optimal strategy, in cases of narrow curved-valley shape of $U(\beta)$, causes a global deceleration of convergence.

One strategy for changing λ , which ensures a global speeding up of convergence,

concerns three typical situations which may appear in the course of minimization of $U(\boldsymbol{\beta})$, [43].

- (a) It is possible to use the direction of linearization \mathbf{L} , i.e. $\lambda = 0$.
- (b) When the criterion function $U(\boldsymbol{\beta})$ has a curved-valley shape, the smallest eigenvalue C_1 of matrix $(\mathbf{J}^T \mathbf{J})^{-1}$ is significantly smaller than the second smallest eigenvalue C_2 , e.g. $10C_1 < C_2$. The direction of this valley is determined by the eigenvector \mathbf{k}_1 corresponding to the smallest eigenvalue. In this case, an increment Δ_i is chosen so that the criterion function $U(\boldsymbol{\beta})$ remains approximately unchanged. A search is then carried out on the hill-side of the valley [43].
- (c) In other cases, the curved valley is not so distorted and the quadratic approximation is inappropriate. For these cases, a suitable strategy for the selection of Δ_i has been described [43].

Meyer and Roth [44] have proposed the method MDLS (modified damped least-squares) for finding the vector $\mathbf{V}(\lambda)$, which involves a one-way minimization in a specific direction.

Generally, Marquardt-type methods are included in standard program libraries because of their robustness.

8.5.2.3 Dog-leg type procedures

The main disadvantages of Marquardt methods include:

- (a) the need for an inversion after every change of parameter λ ;
- (b) the small length of vector $\mathbf{V}(\lambda)$ for a large λ .

Both these disadvantages are removed in the next hybrid methods, in which the optimal directional vector $\mathbf{V}(\mu)$, a convex combination of vectors \mathbf{L} and $-\alpha^* \cdot \mathbf{g}_i$, is searched for. Here α^* is estimated from Eq. (8.132) and $0 \leq \mu \leq 1$. It follows that

$$\mathbf{V}(\mu) = \boldsymbol{\beta}^{(i)} + (1 - \mu)\mathbf{L}_i\alpha_1 - \mu\alpha^*\mathbf{g}_i \quad (8.134)$$

The function $\mathbf{V}(\mu)$ for $\alpha_1 = 1$ and $\alpha_1 < 1$ is shown in Fig. 8.17 as hypotenuses of right angle triangles. The dotted line represents $\alpha_1 < 1$ and the solid line $\alpha_1 = 1$. The classical strategy of the *Powell dog-leg method* estimates an optimal vector $\mathbf{V}_i(\mu)$ on the abscissa, TB, of a triangle defined by the vertices $\mathbf{O} = \mathbf{b}^{(i)}$, $\mathbf{T} = \mathbf{b}^{(i)} + \mathbf{L}_i$, and $\mathbf{B} = \mathbf{b}^{(i)} - \alpha^*\mathbf{g}_i$, where α^* is defined by Eq. (8.132).

It can be seen that for $\mu = 0$, the vector $\mathbf{V}(\mu)$ is identical to the linearization direction \mathbf{L}_i and for $\mu = 1$, the vector $\mathbf{V}(\mu)$ is identical to the direction of the negative gradient $-\mathbf{g}$. The magnitude of the total increment in direction $-\mathbf{g}$ corresponds to the optimal value α^* .

Dennis and Mei [45] use the *double dog-leg* strategy in which, instead of the vector \mathbf{L}_i , the “shorter” vector $\alpha_1 \mathbf{L}_i$ is used. The parameter α_1 is determined in the linearization direction so that the increment corresponds approximately to the Cauchy point [Eq. (8.119)]. Therefore, it can be shown [45] that

$$\alpha_1 = 0.2 + 0.8 \|\mathbf{g}_i\|^4 [\mathbf{g}_i^T (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{g}_i \mathbf{g}_i^T (\mathbf{J}^T \mathbf{J}) \mathbf{g}_i]^{-1}$$

From Fig. 8.17, it is evident that shortening $\alpha_1 \mathbf{L}_i$ leads to a directional vector $\mathbf{V}_1^*(\mu)$

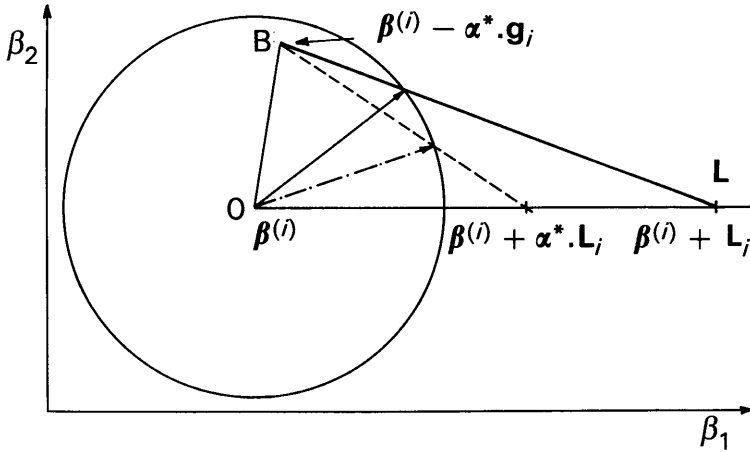


Fig. 8.17—Geometrical illustration of the dog-leg strategy. The circle shows the admissible range of increments. The solid hypotenuse is $V(\mu)$ for $\alpha_1 = 1$ and the dotted hypotenuse is $V(\mu)$ for $\alpha_1 = 1$.

which is nearer to the linearization direction than vector $V(\mu)$, calculated from Eq. (8.134) with $\alpha_1 = 1$. The actual strategy of these techniques differs in the admissible range adopted, in the method of inversion for finding L_i , and in improvement of the Hessian matrix by variable metric methods.

The program MINOPT, described later in this chapter, is based on the Dennis and Mei procedure [46].

8.5.3 Complications in nonlinear regression

In nonlinear regressions, many complications arise that are not found in linear regression models.

- (a) A minimum in $U(\beta)$ exists for some regression models only
- (b) There may be local minima and saddle points in $U(\beta)$.
- (c) Parameters may be inestimable.
- (d) Parameters may be ill-conditioned.

8.5.3.1 Parameter estimability

Complications (c) and (d) may be identified by analysing the sensitivity coefficients g_i defined by Eq. (8.5). For practical purposes the normalized sensitivity coefficients [47]

$$C_{ji} = \beta_j \frac{\delta f(x_i, \beta)}{\delta \beta_j}, \quad \begin{matrix} j = 1, \dots, m; \\ i = 1, \dots, n; \end{matrix} \quad (8.135)$$

are recommended. The ill-conditioning of parameters β_j and β_h is a consequence of the approximate multicollinearity between parameters β_j and β_h . For a visual interpretation of the examination of the conditioning of parameters in a model, the

sensitivity graph is used. The sensitivity graph is a plot of $C_{j(i)}$ and $C_{h(j)}$ vs. x_i , $i = 1, \dots, n$. The dependence of the normalized sensitivity coefficients on the index i may also be plotted.

Figure 8.18a shows several possible ill-conditioned models, where the sensitivity coefficients C_j and C_h are linearly dependent. Figure 8.18b, on the other hand, shows linearly independent sensitivity coefficients. More details of similar cases have been presented in the literature [47]. From Fig. 8.18, it follows that for the examination of the linear dependence of the sensitivity coefficients the location of points for $C = 0$ is important. The situation is more complicated when there are more parameters in the model [47].

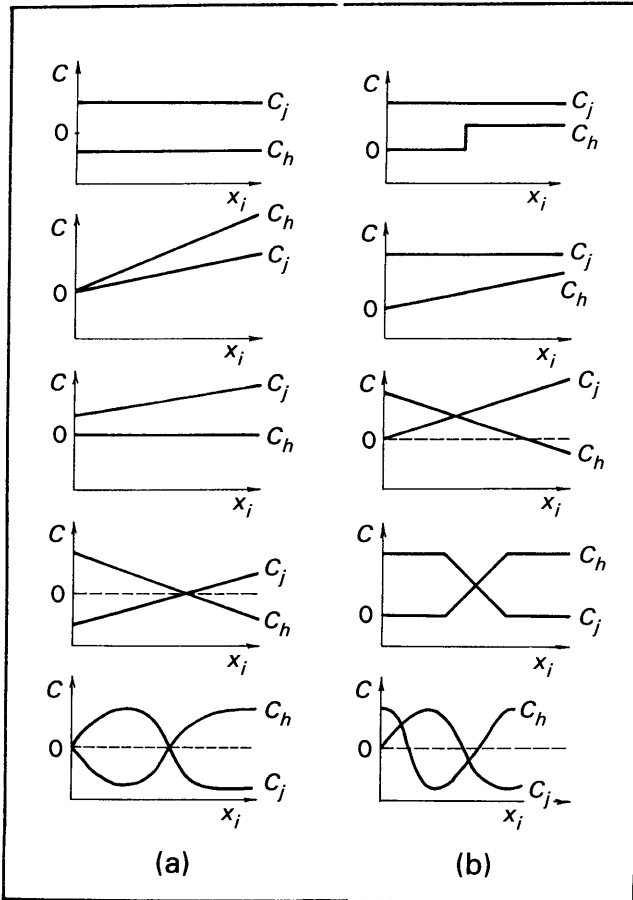


Fig. 8.18—The sensitivity graphs for the sensitivity coefficients where they are (a) linearly dependent, and (b) linearly independent.

To express the sensitivity of a regression model in terms of a change in parameter β_j , the *total sensitivity function* [1] C_{c_j} may be used. This function is defined by

$$C_{cj} = \frac{1}{n} \sum_{i=1}^n \left[\frac{\delta f(x_i, \beta)}{\delta \beta_j} \right]^2 \quad (8.136)$$

This sensitivity function has meaning only for the nonlinear parameters β_j in the regression model $f(x, \beta)$. For the linear parameters β_h in the regression model $f(x, \beta)$, the total sensitivity function C_{ch} reaches a constant value. A graphical illustration shows the *sensitivity graph of parameters*, which plots the dependence of C_{cj} on β_j in the vicinity of $\beta_j^{(0)}$ or \mathbf{b}_j . If the sensitivity graph of parameters is nearly constant, then the regression model has a low sensitivity to changes in the j th parameter, or the regression model $f(x, \beta)$ is linear with respect to the parameter β_j .

8.5.3.2 Existence of a minimum of $U(\beta)$

If, in the vicinity of a minimum, the model function $f(x, \beta)$ reaches an infinite value, then $U(\beta) \rightarrow \infty$. For example, for the model

$$f(x, \beta) = \frac{\beta_0 + \beta_1 x_{1i}}{\beta_2 x_{1i} + \beta_3 x_{2i}}$$

with two independent variables x_{1i} and x_{2i} such that $\beta_2 x_{2i} = -\beta_3 x_{1i}$, then $U(\beta) \rightarrow \infty$.

Gallant [48] has given an example when, with a simulated data set and an exponential model, for one parameter $U(\beta)$ does not have a minimum, but a maximum. Demidenko [49] has proposed a procedure for testing for the existence of a least-squares estimate, and therefore a minimum in $U(\beta)$, in which the regression criterion is restricted (from below) to a limit value. Generally, the existence of a minimum is connected with the problem of identifying the parameters \mathbf{b} of the regression model $f(x, \beta)$ [50].

A given regression model is *globally unidentified on a region Ω* , if, for any parameter vector $\mathbf{b} \in \Omega$, we can find another vector $\mathbf{b}^* \in \Omega$ for which

$$f(x, \mathbf{b}^*) = f(x, \mathbf{b}) \quad (8.137)$$

where the symbol $f(x, \beta)$ means the vector with elements $f(x_i, \beta)$.

When the columns of the Jacobian matrix \mathbf{J} are independent, the unidentifiability is *structural* in nature, i.e. independent of the actual numerical values of parameters. The cause of unidentifiability is symmetry of the parameters [51]. This means that the model $f(x, \beta)$ is invariant to a transformation of points in the parametric space. The condition of invariance of this function, with respect to the continuous transformation corresponding to Lie group, is given by [51]

$$\frac{\delta f(x, \beta)}{\delta \beta} \mathbf{h}(\beta) = 0 \quad (8.138)$$

where $\mathbf{h}(\beta)$ is the tangent vector which unambiguously defines continuous transformations of parameters, and $\delta f(x, \beta)/\delta \beta$ is the vector with components $\delta f(x, \beta)/\delta \beta_j$, $j = 1, \dots, m$. If we compare Eq. (8.138) with Eq. (8.5), we find that they express the same condition. This means that for unidentified models the sensitivity coefficients are linearly dependent.

8.5.3.3 Existence of local minima

The existence of local minima is characteristic of overdetermined models. Local minima exist in various models formed as a sum of partial nonlinear terms, sums of exponentials, etc. Let us suppose, for example, that for the model

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_1 \exp(\beta_2 x)$$

a global minimum $U(\boldsymbol{\beta}^*)$ with estimates, β_1^* and β_2^* , were found by the least-squares method. If, for the same data set, we use the model

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_1 \exp(\beta_2 x) + \beta_2 \exp(\beta_4 x)$$

we will notice that behind a global minimum many local stationary points will satisfy the condition, $\delta U(\boldsymbol{\beta})/\delta \boldsymbol{\beta} = 0$. There are points for which

- (a) $\beta_1 = \beta_3 = 0$ or $\beta_2 = \beta_4 = 0$, i.e. the model is reduced to a simplified form;
- (b) $\beta_4 = \beta_3 = \beta_2^*$ and all combinations of β_2, β_1 for which $\beta_1 + \beta_2 = \beta_1^*$, once again simplifying the model.

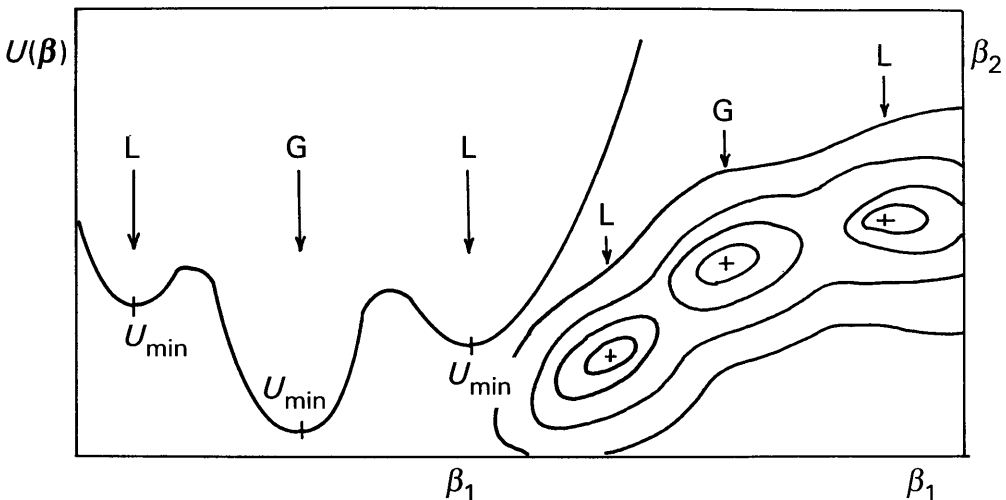


Fig. 8.19—Existence of local minima (L) beside the global one (G).

For higher numbers of exponential terms, the number of local minima increases sharply. Owing to the influence of measurement errors, the global minimum can be determined for a model with fewer terms [52].

A decision about whether a minimum found is global or local may be made on the basis of the Hessian matrix of the criterion function $U(\boldsymbol{\beta})$. If the Hessian matrix \mathbf{H} is positive-definite at point \mathbf{b} , a global minimum has been found [49]. For practical purposes, instead of the matrix \mathbf{H} , we examine the positive-definiteness of matrix $(\mathbf{J}^T \mathbf{J})$. If both matrices \mathbf{H} or $(\mathbf{J}^T \mathbf{J})$, are positive-definite, a global minimum has been found. Positive semi-definiteness indicates an overestimated or unidentified model.

8.5.3.4 Ill-conditioning of parameters

Ill-conditioning of parameters in a model which causes approximately linear dependence in the sensitivity graphs, depends not only on the type of regression model but also on the location and the range of experimental data. Ill-conditioning of parameters is indicated as ill-conditioning of the matrix $(\mathbf{J}^T \mathbf{J})$. Algorithms used for matrix inversion, such as pseudoinversion or the Marquardt type methods, are resistant against this difficulty. Difficulties arise when Gauss–Newton type methods are used. In many cases, ill-conditioning appears on application of numerical differentiation instead of analytical, because there is a loss of precision in the construction of matrix $(\mathbf{J}^T \mathbf{J})$. Numerical differentiation may cause even good algorithms to fail. For numerical differentiation, the method of forward difference is often used.

$$\frac{\delta f(x_i, \boldsymbol{\beta})}{\delta \beta_j} \approx \frac{1}{h_j} [f(x_i, h_j \mathbf{I} + \boldsymbol{\beta}) - f(x_i, \boldsymbol{\beta})] \quad (8.139)$$

where \mathbf{I} is a unit vector, with the j th component equal to 1, and other components equal to zero. The increment, h_j , in many cases determines the quality of numerical differentiation. It can be selected as

$$h_j = |\beta_j| + \sqrt{EP}$$

where EP is the computer precision. For nonlinear regression, it is convenient to use the algorithm written by Brown and Dennis [52], where h_j is evaluated by

$$h_j = \min(U(\boldsymbol{\beta}), \delta_j)$$

and the parameter δ_j is estimated by

$$\delta_j = \begin{cases} \rightarrow 10^{-9} & \text{if } |\beta_j| < 10^{-6} \\ \rightarrow 10^{-3} |\beta_j| & \text{if } |\beta_j| \geq 10^{-6} \end{cases}$$

This technique is also used in our algorithm MINOPT [46].

Reparameterization [53] may cause significant improvement in the shape of the criterion function $U(\boldsymbol{\beta})$ and hence improve the conditioning of the matrix $(\mathbf{J}^T \mathbf{J})$. Suitable reparameterization procedures for some nonlinear models are described by Ratkowsky [53].

8.5.3.5 Small range of experimental data

In many practical problems, difficulties with overdetermination and ill-conditioning of model parameters are partly the consequence of a small range of experimental data.

The application of the least-squares method leads then to a model which gives a good description of experimental dependence, but the parameters have no physical meaning. In these situations, it is possible

- (a) to collect more experimental data;
- (b) to examine the possibility of model simplification;
- (c) to investigate the possibility of estimating some parameters on the basis of other supplementary experiments, previous knowledge, experience, theory, etc.

- (d) to select suitable restrictions to be placed on the parameters so that their estimates have physical meaning.

The actual procedure depends on the problem in question and on the experience of the experimenter.

Problem 8.12. Search for a model of the Kohlrausch equation

The tension relaxation after a jump change of deformation on the value $\varepsilon = 0.4$ was studied for laboratory-made fibres PADt-G. The RETEST apparatus monitored the dependence between tension N_t (MPa) and time t (sec) in range 0–500 sec; and 21 points were recorded. The proposed model was

$$N_t = \beta_1 + (\beta_1 - \beta_2) \exp[-(\beta_4 t)^{\beta_3}] \quad (8.140)$$

where β_1 is the initial tension, β_2 is the equilibrium tension and β_3, β_4 are empirical constants. A graphical illustration of Eq. (8.140) is shown in Fig. 8.20.

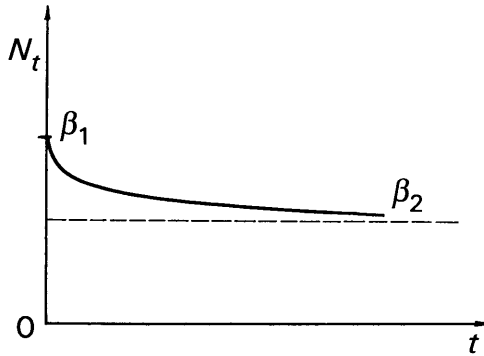


Fig. 8.20—The parameters of the Kohlrausch model.

The task is to estimate the parameters, by using five variants of calculation strategy.

Variant I: All parameters are estimated simultaneously.

Variant II: The model is simplified by assuming that $\beta_2 = 0$.

Variant III: The value for parameter $\beta_1 = 69$ MPa is fixed from previous knowledge.

Variant IV: The value for parameter $\beta_2 = 12$ MPa is fixed from previous knowledge.

Variant V: The values of parameters $\beta_1 = 69$ MPa and $\beta_2 = 12$ MPa are fixed.

Data: $n = 21$

t , sec	N_t , MPa	t , sec	N_t , MPa	t , sec	N_t , MPa
1	62.22	85	52.22	220	49.26
5	58.52	100	52.22	250	48.15
10	58.15	115	51.11	310	47.78
25	55.56	130	50.74	340	47.04
40	54.07	145	50.74	400	46.66
55	53.33	160	50.37	430	46.29
70	52.59	190	49.63	490	45.93

Solution: Variant I: From the initial guesses of the parameters $\beta_1^{(0)} = 65$, $\beta_2^{(0)} = 35$, $\beta_3^{(0)} = 0.118$, $\beta_4^{(0)} = 0.811$, refined values were found (Table 8.1). Parameter β_2 is negative and therefore has no physical meaning. Another initial guess of parameters leads to the refined estimates in column Ib of Table 8.1. The approximate singularity of matrix $\mathbf{J}^T \mathbf{J}$ causes problems.

Table 8.1. Parameter estimates \mathbf{b} and $U(\boldsymbol{\beta})$ with $\hat{\sigma}$ for the various variants of the model

Parameter estimates	Model variants					
	Ia	Ib	II	III	IV	V
b_1 , MPa	71.32	72.11	58.80	69*	58.27	69*
b_2 , MPa	-811.2	-4414	0*	-38.69	12*	12*
b_3	0.163	0.152	0.216	0.202	0.231	0.222
b_4	0.0106	0.0227	0.106	0.0691	0.121	0.131
$U(\mathbf{b})$	1.665	1.775	1.750	1.718	1.792	1.818
$\hat{\sigma}$	0.09797	0.104	0.0972	0.0954	0.0966	0.957

*constant values

Variant II: The refined parameters are in column II of Table 8.1. The simplification of the model leads to an insignificant increase in the minimum of the criterion function $U(\boldsymbol{\beta})$. However, the corresponding residual variance $\hat{\sigma}^2$ shows that this model is better than variant I. The matrix $\mathbf{J}^T \mathbf{J}$ is regular and well-conditioned.

Variant III: Even with the improvement in the model (lower $\hat{\sigma}$), the estimate of parameter β_2 has no physical meaning so that knowledge of parameter β_1 does not help here.

Variant IV: Knowledge of parameter β_2 leads to acceptable estimates of other parameters and the degree of fit is still quite good.

Variant V: With two parameters, β_1 and β_2 , known, the model becomes much simpler.

Conclusion: The goodness-of-fit for various variants of the model was compared by the residual standard deviation $\hat{\sigma} = \sqrt{U(\mathbf{b})/(n - m)}$. For a given data set, the individual variants of the model do not differ significantly. Moreover, it may be concluded that

- The Kohlrausch equation is inappropriate for practical purposes, because the quality of the results depends on the time units. After modification of the exponential term, $\exp[-(\beta_4 t)^{\beta_3}]$, the parameter β_3 is dimensionless and β_4 has dimensions of reciprocal time. This modification does not solve the problem of insufficient sample size.
- With regard to the physical meaning of model and knowledge about the data, the best method is a reduced model with $\beta_2 = 0$. Also, other physically acceptable variants IV and V hardly differ in the estimates of β_3 and β_4 .
- The estimate of parameter β_1 varies markedly depending on the choice of model variant.

From the results, it can be seen that without previous knowledge of the model system, it is not possible to estimate the equilibrium tension β_2 . Also, estimate β_1

becomes loaded by quite a high uncertainty which comes from the variants used.

The aim of Problem 8.11 was to demonstrate that numerical application of linear regression does not guarantee finding a useful model. Statistical analysis of this Problem is then necessary and is detailed in Section 8.6.

8.5.4 Examination of the reliability of the regression algorithm

In the literature, many regression algorithms and packages of programs for nonlinear regression are described [54, 55]. To examine the reliability of a regression algorithm, various test models are used. A good reliable algorithm should estimate correct values of the regression parameters.

For six test models, with their typical data sets, the final results depend on the initial guesses of the parameters. In comparison of the numerical results of these models, no restart or other technique of repeated determination of new initial guesses of parameters (if divergence occurred) was used.

Six test models:

- Model I. $y = \beta_1 + \beta_2 \exp(\beta_3 x)$
- Model II. $y = \exp(\beta_1 x) + \exp(\beta_2 x)$
- Model III. $y = \beta_1 \exp[\beta_2 / (\beta_3 + x)]$
- Model IV. $y = \beta_1 \exp(\beta_3 x) + \beta_2 \exp(\beta_4 x)$
- Model V. $y = \beta_1 x^{\beta_3} + \beta_2 x^{\beta_4}$
- Model VI. $y = \beta_1 [\exp(-\beta_2 x_1) + \exp(\beta_3 x_2)]$

Test data sets $\{x, y\}$ for

Model I: $n = 10$

x	1	5	10	15	20	25	30	35	40	50
y	16.7	16.8	16.9	17.1	17.2	17.4	17.6	17.9	18.1	18.7

Model II: $n = 10$

x	1	2	3	4	5	6	7	8	9	10
y	4	6	8	10	12	14	16	18	20	22

Model III: $n = 16$

x	50	55	60	65	70	75	80	85
y	34780	28610	23650	19630	16370	13720	11540	9744
	90	95	100	105	110	115	120	125
	8261	7030	6005	5147	4427	3820	3307	2872

Model IV: $n = 16$

x	7.448	7.448	7.969	8.176	9.284	9.439	7.552
y	57.544	53.546	19.498	16.444	4.305	3.006	45.290
	7.877	8.522	9.314	7.607	7.847	8.176	8.523
	27.952	11.803	4.764	51.286	31.623	21.777	13.996
	8.903	9.314					
	7.727	4.999					

Model V: $n = 12$

x	12	13	14	15	16	17	18	19	20
y	7.31	7.55	7.80	8.05	8.31	8.57	8.84	9.12	9.40
	21	22	23						
	9.69	9.99	10.30						

Model VI: $n = 23$

x_1	0	0.6	0.6	1.4	2.6	3.2	0.8	1.6	2.6	4.0
x_2	0	0.4	1.0	1.4	1.4	1.6	2.0	2.2	2.2	2.2
y	40	10	5.0	2.5	2.5	2.0	1.0	0.7	0.8	0.7
	1.2	2.0	4.6	3.2	1.6	4.2	4.2	3.2	2.8	
	2.6	2.6	2.8	3.0	3.2	3.4	3.4	3.8	4.2	
	0.4	0.4	0.3	0.22	0.22	0.1	0.05	0.07	0.03	
	4.2	5.4	5.6	3.2						
	4.2	4.4	4.8	5.0						
	0.03	0.03	0.02	0.01						

Table 8.2. Initial guess of parameters estimated for the six test models

Model	$\beta_1^{(0)}$	$\beta_2^{(0)}$	$\beta_3^{(0)}$	$\beta_4^{(0)}$	$U(\beta^{(0)})$
I	20	2	0.5	—	2×10^{23}
II	0.3	0.4	—	—	4×10^3
III	0.02	4000	250	—	1.7×10^9
IV	10^5	10^5	-1.679	-1.31	1.12×10^4
V	100	0.1	2	10	2.68×10^3
VI	12	1.0	25	—	226.9

Table 8.3. Best estimates of parameters of the six test models

Model	b_1	b_2	b_3	b_4	$U(\mathbf{b})$
I	15.67	0.9994	0.0222	—	5.98×10^{-3}
II	0.2578	0.2578	—	—	124.34
III	0.005618	6180	345.2	—	87.9
IV	8.315×10^7	5.088×10^3	-1.95	-0.7786	134
V	3.802	4.141×10^{-3}	0.223	2.061	2.98×10^{-5}
VI	31.5	1.51	19.9	—	1.25

Three regression methods, the method of modified simplex (MSM), the Gauss–Newton method (GN) and program MINOPT were tested. Tables 8.2 and 8.3 show the initial and final parameter estimates, respectively. Table 8.4 shows that program MINOPT works very reliably.

Table 8.4. Number of iterations necessary to reach a minimum $U(\boldsymbol{\beta})$ for the six test models (F means convergence to a false solution, S means a slow convergence)

Model	MSM	GN	MINOPT
I	F	55	22
II	179	S	10
III	2452	15	32
IV	F	S	42
V	F	F	65
VI	387	12	16

Models I–V were used to compare the ability of well known statistical packages to solve these problems. The packages tested were:

BMDP SOLO version 3.1	abbreviation SOLO
BMDP version 1987	abbreviation BMDP
SAS version 6.03	abbreviation SAS
SPSS PC+ version 3.1	abbreviation SPSS
STATGRAPHICS version 4.2	abbreviation STATGR
ASYSTANT + version 1.0	abbreviation ASYST
SYSTAT version 4.0	abbreviation SYSTAT
CHEMSTAT 1.1 (TRILOBYTE Ltd.)	abbreviation CHEMSTAT
MINSQ 3.12 (MICROMATH)	abbreviation MINSQ

For overall comparison, the performance index, PI , is defined as

$$PI = 100 * (\text{number of correct results}) / (T * \text{number of methods})$$

was computed. Here T is number of tests used. The greater PI , the better the package at solving nonlinear regression problems. ‘Number of methods’ means the number of optimization methods available in the packages. The possibility of combining methods (as in MINSQ) was not tested. The results are summarized in Table 8.4a.

The best results were obtained with ADSTAT and SPSS. Program MINSQ is very quick and can use a combination of methods. Other packages were not very satisfactory for these test problems. This comparison will disappoint many users of

Table 8.4b. Comparison of various packages for nonlinear regression

Package	PI (tests I–IV)	Number of methods
BMDP	25	2
SAS	25	4
SYSTAT	37.5	2
STATGR	50	1
ASYST	8.3	4
SPSS	100	1
ADSTAT	100	1
SOLO	20	1
MINSQ	80	1

the standard statistical packages, showing, as it does, that errors due to false optimum location can cause failure of the whole regression analysis.

8.6 STATISTICAL ANALYSIS OF NONLINEAR REGRESSION

Statistical analysis in nonlinear regression depends on the model used, the model of measurement errors and the criterion function. Let us limit ourselves to the method of maximum likelihood when the estimates \mathbf{b} minimize the logarithm of maximum likelihood, $l(\boldsymbol{\beta}) = \ln L(\boldsymbol{\beta})$, as defined by Eq. (8.28).

In the construction of confidence intervals for parameters $\boldsymbol{\beta}$ or in the testing of statistical hypotheses, three main approaches are used.

(a) *The method of linearization* is based on the asymptotic normality of the vector $\sqrt{n}(\mathbf{b} - \boldsymbol{\beta})$ and an estimate of the variance $D(\mathbf{b})$, defined by Eq. (8.29). As was shown in Problem 8.8, for additive and normally distributed errors:

$$D(\mathbf{b}) = \sigma^2(\mathbf{J}^T \mathbf{J})^{-1}$$

This expression corresponds to a *linearization* of the regression model [Eq. (8.66a)]. Here, the same expression is used as in a linear regression but matrix \mathbf{J} replaces matrix \mathbf{X} .

(b) *The method of Lagrange multipliers* is based on the fact that the quadratic form

$$QF = \mathbf{U}^T \mathbf{I}^{-1} \mathbf{U} \quad (8.141)$$

has asymptotically the $\chi^2(m)$ distribution, and the matrix \mathbf{I} is defined either by Eq. (8.30) or Eq. (8.31). The vector \mathbf{U} , of dimension $(m \times 1)$, has components $\delta l(\boldsymbol{\beta}) / \delta \beta_j$.

(c) *The method of the maximum likelihood ratio* is based on:

$$P(\boldsymbol{\beta}) = \frac{L(\boldsymbol{\beta})}{L(\mathbf{b})} \quad (8.142)$$

The random variable, $\{-2 \ln[P(\boldsymbol{\beta})]\}$, has asymptotically the $\chi^2(m)$ distribution.

The method of linearization is, in practice, the most widely used one as it does not involve additional calculation. The method of Lagrange multipliers and the maximum likelihood ratio require a numerical search for the roots of nonlinear functions, but lead to more accurate and more useful results.

Statistical properties of the least-squares method (LS)

The LS method is a special case of the maximum likelihood method for an additive model of measurement and an independent normal distribution of measurement errors. Gallant [5] derived the following equations:

$$\mathbf{b} = \boldsymbol{\beta}^* + (\mathbf{J}_t^T \mathbf{J}_t)^{-1} \mathbf{J}_t^T \boldsymbol{\varepsilon} + \text{term}(1/\sqrt{n}) \quad (8.143)$$

$$\sigma^2 = \frac{\boldsymbol{\varepsilon}^T (\mathbf{E} - \mathbf{J}_t (\mathbf{J}_t^T \mathbf{J}_t)^{-1} \mathbf{J}_t^T) \boldsymbol{\varepsilon}}{n - m} + \text{term}(1/n) \quad (8.144)$$

Here $\boldsymbol{\beta}^*$ is the true value of the parameters in the model, \mathbf{J}_t is the Jacobian matrix evaluated at the theoretical point $\boldsymbol{\beta}^*$ and term $(1/\sqrt{n})$ denotes a random quantity with its mean value equal to $1/\sqrt{n}$. If in Eqs. (8.143) and (8.144), the higher terms are neglected, then from the theory of linear regression, it follows that the distribution of the random quantity \mathbf{b} is m -dimensionally normal

$$\mathbf{b} \sim N[\boldsymbol{\beta}, \sigma^2 (\mathbf{J}^T \mathbf{J})^{-1}] \quad (8.145)$$

The random variable, $(n - m)\hat{\sigma}^2/\sigma$, has a $\chi^2(n - m)$ distribution and \mathbf{b} and $\hat{\sigma}^2$ are independent. In practice, the matrix \mathbf{J}_t is replaced by matrix \mathbf{J} , evaluated at point \mathbf{b} . The asymptotic normality of estimates \mathbf{b} , determined by the least-squares method, does not require normality of errors $\boldsymbol{\varepsilon}$, [5]. Application of Eq. (8.143) requires the following conditions to be fulfilled:

- (1) The regression model must be twice differentiable.
- (2) The regression model must be identifiable; that is, the function

$$U(\boldsymbol{\beta}^{(k)}) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n [f(x_i, \boldsymbol{\beta}^{(k)}) - f(x_i, \boldsymbol{\beta})]^2$$

should have an unambiguous minimum at the point $\boldsymbol{\beta}^{(k)} = \boldsymbol{\beta}$.

- (3) The matrix $\mathbf{Q} = \lim_{n \rightarrow \infty} (1/n) \mathbf{J}_t^T \mathbf{J}_t$ must be asymptotically regular.

Practical examination of all these conditions is rather complicated [5].

For real experimental data, the estimates \mathbf{b} and other statistical characteristics are *biased*, and therefore the application of Eq. (8.145) is limited.

The usefulness of the statistical analysis of nonlinear regression models by least-squares methods depends on the magnitude of the bias, and this depends on the degree of nonlinearity in the regression model.

The covariance matrix of parameter estimates

From Eq. (8.145), it follows that the asymptotic covariance matrix of estimates \mathbf{b} obtained by the LS method is expressed by

$$D(\mathbf{b}) = \sigma^2 (\mathbf{J}^T \mathbf{J})^{-1} \quad (8.146)$$

When errors are independent and identically distributed with constant variance, it is possible to find a more accurate approximation, based on a linearization of estimate $\mathbf{b}(\boldsymbol{\varepsilon})$ as a function of $\boldsymbol{\varepsilon}$:

$$D(\mathbf{b}) = 4\sigma^2 \mathbf{H}^{-1} (\mathbf{J}^T \mathbf{J}) \mathbf{H}^{-1} \quad (8.147)$$

When $(\mathbf{J}^T \mathbf{J}) \approx 0.5 \mathbf{H}$ is simplified, a less accurate approximation is obtained:

$$D(\mathbf{b}) = 2\sigma^2 \mathbf{H}^{-1} \quad (8.148)$$

When the residuals $\hat{\mathbf{e}}$ are small or the elements of matrix \mathbf{B} in Eq. (8.62) are approximately zero, then

$$\mathbf{H}^{-1} = 0.5(\mathbf{J}^T \mathbf{J})^{-1}$$

The effect of the application of Eqs. (8.146), (8.147) and (8.148) on the accuracy of the estimate of the confidence regions of parameters \mathbf{b} has been studied [57], and it was concluded that the two more accurate relations, Eqs. (8.147) and (8.148), are not significant. For practical calculations, the asymptotic formula [Eq. (8.146)] is obviously acceptable.

With a knowledge of the covariance matrix $D(\mathbf{b})$, either the variance of individual parameters $D(b_i)$ or the correlation coefficients r_{ij} between estimates b_i and b_j , may be estimated. From Eq. (8.146), we can write

$$D(b_j) = \sigma^2 V_{jj} \quad (8.149)$$

where V_{jj} are the diagonal elements of the matrix $\mathbf{V} = (\mathbf{J}^T \mathbf{J})^{-1}$. Similarly, the correlation coefficient between the parameters b_i and b_j is

$$r_{ij} = \frac{V_{ij}}{\sqrt{V_{jj} V_{ii}}} \quad (8.150)$$

If the value of r_{ij} is close to unity, the estimates b_i and b_j are linearly dependent and the model is overdetermined or ill-conditioned with respect to parameters b_i and b_j .

8.6.1 Degree of nonlinearity of a regression model

For characterization of nonlinear behaviour in regression models, the intrinsic curvature K_h^N [Eq. (8.75a)], the parameter-effects curvature K_h^P [Eq. (8.75b)] or the maximum intrinsic curvature Γ^N [Eq. (8.75c)] and the maximum parameter-effects curvature Γ^P [Eq. (8.75d)] can be adopted. If Γ^N and Γ^P are sufficiently small [56] for statistical analysis and for construction of confidence intervals, the *linearization of regression model* [Eq. (8.145)] may be used.

From many practical experiments, it has been concluded that

- the influence of nonlinearity may be, in many cases, removed by a suitable reparameterization when Γ^N is small and Γ^P high;
- the efficiency of reparameterization also depends on the data;
- in some cases, even for high values of Γ^N and Γ^P , linearization may be applied. There are cases, however, when even for small values of Γ^N and Γ^P the confidence intervals estimated after linearization are not suitable.

The measures of nonlinearity based on a curvature are by no means universal.

8.6.1.1 Bias of parameter estimates

For expressing the bias of parameter estimates,

$$\mathbf{h} = E(\mathbf{b} - \boldsymbol{\beta}^*)$$

many approximations exist in the literature. For simplicity, we use the following definition of parameter bias [58]

$$\mathbf{h} = (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \mathbf{d} \quad (8.151)$$

where \mathbf{d} is the $(n \times 1)$ vector with the components

$$d_i = \frac{-\sigma^2 \text{tr}[(\mathbf{J}^T \mathbf{J})^{-1} \mathbf{G}_i]}{2} \quad (8.152)$$

where $\text{tr}(\cdot)$ denotes the trace of a matrix and \mathbf{G}_i is the matrix of second derivatives of the model function [Eq. (8.65)]. The vector \mathbf{d} is the expected value of the difference between the linear approximation [Eq. (8.66a)] and the quadratic approximation of the model function. Equation (8.151) enables the bias \mathbf{h} to be found from the coefficients (parameters) of the linear regression model,

$$\mathbf{d} = \mathbf{J} \mathbf{h} \quad (8.153)$$

It is obvious that the bias \mathbf{h} will be small if

(a) the vector \mathbf{d} is perpendicular to the tangent hyperplane, defined by columns of matrix \mathbf{J} , such that

$$\mathbf{J}^T \mathbf{d} = \mathbf{0}$$

(b) the vector elements \mathbf{d} are small, i.e. the increment of the quadratic term will be insignificant and the model $f(\mathbf{x}, \boldsymbol{\beta})$ will be well linearized.

Similarly, the bias of residuals is given by

$$\hat{e}_i = y_i - f(x_i, \mathbf{b}).$$

The mean value of the vector of residuals,

$$\mathbf{E} = E(\hat{\mathbf{e}}) \quad (8.154a)$$

can be rewritten as

$$\mathbf{E} = (\mathbf{E} - \mathbf{P}) \mathbf{d} \quad (8.154b)$$

where \mathbf{P} is the projection matrix [Eq. (8.68)]. The mean value of the residuals is called the *residuals bias* because it is assumed that $E(\mathbf{e}) = \mathbf{0}$.

For practical calculation, we often use the relative bias of the parameter estimates defined by

$$h_{\mathbf{R},j} = \frac{h_j}{b_j} \times 100[\%] \quad (8.154c)$$

The bias of estimates is considered significant if $h_{\mathbf{R},j} > 1\%$. For such biased estimates,

the statistical analysis based on linearization of regression model cannot be legitimately used.

Problem 8.13. *Calculation of the relative bias of parameters*

Calculate the relative bias of parameters $h_{R,j}$ of individual variants of the Kohlrausch model from Problem 8.12.

Data: from Problem 8.12

Solution: The values calculated for the relative bias of parameter estimates for all the variants of Problem 8.11 are given in Table 8.5.

Table 8.5. Relative bias $h_{R,j}$ of the parameters of the Kohlrausch model

Relative bias, %	Variant of model					
	I _a	I _b	II	III	IV	V
$h_{R,1}$	-1.0×10^{-4}	-1.8×10^6	0.18	—*	0.163	—*
$h_{R,2}$	-4.1×10^6	-6.0×10^8	—*	379.3	—*	—*
$h_{R,3}$	8.62×10^4	1.46×10^7	0.058	0.44	0.058	0.007
$h_{R,4}$	4.12×10^6	3.52×10^9	1.49	—22.5	1.402	1.42

Conclusion: The original model (Ia, Ib) is very ill-conditioned and has high values of bias, indicating the inadequacy of the proposed model. For variants II, IV and V, the parameter estimates have physical meaning and the bias values are sufficiently small. Therefore for these three variants, the statistical analysis based on linearization of the model may be used.

To express the *total bias* of parameter estimates, Box [59] proposed the scalar characteristic

$$\hat{M} = \frac{\mathbf{h}^T(\mathbf{J}^T\mathbf{J})\mathbf{h}}{m\hat{\sigma}^2} \quad (8.155)$$

and proved that the scalar \hat{M} is related to the maximum parameter-effect curvature caused by parameters, by following the inequality

$$\hat{M} \leq \frac{(\Gamma^P)^2}{4}$$

Similarly, for the norm of the residual bias, the following inequality applies

$$\|\mathbf{E}\| \leq 0.5\Gamma^N\hat{\sigma}\sqrt{m}$$

The bias of parameters is related to the curvature caused by parameters K_h^P and the residual bias is related to the intrinsic curvature K_h^N .

Problem 8.14. *Kinetic parameters of dyeing*

A modified polyester fibre, Tesil 32, was dyed by a disperse dye, C.I. Dispersion Red 54, and the kinetics of the process were studied. The relative concentration of dye C_i on the fibre was measured as a function of time, t_i . The kinetics of the isothermal dyeing may be expressed by the Cegarra–Puente model

$$C_i = \beta_1 \sqrt{1 - \exp(-\beta_2 t_i)}$$

where β_1 is the relative equilibrium concentration of the dye on a fibre and β_2 is the rate constant. Estimate b_1 and b_2 . Calculate the correlation r_{12} between the parameters β_1 and β_2 , the total bias of parameters \hat{M} and the norm of vector, $\|\mathbf{E}\|$.
Data: $n = 6$

t , min	20.0	40.0	60.0	80.0	100.0	120.0
C , %	43.5	53.6	64.1	66.5	72.0	76.5

Solution: From the initial guesses of parameters $\beta_1^{(0)} = 100$, $\beta_2^{(0)} = 10^{-4}$ with $U(\beta^{(0)}) = 12300$, the minimum found was $U(\mathbf{b}) = 12.88$, and the best estimates were $b_1 = 82.37\%$ and $b_2 = 0.0147 \text{ min}^{-1}$. The correlation coefficient $r_{12} = -0.963$ indicates strong multicollinearity. The values of the relative bias of parameters $h_{R,1} = 0.46\%$ and $h_{R,2} = 0.5\%$ indicate low insignificant bias. The scalar value of the total bias of parameters $\hat{M} = 0.1316$ and the norm of vector $\|\mathbf{E}\| = 3.58$ shows low bias also.

Conclusion: Linearization may be used for the statistical analysis of this model.

Because the bias of parameters is related to the curvature caused by parameters, it may be affected by reparameterization. Let us suppose that, instead of parameters β , the transformed parameters γ defined by Eq. (8.6) are used. Each new parameter, γ_j , is a function of all components of the vector β such that $\gamma_j = l(\beta)$. The bias of the parameter estimates c_j of parameters γ_j , is given by

$$h_j(\mathbf{c}) = \mathbf{l}^T \mathbf{h}(\mathbf{b}) + \frac{\sigma^2}{2} \text{tr}[\mathbf{M}(\mathbf{J}^T \mathbf{J})^{-1}] \quad (8.156)$$

where \mathbf{l} is the vector comprising first derivatives of the transformation with elements

$$l_j = \frac{\delta l(\beta)}{\delta \beta_j} \quad (8.156a)$$

and the matrix \mathbf{M} contains the second derivatives of the transformation with elements

$$M_{jk} = \frac{\delta^2 l(\beta)}{\delta \beta_j \delta \beta_k} \quad (8.156b)$$

The symbol $\mathbf{h}(\mathbf{b})$ denotes the bias, \mathbf{h} , of the parameter estimates, \mathbf{b} , calculated from Eq. (8.151). If the actual reparameterization is known, the change of bias of the parameters may be predicted. Another task is to select the reparameterization that leads to the smallest bias.

Problem 8.15. *Change of bias of parameters after model reparameterization*

Determine the change in the bias of estimate b_1 in the model $f(x, \beta) = \beta_1 \exp(\beta_2 x)$ by reparameterization, $\gamma_1 = \ln \beta_1$.

Solution: The model reparameterization of

$$f(x, \beta) = \beta_1 \exp(\beta_2 x)$$

leads to the function

$$f(x, \gamma_1, \beta_2) = \exp(\gamma_1 + \beta_2 x)$$

The vector \mathbf{l} has elements, $\mathbf{l} = [1/b_1 \quad 0]$ and matrix \mathbf{M} has the form

$$\mathbf{M} = \begin{bmatrix} 1/b_1^2 & 0 \\ 0 & 0 \end{bmatrix}$$

On substitution into Eq. (8.156), we get

$$h_1(\mathbf{c}) = \frac{h_1(b_1)}{b_1} - \frac{\sigma^2 V_{11}}{2 b_1^2} = \frac{h_1(b_1) - 0.5D(b_1)/b_1}{b_1}$$

where V_{11} is the first diagonal element of the matrix $(\mathbf{J}^T \mathbf{J})^{-1}$ and $D(b_1)$ is the variance of b_1 .

Conclusion: The decrease of bias, $h_1(b_1) - h_1(\mathbf{c})$, is bigger for big variances of estimate b_1 , and smaller for small values of b_1 . It is assumed that b_1 is always positive.

8.6.1.2 Asymmetry of parameter estimates

Nonlinearity of the regression model results in an unsymmetrical distribution of estimates \mathbf{b} . The measure of nonlinearity is a measure of the asymmetry of the estimates. Ratkowsky [53] used n points generated as identically distributed quantities, ε_i^* , with mean value equal to zero and variance $\hat{\sigma}^2$. Then, the random dependent variables y_i^+ and y_i^- are generated by

$$y_i^+ = f(x_i, \mathbf{b}) + \varepsilon_i^* \quad (8.157a)$$

$$y_i^- = f(x_i, \mathbf{b}) - \varepsilon_i^* \quad (8.157b)$$

Parameter estimates obtained by the least-squares method when using \mathbf{y}^+ , instead of \mathbf{y} , are denoted \mathbf{b}^+ and similarly, \mathbf{b}^- when using \mathbf{y}^- , instead of \mathbf{y} . Linear models are symmetrical, in that $(b_i^+ - \beta_i^*) = -(b_i^- - \beta_i^*)$. For nonlinear models, there is no such symmetry. A convenient measure of asymmetry is the expression,

$$\psi_i = \frac{1}{2} [(b_i^+ - \beta_i^*) + (b_i^- - \beta_i^*)] \quad (8.158)$$

The mean value, $E(\psi_i)$, is equal to the bias, h_i , and the variance $D(\psi_i)$ is given by

$$D(\psi_i) = 0.25D(b_i^+) + 0.25D(b_i^-) + 0.25 \text{cov}(b_i^+, b_i^-) \quad (8.159)$$

Nonlinearity is indicated by the ratio

$$\lambda_{Ni} = \frac{D(\psi_i)}{D(b_i)}$$

When $\lambda_{Ni} < 0.01$, the distribution of parameter estimates is nearly symmetrical. For $\lambda_{Ni} > 0.01$, the distribution of estimate b_i is strongly asymmetrical and the model,

with respect to this parameter, is strongly nonlinear.

Morton [60] published expressions relating the statistical measure, λ_{Ni} , to the bias and other measures of nonlinearity. The great advantage of the measure λ_{Ni} is that it is based on statistical arguments. For the determination of λ_{Ni} , Ratkowsky [53] generated many estimates \mathbf{b}^+ and \mathbf{b}^- from various vectors of errors, $\mathbf{\varepsilon}^*$.

8.6.2 Interval estimates of parameters

Point estimates \mathbf{b} of regression parameters $\boldsymbol{\beta}$ are, from a statistical point of view, worthless as they do not mention the intervals in which a true value $\boldsymbol{\beta}^*$ may be expected. The estimates \mathbf{b} are random quantities estimated from a sample of size n , $\{x_i, y_i\}$, $i = 1, \dots, n$. The confidence regions, simultaneous confidence regions and confidence intervals for multivariate samples are constructed similarly. For their elucidation, the same rules as for univariate data are applied (Chapter 3).

In nonlinear regression models for construction of confidence regions and intervals, the linearization often used has elliptic confidence regions. However, linearization is useful only when the model is *not* strongly nonlinear and the nonlinearity measures of asymmetry and bias are small. More accurate confidence regions can be found by using Lagrange multipliers or the likelihood ratio; these are non-elliptic and do not have to be continuous.

8.6.2.1 Confidence regions of parameters

From the normality of estimates \mathbf{b} , it follows that the quadratic form

$$Q = (\boldsymbol{\beta}^* - \mathbf{b})^T D(\mathbf{b})^{-1} (\boldsymbol{\beta}^* - \mathbf{b}) \quad (8.159a)$$

has the $\chi^2(m)$ distribution. The corresponding $100(1 - \alpha)\%$ confidence region of parameters $\boldsymbol{\beta}^*$ is the m -dimensional ellipsoid with boundaries expressed by

$$(\boldsymbol{\beta}^* - \mathbf{b})^T D(\mathbf{b})^{-1} (\boldsymbol{\beta}^* - \mathbf{b}) = \chi_{1-\alpha}^2(m) \quad (8.160)$$

where $\chi_{1-\alpha}^2(m)$ is the $100(1 - \alpha)\%$ quantile of $\chi^2(m)$ with m degrees of freedom. The centre of this ellipsoid is at the point \mathbf{b} .

For the least-squares method

$$D(\mathbf{b}) = \frac{U(\mathbf{b})(\mathbf{J}^T \mathbf{J})^{-1}}{n - m}$$

After substitution into Eq. (8.159), the quadratic form can be formulated as

$$\left[\frac{\Delta \boldsymbol{\beta}^T (\mathbf{J}^T \mathbf{J}) \Delta \boldsymbol{\beta}}{\sigma^2 m} \right] \left[\frac{U(\mathbf{b})}{\hat{\sigma}^2 (n - m)} \right]^{-1} = \frac{\Delta \boldsymbol{\beta}^T (\mathbf{J}^T \mathbf{J}) \Delta \boldsymbol{\beta}}{\hat{\sigma}^2 m} \quad (8.161a)$$

where $\Delta \boldsymbol{\beta} = \boldsymbol{\beta}^* - \mathbf{b}$. This quadratic form has the distribution $\chi^2(m)/\chi^2(n - m) = F(m, n - m)$, i.e. the Fisher-Snedecor distribution with m and $(n - m)$ degrees of freedom. The confidence ellipsoid then has the boundary

$$\Delta \boldsymbol{\beta}^T (\mathbf{J}^T \mathbf{J})^{-1} \Delta \boldsymbol{\beta} = m \hat{\sigma}^2 F_{1-\alpha}(m, n - m) \quad (8.161b)$$

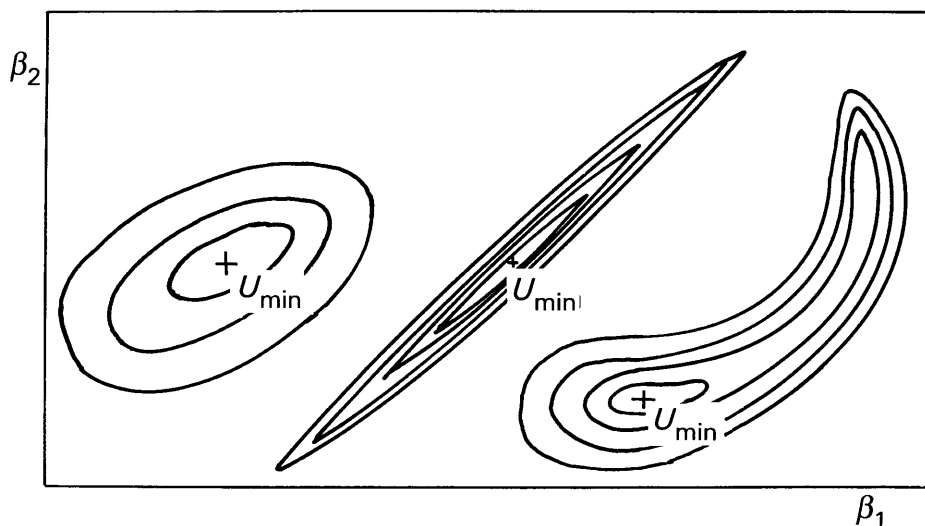


Fig. 8.21—Contours for various hyperparaboloid shapes in the vicinity of the minimum (pit).

In some chemometrics programs, instead of confidence ellipsoids, the boundary of the last contour of the least squares criterion is used. (Fig. 8.21).

When the bias of parameters, \mathbf{h} , is calculated, the correction $\Delta\beta^* = \mathbf{b} - \mathbf{h} - \beta^*$ may be used instead of $\Delta\beta$. To express the geometry of the confidence ellipsoids, the decomposition of the matrix $(\mathbf{J}^T\mathbf{J})^{-1}$, the eigenvalues λ_i and eigenvectors \mathbf{V}_i may be introduced such that

$$(\mathbf{J}^T\mathbf{J})^{-1} = \mathbf{V}\boldsymbol{\lambda}\mathbf{V}^T \quad (8.162a)$$

where \mathbf{V} is the matrix containing the eigenvectors in its columns and the diagonal matrix $\boldsymbol{\lambda}$ contains eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$ on the diagonal. For the corresponding decomposition of the matrix $(\mathbf{J}^T\mathbf{J})$, we have

$$(\mathbf{J}^T\mathbf{J}) = \mathbf{V}^T\boldsymbol{\lambda}^{-1}\mathbf{V} \quad (8.162b)$$

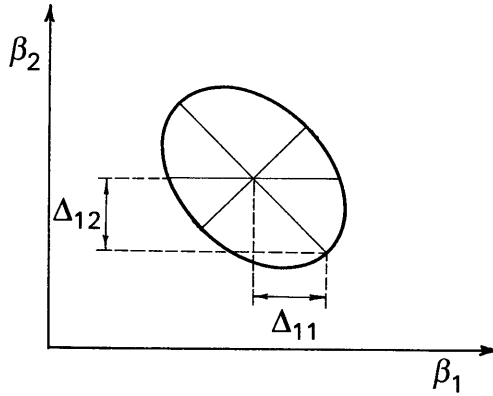
The matrix $\boldsymbol{\lambda}^{-1}$ is, once again, a diagonal matrix with reciprocal values λ_j^{-1} , $j = 1, \dots, m$ as elements. After substitution from Eq. (8.162a) into Eq. (8.161), we have

$$\Delta\beta^T\mathbf{V}^T\boldsymbol{\lambda}^{-1}\mathbf{V}\Delta\beta = \mathbf{Y}^T\boldsymbol{\lambda}^{-1}\mathbf{Y} = \sum_{i=1}^m \frac{1}{\lambda_i} Y_i^2 \quad (8.163)$$

Here $\mathbf{Y} = \mathbf{V}\Delta\beta$, which is the new orthogonal set of co-ordinates having the important property that the axes of the confidence ellipsoid are identical with the axes of the co-ordinate system. If we introduce the notation

$$p^2 = m\hat{\sigma}^2 F_{1-\alpha}(m, n-m)$$

the confidence ellipsoid can be expressed by the simple formula

Fig. 8.22—Graphical illustrations of the projections Δ_{11} and Δ_{12} .

$$\sum_{i=1}^m \frac{Y_i^2}{\lambda_i} = p^2 \quad (8.164)$$

The lengths of the half-axes of the ellipsoid are equal to $p\sqrt{\lambda_i}$. The projection Δ_{jk} of the j th half-axis into the axis of parameter β_k , is given by

$$\Delta_{jk} = p|V_{kj}\sqrt{\lambda_j}| \quad (8.165)$$

where V_{kj} is the k th element of the vector \mathbf{V}_j which is the j th column of matrix \mathbf{V} .

Problem 8.16. *Confidence ellipsoid for the parameters of the Cegarra–Puente model*
Estimate the 95% confidence ellipsoid for parameters β_1 and β_2 of the Cegarra–Puente model given in Problem 8.13.

Data: from Problem 8.13

Solution: With use of Eq. (8.164) the co-ordinates of the confidence region in systems Y_1 and Y_2 were calculated, then, by reverse transformation, β_1 and β_2 . The results are shown in Fig. 8.23.

Conclusion: The elongated shape of the confidence ellipsoid and its orientation prove strong negative correlation between parameters β_1 and β_2 .

When the dimension, m , of the parameter vector is greater than 2, the partial confidence ellipsoid is constructed for only two of the model parameters. Let us assume that the vector $\boldsymbol{\beta}^* = (\beta_1^*, \beta_2^*)^T$ is known, and that the confidence region for parameters $\boldsymbol{\beta}_2^*$ is to be constructed. As for linear regression models, the boundaries of the $100(1 - \alpha)\%$ confidence ellipsoid are found from

$$(\boldsymbol{\beta}_2^* - \mathbf{b}_2)^T \mathbf{D}_2^{-1} (\boldsymbol{\beta}_2^* - \mathbf{b}_2) = q\hat{\sigma}^2 F_{1-\alpha}(q, n - m) \quad (8.166)$$

where q is the dimension of vector $\boldsymbol{\beta}_2^*$. The matrix \mathbf{D}_2 is of dimension $(q \times q)$ and is formed from the matrix $(\mathbf{J}^T \mathbf{J})^{-1}$ by omitting $(m - q)$ rows and $(m - q)$ columns. Before constructing the confidence region, the parameters must be renumbered to make the parameters for which the ellipsoid is constructed the last ones.

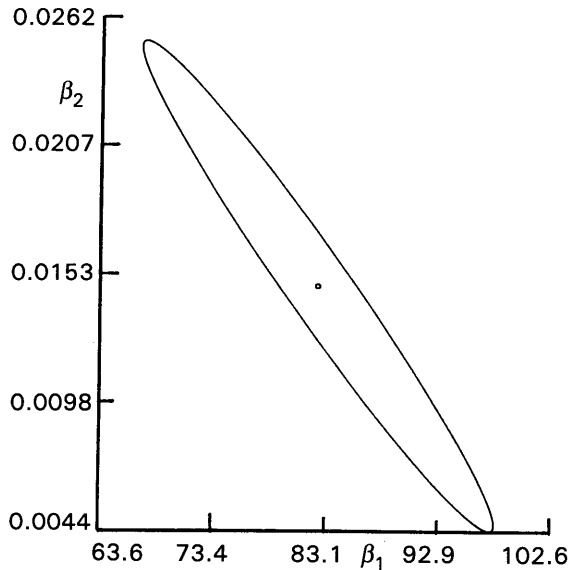


Fig. 8.23—The simultaneous confidence region of parameters β_1 and β_2 for the Cegarra-Puente model.

For the method of Lagrange multipliers, the quadratic QF from Eq. (8.141) is directly applied.

In the general case of the maximal likelihood method, the boundary of the $100(1 - \alpha)\%$ confidence region is defined by

$$\mathbf{U}^T \mathbf{I}^{-1} \mathbf{U} = \chi^2_{1-\alpha}(m) \quad (8.167a)$$

Various formulae can be derived, depending on the actual likelihood function. For example, in the case of the least-squares method, the matrix $\mathbf{I}^{-1} \approx \sigma^2(\mathbf{J}^T \mathbf{J})^{-1}$ and for vector \mathbf{U} is given by

$$\mathbf{U} = \mathbf{J}^T \hat{\mathbf{e}} \quad (8.167b)$$

where $\hat{\mathbf{e}}$ is the vector with elements $\hat{e}_i = y_i - f(x_i, \boldsymbol{\beta}^*)$. When σ^2 is estimated, the boundary of the $100(1 - \alpha)\%$ confidence region has the form

$$\hat{\mathbf{e}}^T \mathbf{J}(\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J} \hat{\mathbf{e}} = m \hat{\sigma}^2 F_{1-\alpha}(m, n - m) \quad (8.167c)$$

The boundary of the confidence region is determined as a set of vectors of parameters $\boldsymbol{\beta}$ which fulfil Eq. (8.167c). These confidence regions [Eq. (8.167c)] do not always have to be elliptic.

According to the *likelihood ratio* method, the boundaries of the $100(1 - \alpha)\%$ confidence region may be defined as

$$\ln L(\mathbf{b}) - \ln L(\boldsymbol{\beta}^*) = \chi^2_{1-\alpha}(m)/2 \quad (8.168a)$$

In the least-squares method, when both the parameter vector \mathbf{b} and the variance $\hat{\sigma}^2$ are estimated, Eq. (8.168a) may be expressed in the form

$$U(\mathbf{\beta}^*) - U(\mathbf{b}) = p^2 \quad (8.168b)$$

where p^2 is defined by Eq. (8.164). For determination of the confidence region boundaries, a numerical method should be used.

Problem 8.17. *Confidence regions for a model describing the activity of sea-weed as a function of temperature*

To express the empirical dependence of the activity P , of sea-weed on the temperature T at illumination level of $96 \text{ W} \cdot \text{m}^{-2}$, the model proposed was

$$P = \beta_2 \left[\frac{\beta_3 - T}{\beta_3 - \beta_1} \right]^Z \exp \left[Z \left[1 - \frac{\beta_3 - T}{\beta_3 - \beta_1} \right] \right] \quad (8.169)$$

where $Z = B^2(1 + \sqrt{1 + 40/R})^2/400$, $B = \ln(\beta_4)(\beta_3 - \beta_1)$, β_1 corresponds to the temperature with the maximum activity of sea-weed, $P_{\max} = \beta_2$, β_3 is the temperature at which the activity of the sea-weed is zero and β_4 is the shape factor. A graphical illustration of the proposed model $P = f(T, \beta_1, \beta_2, \beta_3, \beta_4)$ is shown in Fig. 8.24. Estimate all four parameters and determine the confidence regions for pairs of parameters, $\beta_1^* - \beta_2^*$, $\beta_2^* - \beta_4^*$ and $\beta_1^* - \beta_3^*$.

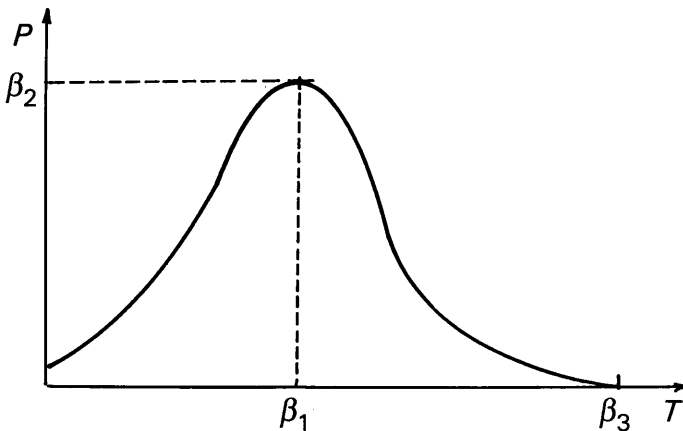


Fig. 8.24—Graphical illustration of a model [Eq. (8.169)] and the individual parameters.

Data: $n = 7$

$T, ^\circ\text{C}$	5	10	15	20	25	30	35
P	2.18	3.16	4.11	6.67	9.04	5.83	1.13

Solution: From the initial guesses of parameters $\beta_1^{(0)} = 25$, $\beta_2^{(0)} = 9$, $\beta_3^{(0)} = 36$ and $\beta_4^{(0)} = 1.8$, the corresponding value of the sum of squares function $U(\mathbf{\beta}^{(0)}) = 10.52$. The minimization process terminated at $U(\mathbf{b}) = 1.782$. The parameter estimates \mathbf{b}

Table 8.6. Parameter estimates and their relative bias, for the model in Eq. (8.169)

Parameter	Best estimate b_j	Relative bias $h_{R,j}$, %
β_1	25.06	0.0114
β_2	8.296	0.499
β_3	37.23	2.551
β_4	2.201	2.284

with their relative bias are listed in Table 8.6.

Figures 8.25, 8.26 and 8.27 show the 95% confidence region of pairs of parameters $\beta_1^* - \beta_2^*$, $\beta_1^* - \beta_3^*$ and $\beta_2^* - \beta_4^*$. The solid curve refers to Eq. (8.169) and the dotted curve to Eq. (8.168).

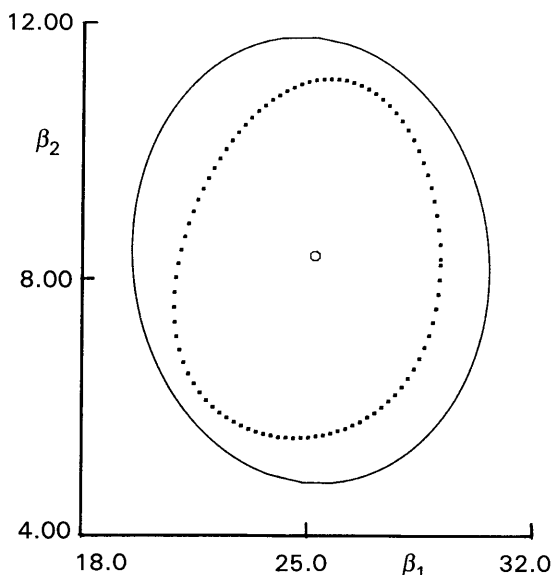


Fig. 8.25—The confidence region for parameters $\beta_1^* - \beta_2^*$ for the proposed model, Eq. (8.169), (solid curve) and with the use of Eq. (8.168) (dotted curve).

Conclusion: From the figures, it is obvious that for even a small bias, the differences between the confidence ellipsoids and the more accurate confidence regions are highly significant. The smallest difference is for the pair of parameters, $\beta_1^* - \beta_2^*$, where the values of the relative bias are smaller than 1%.

For the construction of the confidence regions, a reparameterization limiting the bias followed by a reverse parameterization, may also be used [61]. If a nonparametric technique is required, the Jack-knife or the Bootstrap methods are often used. The principle involved is the same as for univariate samples (Chapter 3).

8.6.2.2 Confidence intervals of parameters

With the use of Eq. (8.161), the $100(1 - \alpha)\%$ confidence interval of parameter β_j in the form

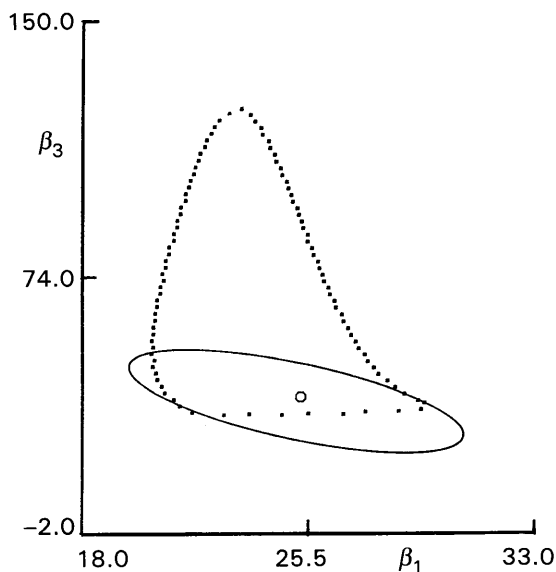


Fig. 8.26—The confidence region for parameters $\beta_1^* - \beta_3^*$ for the proposed model, Eq. (8.169), (solid curve) and with the use of Eq. (8.168) (dotted curve).

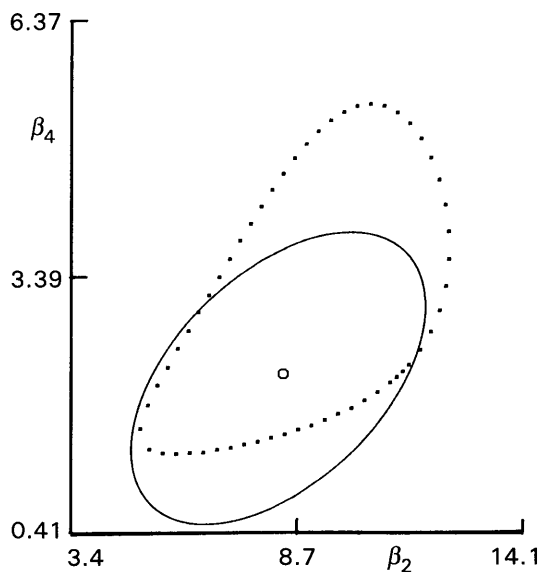


Fig. 8.27—The confidence region for parameters $\beta_2^* - \beta_4^*$ for the proposed model, Eq. (8.169), (solid curve) and with the use of Eq. (8.168) (dotted curve).

$$b_j - \hat{\sigma} \sqrt{V_{jj}} t_{1-\alpha/2}(n-m) \leq \beta_j^* \leq b_j + \hat{\sigma} \sqrt{V_{jj}} t_{1-\alpha/2}(n-m) \quad (8.170)$$

is a direct analogy of the confidence intervals of the parameters of linear models. The

influence of other parameters is neglected. When all off-diagonal elements of the matrix $(\mathbf{J}^T \mathbf{J})^{-1}$ are zero, Equation (8.170) may be used. However, the elements of the vector \mathbf{b} are often mutually correlated, so that the intervals of Eq. (8.170) are underestimated, i.e. they are too narrow.

A more suitable determination of the confidence interval of parameter β_k^* is on the basis of the maximal length Δ_k of the projection Δ_{kj} onto the parameter axis β_k . In the program LETAGROP, and the related system ABLET, the estimate of the standard deviation of the k th parameter, β_k^* , is calculated from

$$\Delta_k = \max_j (\Delta_{kj}) \quad (8.171a)$$

and the confidence interval of the parameter β_k is estimated from

$$b_k - \Delta_k \leq \beta_k^* \leq b_k + \Delta_k \quad (8.171b)$$

Instead of projections it is simpler to search directly for the co-ordinates of the extreme points on the confidence ellipsoid in the directions of the individual parameter axes [62]. The confidence interval of the parameter β_k^* is given by

$$b_k - p\sqrt{V_{kk}} \leq \beta_k^* \leq b_k + p\sqrt{V_{kk}} \quad (8.172)$$

For $m = 1$, these confidence intervals are identical. When the number of regression parameters m is increased, the confidence intervals in Eqs. (8.171) and (8.172) become broader than those in Eq. (8.170). All confidence intervals are symmetrical.

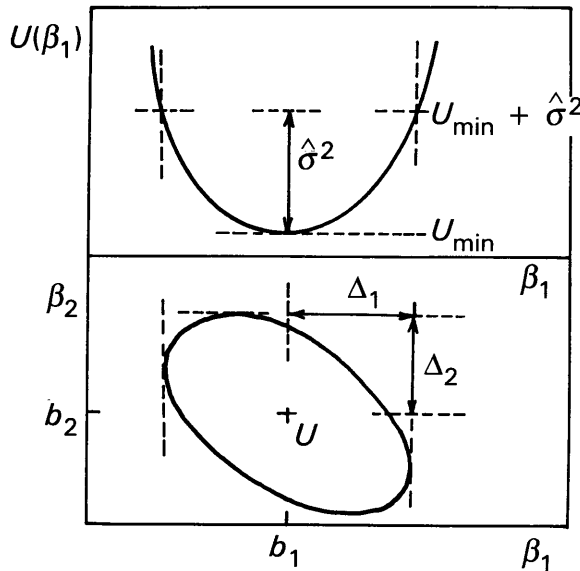


Fig. 8.28—The estimates of the parameter standard deviations Δ_1 and Δ_2 by programs LETAGROP and ABLET.

Problem 8.18. *The confidence intervals for the sea-weed activity model*

Estimate the half-length of the 95% confidence intervals for four parameters of the proposed model for the dependence of activity of sea-weed on temperature, from Problem 8.17.

Data: from Problem 8.17

Solution: The half-length of the confidence intervals of the four parameters of the model (8.169) calculated by three different approaches are shown in Table 8.7.

Table 8.7. Half-lengths of the 95% confidence interval of four parameters of model (8.169).

Parameter	b_k	Δ_k (8.170)	Δ_k (8.171)	Δ_k (8.172)
β_1	25.06	2.931	3.922	5.541
β_2	8.296	1.845	3.466	3.488
β_3	37.23	8.199	15.470	15.500
β_4	2.201	0.8984	1.160	1.690

Conclusion: Equation (8.170) leads confidence intervals that are false and too narrow. The broad confidence intervals are the consequence of too small a sample size, $n = 7$, relative to the number of unknown parameters, $m = 4$.

Asymmetrical confidence intervals of parameter estimates may be obtained when Eqs. (8.167) and (8.168) are solved numerically with respect to parameter β_k^* , when estimates \mathbf{b} are supplied for the other components of vector $\boldsymbol{\beta}^*$.

8.6.2.3 Confidence intervals of prediction

If the regression model can be linearized, the $100(1 - \alpha)\%$ confidence interval of a prediction $f(x^*, \mathbf{b})$ at the point x^* may be calculated. It then follows that

$$f(x^*, \mathbf{b}) - t_{1-\alpha/2}(n-m)\hat{\sigma}_p(x^*) \leq f(x^*, \boldsymbol{\beta}) \leq f(x^*, \mathbf{b}) + t_{1-\alpha/2}(n-m)\hat{\sigma}_p(x^*) \quad (8.173)$$

where $\hat{\sigma}_p^2(x^*)$ is the estimate of the prediction variance for which

$$\hat{\sigma}_p^2(x^*) = \mathbf{J}^T D(\mathbf{b}) \mathbf{J} \quad (8.173a)$$

The symbol \mathbf{J} denotes the vector of derivatives of a model function at the point x^* with elements

$$J_j = \frac{\delta f(x^*, \boldsymbol{\beta})}{\delta \beta_j} \quad (8.173b)$$

The confidence intervals of prediction calculated for the whole range of the independent variable x (if scalar) form the *confidence bands*. Accurate confidence bands may be constructed with the aid of a suitable reparameterization [61].

Problem 8.19. *Confidence bands of prediction for the model of the effect of temperature on sea-weed activity*

Calculate the 95% confidence bands of prediction for the model of activity of sea-weed vs. temperature, Eq. (8.169).

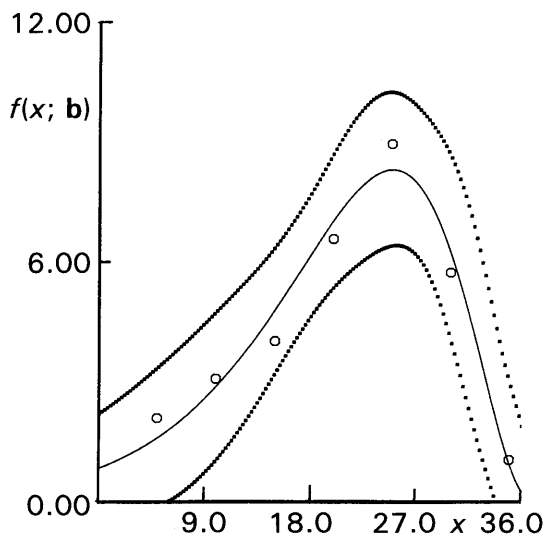


Fig. 8.29—Confidence bands of prediction (dotted curve) with calculated regression line for the model, Eq. (8.169) (solid curve) with experimental points.

Data: from Problem 8.17

Solution: Figure 8.29 shows the regression line (solid curve) and the 95% bands of prediction (dotted curves).

Conclusion: Here, the width of the confidence band is negatively affected by the small sample size.

8.6.3 Hypothesis tests about parameter estimates

Hypothesis testing is closely related to construction of confidence bands. If parameters β_0 lie in the 95% confidence range around \mathbf{b} , the differences $(\beta^* - \beta_0)$ may be considered as statistically insignificant at the significance level $\alpha = 0.05$. (The principle of testing is described in Chapter 3.) We restrict ourselves here to the main tests.

To examine the regression parameters, the null hypothesis, $H_0: \beta = \beta_0$ is often tested against the alternative $H_A: \beta \approx \beta_0$, where β_0 is a given parameter vector. If a regression model can be linearized, Eq. (8.161) leads to the test characteristic

$$T = \frac{(\mathbf{b} - \beta_0)^T (\mathbf{J}^T \mathbf{J}) (\mathbf{b} - \beta_0)}{m \hat{\sigma}^2} \quad (8.174)$$

which, if the null hypothesis is valid, has the Fisher–Snedecor F -distribution with m and $(n - m)$ degrees of freedom. If $T > F_{1-\alpha}(m, n - m)$ the null hypothesis H_0 , about the equality of β and β_0 , is rejected.

When the null hypothesis about one parameter, $H_0: \beta_j = \beta_0$, is tested against the alternative, $H_A: \beta_j \neq \beta_0$, the criterion

$$T_j = \frac{|b_j - \beta_0|}{\hat{\sigma} \sqrt{V_{jj}}} \quad (8.175)$$

may be used. When the null hypothesis is valid, T_j has the Student distribution with $(n - m)$ degrees of freedom. If $T_j > t_{1-\alpha/2}(n - m)$, the null hypothesis about a parameter identity is rejected at significance level α .

When test criteria T and T_j are used, the same restrictions as for the confidence regions hold. If the error distribution ε does not differ from normality, the distribution of the test criterion, T_j , even for strongly nonlinear models. As in the construction of confidence regions for parameter subsets β_2 , tests can be constructed for parameter subsets β_2 . When $\beta_0 = 0$ is selected, the classical tests of significance of the regression parameters result.

Problem 8.20. Tests of parameters

For the kinetic model in Problem 8.13, examine the following null hypotheses:

- (1) $H_0: \beta = 0$ (which implies that the model is insignificant) against $H_A: \beta \neq 0$;
- (2) $H_0: \beta_1 = 80$ against $H_A: \beta_1 \neq 80$.

Data: from Problem 8.13

Solution: The parameter estimates obtained by the least-squares method are $b_1 = 82.36$, $b_2 = 0.0148$, the variance estimate $\hat{\sigma}^2 = 3.22$. The matrix $(\mathbf{J}^T \mathbf{J})$ has the form

$$\mathbf{J}^T \mathbf{J} = \begin{bmatrix} 3.59 & 5330 \\ 5330 & 8.55 \times 10^6 \end{bmatrix}$$

and the matrix $\hat{\sigma}^2(\mathbf{J}^T \mathbf{J})^{-1} = D(\mathbf{b})$ is equal to

$$D(\mathbf{b}) = \begin{bmatrix} 12.18 & -0.00758 \\ -0.00758 & 5.11 \times 10^6 \end{bmatrix}$$

For testing, we select the significance level $\alpha = 0.05$.

(1) From Eq. (8.174), we estimate $\mathbf{b}^T(\mathbf{J}^T \mathbf{J})\mathbf{b} = 3.91 \times 10^4$. The test criterion $T = 5081.9$ is significantly higher than the quantile $F_{0.95}(2, 4) = 6.94$ and therefore the null hypothesis about model insignificance, $\beta = 0$, is rejected.

(2) From Eq. (8.175), the test criterion is

$$T_1 = \frac{|82.36 - 80|}{\sqrt{12.18}} = 0.0675$$

In comparison with the quantile $t_{0.975}(4) = 2.776$, the T_1 criterion is significantly lower, and therefore the null hypothesis $H_0: \beta_1 = 80$ cannot be rejected.

Conclusion: The statistical tests described are quite simple and do not require complicated calculations.

For testing general parametric hypotheses, the tests of the likelihood ratio and the Lagrange multipliers may be used. Any parametric hypothesis may be expressed as $H_0: \beta \in \omega$ against $H_A: \beta \in \Omega - \omega$, where Ω is an admissible parameter space and ω is its subspace. Often the null hypothesis, expressing q relationships between regression parameters, $H_0: f_1(\beta) = 0, f_2(\beta) = 0, \dots, f_q(\beta) = 0$, is tested. Then ω is given by restriction conditions of the type $f_j(\beta) = 0, j = 1, \dots, q$.

For a test of the null hypothesis H_0 , the likelihood ratio has the form

$$l(\mathbf{b}_\omega) = \frac{\max_{\beta \rightarrow \omega} L(\beta)}{\max_{\beta \rightarrow \Omega} L(\beta)} = \frac{L(\mathbf{b}_\omega)}{L(\mathbf{b})} \quad (8.176)$$

where \mathbf{b}_ω is the maximum likelihood estimate of parameters β , with the restriction that it must be in a range ω . The test uses the fact that $-2 \ln l(\mathbf{b}_\omega)$ has the $\chi^2(q)$ distribution. Especially in the case of the least-squares method, when the residual variance is calculated from Eq. (8.176), the test criterion has the form

$$TL = \frac{[U(\mathbf{b}_\omega) - U(\mathbf{b})](n - m)}{qU(\mathbf{b})} \quad (8.177)$$

This statistic has, if the null hypothesis is valid, the Fisher–Snedecor F -distribution with q and $(n - m)$ degrees of freedom.

Problem 8.21. Tests of parametric hypotheses

For the Cegarra–Puente kinetic model from Problem 8.13, examine the following hypotheses:

- (a) $H_0: \beta_1 = 80$ vs. $H_A: \beta_1 \neq 80$,
- (b) $H_0: \beta_1 = 80$ and $\beta_2 = 0.01$ vs. $H_A: \beta_1 \neq 80$ and $\beta_2 \neq 0.01$.

Use the TL test criterion.

Data: from Problem 8.13

Solution: (1) For the model, $y = 80\sqrt{(1 - \exp(-\beta_2 x))}$, the minimum $U(\beta_2)$ is reached when $U(\beta_2) = 14.51$. In Problem 8.13, $U(\beta) = 12.88$ was achieved. Putting these values into the test criterion, we get $TL = (14.51 - 12.88) \times 4/12.88 = 0.506$, which is smaller than the quantile $F_{0.95}(1, 4) = 7.7$. Therefore the null hypothesis $H_0: \beta_1 = 80$ cannot be rejected at the significance level $\alpha = 0.05$.

(2) The criterion function $U(\mathbf{b}_\omega) = 469.3$ for $\beta_1 = 80$ and $\beta_2 = 0.01$, for the model $y = 80\sqrt{(1 - \exp(-0.01x))}$. Equation (8.177) then gives the test statistic $TL = 70.87$. This value is significantly higher than the quantile $F_{0.95}(2, 4) = 5.94$, and therefore the null hypothesis $H_0: \beta_1 = 80$ and $\beta_2 = 0.01$ is rejected.

Conclusion: We have shown that application of the TL criterion requires two minimizations, except when all values of parameters are known or assumed.

In some cases of a search for a minimum, overdetermined models are formed. Galant [5] has proposed a special procedure for these.

8.6.4 Goodness-of-fit tests

The examination of residuals is useful, not only for the linear regression model (Section 6.5.2.1), but also for nonlinear regression models and analysis of variance models.

Residuals are defined as the differences

$$\hat{e}_i = y_i - \hat{y}_{P,i} = y_i - f(x_i, \mathbf{b}), \quad i = 1, \dots, n \quad (8.178)$$

where y_i is an observation and $\hat{y}_{P,i} = f(x_i, \mathbf{b})$ is the calculated value, a “prediction”,

i.e. the value found from the equation fitted.

We now use graphical and analytical methods for examining the residuals, in order to check the quality of nonlinear models.

8.6.4.1 Graphical analysis of residuals

The principle ways of plotting the residuals \hat{e}_i have already been described in regression diagnostics for linear regression models. The following plots are often used in the examination of nonlinear models:

- (1) The *overall diagram* gives an initial impression of the residuals. If the model is correct, the residuals should resemble observations from a normal distribution with zero mean.
- (2) Plot type I (the *index plot*) is a plot of residuals \hat{e}_i against the index i in time order.
- (3) Plot type II (the plot *vs. the independent variable*) is a plot of residuals \hat{e}_i against the independent variable x_j , $j = 1, \dots, m$.
- (4) Plot type III (the plot *vs. the prediction*) is a plot of residuals against the predicted value $\hat{y}_{P,i}$.

Plot type II is usually adopted as the standard plot (Fig. 8.30). If the proposed model represents the data adequately, the residuals should form a random pattern. Systematic departures from randomness indicate that the model is not satisfactory. To examine the normality of a residual distribution, the rankit plot, used in regression diagnostics for linear models, may be applied.

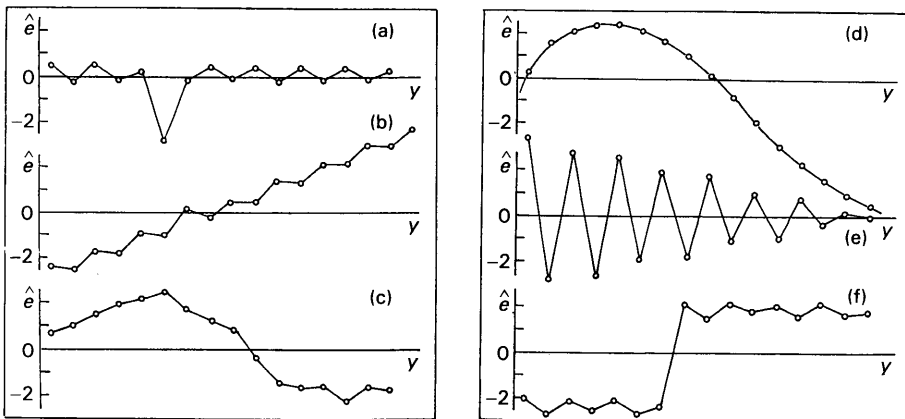


Fig. 8.30—Plot II (residuals vs. the independent variable x): (a) detection of an outlier, (b) detection of a trend in residuals, (c) detection of sign changes, (d) detection of a false model, (e) detection of heteroscedasticity, and (f) detection of an abrupt shift in level of the experiment.

8.6.4.2 Statistical analysis of residuals

The plots that have been recommended in the previous sections are visual techniques for easy checking of some of the basic assumptions of the least-squares method and

of the proposed model. Certain statistics provide a numerical measure for some of the discrepancies previously described.

In many regression programs used in the chemical laboratory, the statistical analysis of residuals represents the main diagnostic tool used to search for the “best” model when more than one are possible or proposed. The goodness-of-fit test analyses the set of residual and examines the following criteria.

(1) The arithmetic mean of residuals known as the *residual bias*, $E(\hat{\epsilon})$, should be equal to zero.

(2) The mean of the absolute values of residuals, $|\bar{\epsilon}|$, and the square-root of the residual variance (the estimate of the *residual standard deviation*), $s(\hat{\epsilon})$, should both be of the same magnitude as the (instrumental) error of the dependent variable (observation, measured quantity y), $s_{\text{inst}}(y)$, i.e. $|\bar{\epsilon}| \approx s_{\text{inst}}(y)$ and $s(\hat{\epsilon}) \approx s_{\text{inst}}(y)$.

(3) The *residual skewness*, $g_1(\hat{\epsilon})$, for a Gaussian normal distribution should be equal to zero.

(4) The *residual kurtosis*, $g_2(\hat{\epsilon})$, for a Gaussian normal distribution should be equal to 3.

(5) The *residual variance* is calculated from the residual sum of squares $\hat{\sigma}^2 = U(\mathbf{b})/(n - m)$.

(6) The *coefficient of determination*, D^2 , is calculated from:

$$D^2 = 1 - \frac{U(\mathbf{b})}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (8.179)$$

where $\bar{y} = (1/n) \sum_{i=1}^n y_i$. The coefficient of determination is equal to the square of the correlation coefficient, for linear models.

(7) In chemometrics, we often use the Hamilton R -factor:

$$R = \sqrt{\frac{U(\mathbf{b})}{\sum_i y_i^2}} \quad (8.180)$$

When $\bar{y} = 0$, $R^2 = 1 - D^2$, so Eq. (8.180) may be expressed as

$$R = \sqrt{(1 - D^2) - \frac{(1 - D^2)n\bar{y}^2}{\sum_i y_i^2}} \quad (8.181)$$

The Hamilton R -factor illustrates the difference between the two models, $y = f(x, \beta)$ and $y = 0$. This rule is not correct for models with intercept terms and the values of the Hamilton R -factor are incorrectly low. It should be noted that D^2 and the R -factor are continuous functions of the number of parameters. D^2 is an increasing function of the number of parameters, whereas the R -factor is decreasing function of this number. Therefore, both of these statistics are unsuitable as resolution diagnostics for comparing models with different numbers of parameters.

(8) To distinguish between models, the *Akaike information criterion*, AIC, is more suitable:

$$\text{AIC} = -L(\mathbf{b}) + 2m \quad (8.182)$$

The “best” model is considered to be that for which this criterion reaches a minimum value. For least-squares models and models that do not belong to the same class, the AIC criterion may be expressed as

$$\text{AIC} = n \times \ln \left[\frac{U(\mathbf{b})}{n - m} \right] + 2m \quad (8.183)$$

It should be noted that the diagnostic use of classical residuals is not rigorous but rather approximate. Classical residuals do not have zero mean, they are biased and they are a linear combination of errors ε . Moreover, they depend on the true values of parameters β^* , and these are unknown. Therefore, it has been suggested [58] that the projections of residuals are used, because this partly limits all these disadvantages. Lyoness [63] proposed various approximate expressions for the determination of different types of residuals for nonlinear models, for example, Jack-knife residuals, recursive residuals and partial residuals, with applications similar to those for linear models.

For more objective examination of residuals, all the statistical regression diagnostics for linear models may also be used for nonlinear models. A difficulty may arise from the distributions of some test criteria, which are affected by the nonlinearity of the model. Some of the test criteria are derived from a general criterion

$$T_{p,q} = \sum_{i=1}^n \hat{\varepsilon}_i^p [f(x_i, \mathbf{b})]^q \quad (8.184)$$

As for linear regression models, the following conditions should be met:

- (a) The test criterion $T_{1,1}$ should be approximately equal to zero since

$$\hat{\varepsilon}^T f(x_i, \mathbf{b}) = 0$$

A high value of $T_{1,1}$ indicates a false minimum or a false model.

- (b) The test criterion $T_{2,1}$ indicates heteroscedasticity.

- (c) The test criterion $T_{1,2}$ indicates that a false model has been proposed.

- (d) The test criterion $T_{1,0}$ should be approximately equal to zero. From Eq. (8.154) it follows that the mean of the residuals $E(\hat{\varepsilon})$ is not equal to zero. Therefore:

$$T_{1,0} = n \sum_{i=1}^n \left[d_i - \sum_{k=1}^q P_{ik} d_k \right]$$

Expressions for the variance and covariance matrix have been proposed [60].

The predictive ability of a proposed model may be examined by the *mean quadratic error of prediction*, defined by

$$\text{MEP} = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i, \mathbf{b}_{(i)}))^2 \quad (8.185)$$

The one-step approximation $\mathbf{b}_{(i)}^1$, defined by Eq. (8.188), may be used instead of the parameter estimate, $\mathbf{b}_{(i)}$. The lower the values of MEP, the better the predictive ability

of the proposed model.

The correctness of the proposed model can be examined by the White test [67]. The coefficient C , defined as

$$C = (n - 1) \sum_{i=1}^n \frac{\delta f(x_i, \mathbf{b})}{\delta \beta_j} \times \frac{\delta f(x_i, \mathbf{b})}{\delta \beta_k} \left[\hat{e}_i - \frac{U(\mathbf{b})}{n} \right] \quad (8.186)$$

should be equal to zero for a correct model only. The test procedure calculates $(m(m + 1)/2 - n)$ variables, and

$$w_s = \frac{\delta f(x_i, \boldsymbol{\beta})}{\delta \beta_j} \times \frac{\delta f(x_i, \boldsymbol{\beta})}{\delta \beta_k} \quad (8.187)$$

for $s = 1, \dots, m(m + 1)/2$. The test criterion represents the correlation coefficient of regression of the variable \hat{e}^2 on the vector of variables \mathbf{w} . In the case of a correct model, this test criterion has the $\chi^2_{1-\alpha}(m(m + 1)/2)$ distribution.

8.6.4.3 Identification of influential points

For linear regression models (Chapter 6), all characteristics which aid the identification of influential points are functions of residuals \hat{e}_i and diagonal elements H_{ii} of the projection matrix, $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$. Some diagnostics of influential points use estimate $\mathbf{b}_{(i)}$ calculated from all points except the i th one. For linear regression models, the estimate $\mathbf{b}_{(i)}$ is easily obtained from information on the matrix $\mathbf{X}^T \mathbf{X}^{-1}$ and quantities \hat{e}_i and H_{ii} .

For nonlinear regression models, the situation is rather more complicated as the parameter estimates and residuals cannot be expressed simply as a linear combination of experimental data. When a Taylor-type linearization of the original nonlinear model is used, all methods for identification of influential points in linear models can be used. This starts with a one-step approximation of the parameter estimate

$$\mathbf{b}_{(i)}^1 = \mathbf{b} - (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}_i \frac{\hat{e}_i}{1 - P_{ii}} \quad (8.188)$$

where P_{ii} are elements of the projection matrix [Eq. (8.68)]. With the use of Eq. (8.188), the test criterion DFS_{ij} may be written as

$$DFS_{ij} = \frac{b_j - b_{j(i)}^1}{\hat{s}_{(i)}^2 V_{ii}} \quad (8.189)$$

This criterion expresses the influence of the i th point on the estimate of the j th parameter. Quality, $\hat{s}_{(i)}^2$, is the variance estimate calculated when the i th point is omitted, i.e.

$$\hat{s}_{(i)}^2 = \frac{U(\mathbf{b}) - \frac{\hat{e}_i^2}{1 - P_{ii}}}{n - m - 1} \quad (8.190)$$

The symbol V_{ii} , in Eq. (8.189), denotes elements of the matrix $\mathbf{V} = (\mathbf{J}^T \mathbf{J})^{-1}$. Applying

the DFS_{ij} criterion, the i th point is considered to be influential if $DSF_{ij} > 2/\sqrt{n}$.

Influential points may be identified readily on the basis of a one-step approximation of the Jack-knife residuals \hat{e}_{ji} expressed as,

$$\hat{e}_{ji} = \frac{\hat{e}_i}{\hat{s}_{(i)}\sqrt{1 - P_{ii}}} \quad (8.191)$$

To express the influence of individual points on parameter estimates, the quadratic expansion of a regression model may be used. Often an examination of either the changes of the vector of bias with the omission of the i th point, $\mathbf{h}_{(i)}$, or changes of the mean value of the i th residual with the i th point omitted, is suitable [66].

A nonlinear measure of the influence of the i th point on the parameter estimates is represented by the likelihood distance,

$$LD_i = 2[\ln L(\mathbf{b}) - \ln L(\mathbf{b}_{(i)})] \quad (8.192)$$

In the case of least-squares, the likelihood distance is expressed by

$$LD_i = n \times \ln \left[\frac{U(\mathbf{b}_{(i)})}{U(\mathbf{b})} \right] \quad (8.193)$$

In Eqs. (8.192) and (8.193), the estimates $\mathbf{b}_{(i)}$ calculated by nonlinear regression, when the i th point was left out, or the one-step approximation $\mathbf{b}_{(i)}^1$ of parameter estimates may be used. When $LD_i > \chi^2_{1-\alpha}(2)$, the i th point is said to be strongly influential. The significance level, α , is usually chosen as 0.05.

Some diagnostics for identification of influential points were compared [66], and the following conclusions were reached.

- Influential points affect not only the parameter estimates but also the relative bias \mathbf{h}_R , which is rather sensitive to the presence of influential points.
- Diagnostics based on linearization or quadratic expansion of a nonlinear model do not always indicate the presence of influential points. They are not suitable for strongly nonlinear models.
- The best identification of influential points is given by the likelihood distance LD_i . In some cases, groups of influential points can cause masking effects.
- For practical calculations, the approximation of LDS_i is sufficient when the quantity $\mathbf{b}_{(i)}^1$ from Eq. (8.188) is used in Eq. (8.185) instead of $\mathbf{b}_{(i)}$.

Problem 8.22. *Identification of influential points in the dependence of the activity of sea-weed on temperature*

Determine the influential points in the data in Problem 8.15. Calculate the mean values of residuals $E(e_i)$, the Jack-knife residuals \hat{e}_{ji} and the likelihood distance LD_i and LDS_i .

Data: from Problem 8.15

Table 8.8. The mean values of residuals and other measures of influential points for the data in Problem 8.15

Point	E_i	\hat{e}_{ji}	LDS_i	LD_i
1	0.61	0.81	1.05	0.277
2	0.304	0.34	0.82	0.085
3	-0.661	-0.63	1.99	0.376
4	-0.343	-0.35	0.98	0.121
5	0.745	0.89	1.22	4.75
6	-0.419	-0.89	0.93	136
7	0.181	0.415	0.0216	1968

Solution: Diagnostics for the identification of influential points are listed in Table 8.8.

Conclusion: When the diagnostics for identification of influential points are compared, only the nonlinear measure LD_i lead to a conclusion that two points, 6 and 7, are strongly influential, because they control the decreasing part of the curve. The other diagnostics do not indicate obviously influential points. The mean values of residuals E_i show their bias directly.

8.7 PROCEDURE FOR BUILDING AND TESTING A NONLINEAR MODEL

The quality of a proposed nonlinear model is examined in the same way as for linear models, using the following criteria:

(1) The quality of parameter estimates

The quality of parameter estimates obtained is considered according to their confidence intervals or their variances $D(b_j)$. The empirical rule that is often used is that a parameter is considered to be significant when its estimate is greater than 3 standard deviations, i.e. $3\sqrt{D(b_j)} < |b_j|$. High parameter variances are often caused either by termination of the minimization process before a minimum is reached, by inaccuracy of determination of matrix \mathbf{J} , or high nonlinearity of the regression model.

(2) The quality of the curve fitting

Agreement of the proposed model with the experimental data is examined by the goodness-of-fit test based on the statistical analysis of residuals. The following statistical characteristics for a set of classical residuals are calculated: from the residual square-sum $U(\mathbf{b})$ reached at a minimum, the estimate of residual variance $\hat{\sigma}^2$ and estimates of the determination coefficient D^2 , the regression rabat $100D^2[\%]$, the arithmetic mean of the residuals $E(\hat{e})$, the mean of the absolute values of the residuals

$|\bar{e}|$, the mean of the relative residuals \bar{e}_{rel} , the residual standard deviation $s(\hat{e})$, the residual skewness $g_1(\hat{e})$, the residual kurtosis $g_2(\hat{e})$, and the Pearson χ^2 -test of normality of the residual distribution are carried out. In addition, the four test criteria $T_{1,1}$, $T_{2,1}$, $T_{1,2}$ and $T_{1,0}$ are calculated, for more objective residual analysis.

(3) *The predictive ability of the proposed model*

The predictive ability of a model is classified by the following procedure. Data are divided into two groups, M_1 with indices $i = 1, \dots, \text{int}(n/2)$ and M_2 with indices $i = \text{int}(n/2) + 1, \dots, n$. Estimates of the parameters are calculated from points in the subgroup M_1 as $\mathbf{b}(M_1)$. The predictive ability of the model is expressed by

$$K = \frac{U(\mathbf{b})}{\sum_{i \in M_1} [y_i - f(x_i, \mathbf{b}(M_2))]^2 + \sum_{i \in M_2} [y_i - f(x_i, \mathbf{b}(M_1))]^2} \quad (8.194)$$

The predictive ability of the model is higher as the criterion, K , tends towards a value of 1.

The mean quadratic error of prediction MEP (8.185) is then calculated. The lower the value of MEP , the better the predictive ability of the proposed model.

(4) *The quality of the experimental data*

For the examination of the quality of the experimental data, influential points are identified by regression diagnostics. The most important diagnostics are the likelihood distances LD_i and LDS_i . The test criterion DFS_{ij} and the mean value of each residual $E(\hat{e}_i)$ are also useful.

(5) *The correctness of the model proposed*

The White test calculates the coefficient C to prove that proposed model is correct. Some other tests of accuracy of the proposed model have been proposed [68].

(6) *The physical meaning of parameter estimates*

In chemometrics models, there are often restrictions from the physical meaning of the parameters. For example, concentrations or molar absorptivities must be positive numbers.

8.8 ADDITIONAL PROBLEMS

Problem 8.23. *Estimation of the parameters of the extended Debye–Hückel equation*

Estimate the thermodynamic dissociation constant pK_a^T , the effective ion-size parameter \hat{a} and the salting-out parameter C from the dependence of the mixed dissociation constant pK_a on the ionic strength, according to the extended Debye–Hückel law [65]. The dissociation $HL^z = L^{(z-1)} + H^+$ is described by the thermodynamic dissociation constant

$$K_a^T = a_{H^+} a_L^{(z-1)} / a_{HL}^z$$

or the mixed dissociation constant

$$K_a = a_{H^+} [L^{(z-1)}] / [HL^z].$$

If the two ions $L^{(z-1)}$ and HL^z have roughly the same ion-size \bar{a} (in 10^{-10} m) and the overall salting out coefficient is $C = C_{HL^z}^* - C_{L^{(z-1)}}^*$, then the extended Debye-Hückel law is expressed in the form

$$pK_{a,i} = pK_a^T - \frac{A\sqrt{I_i}(1-2z)}{1+B\bar{a}\sqrt{I_i}} + CI_i \quad (8.192)$$

where $pK_{a,i}$ is the dependent variable, I_i is the independent variable; pK_a^T , \bar{a} and C are three unknown parameters to be estimated; and two known numerical constants are $A = 0.5112 \text{ mol}^{-1/2} \text{ l}^{1/2} \text{ K}^{3/2}$, $B = 3.291 \times 10^{10} \text{ mol}^{-1/2} \text{ m}^{-1} \text{ l}^{1/2} \text{ K}^{1/2}$ for aqueous solution at 298.16K. We assume an additive model of measurements, and normality of errors of the dependent variable pK_a . The independent variable I , has a significantly smaller experimental error.

Data: 20 points of the dependence, $pK_a = f(I)$ were generated for pre-selected parameters $pK_a^T = 5.000$, $\bar{a} = 0.45$ and $C = 0.300$ and an instrumental error of the dependent variable of $s_{\text{inst}}(pK_a) = 0.005$. The simulated data set $\{I, pK_a\}$ is:

I	0.01	0.04	0.09	0.16	0.25	0.36	0.49
pK_a	4.8646	4.7752	4.7019	4.6661	4.6407	4.6145	4.6084
	0.64	0.81	1.00	1.21	1.44	1.69	1.96
	4.6318	4.6484	4.6726	4.7179	4.7769	4.8213	4.8896
	2.25	2.56	2.89	3.24	3.61	4.00	
	4.9522	5.0424	5.1242	5.2178	5.3129	5.4196	

Solution: With initial guesses of parameters $(pK_a^T)^{(0)} = 1$, $(\bar{a})^{(0)} = 0$, $C^{(0)} = 1$, the sum of squares $U(\mathbf{b}^{(0)}) = 325.7$. With MINOPT, a minimum is reached at $U(\mathbf{b}) = 4.95 \times 10^{-4}$. The best values of the parameter estimates \mathbf{b} , their half lengths of confidence interval Δ_j [Eq. (8.170)] and $\Delta_{R,j}$ [Eq. (8.172)] and the relative bias of parameter estimates $h_{R,j}$ are given in Table 8.9.

Table 8.9. Estimates of the parameters of the Debye-Hückel law with their statistical characteristics

Parameter	Estimate, b_j	Δ_j	$\Delta_{R,j}$	$h_{R,j}$, %
pK_a^T	4.997	0.0073	0.0106	-0.0002
\bar{a}	0.452	0.0167	0.0245	0.0174
C	0.299	0.0051	0.0075	0.0036

Since the bias of the parameter estimates is very small, the confidence intervals of the parameters are obtained either from the linearization $b_j \pm \Delta_j$ or $b_j \pm \Delta_{R,j}$. The residual standard deviation, $\hat{\sigma} = 0.0054$, is close to the instrumental error

$s_{\text{inst}}(\text{p}K_a) = 0.005$. Good curve-fitting of the calculated regression curve through the points is evident from the high value of the regression rabat, $D = 99.96\%$, and the low values of the Hamilton R -factor, $R = 0.13\%$, the mean of the absolute residuals $E(\hat{e}) = 0.0042$ and the mean of relative residuals $R(\hat{e}_R) = 0.087\%$. The correlation matrix of parameter estimates shows some correlation between \hat{a} and $\text{p}K_a^T$ or \hat{a} and C , as shown in Table 8.10.

Table 8.10. Correlation matrix of parameter estimates.

	$\text{p}K_a^T$	\hat{a}	C
$\text{p}K_a^T$	1	-0.84	0.612
\hat{a}		1	-0.92
C			1

Conclusion: Because of the high precision of the data [$s_{\text{inst}}(\text{p}K_a) = 0.005$], the asymptotic expression coming from the linearization of the model function may be used in statistical analysis.

Problem 8.24. Estimation of dissociation constants of two overlapping protonation equilibria by analysis of the A–pH graph

Estimate three dissociation constants and four molar absorptivities for a dissociating acid, H_3L , by regression analysis of the A–pH curve [65, 91].

Data: The absorbance–pH graph for $4.18 \times 10^{-5} \text{ M}$ 3-CAPAZOXS was measured at 470 nm in a cuvette of length 1.000 cm by a glass electrode–SCE cell (59.16 mV/pH slope) at 25°C. The additive measurement model and normality of errors are assumed. The measured values of pH and absorbance, A are:

pH	1.565	1.750	1.817	2.000	2.058	2.224	2.500	2.550
A	0.660	0.666	0.653	0.640	0.631	0.593	0.547	2.530
	2.750	2.788	2.956	3.000	3.185	3.250	3.364	3.518
	0.483	0.471	0.428	0.420	0.382	0.366	0.350	0.327
	3.772	4.000	4.082	4.369	4.872	5.569	6.266	6.691
	0.295	0.273	0.272	0.257	0.245	0.238	0.232	0.222
	6.750	7.056	7.295	7.500	7.740	8.000	8.072	8.250
	0.220	0.205	0.192	0.180	0.164	0.155	0.153	0.147
	8.464	8.915	9.316	9.855				
	0.144	0.140	0.140	0.140				

Solution: At different pH values, four species L^{3-} , HL^{2-} , H_2L^- and H_3L exist in an aqueous solution of 3-CAPAZOXS, and can be described by the following mixed

dissociation constants

$$K_{a1} = \frac{a_{H^+} [L^{3-}]}{[HL^{2-}]}$$

$$K_{a2} = \frac{a_{H^+} [HL^{2-}]}{[H_2L^-]}$$

and

$$K_{a3} = \frac{a_{H^+} [H_2L^-]}{[H_3L]}$$

When all four species absorb light at a given wavelength, the absorbance A in a cuvette of length d cm is expressed by

$$A = dc_{H_3L} \frac{\varepsilon_L + \sum_{j=1}^3 \varepsilon_{H_jL} 10^{\sum_{i=1}^j pK_{a,i} - j \text{ pH}}}{1 + \sum_{j=1}^3 10^{\sum_{i=1}^j pK_{a,i} - j \text{ pH}}}$$

where c_{H_3L} is the analytical concentration of 3-CAPAZOXS and the parameters estimated are pK_{a1} , pK_{a2} , pK_{a3} , ε_L , ε_{HL} , ε_{H_2L} and ε_{H_3L} (the charges on the ions are omitted for the sake of simplicity).

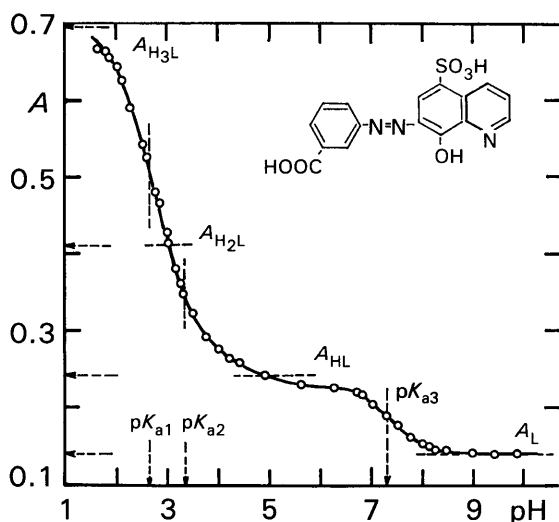


Fig. 8.31—Absorbance-pH graph for 3-CAPAZOXS.

With the initial guess of parameters of $pK_{a1}^{(0)} = 7.3$, $pK_{a2}^{(0)} = 3.3$, $pK_{a3}^{(0)} = 2.55$, $\varepsilon_L^{(0)} = 3300$, $\varepsilon_{HL}^{(0)} = 5800$, $\varepsilon_{H_2L}^{(0)} = 9800$, $\varepsilon_{H_3L}^{(0)} = 16800$, the sum of squares reaches a

minimum value of $U(\mathbf{b}^{(0)}) = 5.04 \times 10^{-4}$. By the least-squares method MINOPT, a minimum value of $U(\mathbf{b}) = 3.95 \times 10^{-4}$ was found. The estimates of the seven parameters with their half lengths of the 95% confidence interval [Eqs. (8.170), (8.171) and (8.172)] and the relative bias of parameters are listed in Table 8.11.

Table 8.11. Estimates of seven parameters and their statistical characteristics

Parameter	Estimate	Half-length of confidence interval by			$h_{R,j}$ %
		Eq. (8.170)	Eq. (8.171)	Eq. (8.172)	
ε_L	3290	89	153	177	-0.41
ε_{HL}	5731	126	177	249	-0.61
ε_{H_2L}	7154	1754	3471	3471	-6.84
ε_{H_3L}	16702	152	237	300	-0.25
pK_{a1}	7.338	0.101	0.124	0.199	0.33
pK_{a2}	3.925	0.765	1.502	1.513	8.20
pK_{a3}	2.720	0.111	0.214	0.219	1.36

It can be seen that the estimates of ε_{H_2L} , pK_{a2} and pK_{a3} are significantly biased. For identification of any mutual linear association between estimates of individual parameters, the correlation matrix of estimates was calculated.

	ε_L	ε_{HL}	ε_{H_2L}	ε_{H_3L}	pK_{a1}	pK_{a2}	pK_{a3}
ε_L	1	0.229	0.102	0.025	-0.609	-0.115	-0.086
ε_{HL}		1	0.469	0.120	-0.681	-0.528	-0.401
ε_{H_2L}			1	0.522	-0.309	-0.992	-0.976
ε_{H_3L}				1	-0.077	-0.476	-0.659
pK_{a1}					1	0.349	0.263
pK_{a2}						1	0.950
pK_{a3}							1

A strong correlation exists between the following pairs of parameters: ε_{H_2L} and pK_{a2} , ε_{H_2L} and pK_{a3} and pK_{a2} and pK_{a3} .

The ratio of the maximum eigenvalue, with a value $\lambda_{\max} = 5.408 \times 10^{10}$, and the minimum eigenvalue, $\lambda_{\min} = 2.195$, of the matrix $(\mathbf{J}^T \mathbf{J})^{-1}$, is $\lambda_{\max}/\lambda_{\min} = 2.5 \times 10^{10}$, which implies very bad conditioning in the model.

Figure 8.32 illustrates the dependence of the sensitivity function $C_{j(i)}$, $j = 1, \dots, 3$, on pH. It is clear that three out of the seven parameters i.e. ε_{H_2L} , pK_{a2} and pK_{a3} , are strongly ill-conditioned in the model and therefore their estimation is rather uncertain.

Although three parameters are ill-conditioned, the model closely resembles the experimental data. This is obvious from the A -pH curve fitting with the 95% confidence intervals (Fig. 8.33).

The goodness-of-fit is proved by the low value of the residual standard deviation $s(A) = 0.037$, which is close to the instrumental error of the spectrophotometer used, $s_{\text{inst}}(A) = 0.003$, at 470 nm, and regression rabat $D = 99.96\%$ with the Hamilton R -factor = 0.83%.

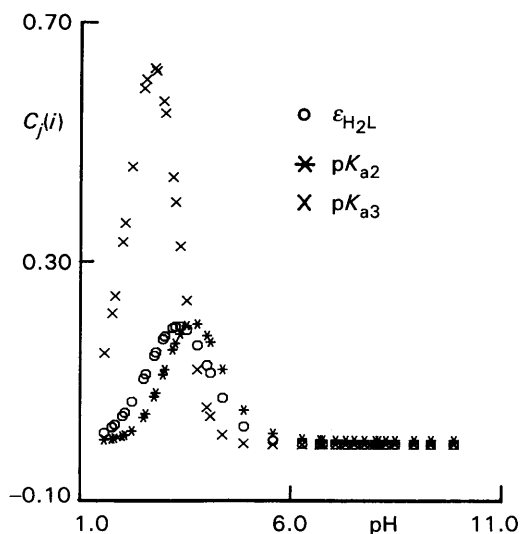


Fig. 8.32—The sensitivity function for the three chosen parameters ϵ_{H_2L} , pK_{a2} and pK_{a3} proves their ill-conditioning in the model.

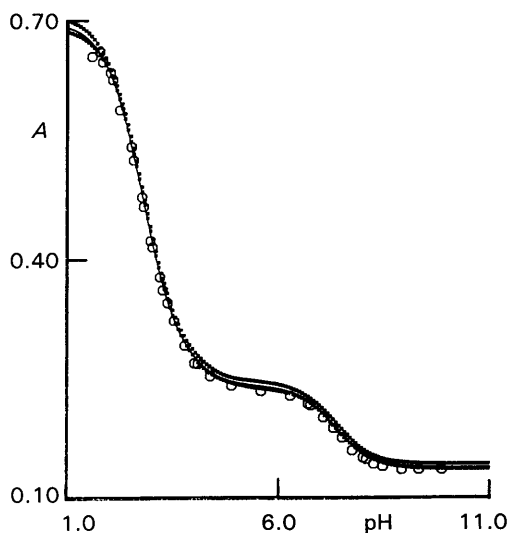


Fig. 8.33—Regression curve for A -pH with the 95% confidence interval.

Conclusion: Even ill-conditioned models can from a numerical point-of-view (goodness-of-fit of the regression curve through experimental points) be quite acceptable. The ill-conditioning of the three parameters ϵ_{H_2L} , pK_{a2} and pK_{a3} is caused by overlapping of the protonation steps, $H_3L \rightleftharpoons H_2L^- \rightleftharpoons HL^{2-}$.

Problem 8.25. *Estimation of parameters of the Freundlich adsorption isotherm*

The relationship between the equilibrium amount of sodium hydroxide, y , g, adsorbed on an active carbon and the concentration of NaOH, x , mmol dm⁻³, was investigated. The additive model and normality of errors are assumed. For adsorption, the model proposed is the Freundlich adsorption isotherm

$$y = \beta_1 x^{\beta_2}$$

Estimate the two parameters, β_1 and β_2 , and their correlation coefficient.

Data: $n = 6$

x	774	381.3	184.8	87.24	43.1	20.46
y	26.0	19.2	15.0	11.6	7.74	5.0

Solution: From the initial guesses of parameters $b_1^{(0)} = 1$, $b_2^{(0)} = 0.1$ and the sum of squares $U(\mathbf{b}^{(0)}) = 1212$, the minimum $U(\mathbf{b}) = 2.072$ was found. The parameter estimates with their half-lengths of the 95% confidence interval Δ_j (8.170), $\Delta_{R,j}$ (8.172) and the relative bias $h_{R,j}$, % are given in Table 8.12.

Table 8.12. Estimates of the parameters of the Freundlich adsorption isotherm with their statistical characteristics

Parameter	Estimate	Δ_j	$\Delta_{R,j}$	$h_{R,j}$, %
β_1	1.724	0.568	0.766	0.534
β_2	0.4084	0.0552	0.0744	0.0572

Good agreement of the model with the experimental data is evident from the high value of the regression rabat $D = 99.3\%$, the mean of absolute residuals $E(\hat{\epsilon}) = 0.498$, the mean of relative residuals $E(\mathbf{e}_R) = 5.44\%$ and the residual standard deviation $s(y) = 0.72$. The correlation coefficient, R_{12} , between β_1 and β_2 , is -0.988 . This implies a strong correlation, so it is not possible to make an interpretation of the two parameters separately.

Figure 8.34 shows the course of the regression curve through the experimental points, together with the 95% confidence interval. The 95% confidence ellipsoid of parameters β_1 and β_2 is shown in Fig. 8.35. Some regression diagnostics for identification of influential points such as Jack-knife residuals \hat{e}_j , likelihood distance LD and LDS are listed in Table 8.13. Two suspicious points are discovered, 4 and 6. *Conclusion:* The Freundlich adsorption isotherm fits the experimental data. For the small number of points, it is difficult to decide whether point 4 is an outlier or not.

Problem 8.26. *Change of estimates of parameters of the Freundlich isotherm after correction of one outlier in data*

Point 4 of the measured data in Problem 8.25 is an outlier and the true value is 11.16 g instead of 11.60 g. Find out whether the change in one point causes any significant change in the parameter estimates and their statistical characteristics.

Data: from Problem 8.25

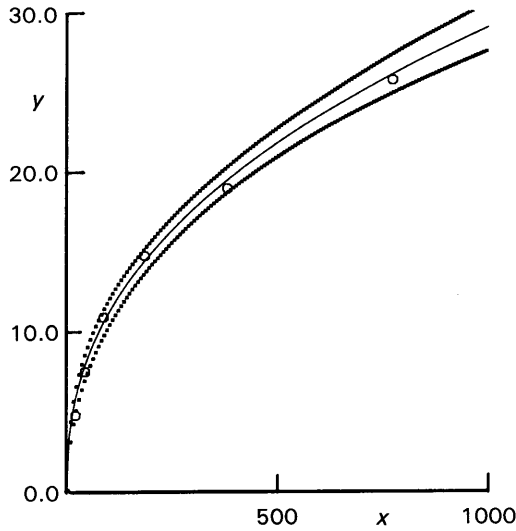


Fig. 8.34—Model of the Freundlich adsorption isotherm, together with the 95% confidence interval.

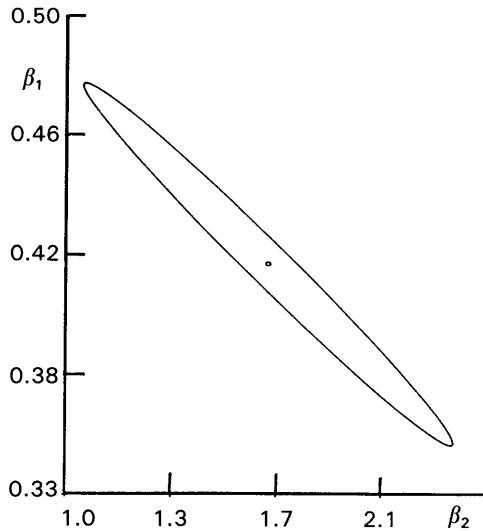


Fig. 8.35—Simultaneous 95% confidence interval for parameters β_1 and β_2 .

Solution: The residual sum-of-squares is $U(\mathbf{b}) = 1.422$. Parameter estimates with their half-length of the 95% confidence interval Δ_j (8.170), $\Delta_{R,j}$ (8.172) and the relative bias $h_{R,j}$, % are given in Table 8.14.

The confidence intervals are narrower and the parameter bias is smaller. The fitting of the regression curve through experimental points is better than in Problem 8.25. The residual standard deviation $s(y) = 0.596$, the regression rabat $D = 99.53\%$, the

Table 8.13. Regression diagnostics indicating influential points (* means a suspicious point)

Point	\hat{e}_{ji}	LD	LDS
1	-0.258	0.051	0.015
2	-0.73	0.028	0.013
3	3.34	0.50	0.016
4	1.308*	0.036*	0.26*
5	-0.62	0.024	0.012
6	1.72*	0.53*	0.23*

Table 8.14. Estimates of the parameters of the Freundlich adsorption isotherm with their statistical characteristics after correction of one outlier

Parameter	Estimate	Δ_j	$\Delta_{R,j}$	$h_{R,j}$, %
β_1	1.669	0.46	0.62	0.374
β_2	0.414	0.046	0.062	0.0397

mean absolute residual $E|\hat{e}| = 0.417$ and the mean relative residual $E(\hat{e}_R) = 4.54\%$. The curve shape, the confidence interval and the parameter estimates are nearly the same as in Problem 8.25. The regression diagnostics \hat{e}_j , and LD , for influential points, are given in Table 8.15. The value of \hat{e}_j does not indicate any influential point in this case.

Table 8.15. Regression diagnostics for influential points

Point	\hat{e}_{ji}	LD_i
1	-0.00528	0.095
2	-0.0189	0.028
3	0.0518	0.102
4	0.074	0.144
5	-0.0304	0.019
6	-0.21	0.81

Conclusion: The correction of one outlier caused a significant change in the statistical characteristics.

REFERENCES

- [1] L. Endrenyi, ed., *Kinetic Data Analysis*, p. 47, Plenum Press, New York, 1983.
- [2] R. C. Magel and D. Hertsgaard, *Commun. Stat.*, 1987, **16**, 85.
- [3] J. M. Criado, M. Gonzalez, A. Ortega and C. Real, *J. Therm. Anal.*, 1984, **29**, 243.
- [4] F. J. Anscombe, *J. Royal Stat. Soc.*, 1969, **B29**, 1.
- [5] A. R. Gallant, *Nonlinear Statistical Models*, Wiley, New York, 1987.
- [6] Y. Bard, *Nonlinear Parameter Estimation*, Academic Press, New York, 1974.
- [7] D. M. Bates and D. G. Watts, *J. Roy. Stat. Soc.*, 1980, **B24**, 1.
- [8] J. C. Nash, *J. Inst. Math. Applics.*, 1977, **19**, 321.
- [9] K. Hiebert, *ACM Trans. Math. Software*, 1981, **7**, 1.
- [10] J. L. Kuester and J. N. Mize, *Optimization Techniques in FORTRAN*, McGraw Hill, New York, 1973.

- [11] M. A. Wolfe, *Numerical Methods for Unconstrained Optimization*, Van Nostrand, New York, 1978.
- [12] P. E. Gill, W. Murray and M. M. Wright, *Practical Optimization*, Academic Press, London, 1981.
- [13] R. Schmidt, *Advances in Nonlinear Parameter Optimization*, Springer, Berlin, 1982.
- [14] T. Nakagawa and Y. Oyanagi, *Program System SALS for Nonlinear Least Squares Fitting*, ISE-TR-13, University of Tsukuba, Japan, 1980.
- [15] M. M. Rosenbrock and C. Storey, *Computational Techniques for Engineers*, Pergamon, Oxford, 1966.
- [16] M. D. J. Powell, *J. Comput.*, 1964, **7**, 155.
- [17] W. Spendley, G. R. Hext and F. R. Himworth, *Technometrics*, 1962, **4**, 441.
- [18] J. A. Nelder and R. Mead, *J. Comput.*, 1965, **7**, 308.
- [19] M. W. Routh, P. A. Swartz and M. B. Denton, *Anal. Chem.*, 1977, **49**, 1422.
- [20] P. B. Ryan, P. L. Barr and M. D. Tod, *Anal. Chem.*, 1977, **49**, 1460.
- [21] S. Marsili-Libelli and M. Castelli, *Appl. Mathematics and Comput.*, 1987, **23**, 341.
- [22] I. A. Volkov, P. I. Grabov and A. B. Potapov, *Zavod. Lab.*, 1985, No. 5, 60.
- [23] F. T. Lindstrom, *Am. Statist.*, 1980, **34**, 183.
- [24] W. Spendley, in *Optimization*, R. Fletcher (ed.), Academic Press, London, 1969.
- [25] W. L. Price, *J. Opt. Theor. Appl.*, 1983, **40**, 333.
- [26] I. O. Bohachevsky, M. E. Johnson and M. L. Stein, *Technometrics*, 1986, **28**, 209.
- [27] M. V. Henckroth, *AIChE Journal*, 1976, **22**, 744.
- [28] L. Pronzato, E. Walter, A. Venot and J. F. Lebruchec, *Math. Comput. Simulation*, 1984, **26**, 412.
- [29] L. G. Sillen and N. Ingri, *Acta Chem. Scand.*, 1962, **16**, 173.
- [30] M. Meloun and J. Cermak, *Talanta*, 1984, **31**, 947.
- [31] G. Peckham, *J. Comput.*, 1970, **13**, 418.
- [32] M. L. Ralston and R. I. Jennrich, *Technometrics*, 1978, **20**, 7.
- [33] J. E. Dennis, D. M. Gay and R. E. Welsch, *ACM Trans. Math. Software*, 1981, **7**, 348.
- [34] H. Ramsin and P. Wedin, *BIT*, 1977, **17**, 72.
- [35] R. I. Jennrich and P. F. Sampson, *Technometrics*, 1968, **10**, 63.
- [36] R. Schmidt, *Advances in Nonlinear Parameter Optimization*, Springer, Berlin, 1982.
- [37] J. E. Dennis and R. E. Welsch, *Commun. Statist.*, 1978, **B7**, 345.
- [38] P. E. Gill and W. Murray, *SIAM J. Num. Anal.*, 1978, **15**, 977.
- [39] J. M. Chambers, *Biometrika*, 1973, **60**, 1.
- [40] J. C. Nash, *Compact Numerical Methods for Computers*, Adam Hilger Ltd., Bristol, 1979.
- [41] N. Wharton and D. K. Olson, *A Generalized Nonlinear Least-Squares Fitting*, Program Rept. ORNL ITM-6545, Oak Ridge Natl. Lab., 1978.
- [42] J. J. More, in *Lecture Notes in Mathematics*, D. Watson (eds.), No. 630, Springer Verlag, Berlin, 1978.
- [43] S. G. Linquist, *Proc. Conf. COMPSTAT 80*, Physica Verlag, Wien, 1980.
- [44] R. R. Meyer and D. M. Roth, *J. Inst. Math. Applics*, 1973, **9**, 218.
- [45] J. E. Dennis and H. H. W. Mei, *J. Opt. Theor. Appl.*, 1979, **28**, 453.
- [46] J. Militky and J. Cap, *Proc. Conf. CEF 87*, Taormina, Sicilia, May, 1987.
- [47] J. V. Beck and K. J. Arnold, *Parameter Estimation in Engineering and Science*, Wiley, New York, 1977.
- [48] A. R. Gallant, *J. Am. Statist. Assoc.*, 1977, **72**, 523.
- [49] E. Z. Demidenko, *Linejnaja i nelinejnaja regresija*, Finansy i Statistika., Moskva, 1981.
- [50] G. M. Stanley and R. S. H. Mah, *Chem. Eng. Sci.*, 1981, **36**, 259.
- [51] V. G. Gorskij, *Zavod. Lab.*, 1987, No. 1, 50.
- [52] K. M. Brown and J. E. Dennis, *Num. Math.*, 1972, **18**, 289; and A. J. Miller in *Interactive Statistics*, D. McNeil (eds.), North Holland, Amsterdam, 1979.
- [53] D. A. Ratkowsky, *Nonlinear Regression Modelling*, Dekker, New York, 1983.
- [54] L. Luksan, *SPONA – Package for Optimization and Nonlinear Approximation*, Resp. Rept. V-4, Central Computing Center of Academy of Sciences, Prague, 1976.
- [55] F. James and M. Ross, *Comp. Phys. Commun.*, 1976, **10**, 343.
- [56] D. M. Bates and D. G. Watts, *J. Roy. Stat. Soc.*, 1986, **B42**, 1.
- [57] J. R. Donaldson and R. B. Schnabel, *Technometrics*, 1987, **29**, 67.
- [58] R. D. Cook E. C.-L. Tsai and B. C. Wei, *Biometrika*, 1986, **73**, 615.
- [59] M. J. Box, *J. Roy. Stat. Soc.*, 1971, **B32**, 171.
- [60] R. Morton, *Biometrika*, 1987, **74**, 679.
- [61] G. P. Clarke, *J. Am. Statist. Assoc.*, 1987, **82**, 221.
- [62] L. Schwartz, *Anal. Chim. Acta*, 1980, **122**, 291.
- [63] R. M. Lyoness, *Commun. Statist.*, 1987, **16**, 997.

- [64] D. M. Himmelblau, *Process Analysis by Statistical Methods*, Wiley, New York, 1970.
- [65] M. Meloun, J. Havel and E. Högfeldt, *Computation of Solution Equilibria*, Ellis Horwood, Chichester, 1988.
- [66] W. C. Hamilton, *Statistical in Physical Science*, Ronald Press, New York, 1964.
- [67] J. Militky, *Proc. 2nd Int. Statist. Conference*, Tampere University Press, 1987.
- [68] H. White and I. Dorniwotz, *Econometrica*, 1984, **52**, 143.
- [69] M. R. Hestena and E. Stiefel, *J. Res. NBS*, 1952, **49**, 409.
- [70] L. C. W. Dixon, *J. Inst. Math. Appl.*, 1975, **15**, 9.
- [71] R. Fletcher, *Comput. J.*, 1970, **13**, 317.
- [72] P. E. Gill and W. Murray, *J. Inst. Maths. Appl.*, 1972, **9**, 91.
- [73] S. S. Oren and D. G. Luenberger, *Management Sci.*, 1974, **20**, 845.
- [74] S. S. Oren, *Management Sci.*, 1974, **20**, 863.
- [75] H. Y. Huang, J. P. Chambliss, *J. Optimization Theory Appl.*, 1974, **13**, 620.
- [76] R. Bass, *Math. of Comp.*, 1972, **26**, 129.
- [77] D. H. Jacobson and W. Oksma, *J. Math. Anal. Appl.*, 1972, **38**, 535.
- [78] E. J. Davison and P. Wong, *Automatica*, 1975, **11**, 197.
- [79] K. Ritter, *Computing*, 1975, **14**, 79.
- [80] W. C. Davidon, *Math. Programming*, 1975, **9**, 1.
- [81] W. H. Swann, *Central Instrument Laboratory Research Note*, 1964, **64**, 3.
- [82] M. J. D. Powell, *J. Comput.*, 1964, **7**, 155.
- [83] W. I. Zangwill, *J. Comput.*, 1967, **10**, 293.
- [84] K. W. Brodlic, *J. Inst. Maths. Appl.*, 1975, **15**, 385.
- [85] L. Nazareth, *Res. Report LBL2692*, University of California 1973.
- [86] R. Mifflin, *Math. Programming*, 1975, **9**, 100.
- [87] R. F. Denneweyer and E. H. Mookini, *J. Optimization Theory Appl.*, 1975, **16**, 67.
- [88] P. E. Gill, Murray, *Math. Programming*, 1974, **7**, 311.
- [89] W. E. Bosarge and P. L. Fabi, *J. Optimization Theory Appl.*, 1969, **4**, 156.
- [90] R. Fletcher, *Res. Report R-6799*, AERE Harwell, 1971.
- [91] M. Meloun and M. Javurek, *Talanta*, 1985, **32**, 973.