

Výstavba regresního modelu regresním tripletem

Prof. RNDr. Milan Meloun, DrSc.,
Katedra analytické chemie, Univerzita Pardubice, 532 10 Pardubice

Souhrn: Postup hledání regresního modelu je popsán obecně a dokumentován na 3 úlohách analytické laboratoře. Skládá se z těchto kroků: 1. Návrh modelu začíná vždy od nejjednoduššího modelu, lineárního. 2. Předběžná analýza dat sleduje proměnlivost proměnných na rozptylových diagramech, indexových grafech. Vyšetruje se multikolinearita, heteroskedasticita, autokorelace a vlivné body. 3. Odhadování parametrů se provádí klasickou metodou nejmenších čtverců, následuje testování významnosti parametrů Studentovým *t*-testem. Střední kvadratická chyba predikce MEP a Akaikovo informační kritérium AIC jsou rozhodčí kritéria při hledání modelu. 4. Regresní diagnostika provádí identifikaci vlivných bodů a ověření předpokladů metody nejmenších čtverců. V případě více vysvětlujících proměnných se posoudí vhodnost proměnných pomocí parciálních regresních grafů a parciálních reziduálních grafů. 5. Konstrukce zpřesněného modelu: parametry zpřesněného modelu jsou odhadovány s využitím (a) metody vážených nejmenších čtverců (MVNC) při nekonstantnosti rozptylu, (b) metody zobecněných nejmenších čtverců (MZNČ) při autokorelaci, (c) metody podmíkových nejmenších čtverců (MPNC) při omezení kladených na parametry, (d) metody racionálních hodnot u multikolinearity, (e) metody rozšířených nejmenších čtverců (MRNČ) pro případ, že všechny proměnné jsou zatížené náhodnými chybami, a konečně (f) robustních metod pro jiná rozdělení než normální a data s vybočujícími hodnotami a extrémy.

Při výstavbě regresních modelů se běžně užívá metody nejmenších čtverců. Metoda nejmenších čtverců poskytuje postačující odhady parametrů jenom při současném splnění všech předpokladů o datech a o regresním modelu. Pokud tyto předpoklady nejsou splněny, ztrácí metoda nejmenších čtverců své vlastnosti.

Základní předpoklady metody nejmenších čtverců (MNČ): Statistické vlastnosti odhadů $\hat{\mathbf{y}}_p, \hat{\boldsymbol{\epsilon}}, \mathbf{b}$ závisí na splnění jistých předpokladů. Pokud platí předpoklady I až IV, jsou odhady \mathbf{b} parametrů β nejlepší, nestranné a lineární (NNLO). Navíc mají asymptoticky normální rozdělení. Pokud platí ještě předpoklad VII, mají odhady \mathbf{b} normální rozdělení i pro konečné výběry.

- I. Regresní parametry β mohou nabývat libovolných hodnot. V praxi však často existují omezení parametrů, která vycházejí z jejich fyzikálního smyslu.
- II. Regresní model je lineární v parametrech a platí aditivní model měření.
- III. Matice nenáhodných, nastavovaných hodnot vysvětlujících proměnných \mathbf{X} má hodnost rovnou právě m . To znamená, že žádné její dva sloupce $\mathbf{x}_j, \mathbf{x}_k$ nejsou kolineární, tj. rovnoběžné vektory. Tomu odpovídá i formulace, že matice $\mathbf{X}^T \mathbf{X}$ je symetrická regulární matice, ke které existuje inverzní matice a jejíž determinant je větší než nula.
- IV. Náhodné chyby $\boldsymbol{\epsilon}_i$ mají nulovou střední hodnotu $E(\boldsymbol{\epsilon}_i) = 0$. To musí u korelačních modelů platit vždy. U regresních modelů se může stát, že $E(\boldsymbol{\epsilon}_i) = K, i = 1, \dots, n$, což znamená, že model neobsahuje absolutní člen. Po jeho zavedení bude $E(\boldsymbol{\epsilon}'_i) = 0$, kde $\boldsymbol{\epsilon}'_i = y_i - \hat{y}_{p,i} - K$.
- V. Náhodné chyby $\boldsymbol{\epsilon}_i$ mají konstantní a konečný rozptyl $E(\boldsymbol{\epsilon}_i^2) = \sigma^2$. Také podmíněný rozptyl $D(y/x) = \sigma^2$ je konstantní a jde o homoskedastický případ.

VI. Náhodné chyby ε_i jsou vzájemně nekorelované a platí $\text{cov}(\varepsilon_i \varepsilon_j) = E(\varepsilon_i \varepsilon_j) = 0$. Pokud mají chyby normální rozdělení, jsou nezávislé. Tento požadavek odpovídá požadavku nezávislosti měřených veličin y .

VII. Chyby ε_i mají normální rozdělení $N(0, \sigma^2)$. Vektor y má pak vícerozměrné normální rozdělení se střední hodnotou $X\beta$ a kovarianční maticí $\sigma^2 E$, kde E je jednotková matice.

Regresní diagnostika

Metoda nejmenších čtverců nezajišťuje obecně nalezení přijatelného modelu, a to jak ze statistického, tak i z fyzikálního hlediska. Musí být splněny podmínky, odpovídající složkám tzv. *regresního tripletu* [data, model, metoda odhadu].

Regresní diagnostika obsahuje postupy k identifikaci

- a) vhodnosti dat pro navržený regresní model (složka *data*),
- b) vhodnosti modelu pro daná data (složka *model*),
- c) splnění základních předpokladů MNČ (složka *metoda*).

Základní rozdíl mezi regresní diagnostikou a klasickými testy spočívá v tom, že u regresní diagnostiky není třeba přesně formulovat alternativní hypotézu. Tímto pojetím se regresní diagnostika blíží spíše k *exploratorní regresní analýze*, která vychází z faktu, že "uživatel ví o analyzovaných datech přece jenom více než počítač". Počítač slouží jako nástroj analýzy dat, modelu a metody odhadu. Model je navrhován v interakci uživatele s programem. Tím by měl být omezen vznik formálních regresních modelů, které nemají fyzikální smysl a jsou v technické praxi obvykle jen omezeně použitelné.

1. Data: mezi základní techniky diagnostiky patří stanovení rozmezí dat, jejich variability a přítomnosti vybočujících pozorování. K tomu lze využít grafů rozptýlení s kvantily a řady postupů průzkumové analýzy jednorozměrných dat. Přes svoji jednoduchost umožňuje diagnostika identifikovat ještě před vlastní regresní analýzou

- a) *nevhodnost dat* (malé rozmezí nebo přítomnost vybočujících bodů),
- b) *nesprávnost navrženého modelu* (skryté proměnné),
- c) *multikolinearitu*,
- d) *nenormalitu* v případě, kdy jsou vysvětlující proměnné náhodné veličiny.

Kvalita dat úzce souvisí s užitým regresním modelem. Při posuzování se sleduje především výskyt *vlivných bodů* (VB), které mohou být hlavním zdrojem řady problémů, jako je zkreslení odhadů a růst rozptylů až k naprosté nepoužitelnosti regresních modelů. Podle toho, kde se vlivné body vyskytují, lze provést dělení na

1. *Vybočující pozorování* (outliers), které se liší v hodnotách vysvětlované (závisle) proměnné y od ostatních, a

2. *Extrémy* (high leverage points), které se liší v hodnotách vysvětlujujcích (nezávisle) proměnných x nebo v jejich kombinaci (v případě multikolinearity) od ostatních bodů.

Vyskytují se však i body, které jsou jak vybočující, tak i extrémní. K identifikaci vlivných bodů typu vybočujícího pozorování se využívá zejména různých typů reziduí a k identifikaci extrémů pak diagonálních prvků H_{ii} projekční matice H .

2. Model: kvalitu regresního modelu lze posoudit v případě jedné vysvětlující proměnné x přímo z rozptylového grafu závislosti y na x . V případě více vysvětlujících proměnných a multikolinearity mohou však rozptylové grafy *mylně indikovat* nelineární trend i u lineárního modelu. Z řady různých grafů k posouzení vztahu y a x_j se omezíme na a) parciální regresní grafy, a b) parciální reziduální grafy.

Parciální regresní grafy byly Belseyem zařazeny mezi základní nástroje počítačové interaktivní analýzy regresních modelů. Umožňují nejenom posouzení kvality navrženého regresního modelu, ale indikují i přítomnost vlivných bodů a nesplnění předpokladů klasické metody nejmenších čtverců. Parciální regresní graf pro posouzení vztahu mezi y a i -tou vysvětlující proměnnou x_i je závislost *reziduů v regrese* y na sloupcích matice $\mathbf{X}_{(i)}$ a reziduů u regrese x_i na sloupcích matice $\mathbf{X}_{(i)}$. Přitom matice $\mathbf{X}_{(i)}$ vznikne z matice \mathbf{X} vynecháním i -tého sloupce x_i , odpovídajícího i -té vysvětlující proměnné. Parciální regresní grafy mají tyto vlastnosti:

a) Směrnice přímky v parciálním regresním grafu je stejná jako odhad b_j v neděleném modelu a úsek je roven nule. Tato lineární závislost platí pouze v případě, že navržený model je správný.

b) Korelační koeficient mezi oběma proměnnými parciálního regresního grafu odpovídá parciálnímu korelačnímu koeficientu $R_{yx(x)}$.

Parciální reziduální grafy se označují také jako grafy "komponenta + reziduum". Parciální reziduální grafy však poskytují poněkud odlišné informace než parciální regresní grafy. Směrnice lineární závislosti je rovna b_j a úsek je nulový. Lineární závislost pak ukazuje na vhodnost navržené proměnné x_j v modelu.

Parciální reziduální grafy se doporučují především k indikaci rozličných typů nelinearity v případě nesprávně navrženého regresního modelu.

3. Metoda: V praxi bývají některé předpoklady MNČ porušeny, což vede k použití jiných kritérií. K porušení předpokladů dochází v těchto základních případech:

a) Na parametry jsou kladena omezení, což vede na užití *metody podmínkových nejmenších čtverců (MPNČ)*.

b) Kovarianční matice chyb není diagonální (autokorelace), příp. data nemají stejný rozptyl (heteroskedasticita), což vede na užití *metody zobecněných nejmenších čtverců (MZNČ)*, resp. *metody vážených nejmenších čtverců (MVNČ)*.

c) Rozdělení dat nelze považovat za normální nebo se v datech vyskytují vlivné body. V takovém případě se místo kritéria metody nejmenších čtverců užije *robustního kritéria*, které je na porušení předpokladu o rozdělení chyb a na vlivné body málo citlivé. Z robustních kritérií jsou nejznámější *M-odhad*. Jedná se o maximálně věrohodné odhady pro vhodnou hustotu pravděpodobnosti chyb. Pro odhad parametrů \mathbf{b} se užívá *iterační metody vážených nejmenších čtverců (IVNČ)*.

d) Také proměnné x mohou být zatížené náhodnými chybami, což vede na užití *metody rozšířených nejmenších čtverců (MRNČ)*. Pro případ regresní přímky je použití metody rozšířených nejmenších čtverců velmi jednoduché. Postačuje znalost poměru rozptylu σ_y^2 (vysvětlovaná proměnná) a σ_x^2 (vysvětlující proměnné), $K = \sigma_y^2/\sigma_x^2$. Pro odhad směrnice regresní přímky $y = a x + b$ pak platí

$$a = L + \text{sign}(S_{yx}) \sqrt{K + L^2}$$

kde

$$L = \frac{S_{yx} - K S_x}{2S_x}$$

a $\text{sign } S_{yx}$ je znaménková funkce. Symboly S označují součty čtverců, odpovídajících proměnných

$$S_x = \sum_{i=1}^n (x_i - \bar{x})^2, \quad S_y = \sum_{i=1}^n (y_i - \bar{y})^2, \quad S_{yx} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

Při znalosti odhadu směrnice \hat{a} se snadno určí odhad úseku \hat{b} ze vztahu

$$\hat{b} = \bar{y} - \hat{a} \bar{x}$$

Pro případ stejných rozptylů, tj. $K = 1$ vede dosazení do výše uvedených vztahů k odhadům minimalizujícím kolmé vzdálenosti (*orthogonální regrese*). Pro odhady rozptylů odhadů \hat{a}, \hat{b} se pak používá speciálních vztahů.

e) Pro špatně podmíněné matice $X^T X$ se používá *metoda racionálních hodností*, vedoucí k systému vychýlených odhadů, kde vychýlení je řízeno jedním parametrem.

Postup výstavby lineárního regresního modelu:

1. Návrh modelu: začíná se vždy od nejjednoduššího modelu, u kterého vystupují jednotlivé vysvětlující proměnné v prvních mocninách a nevyskytují se žádné interakční členy typu $x_j x_k$.

2. Předběžná analýza dat: sleduje se proměnlivost jednotlivých proměnných a možné párové vztahy. Užívá se proto rozptylových diagramů závislosti x_j na x_k nebo indexových grafů závislosti x_j na j . Posuzuje se významnost proměnných s ohledem na jejich proměnlivost a přítomnost multikolinearity. Přibližně lineární vztah mezi proměnnými v rozptylových grafech závislosti x_j na x_k indikuje multikolinearitu. Lze rovněž odhalit vlivné body, které způsobují multikolinearitu.

Podle volby uživatele se provedou požadované transformace původních proměnných. Zadává se, zda model obsahuje absolutní člen. Uživatel může volit polynomickou transformaci zadáním stupně polynomu.

Provádí se sestavení korelační matice R a její rozklad na vlastní čísla a vlastní vektory. Jsou vypočteny VIF k indikaci multikolinearity a tisknuta setříděná vlastní čísla. K určení inverzní matice R^{-1} se užívá metoda racionálních hodností pro standardně zadávané vychýlení $P = 10^{-15}$. Uživatel může zadat jinou hodnotu parametru vychýlení P , což však vede pro vyšší hodnoty P k vychýleným odhadům. Bývá proto vhodné volit P z tohoto intervalu $10^{-5} \leq P \leq 10^{-3}$.

3. Odhadování parametrů: odhadování parametrů modelu se provádí metodou racionálních hodností s volbou $P = 10^{-5}$. Ze zobecněné inverzní matice R^{-1} jsou určovány odhady parametrů b , jejich směrodatné odchyly $\sqrt{D(b_j)}$ a velikosti testačních statistik Studentova t -testu významnosti pro $\beta_j = 0$. Dále jsou provedeny testy významnosti odhadů b_j , vícenásobného korelačního koeficientu R a koeficientu determinace D . Je vhodné sledovat souhrnné charakteristiky regrese jako je střední kvadratická chyba predikce MEP a Akaikovo informační kritérium AIC , případně posoudit linearitu modelu.

4. Regresní diagnostika: s využitím pěti rozličných grafů je prováděna identifikace vlivných bodů, a to *grafy Williamsovým, Pregibonovým, McCulloh-Meeterovým, L-R*, a *grafem predikovaných reziduí*. Dále pak ověření splnění předpokladů metody nejmenších čtverců jako je homoskedasticita, nepřítomnost autokorelace a normalita rozdělení chyb. Pokud dojde k úpravě dat, je třeba provést znovu regresní diagnostiku se zaměřením na porušení předpokladů metody nejmenších čtverců a posouzení vlivu multikolinearity. V případě více vysvětlujících proměnných se posoudí vhodnost jednotlivých proměnných a jejich funkcí s využitím parciálních regresních grafů nebo grafů "komponenta + reziduum". *Tabulka reziduí*

obsahuje klasická rezidua \hat{e}_i , normovaná rezidua \hat{e}_{Ni} , standardizovaná rezidua \hat{e}_{Si} a Jackknife rezidua \hat{e}_{Ji} . Je uveden odhad autokorelačního koeficientu reziduí prvního řádu $\hat{\rho}_1$. Tabulka vlivných bodů obsahuje veličiny H_{ii} , H_{ii}^* , D_i , A_i , DF_i , $LD_i(\mathbf{b})$, $LD_i(\hat{\sigma}^2)$ a $LD_i(\mathbf{b}, \hat{\sigma}^2)$. Hvězdičkou jsou označeny hodnoty silně vlivných bodů.

5. Konstrukce zpřesněného modelu:

- s využitím
- metody vážených nejmenších čtverců (MVNC) při nekonstantnosti rozptylů,
 - metody zobecněných nejmenších čtverců (MZNČ) při autokorelací,
 - metody podmínkových nejmenších čtverců (MPNČ) při omezeních na parametry,
 - metody racionálních hodností RH u multikolinearity,
 - metody rozšířených nejmenších čtverců (MRNČ) pro případ, že všechny proměnné jsou zatížené náhodnými chybami,
 - robustní metody pro jiná rozdělení dat než normální a data s vybočujícími hodnotami a extrémy jsou odhadovány parametry zpřesněného modelu.

6. Zhodnocení kvality modelu: s využitím klasických testů, postupu regresní diagnostiky a doplnkových informací o modelované soustavě se provede posouzení kvality navrženého lineárního regresního modelu.

Vzorová úloha: Model teplotní závislosti přechodového tlaku bismutu (J6.01)

Ukážeme postup analýzy jednorozměrného lineárního regresního modelu. Byl studován přechodový tlak bismutu I - II p jako funkce teploty t . Nalezněte lineární regresní model, který bude adekvátní daným datům. Vyšetřete regresní triplet a indikujte vlivné body.

Data: Teplota t [$^{\circ}\text{C}$], tlak p [bar]:

20.8	25276,	20.9	25256,	21.0	25216,	21.9	25187,	22.1	25217,
22.1	25187,	22.4	25177,	22.5	25177,	24.8	25098,	24.8	25093,
25.0	25088,	34.0	24711,	34.0	24701,	34.1	24716,	42.7	24374,
42.7	24394,	42.7	24384,	49.9	24067,	50.1	24057,	50.1	24057,
22.5	25147,	23.1	25107,	23.0	25077				

Řešení:

1. Odhadování parametrů: klasickou metodou nejmenších čtverců (MNČ) byly nalezeny nejlepší odhady úseku β_0 a směrnice β_1 . Studentův t -test ukázal, že úsek (absolutní člen) β_0 je statisticky významný a směrnice β_1 je statisticky významná.

Odhad	Směrodatná odchylka	Test H0: $B[j] = 0$ vs. HA: $B[j] \neq 0$	t-kriterium	hypoteza H0 je	Hlad. význam.
B[0]	2.6068E+04	1.6169E+01	1.6122E+03	Zamítnuta	0.000
B[1]	-3.9874E+01	5.0419E-01	-7.9084E+01	Zamítnuta	0.000

2. Regresní diagnostika: absolutní hodnota párového korelačního koeficientu R ukazuje, že navržený lineární regresní model je statisticky významný. Vysoká hodnota koeficientu determinace $D = R^2$ (99.67%), představuje procento variability, vysvětlené modelem. Predikovaný koeficient determinace R_p^2 ukazuje na predikční schopnost modelu, je však vyčíslen jinak než R^2 , místo RSC se ve vztahu užije MEP. Střední kvadratická chyba predikce MEP a Akaikovo informační kritérium AIC se užívají k rozlišení mezi několika navrženými modely. Za optimální se považuje model, pro který dosahuje MEP a AIC minimální hodnotu.

Vícenásobný korelační koeficient, R	: 9.9833E-01
Koeficient determinace, D	: 9.9665E-01
Predikovaný koeficient determinace, Rp^2	: 9.9804E-01

Střední kvadratická chyba predikce, *MEP*
Akaikeho informační kritérium, *AIC*

: 6.8546E+02
: 1.5054E+02

3. Konstrukce zpřesněného modelu:

(a) Po odstranění bodů č. 23 (*kritika dat*) byly nalezeny nové odhady parametrů zpřesněného modelu. Zpřesněný model (v závorce je uveden vždy odhad směrodatné odchylky parametru) $y = 26\ 078\ (13) - 40.1\ (0.4)\ x_1$ je doložen statistickými charakteristikami: *párový korelační koeficient R = 0.9990, koeficient determinace D = 99.808% a predikovaný korelační koeficient R_p = 0.99885* dosáhly vesměs vysokých hodnot. *Střední kvadratická chyba predikce MEP = 414.22 a Akaikeho informační kritérium AIC = 132.62* dosáhly nižších hodnot než u předešlého modelu, což dokazuje, že zpřesněný model je lepší. Rezidua nyní vykazují normální rozdělení a nevykazují trend, stále však vykazují heteroskedasticitu, a proto lze doporučit použít metodu vážených nejmenších čtverců.

(b) Užitím statistické váhy ($w_i = 1/y_i^2$) kompenzujeme heteroskedasticitu v datech. Obdržíme nové odhady parametrů. Opravený model má tvar, (v závorce je uveden odhad směrodatné odchylky parametru) $y = 26\ 079\ (13) - 40.1\ (0.4)\ x_1$. Jelikož došlo ke snížení rozhodujících kritérií, tj. *střední kvadratické chyby predikce MEP = 410.29 a Akaikeho informačního kritéria AIC = 132.39*, lze považovat tyto odhady za lepší než předešlé.

4. Zhodnocení kvality modelu: porovnáním hodnot regresní diagnostiky lze snadno provést zhodnocení *regresního tripletu* u dosaženého lineárního regresního modelu pro upravená data, zbavená odlehčích hodnot a metodou vážených nejmenších čtverců. Nalezený a prokázaný model teplotní závislosti přechodového tlaku bizmutu má tvar, (v závorce je vždy uveden odhad směrodatné odchylky parametru)

$$y = 26\ 079\ (13) - 40.1\ (0.4)\ x_1.$$

Poděkování: Práce vznikla za podpory grantu Ministerstva zdravotnictví NS9831-4/2008 a vědeckých záměrů MSMT0021627502.

Doporučená literatura

- [1] Meloun M., Militký J., *Statistické zpracování experimentálních dat*, Plus Praha 1994, Academia Praha 2004..
- [2] Meloun M., Militký J., Hill M., *Počítačová analýza vícerozměrných dat v příkladech*, Academia Praha 2005.
- [3] Meloun M., Militký J., *Kompendium statistického zpracování experimentálních dat*, Academia Praha 2002, 2006.

Výstavba regresního modelu regresním tripletem

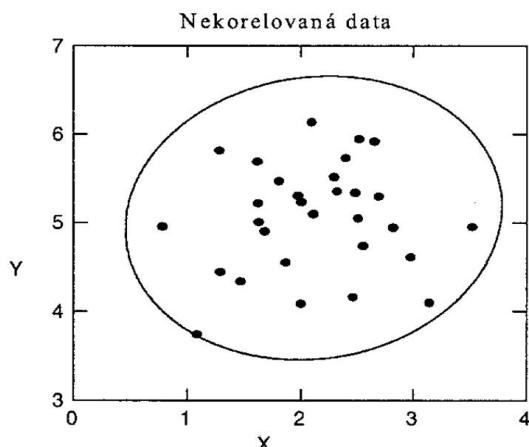
Prof. RNDr. Milan Meloun, DrSc.,
Katedra analytické chemie,
Univerzita Pardubice, 532 10 Pardubice

<http://meloun.upce.cz>

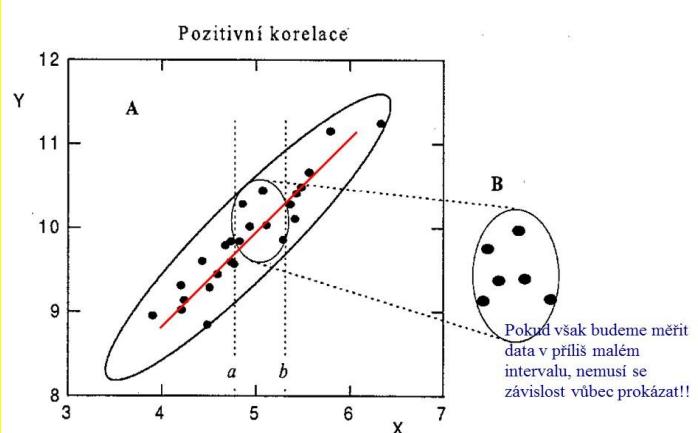
Definice regresního modelu

6

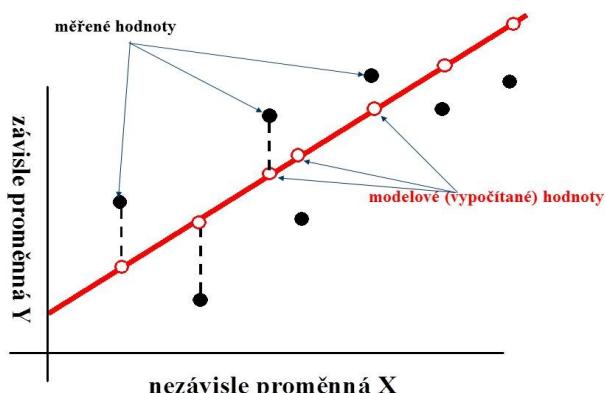
STATISTICKÁ ZÁVISLOST



STATISTICKÁ ZÁVISLOST



Grafické vysvětlení:



Definice regresního modelu

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1j} & \cdots & x_{1m} \\ x_{21} & x_{22} & \cdots & x_{2j} & \cdots & x_{2m} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{i1} & x_{i2} & \cdots & x_{ij} & \cdots & x_{im} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nj} & \cdots & x_{nm} \end{bmatrix} \cdot \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_j \\ \vdots \\ \beta_m \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_i \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

\mathbf{y} závisle proměnná \mathbf{X} nezávisle proměnná $\boldsymbol{\beta}$ regresní parametry $\boldsymbol{\varepsilon}$ náhodná chyba

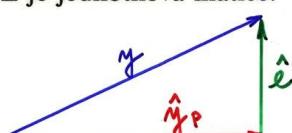
$$\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

Vlastnosti projekční matice P :

Projekční matice P pro kolmou projekci do nadroviny L^\perp , kolmé na nadrovinu L , má tvar

$$\mathbf{P} = \mathbf{E} - \mathbf{H}$$

kde \mathbf{E} je jednotková matice.



Závěr: Rozklad vektoru y do dvou složek

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

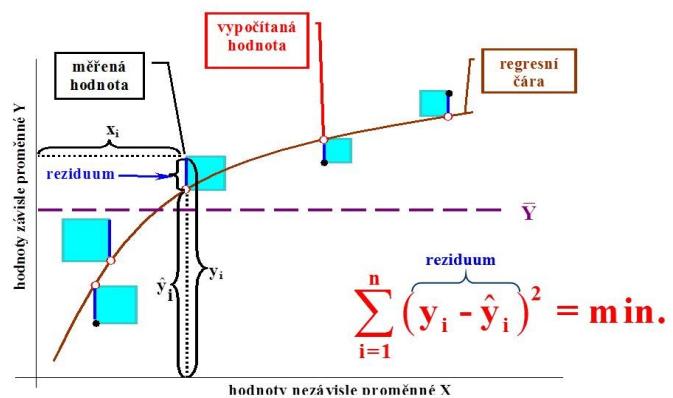
Geometricky: vektor y byl rozložen na dva kolmé vektory

Předpoklady metody nejmenších čtverců

- Parametry β mohou nabývat libovolných hodnot. Omezení jsou pouze fyzikálního smyslu.
- Model je lineární v parametrech β ; platí přitom additivní model měření $y = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}$.
- Matrice \mathbf{X} je nenáhodná, nastavitelných hodnot nezávisle promenných. Má hodnost m :
 - Žádné dva sloupce x_j a x_k nejsou kolineární (čili paralelní).
 - $\mathbf{X}^T \mathbf{X}$ je pak symetrická regulární matice.
 - Rovina L je m rozměrná a vektory $\mathbf{X} \mathbf{b}$ jsou jednoznačné.

4. Náhodné chyby ϵ_i mají nulovou střední hodnotu $E(\epsilon_i) = 0$.
Je-li $E(\epsilon_i) = K$, je nutno zavést absolutní člen a pak bude $E(\epsilon_i) = 0$.
5. Náhodné chyby ϵ_i mají konstantní rozptyl, homoskedasticita, $E(\epsilon_i^2) = \sigma^2$.
6. Náhodné chyby ϵ_i jsou vzájemně nekorelované, $\text{cov}(\epsilon_i \cdot \epsilon_j) = E(\epsilon_i \cdot \epsilon_j) = 0$.
7. Náhodné chyby mají normální rozdělení $\epsilon \sim N(0, \sigma^2)$.

VÝČÍSLENÍ ODHADŮ PARAMETRŮ REGRESNÍHO MODELU METODOU NEJMENŠÍCH ČTVERCŮ (MNČ)



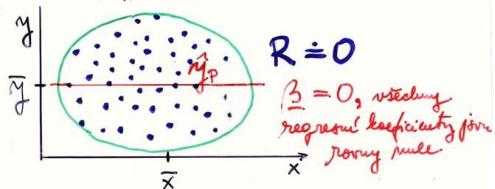
Využití vícenásobného korelačního koeficientu R

$$R = \sqrt{\frac{TSC}{CSC}} = \sqrt{1 - \frac{RSC}{CSC}} = \cos\alpha$$

$$= \sqrt{1 - \frac{\sum_{i=1}^n e_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

DISKUSE: (a) $\hat{R} \rightarrow 0$, pak platí $H_0: \hat{y}_p = \bar{y} \cdot \bar{y}$

čili $H_0: \beta_2 = 0$

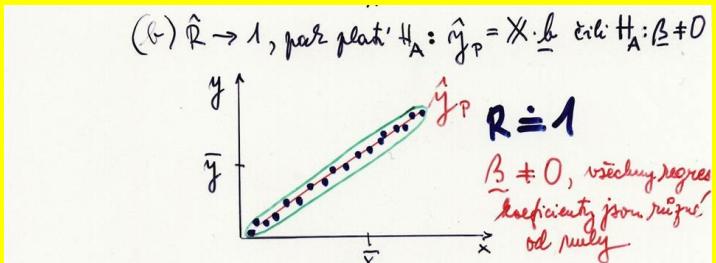
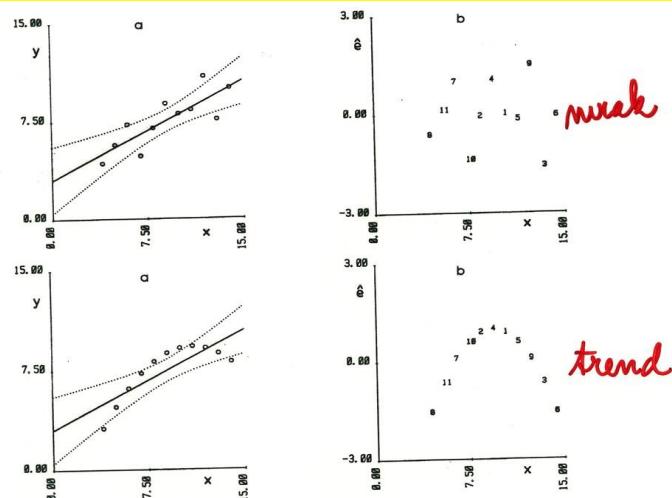


Příklad 6.8 Omezení klasické analýzy lineárního modelu

Anscomb⁵ uvádí testační data pro čtyři simulované výběry. Testujte statistickou významnost obou parametrů β_1 a β_2 a provedte grafickou analýzu reziduí.

Data: čtyři simulované výběry vykazují stejné charakteristiky
 $b_1 = 0.5$, $b_2 = 3.0$, $D(b_1) = 0.0139$ a $D(b_2) = 1.2656$.

Výběr	A	B	C	D		
Proměnná	x	y	y	y	x	y
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



TESTOVÁNÍ R: $H_0: \beta_2 = 0$ čili $\hat{R}^2 = 0$ vs. $H_A: \beta_2 \neq 0$ čili $\hat{R}^2 \neq 0$

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

~ porovnat s $F_{1-\alpha}(m-1, m-n)$

Řešení:
Lineární regresní model $E(y/x) = \beta_1 x + \beta_2$ má pro všechny výběry

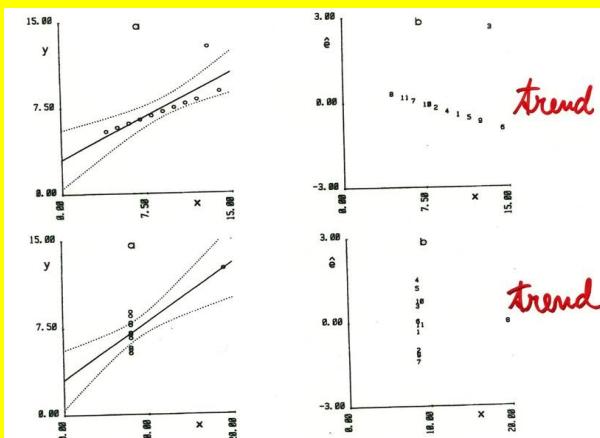
a) Stejné odhad parametrů:

$$b_1 = 0.5, b_2 = 3.0, D(b_1) = 0.0139, D(b_2) = 1.2656.$$

b) Stejné testační statistiky významnosti parametrů:
 $T_1 = 2.667$ a $T_2 = 4.241$.

c) Stejné testační statistiky: $F_R = 17.97$, $\hat{R}^2 = 0.66$, $\hat{\sigma} = 1.237$ ukazují,

že β_1 a β_2 jsou významně odlišné od nuly.



Závěr: Neshodu modelu s daty indikuje grafická analýza reziduí

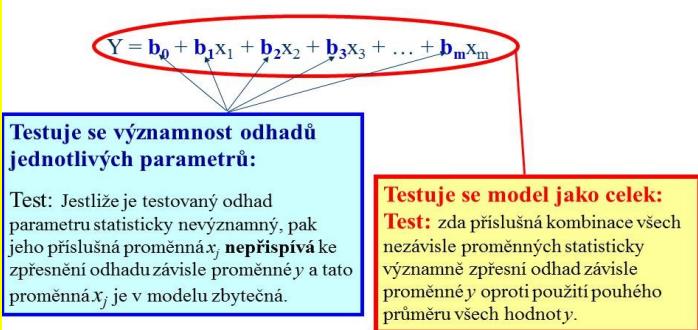
Testy hypotéz a testační kritéria

PODSTATNÉ TESTY VÝZNAMNOSTI V KORELAČNÍ A REGRESNÍ ANALÝZE

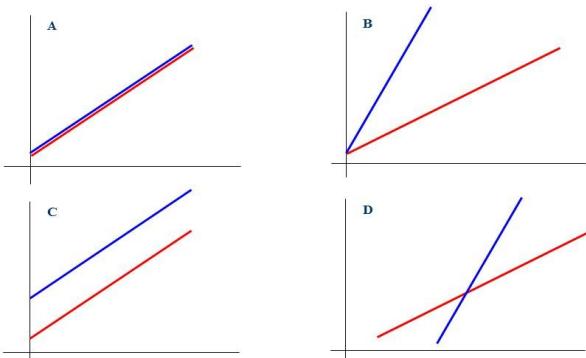
- ◆ test významnosti korelačního koeficientu
- ◆ test významnosti modelu jako celku
- ◆ test významnosti jednotlivých regresních parametrů
- ◆ test shody lineárních regresních modelů
- a mnoho dalších testů.....

21

TEST VÝZNAMNOSTI REGRESNÍHO MODELU – co testujeme?



Porovnání regresních přímek



Porovnání regresních přímek

Porovnání M navržených regresních modelů

$$y_{ij} = \beta_{2j} + \beta_{1j} x_{ij} + \epsilon_{ij} \quad j = 1, \dots, M \quad i = 1, \dots, n_j$$

Postup:

- 1. krok:** vyčíslení β_{2j} , β_{1j} , $\hat{\sigma}_j^2$, $j = 1, \dots, M$
- 2. krok:** test homoskedasticity $\hat{\sigma}_j^2 = \hat{\sigma}^2$, $j = 1, \dots, M$
- 3. krok:** testování, zda-li
 - a) regresní přímky mají společný průsečík,
 - b) regresní přímky mají společnou směrnici,
 - c) regresní přímky jsou totožné.

Test shody dvou lineárních modelů

test shody parametrů β_1 a β_2 dvou lineárních modelů

$$\mathbf{y}_1 = \mathbf{X}_1 \boldsymbol{\beta}_1 + \boldsymbol{\epsilon}_1 \quad \dots \text{RSC}_1$$

$$\mathbf{y}_2 = \mathbf{X}_2 \boldsymbol{\beta}_2 + \boldsymbol{\epsilon}_2 \quad \dots \text{RSC}_2$$

kde \mathbf{X}_1 rozměru $(n_1 \times m)$, \mathbf{y}_1 rozměru $(n_1 \times 1)$, \mathbf{X}_2 rozměru $(n_2 \times m)$, \mathbf{y}_2 rozměru $(n_2 \times 1)$.

$$\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix} \boldsymbol{\beta} + \begin{bmatrix} \boldsymbol{\epsilon}_1 \\ \boldsymbol{\epsilon}_2 \end{bmatrix} \quad \dots \text{RSC}$$

Chowův test shody dvou lineárních modelů

$$H_0: \beta_1 = \beta_2 \text{ vs. } H_A: \beta_1 \neq \beta_2$$

Testační kritérium

$$F_c = \frac{(RSC - RSC_1 - RSC_2)(n - 2m)}{(RSC_1 + RSC_2)m}$$

kde $n = n_1 + n_2$

Testování:

- (a) Při homoskedasticitě $\sigma_1^2 = \sigma_2^2$: F_c vs. $F(m, n - 2m)$,
 (b) Při heteroskedasticitě $\sigma_1^2 \neq \sigma_2^2$: F_c vs. $F(m, r)$,

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

Příklad 6.14 Porovnání výsledků měření ze dvou laboratoří. Stanovení volné enthalpie $-\Delta G$ par oxida boritého v závislosti na teplotě T [Kelvin] bylo provedeno paralelně ve dvou laboratořích. Je možné hodnoty naměřené v obou laboratořích považovat za shodné?

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

T [K]	Laboratoř A	Laboratoř 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

Řešení: Lineární regresní model pro obě skupiny dat

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

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

Chowův test: $H_0: \beta_A = \beta_B$ proti $H_A: \beta_A \neq \beta_B$,
kde $\beta_A = (\beta_{1,A}, \beta_{2,A})^T$ a $\beta_B = (\beta_{1,B}, \beta_{2,B})^T$.

Laboratoř A: $b_{1,A} = -2.768 \cdot 10^{-2} (\pm 5.25 \cdot 10^{-3})$
 $b_{2,A} = 73.73 (\pm 7.865)$, $\hat{\sigma} = 0.916$
 $RSC_A = 3.358$

Laboratoř B: $b_{1,B} = -2.776 \cdot 10^{-2} (\pm 1.57 \cdot 10^{-4})$
 $b_{2,B} = 73.82 (\pm 0.235)$, $\hat{\sigma} = 0.0274$
 $RSC_B = 2.992 \cdot 10^{-3}$

Laboratoř A + B: (označení C = A + B)
 $b_{1,C} = -2.772 \cdot 10^{-2} (\pm 2.35 \cdot 10^{-3})$
 $b_{2,C} = 73.77 (\pm 3.521)$, $\hat{\sigma} = 0.58$
 $RSC_C = 3.364$

Dosazením za $RSC = RSC_C$, $RSC_1 = RSC_A$ a $RSC_2 = RSC_B$ do testačního kritéria Chowova testu

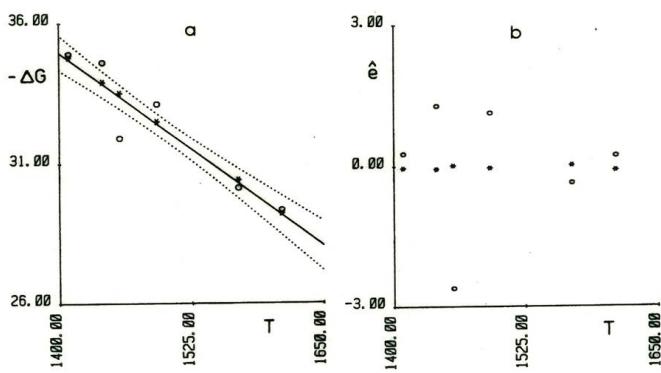
$$F_c = \frac{(RSC - RSC_1 - RSC_2)(n - 2m)}{(RSC_1 + RSC_2)m}$$

bude

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

Data A se od B liší v rozptylu, budou stupně volnosti dle

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



Test: Jelikož $F_{0.95}(2, 4) = 6.94 > F_c$, je $H_0: \beta_A = \beta_B$ přijata.

Závěr: Chowovým testem je prokázáno, že výsledky v obou laboratořích A a B jsou shodné.

HODNOCENÍ KVALITY REGRESNÍHO MODELU

Střední kvadratická chyba predikce (MEP)

$$MEP = \frac{1}{n} \sum_{i=1}^n \frac{e_i^2}{H_{ii}} \quad e_i^2 \text{ čtverec reziduí modelu} \\ H_{ii} \text{ } i\text{-tý diagonální prvek projekční matici } H$$

Akaikovo informační kritérium (AIC)

$$AIC = n \cdot \ln\left(\frac{RSC}{n}\right) + 2m \quad RSC \text{ reziduální součet čtverců} \\ m \text{ počet parametrů}$$

Čím je AIC (MEP) menší, tím je model vhodnější.

Kritika dat regresní diagnostikou

REGRESNÍ DIAGNOSTIKA

Vyšetruje **regresní triplet**, což představuje

- ◆ **Kritiku dat** (zkoumá kvalitu dat pro navržený model)
- ◆ **Kritiku modelu** (zkoumá kvalitu modelu pro daná data)
- ◆ **Kritiku metody odhadu** (prověřuje splnění všech předpokladů požadovaných metodou MNČ)

33

Kritika dat

Vyšetření vlivných bodů:

- (1) Odlehlé body
- (2) Extrémy

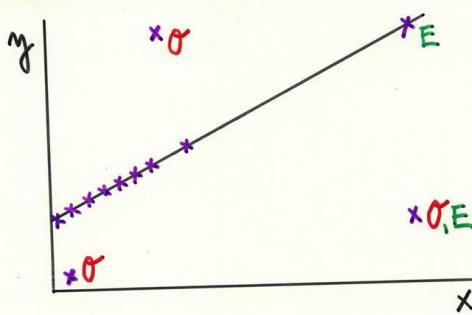
Metoda:

Grafy vlivných bodů

Vybrané diagnostiky vlivných bodů

Dělení vlivných bodů dle výskytu:

1. vybočující pozorování (outliers, O): na y se liší,
2. extrémy (high leverage points, E): liší se na ose x



Statistická analýza reziduí

1. Klasická rezidua: $\hat{e}_i = y_i - \mathbf{x}_i \mathbf{b}$,

Nesprávné představy o reziduích:

1. Rozdělení reziduí je stejné jako rozdělení chyb,
2. Vlastnosti reziduí jsou shodné s vlastnostmi chyb,
2. Čím je reziduum \hat{e}_i větší, tím je bod vlivnější, a tím spíše by se měl z dat vyloučit.

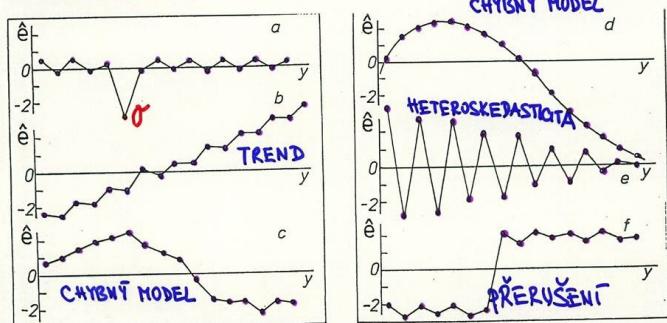
$$\hat{\mathbf{e}} = \mathbf{P} \mathbf{y} = \mathbf{P}(\mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \mathbf{P} \boldsymbol{\varepsilon} = (\mathbf{E} - \mathbf{H}) \boldsymbol{\varepsilon}$$

2. Numerická analýza:

- (a) Střední hodnota reziduí $E(\hat{e}) = 0$,
- (b) Průměrné reziduum $|\bar{e}| = \frac{1}{n} \sum |\hat{e}_i| \approx \epsilon$,
- (c) Směrodatná odchylka střední hodnoty reziduí $s(\hat{e}) \approx \epsilon$,
- (d) Koeficient šikmosti souboru reziduí $g_1(\hat{e}) \approx 0$,
- (e) Koeficient špičatosti souboru reziduí $g_2(\hat{e}) \approx 3$,
- (f) Pearsonův χ^2 -test dobré shody: χ^2_{exp} vs. $\chi^2_{\text{crit.}}$,
- (g) Hamiltonův R-faktor relativní těsnosti: $R \approx 0.5 \%$.

Statistická analýza klasických reziduí

1. Grafická analýza



2. Normovaná rezidua $\hat{\sigma}_N$

$$\hat{e}_{Ni} = \hat{e}_i / \hat{\sigma}$$

normálně rozdělené veličiny $\hat{e}_{Ni} \sim N(0, 1)$.

Diagnostika: pravidlo 3σ , tj. rezidua větší než $\pm 3\hat{\sigma}$ indikují vybočující.

3. Standardizovaná rezidua $\hat{\sigma}_S$

$$\hat{e}_{Si} = \frac{\hat{e}_i}{\hat{\sigma} \sqrt{1 - H_{ii}}}$$

Mají konstantní rozptyl. Maximální hodnota \hat{e}_{Si} je ohraničena velikostí $\sqrt{n - m}$.

Diagnostika: k indikaci heteroskedasticity.

4. Jackknife rezidua $\hat{\sigma}_J$ ("plně studentizovaná")

užijeme místo $\hat{\sigma}$ odhadu směrodatné odchylky $\hat{\sigma}_{(-i)}$,

$$\hat{e}_{Ji} = \hat{e}_{Si} \sqrt{\frac{n - m - 1}{n - m - \hat{e}_{Si}^2}} = \sqrt{n - m} \cotg \Theta_i$$

mají Student. rozdělení s $(n - m - 1)$ stupni volnosti.

Diagnostika: k identifikaci vybočujících bodů (outliers).

5. Predikovaná rezidua $\hat{\sigma}_P$

$$\hat{e}_{Pi} = y_i - \mathbf{x}_i \mathbf{b}_{(i)} = \frac{\hat{e}_i}{1 - H_{ii}}$$

kde $\mathbf{b}_{(i)}$ jsou MNČ odhadové ze všech bodů kromě i-tého.

Diagnostika: indikace vybočujících hodnot (outliers).

6. Rekurzívní rezidua $\hat{\sigma}_R$

Dopředná rekurzívní rezidua jsou definovány vztahy

$$\hat{e}_{Ri} = 0, \quad i = 1, \dots, m$$

$$\hat{e}_{Ri} = \frac{y_i - \mathbf{x}_i \mathbf{b}_{i-1}}{\sqrt{1 + \mathbf{x}_i (\mathbf{X}_{i-1}^{-T} \mathbf{X}_{i-1})^{-1} \mathbf{x}_i^T}} \quad i = m + 1, \dots, n$$

kde \mathbf{b}_{i-1} jsou odhadové získané z prvních $(i - 1)$ bodů.

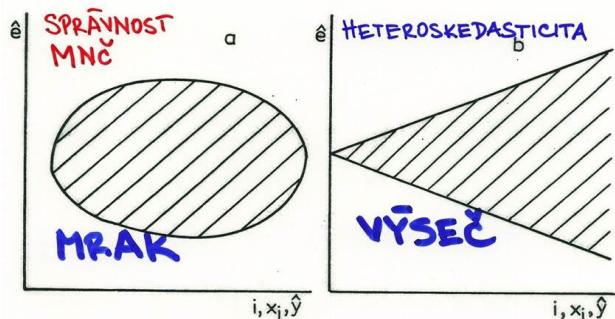
Diagnostika: umožňují identifikovat autokorelací čili nestabilitu modelu, např. v čase.

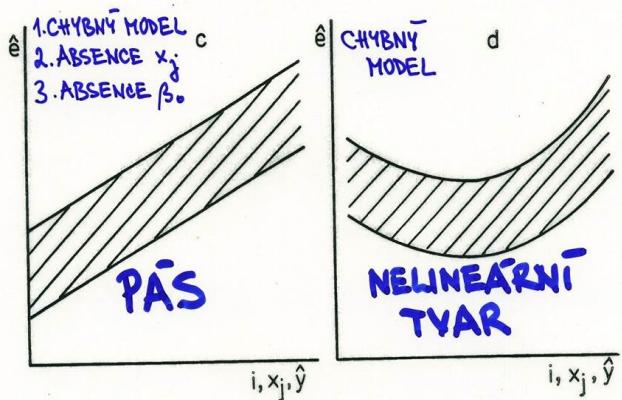
Obrazce v diagnostických grafech

Typ I: Graf reziduů \hat{e}_i proti indexu i.

Typ II: Graf reziduů \hat{e}_i proti proměnné x_j .

Typ III: Graf reziduů \hat{e}_i proti predikci \hat{y}_i .

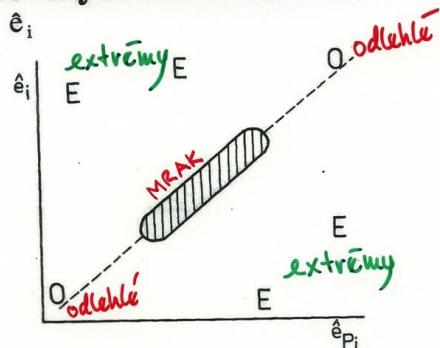




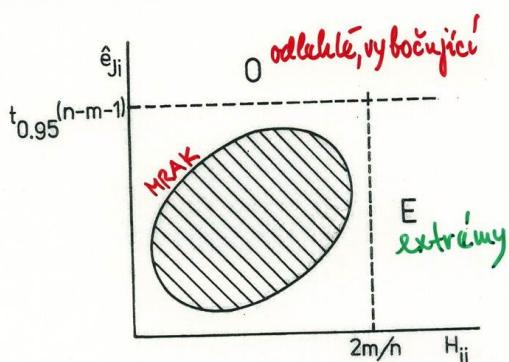
(a) tvar mraku, (b) tvar výseče, (c) tvar pásu a (d) nelineární tvar

Grafy identifikace vlivných bodů

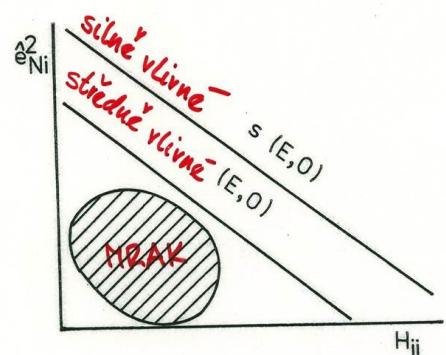
1. Graf predikovaných reziduí (GPR),
osa x: \hat{e}_{pi} , osa y: \hat{e}_i



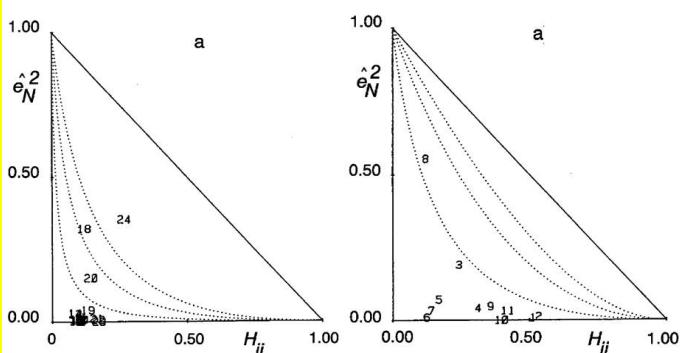
2. Williamsův graf (WG),
osa x: prvky H_{ii} , osa y: \hat{e}_{ji}



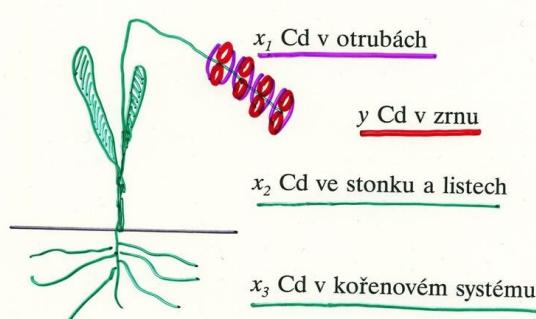
3. Pregibonův graf (PG),
osa x: prvky H_{ii} , osa y: \hat{e}_{Ni}^2



4. L-R graf



Úloha M619. Vliv tří parametrů na obsah kadmia v potravinářské pšenici
Obsah kadmia v zrnu y [mg/l] v závislosti na obsahu kadmia v otrubách x_1 [mg/l], ve stonku s listy x_2 [mg/l] a v kořenovém systému x_3 [mg/l]. Vyšetřete regresní triplet (data, model, metoda) a nalezněte lineární regresní model.



Vzorová úloha výstavby regresního modelu

50

Výstavba lineárního regresního modelu:

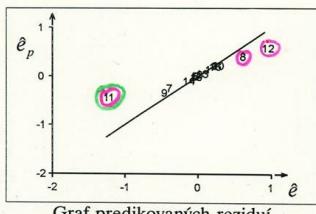
$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

Kritika dat

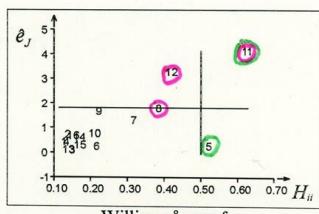
Kritika modelu

Kritika metody

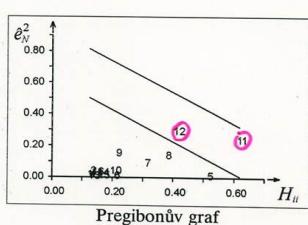
Grafy vlivných bodů (odlehlých O a extrémů E)



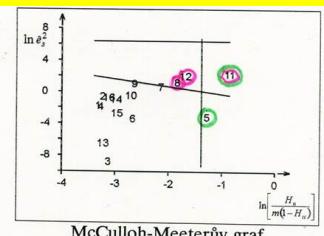
Graf predikovaných reziduí



Williamsův graf



Pregibonův graf



McCulloch-Meeterův graf

VÝSTUP

(1) PŘEDBEŽNÁ STATISTICKÁ ANALÝZA:

	Proměnná	Průměr	Směrodatná odchylka	Párový korelační koeficient	Spočtená hladina výz.
y	6.0125E+00	4.8734E+00	1.0000	----	
x1 OTRUBY	4.8937E+00	3.5692E+00	0.9837	0.000	
x2 LÍSTY	5.7812E+00	4.5296E+00	0.9935	0.000	
x3 KOLEN	5.0812E+00	3.8782E+00	0.9948	0.000	
Párové korelační koeficienty mezi dvojicemi vysvětlujících proměnných					
x1 versus x2 :	9.9344E-01		0.000		
x1 versus x3 :	9.8693E-01		0.000		
x2 versus x3 :	9.8847E-01		0.000		

(2) INDIKACE MULTIKOLINEARITY:

C	Vlastní čísla	Císla podmínky	Variance inflation factor VIF[jj]	Vícenás.korel. koef pro X[jj]
[jj]	korel. matice I[jj]	něnosti K[jj]		
1	6.4568E-03	4.6141E+02	8.3272E+01	0.9940
2	1.4307E-02	2.0823E+02	9.4324E+01	0.9947
3	2.9792E+00	1.0000E+00	4.7508E+01	0.9894
Maximální číslo podmíněnosti K	:	4.6141E+02	(VIF[jj], K > 1000 indikuje silnou multikolinearitu)	
(VIF[jj] > 10 indikuje silnou multikolinearitu)				

Kritika dat: původní data obsahují vlivné body 8, 11, 12

Kritika modelu: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$; technika MEP, AIC, s(e)) našla, že absolutní člen β_0 statisticky nevýznamný

Odhad parametru	Směrodatná odchylka	Student t-test	Parametr je významný	Hladina významnosti
		t _{0.95(16-4)} = 2.18		

A. Původní data a původní model vedou k závěru

$$y = -0.07(0.14) - 0.69(0.19)x_1 + 0.90(0.16)x_2 + 0.84(0.13)x_3$$

(LS: t_{1-0.05/2(14-4)} = 2.179, D = 99.72%, MEP = 0.21014, AIC = -36.18, s(e) = 0.290)

β_0	-0.073	0.138	-0.5269	Nevýznamný	0.608
β_1	-0.685	0.192	-3.5746	Významný	0.004
β_2	0.896	0.161	5.5761	Významný	0.000
β_3	0.838	0.133	6.2879	Významný	0.000

B. Data bez 8, 11, 12 a původní model vedou k závěru

$$y = -0.02(0.10) - 0.80(0.50)x_1 + 0.92(0.34)x_2 + 0.91(0.15)x_3$$

(LS: t_{1-0.05/2(13-4)} = 2.262, D = 99.83%, MEP = 0.06118, AIC = -41.58, s(e) = 0.178)

β_0	-0.017	0.103	-0.1637	Nevýznamný	0.874
β_1	-0.802	0.504	-1.5919	Nevýznamný	0.146
β_2	0.920	0.338	2.7193	Významný	0.024
β_3	0.906	0.152	5.9698	Významný	0.000

Kritika modelu regresní diagnostikou

C. Data bez 11, 12 a opravený model vedou k závěru

$$y = -1.18(0.38)x_1 + 1.25(0.24)x_2 + 0.92(0.15)x_3$$

(LS: t_{1-0.05/2(14-3)} = 2.201, D = 99.80%, MEP = 0.06078, AIC = -43.47, s(e) = 0.193)

β_0	0.000	---	---	---	---
β_1	-1.181	0.383	-3.0854	Významný	0.010
β_2	1.245	0.236	5.2751	Významný	0.000
β_3	0.917	0.151	6.091	Významný	0.000

D. Data bez 8, 11, 12 a opravený model vedou k závěru

$$y = -0.86(0.37)x_1 + 0.95(0.25)x_2 + 0.92(0.13)x_3$$

(LS: t_{1-0.05/2(13-3)} = 2.228, D = 99.83%, MEP = 0.05101, AIC = -43.55, s(e) = 0.170)

β_0	0.000	---	---	---	---
β_1	-0.855	0.372	-2.3002	Významný	0.044
β_2	0.954	0.251	3.8038	Významný	0.003
β_3	0.916	0.132	6.9263	Významný	0.000

Model: $y = -0.86(0.37)x_1 + 0.95(0.25)x_2 + 0.92(0.13)x_3$

REGRESNÍ DIAGNOSTIKA – kvalita modelu

1) Graf reziduí

2) Parciální regresní grafy

Vyjadřuje závislost **mezi vysvětlovanou proměnnou** (vektorem y)

a jednou **vysvětlující proměnnou x_j** při statisticky neměnném vlivu ostatních vysvětlujících proměnných, které tvoří matici $X_{(j)}$ (kdy je vynechaná j -tá proměnná).

Jde zde o určitou grafickou obdobu parciálního korelačního koeficientu u korelačních modelů.

HODNOCENÍ REGRESNÍHO MODELU

1. Kvalita nalezených odhadů parameterů

a) Podle intervalů spolehlivosti (čím menší interval spolehlivosti, tím lépe)

$$\beta_j = b_j \pm \sqrt{C_{mm} \cdot m \cdot s^2 \cdot F_{1-\alpha; m; n-m}}$$

b) Podle rozptylu parameterů, kde by pro statisticky významný odhad mělo platit pravidlo Sillénovy podmínky

$$2 \cdot \sqrt{D(b_j)} < |b_j|$$

HODNOCENÍ REGRESNÍHO MODELU

2. Kvalita dosažené těsnosti proložení

- a) podle reziduálního rozptylu.
- b) podle regresního rabatu, což je v procentech vyjádřený koeficient determinace (čím více se blíží 100 %, tím lepší je proložení modelem).

3. Vhodnost navrženého modelu

- a) Akaikovo informační kritérium (AIC): čím je AIC menší, tím je model vhodnější.
- b) Střední kvadratická chyba predikce (MEP): čím je MEP menší, tím je predikční schopnost modelu lepší.

HODNOCENÍ REGRESNÍHO MODELU

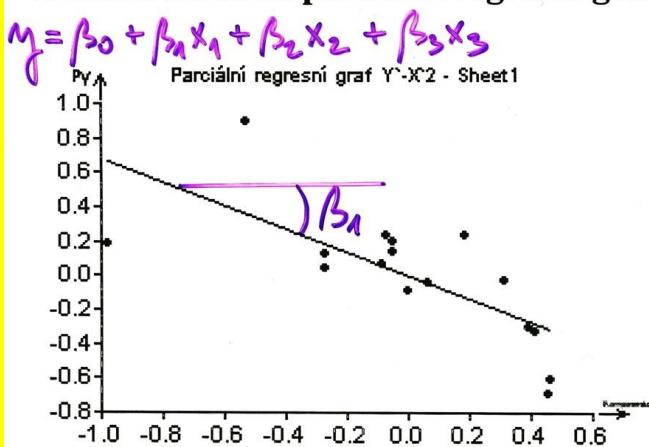
4. Predikční schopnost modelu

- a) na základě analýzy reziduí

5. Kvalita experimentálních dat

- a) na základě analýzy vlivných bodů (podle Jackknife reziduí, Cookovy vzdálenosti, diagonální prvky projekční matice a věrohodnostní vzdálenost).

Kritika modelu: parciální regresní graf



Kritika metody regresní diagnostikou

64

Kritika metody

Ověření předpokladů metody nejmenších čtverců

Metoda:

Vyšetření heteroskedasticity

Vyšetření autokorelace

Vyšetření multikolinearity

Vyšetření normality náhodných chyb

Vyšetření omezení parametrů

Vyšetření trendů reziduí

REGRESNÍ DIAGNOSTIKA

testy vybraných předpokladů klasické MNČ

Multikolinearita: VIF diagnostika indikuje

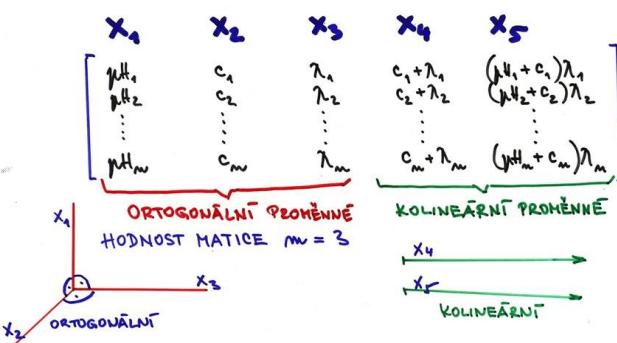
Heteroskedasticita: testy heteroskedasticity (např. Cook-Weisberg)

Autokorelace reziduí: test významnosti autokorelačního koeficientu

Normalita reziduí: testy normality

Test multikolinearity

Paradoxní situace: F-test je významný a všechny t-testy jsou nevýznamné, protože je silná multikolinearita mezi sloupcí matici X , čili existuje rovnoběžnost vektorů x_j a x_k , $j \neq k$, sloupců matici X .



Statistické obtíže:

1. Nestabilita odhadů je způsobená citlivostí odhadů na malé změny v datech. Odhady mívají často nesprávné znaménko, což znemožňuje jejich věcnou (fyzikální) interpretaci a jsou co do absolutních hodnot příliš veliké.
2. Velké rozptyly $D(b_i)$ jednotlivých odhadů způsobují, že t-testy indikují statistickou nevýznamnost β_i .
3. Silná korelovanost mezi prvky vektoru odhadů b způsobuje, že odhady b_j nelze interpretovat odděleně.
4. Koeficient determinace vysoký a regresní model může dobře popisovat experimentální data.

2. Numerická kritéria:

a) Determinant matice \mathbf{R} , $\det(\mathbf{R}) = \prod_{j=1}^m \lambda_j$, kde λ_j jsou vlastní čísla matice \mathbf{R} . Je-li determinant $\det(\mathbf{R})$ příliš malý, tj. menší než 10^{-3} , jde o silnou multikolinearitu.

b) Číslo podmíněnosti K = $\frac{\lambda_{\max}}{\lambda_{\min}}$, kde $\lambda_{\max}, \lambda_{\min}$ jsou maximální a minimální vlastní číslo matice \mathbf{R} . Je-li číslo podmíněnosti $K > 10^3$, jde o silnou multikolinearitu.

c) VIF-faktor (Variance Inflation Factor) je $VIF_j = \tilde{R}_{jj}$, kde

\tilde{R}_{jj} je j-tý diagonální prvek matice \mathbf{R}^{-1} . Platí vztah

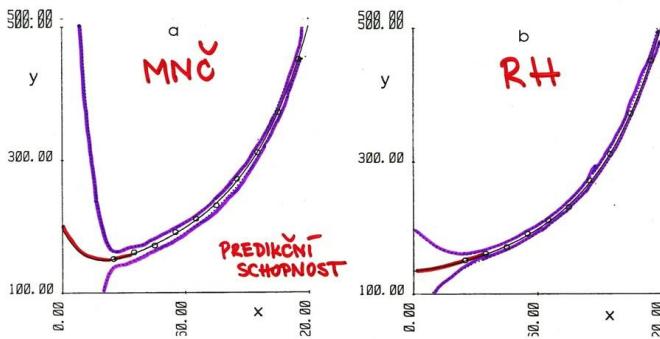
$$VIF_j = \frac{1}{1 - \tilde{R}_{jj}^2}. \text{ Je-li } VIF_j > 10, \text{ jde o silnou multikolinearitu.}$$

MULTIKOLINEARITA - testování

VIF – variance inflation factor – diagonální prvky inverzní matici ke korelační matici nezávisle proměnných (**diag(R⁻¹)**)

The screenshot shows a Microsoft Excel spreadsheet with several tables:

- korrelační matice R** : A 6x6 matrix labeled X1 through X5. The diagonal elements are highlighted in green (2.25, 2.15, 1.51, 0.95, 0.98). The off-diagonal elements are highlighted in orange (e.g., 1.28, -0.88, -10.2, 9.44).
- =INVERZE(B2..F6)**: Formula used to calculate the inverse matrix.
- inverzní matice R^{-1}** : The inverse matrix where the diagonal elements are highlighted in green (2.25, 2.15, 1.51, 0.95, 0.98) and the off-diagonal elements are highlighted in orange (e.g., -1.28, -0.88, -10.2, 9.44).
- kriticky vysoké hodnoty VIF**: A column of red numbers representing the VIF values: 2.25, 2.15, 1.51, 0.95, 0.98.

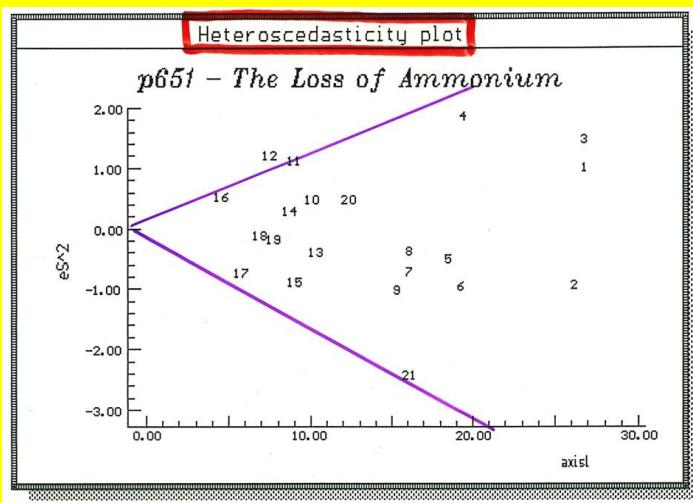
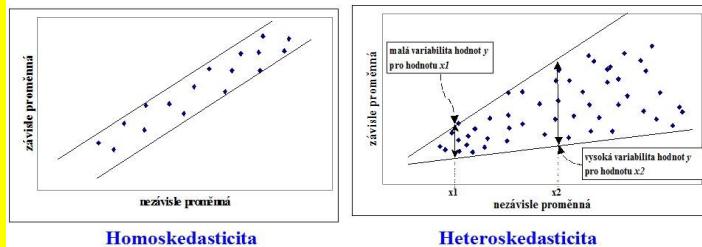


Závěr:

1. RH poskytne odhady parametrů, které zajišťují průběh modelu odpovídající trendům dat a nemá nadbytečné extrémy či inflexy.
2. Při použití klasické MNČ by úloha byla řešitelná pouze při zavedení omezení na regresní parametry.

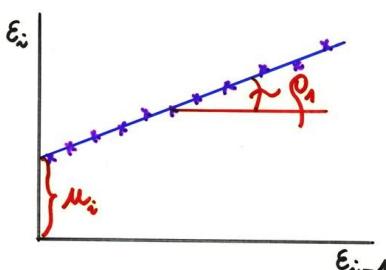
HOMOSKEDASTICITA vs. ETEROSKEDASTICITA Grafické vysvětlení principu

Homoskedasticita znamená, že hodnoty závisle proměnné y mají pro všechny hodnoty nezávisle proměnné x **konstantní rozptyl** (variabilitu, proměnlivost).



Test:

1. Grafická indikace autokorelace:



2. Waldův test pro ρ_1 : $H_0: \rho_1 = 0$ vs. $H_A: \rho_1 \neq 0$

Je-li Waldovo kritérium $W_a = \frac{n \hat{\rho}_1^2}{1 - \hat{\rho}_1^2} < \chi^2(1)$, H_0 je přijata.

2. Autokorelace

Data časových řad mají chyby ϵ_i vzájemně korelované.

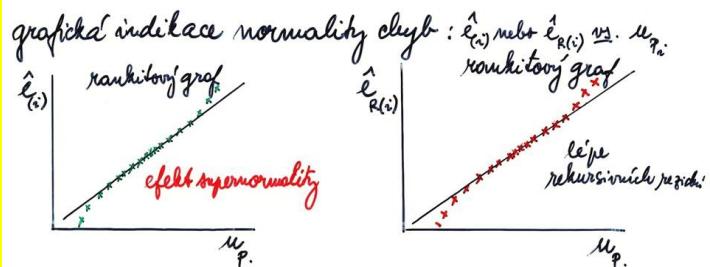
Nejčastější je případ autokorelace prvního řádu

$$\epsilon_i = \rho_1 \epsilon_{i-1} + u_i$$

kde $u_i \sim N(0, \sigma^2)$.

- a) Pro $\rho_1 = 1$ případ kumulativních chyb, který se v chemii vyskytuje často.
- b) Pro $\rho_1 \leq 1$ jde o autokorelační koeficient 1. řádu.

NORMALITA CHYB



Výstavba nelineárního regresního modelu

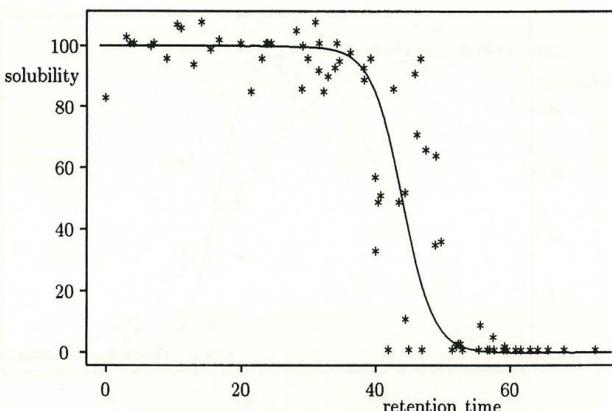


FIGURE 3.7. Peptides example: observed and adjusted solubilities of 75 peptides versus their RP-HPLC retention time

Model: rozšířený Debye-Hückelův zákon

$$pK_{a,\text{smiš}} = pK_{a,T} - \frac{0.5115 \sqrt{I}}{1 + 3.29 \times 10^{10} \text{ Å} \sqrt{I}} + CI$$

Proměnné:

Závisle proměnná y : $pK_{a,\text{smiš}}$, Nezávisle proměnná x : I ,

Neznámé parametry: $\beta_1: pK_{a,T}$, $\beta_2: \text{Å}$, $\beta_3: C$

NELINEÁRNÍ REGRESNÍ MODELY

Základní úlohy:

1. Konstrukce *kalibračních modelů*,
2. Ověření *teoretických modelů* fyzikálně-chemické zákonitosti,
3. Tvorba *empirických modelů*,

Tvorba regresního modelu $f(x, \beta)$ čili funkce

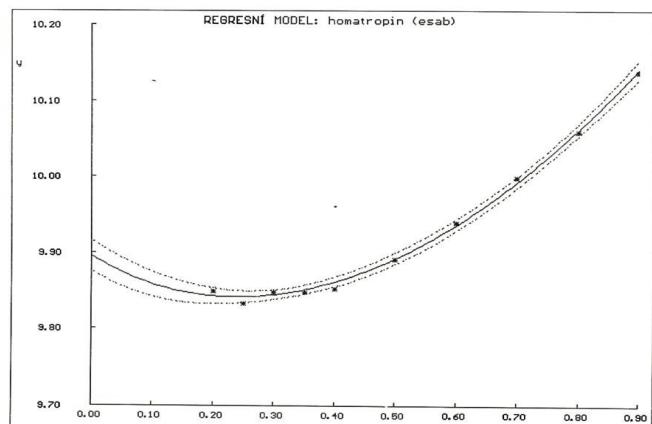
- a) vektoru nastavovaných proměnných (deterministických, kontrolovaných, vysvětlujících, nezávislých) x , tj. bodů $\{\mathbf{x}_i^T, y_i\}$, $i=1, \dots, n$,
- b) vektoru parametrů β o rozměru $(m \times 1)$, $\beta = (\beta_1, \dots, \beta_m)^T$.
- c) y je vysvětlovaná proměnná (závisle p., odezva, měření, pozorování) na zvolenou kombinaci nastavovaných veličin \mathbf{x}_i .

Tvorba se formuluje s ohledem na **regresní triplet**:

1. zadaná data,
2. navržený model,
3. kritérium regrese.

Dělení regresních modelů:

1. **Lineární modely:** parametry nemají fyzikální smysl (koeficienty),
2. **Nelineární modely:** parametry mají přesný fyzikální význam.
(nepř. rovnovážné konstanty (disociační, stability, součiny rozpustnosti) reakčních produktů, rychlostní konstanty u kinetických modelů, neznámé koncentrace, atd.)



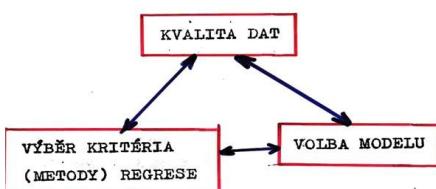
Závislost smíšené disoc. konstanty $pK_{a,\text{smiš}}$ homatropinu na iontové sile:

Nulté přiblížení: $pK_a^T = 1$, $\text{Å} = 1 \text{ Å}$, $C = 1$

Nalezeno: $pK_a^T = 9.90(1)$, $\text{Å} = 6(2) \text{ Å}$ a $C = 0.51(3)$

POSTUP A DIAGNOSTIKA REGRESE:

Regresní triplet:



$$U = \sum_{i=1}^n w_i (y_{\text{exp},i} - y_{\text{vyp},i})^2 = \text{minimum}$$

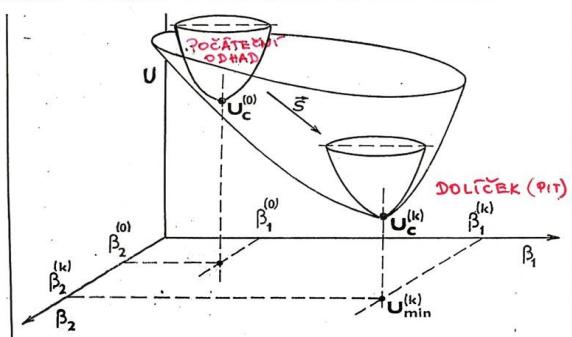
$$\text{kde } y_{\text{vyp},i} = f(\mathbf{x}_i; b_1, \dots, b_m)$$

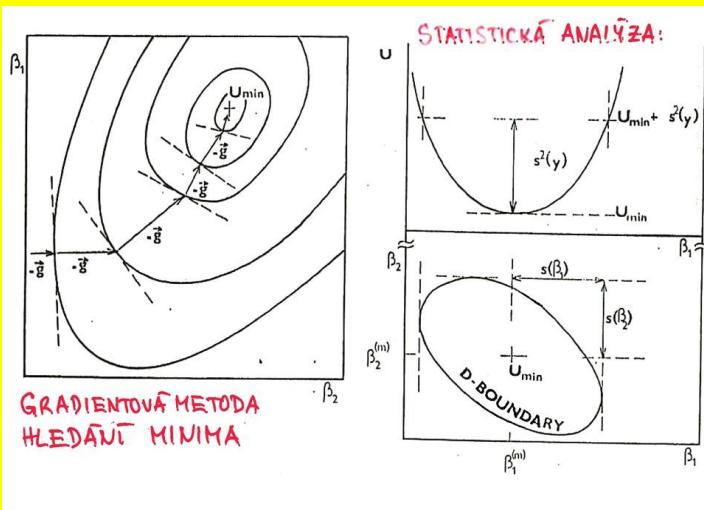
GEOMETRICKÉ ZNÁZORNĚNÍ KRITERIA U

$$U = \sum_{i=1}^m w_i (y_{\text{exp},i} - y_{\text{vyp},i})^2 \approx \text{minimum}$$

$$y_{\text{vyp},i} = f(\mathbf{x}_i; \beta_1, \dots, \beta_m; b_1, \dots, b_m)$$

OPTIMALIZAČNÍ PROBLÉM V $(m+1)$ ROZMĚRNÉM PROSTORU



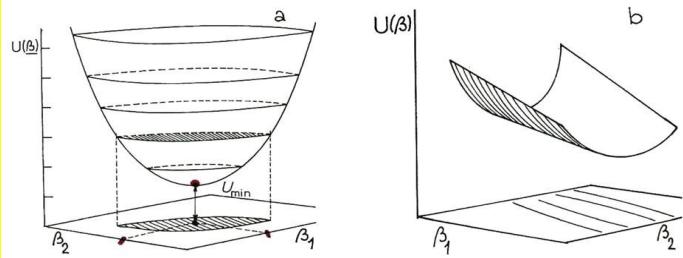


Geometrie nelineární regrese

Kritérium regrese $U(\beta)$ vektorovým zápisem

$$U(\beta) = \|\mathbf{y} - \mathbf{f}\|^2$$

kde $\mathbf{y} = (y_1, \dots, y_n)^T$, $\mathbf{f} = (f(x_1, \beta), \dots, f(x_n, \beta))^T$, a symbol $\|\mathbf{x}\| = \sqrt{\mathbf{x}^T \mathbf{x}}$ označuje euklidovskou normu.



Minimum účelové funkce $U(\beta)$ pro (a) lineární modely, (b) nelineární modely

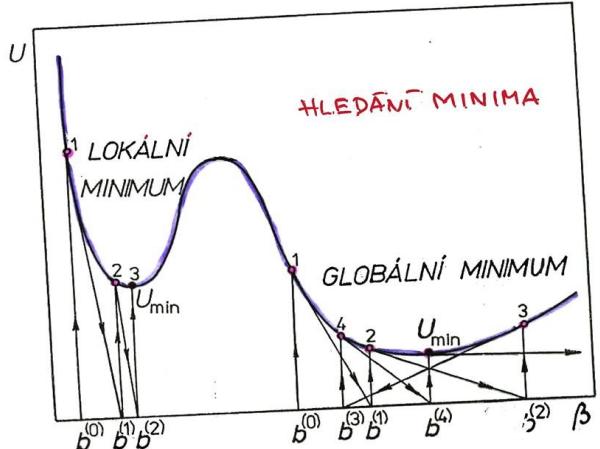
Pro lineární regresní modely:

- Útvary kritéria $U(\beta)$ v prostoru odhadů je elliptický hyperparaboloid se středem v bodě $[\mathbf{b}, U(\mathbf{b})]$, kde nabývá minimální hodnoty.
- Účelová funkce $U(\beta)$ je vzhledem k β kvadratickou formou $\beta^T (\mathbf{X}^T \mathbf{X}) \beta$ a matice $\mathbf{X}^T \mathbf{X}$ je pozitivně definitní.

Pro nelineární regresní modely:

- Složitější útvary vznikají u nelineárních modelů v závislosti na nelinearity funkce $f(x, \beta)$, počtu extrémů a sedlových bodů.

MINIMALIZAČNÍ PROCES:



Numerické postupy

Nelineární minimalizace: model $f(x, \beta)$ je nelineární vzhledem k parametru β_r ,

Nelineární maximalizace: použití maximální věrohodnosti

Extremalizace: užití libovolného kritéria regrese, kde "proměnné" jsou regresní parametry β

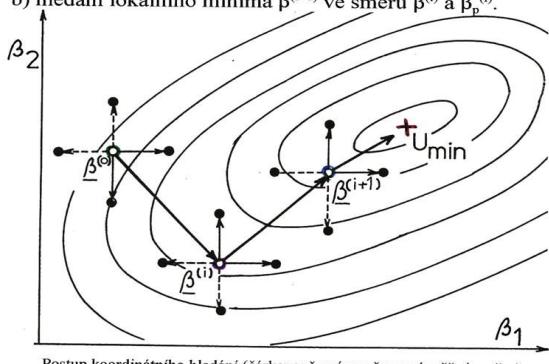
Optimalizační metody:

- hledání volného extrému, pokud nejsou na regresní parametry kladena žádná omezení,
- hledání vázaného extrému, jestliže regresní parametry musí splňovat jisté omezující podmínky.

1. Metody přímého hledání

Hookův-Jeevův algoritmus:

- krokové posuny a nalezení zlepšeného odhadu $\beta_p^{(i)}$, pro který je $G(\beta_p^{(i)}) < G(\beta^{(i)})$,
- hledání lokálního minima $\beta^{(i+1)}$ ve směru $\beta^{(i)}$ a $\beta_p^{(i)}$.



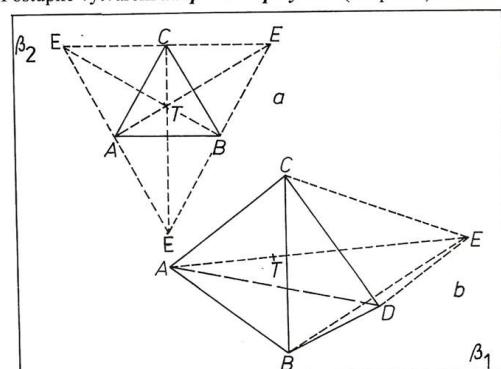
Nederivační metody

Umožňují nalezení extrému $G(\beta)$ vzhledem k β :

- metody přímého hledání,
- simplexové metody,
- metody využívající náhodných čísel,
- postupy speciálně vhodné pro MNČ.

2. Simplexové metody

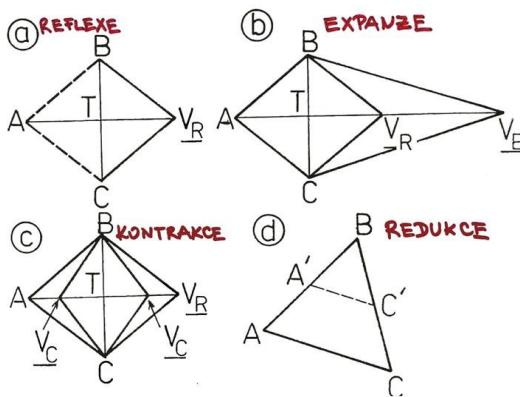
Postupné vytváření adaptivních polyedrů (simplexů):



Simplex pro (a) $m = 2$, a (b) $m = 3$ parametry.
Simplex A, B, C lze převrátit do tří poloh CBE, ABE, ACE

2. Iterativní postup k minimu

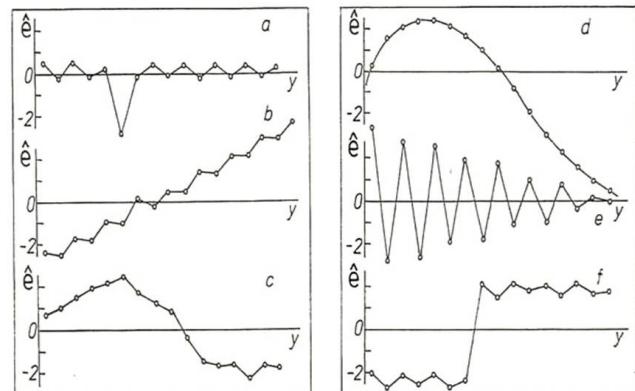
Na spojnici mezi V_H a jeho zrcadlovým obrazem pět operací: reflexe, expanze, kontrakce, redukce a přenesení.



1. Grafická analýza reziduů

- odlehlé (extrémní) hodnoty v souboru reziduů,
- trend v reziduích,
- nedostatečné střídání znaménka u reziduů,
- chybný model nebo vzájemnou závislost reziduů,
- heteroskedasticita (nekonstantnost rozptylu) závisle proměnné (měřené) veličiny y ,
- náhlou změnu podmínek při měření hodnot y .

Těsnost proložení v nelineární regresi



a) odlehlá hodnota v datech;
b) trend v reziduích;
c) nedostatečné střídání znaménka reziduů;
d) chybný model;
e) heteroskedasticita;
f) náhlá změna podmínek

2. Numerická analýza reziduů

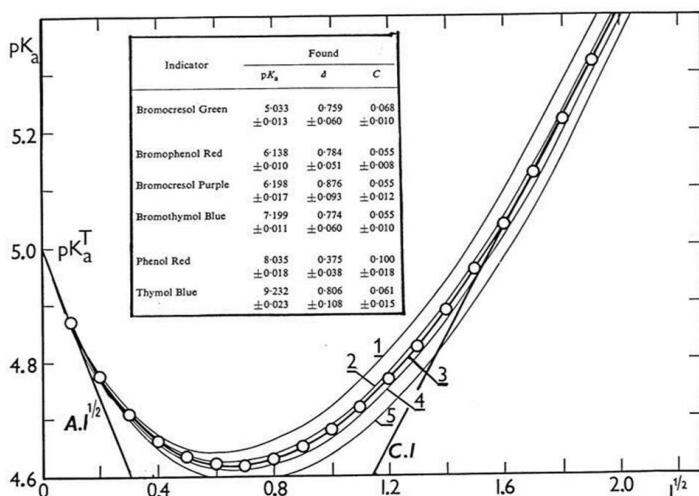
- Střední hodnota reziduů**, $E(\hat{e})$, by se měla rovnat nule,
- Průměrné reziduum** $|\bar{e}| = \frac{1}{n} \sum_{i=1}^n |\hat{e}_i|$ by se mělo rovnat náhodné chybě.
- Směrodatná odchylka střední hodnoty reziduů** $s(\hat{e})$ by se měla rovnat náhodné chybě.
- Koefficient šikmosti** $g_1(\hat{e})$ se pro Gaussovo rozdělení rovná nule.
- Koefficient špičatosti** $g_2(\hat{e})$ se pro Gaussovo rozdělení rovná třem.

DHLET

Model: rozšířený Debye-Hückelův zákon

$$pK_i = pK_0 - \frac{Az_i^2 \sqrt{I_i}}{1 + B\ddot{a}\sqrt{I_i}} + CI_i$$

Parametry: pK_0 , \ddot{a} , C



Úloha 8.1 Odhad parametrů Debyova-Hückelova vztahu

Stanovte termodynamickou disociaci konsantu pK_a^T (parametr β_1), efektivní průměr iontů \ddot{a} (parametr β_2) a vysolovací konsantu C (parametr β_3) závislosti smíšené disociaci konsanty y na iontové síle x podle rozšířeného Debyova-Hückelova vztahu pro vybrané sulfonftaleinu. Mají-li oba ionty L^{Z-1} a HL^Z zhruhu stejnou velikost \ddot{a} [10^{-10} mol] a je-li celkový vysolovací koeficient $C = C_{HL}^{Z-1} - C_{L}^{Z-1}$, lze formulovat Debyev-

$$\text{Hückelův vztah tvarem } y = \beta_1 - \frac{(1 - 2Z) A \sqrt{x}}{(1 + B \beta_2 \sqrt{x})} + \beta_3 x, \text{ kde } A = 0.5112$$

$\text{mol}^{1/2}, \text{l}^{1/2}, \text{K}^{3/2}, B = 0.3291 \text{ mol}^{1/2}, \text{m}^4, \text{l}^{1/2}, \text{K}^{1/2}$ jsou pro 25°C . Předpokládejte aditivní model měření a normalitu chyb závisle proměnné y .

Data: Bromkrezolová zeleň: $Z = -1, \{x, y\}$.

0.0104.901	0.0224.871	0.0404.834	0.0604.808	0.1164.765	0.232 4.709
0.3924.691	0.5944.677	0.9234.664	1.3304.662	2.0504.686	3.720 4.785

Řešení:

1. Návrh modelu:

označme
parametr pK_a^T jako první parametr β_1 s jeho odhadem b_1 ,
parametr \bar{a} jako druhý parametr β_2 s jeho odhadem b_2 a
parametr C třetí parametr β_3 se odhadem b_3 .

Pro nulté přiblížení odhadovaných parametrů je voleno
 $b_1^{(0)} = 1.0, b_2^{(0)} = 1.0, b_3^{(0)} = 1.0$.

2. Odhadování parametrů:

Bodové odhady parametrů:

Parametr	Bodový odhad b_j	Směrodatná odchylka $s(b_j)$	Absolutní vychýlení h_j	Relativní vychýlení $h_{R,j} [\%]$
β_1	5.0336E+00	4.2468E-03	-1.9583E-05	-3.8905E-04
β_2	7.6010E+00	2.1432E-01	4.2364E-03	5.5735E-02
β_3	6.8337E-02	3.4380E-03	1.7995E-05	2.6332E-02

Intervalové odhady parametrů:

Parametr	Bodový odhad b_j	Poloviční délka intervalu spohlednosti spočtená z délky poloos maxim
β_1	5.0336E+00	+/-1.1915E-02
β_2	7.6010E+00	+/-7.2957E-01
β_3	6.8337E-02	+/-9.9546E-03

5. Základní statistické charakteristiky:

Bylo dosaženo výtečného proložení, regresní rabat 99.49% ukazuje, že vysoké procento bodů vyhovuje navrženému modelu Debyeova-Hückelovy závislosti.

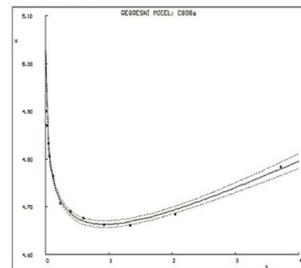
Regresní rabat $D [\%]$: 99.487
Akaikeho informační kriterium AIC	: -117.59
Střední kvadratická chyba predikce MEP	: 1.5516E-04
První statistický moment residu m_1	: -1.9709E-07
Třetí statistický moment residu m_3	: 1.1306E-07
Čtvrtý statistický moment residu m_4	: 2.0362E-09
Parametr Γ^T	: 8.4501E-02
Parametr Γ^N	: 8.6108E+02

B. Tabulka vlivných bodů:

Obsahuje základní charakteristiky k odhalení vlivných bodů.

Bod	Jackknife reziduum	Cookova vzdálenost	Diagonální prvky	Normalizovaná vzdálenost	Věrohodnostní vzdálenost
<i>i</i>	$e_{j,i}$	D_i	H_{ii}	FDA	LDA
1	-2.2875E+00	5.2536E-01	3.0692E-01	1.2144E-01	2.6567E-01
2	3.1245E-01	1.1668E-02	2.4392E-01	3.8138E-03	2.8886E-03
3	3.5252E-01	1.0641E-02	1.8824E-01	4.6230E-03	7.9056E-03
4	5.2560E-01	1.8108E-02	1.5314E-01	9.3211E-03	8.6827E-03
5	8.9337E-01	3.6124E-02	1.1718E-01	2.1943E-02	1.4799E-02
6	-8.1923E-01	3.4006E-02	1.2774E-01	1.962E-02	1.3011E-02
7	8.0345E-01	4.4839E-02	1.6679E-01	2.1813E-02	1.4450E-02
8	1.1447E+00	1.0565E-01	2.0015E-01	4.4082E-02	3.4816E-02
9	8.3151E-02	7.1250E-04	2.1571E-01	2.7009E-04	7.3859E-03
10	-1.1908E+00	1.1815E-01	2.0734E-01	4.5560E-02	3.7563E-02
11	-1.4634E+00	1.8341E-01	2.2452E-01	6.5850E-02	7.1201E-02
12	2.6179E+00	7.7430E+00	8.4834E-01	4.1027E-03	2.7646E-03

3. Graf regresní křivky:



Obr. 8-1a Rozptylový graf, ADSTAT.

4. Korelační matici parametrů:

Ukazuje, že korelace mezi parametry je výrazná.

	b_1	b_2	b_3
$x[1, i]$	1.00000	-0.82415	0.52400
$x[2, i]$	-0.82415	1.00000	-0.85058
$x[3, i]$	0.52400	-0.85058	1.00000

6. Regresní diagnostika:

Obsahuje pomůcky pro kritiku dat, kritiku modelu a kritiku metody.

A. Analýza klasických reziduí:

Směrodatná odchylka rezidui dosahuje hodnoty stejně velikosti, jako je odhad náhodných chyb (sumu) proměny y , tj. $\epsilon(pK_{a^T}) \approx 0.01$. Rozdělení rezidui je mírně asymetrické, sešímkem k nižším hodnotám, protože odhad šírkosti dosahuje záporné hodnoty. Rozdělení se blíží rovnoramennému, protože odhad šípkosti je blízký hodnotě 1.80. S ohledem na malý počet dat nelze z výběrové šírkosti a šípkosti usuzovat na nenormalitu.

Bod	Mířená hodnota	Predikovaná hodnota	Směrodatná odchylka	Vychýlení	Klasické reziduum
<i>i</i>	y_i	$\hat{y}_{pred,i}$	$s(\hat{y}_{pred,i})$	h_i	ϵ_i
1	4.9010	4.9115	3.7122E-03	-9.6193E-06	-1.0524E-02
2	4.8710	4.8691	3.3094E-03	-3.8063E-06	1.9192E-03
3	4.8340	4.8318	2.9072E-03	1.0885E-06	2.2400E-03
4	4.8080	4.8046	2.6222E-03	4.1321E-06	3.3797E-03
5	4.7650	4.7593	2.2938E-03	7.4433E-06	5.6887E-03
6	4.7090	4.7142	2.3949E-03	7.2323E-06	-5.2232E-03
7	4.6910	4.6860	2.7365E-03	4.1517E-06	5.0140E-03
8	4.6770	4.6703	2.9978E-03	3.8001E-07	6.7445E-03
9	4.6640	4.6635	3.1121E-03	-3.8612E-06	5.2314E-04
10	4.6620	4.6689	3.0511E-03	-6.3855E-06	-6.9445E-03
11	4.6860	4.6941	3.1750E-03	-6.2046E-06	-8.1348E-03
12	4.7850	4.7797	6.1717E-03	5.4489E-06	5.3175E-03

Reziduální součet čtvereců RSC : 4.0410E-04
Směrodatná odchylka rezidui $s(\epsilon)$: 5.8030E-03
Odhad šírkosti g_i : -0.579
Odhad šípkosti g_i : 1.796
Hamiltonov R-faktor [%] : 0.122

7. Mapa citlivostní funkce: Citlivostní funkce vyjadřuje změnu regresního modelu při změně parametru o $\pm 5\%$ ukazují, že parametry b_1 a b_3 jsou dobré podmíněny v modelu, jejich změna způsobí změnu účelové funkce 8 až 9 řádu. Parametr b_2 je ve srovnání s předchozími parametry méně citlivý, hůře podmíněný v modelu, změna je podstatně menší.

Parametr	Relativní změna	Souhrnná citlivost	Relativní změna
j	$C_{jR}(-5\%), [\%]$	$C_{jR} [\%]$	$C_{jR}(+5\%), [\%]$
1	8.0486E-09	1.0000	-7.2821E-09
2	1.5937E+01	1.0301E-03	-1.3244E+01
3	-3.3799E-07	1.7701	9.5080E-08

8. Predikční schopnost modelu:

Pro $n = 12$ bude $M_1 = 1$ až 6, $M_2 = 7$ až 12, $RSC(M_1) = 2.8534E-05$, $RSC(M_2) = 5.3549E-05$, $U(b) = 40.410E-05$, a proto $K = 4.923$.
Jelikož se K neblíží k 1, je predikční schopnost modelu slabší.

9. Souhlas s požadavky fyzikálního smyslu:

Termodynamickou disociační konstantu $pK_a^T = 5.034$ a vysolovací konstanta $C = 0.068$ mají fyzikální smysl, efektivní průměr iontu $\bar{a} = 7.6 \times 10^{-10}$ m je v souladu s hodnotami Kielandových tabulek.

