

# Výstavba regresního modelu v analytické laboratoři

Prof. RNDr. Milan Meloun, DrSc.,

Katedra analytické chemie, Univerzita Pardubice, 532 10 Pardubice

a

Prof. Ing. Jiří Militký, CSc.,

Katedra textilních materiálů, Technická univerzita Liberec, 461 17 Liberec

**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. Vyšetřuje 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 Akaikeovo 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 (MZNC) při autokorelaci, (c) metody podmínkových nejmenších čtverců (MPNC) při omezení kladených na parametry, (d) metody racionálních hodnot u multikolinearity, (e) metody rozšířených nejmenších čtverců (MRNC) 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  $\hat{y}_P, \hat{e}, \hat{b}$  závisí na splnění jistých předpokladů. Pokud platí předpoklady I až IV, jsou odhady  $b$  parametrů  $\beta$  nejlepší, nestranné a lineární (NNLO). Navíc mají asymptoticky normální rozdělení. Pokud platí ještě předpoklad VII, mají odhady  $b$  normální rozdělení i pro konečné výběry.

I. *Regresní parametry  $\beta$  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  $X$  má hodnost rovnou právě  $m$ .* To znamená, že žádné její dva sloupce  $x_j, x_k$  nejsou kolineární, tj. rovnoběžné vektory. Tomu odpovídá i formulace, že matice  $X^T X$  je symetrická regulární matice, ke které existuje inverzní matice a jejíž determinant je větší než nula.

IV. *Náhodné chyby  $\varepsilon_i$  mají nulovou střední hodnotu  $E(\varepsilon_i) = 0$ .* To musí u korelačních modelů platit vždy. U regresních modelů se může stát, že  $E(\varepsilon_i) = K, i = 1, \dots, n$ , což znamená, že model neobsahuje absolutní člen. Po jeho zavedení bude  $E(\varepsilon'_i) = 0$ , kde  $\varepsilon'_i = y_i - \hat{y}_{P,i} - K$ .

V. Náhodné chyby  $\varepsilon_i$  mají konstantní a konečný rozptyl  $E(\varepsilon_i^2) = \sigma^2$ . Také podmíněný rozptyl  $D(y/x) = \sigma^2$  je konstantní a jde o homoskedastický případ.

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

VII. Chyby  $\varepsilon_i$  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 rozptylení s kvantily a řady postupů průzkumové analýzy jednorozměrných dat. Přes svoji jednoduchost umožňuje diagnostika identifikovat ještě před vlastní regresní analýzou

- a) *nevzhodnost 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 užitym 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  $x_i$  se omezíme na a) parciální regresní grafy, a b) parciální reziduální grafy.

*Parciální regresní grafy* byly Belseyem zařazeny mezi základní nástroje počítacové 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  $X_{(i)}$  a reziduů  $u$  regrese  $x_i$  na sloupcích matice  $X_{(i)}$ . Přitom matice  $X_{(i)}$  vznikne z matice  $X$  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  $\hat{R}_{yx(x)}$ .

*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  $b$  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  $\sigma_y^2$  (vysvětlovaná proměnná) a  $\sigma_x^2$  (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

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

a 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 matici  $R$  a její rozklad na vlastní čísla a vlastní vektory. Jsou vypočteny VIF k indikaci multikolinearity a tisknuta setříděná vlastní čísla. K určení inverzní matici  $R^{-1}$  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 matici  $R^{-1}$  jsou určovány odhady parametrů  $b_j$ , 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 Wiliamsovým*, *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 reziduů 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 reziduů prvního řádu  $\hat{\rho}_1$ . Tabulka vlivných bodů obsahuje veličiny  $H_{ii}$ ,  $H_{ii}^*$ ,  $D$ ,  $A$ ,  $DF$ ,  $LD_i(\mathbf{b})$ ,  $LD_i(\hat{\sigma}^2)$  a  $LD_i(\mathbf{b}, \sigma^2)$ . Hvězdičkou jsou označeny hodnoty silně vlivných bodů.

### 5. Konstrukce zpřesněného modelu: s využitím

- a) metody vážených nejmenších čtverců (MVNC) při nekonstantnosti rozptylů,
- b) metody zobecněných nejmenších čtverců (MZNC) při autokorelací,
- c) metody podmínkových nejmenších čtverců (MPNC) při omezeních na parametry,
- d) metody racionálních hodností RH u multikolinearity,
- e) metody rozšířených nejmenších čtverců (MRNC) pro případ, že všechny proměnné jsou zatížené náhodnými chybami,
- f) 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 doplňkových informací o modelované soustavě se provede posouzení kvality navrženého lineárního regresního modelu.

## 1. 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  $\beta_0$  a směrnice  $\beta_1$ . Studentův  $t$ -test ukázal, že úsek (absolutní člen)  $\beta_0$  je statisticky významný a směrnice  $\beta_1$  je statisticky významná.

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

Hlad. význam.

0.000

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, $R_p^2$	: 9.9804E-01
Střední kvadratická chyba predikce, $MEP$	: 6.8546E+02
Akaikeho informační kritérium, $AIC$	: 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<sub>p</sub> = 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 odlehlých 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)

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

## 2. Vzorová úloha: Validace analytické metody stanovení formaldehydu (V6.16)

Ukážeme postup validace nové analytické metody: obsah formaldehydu ve vzorcích fenolových vod je stanoven polarografickou metodou  $x$ . Laboratoř fenoplastů navrhla používat jednodušší metodu  $y$ , redox-titraci. Rozptyl obou metod je prakticky stejný. Rozhodněte, zda tato metoda bude poskytovat správné a reprodukovatelné výsledky. Jsou v datech odlehlé hodnoty? Jsou výsledky nové metody zatíženy chybou? Aplikujte Studentův  $t$ -test úseku  $b_0$ , (má být  $\beta_0 = 0$ ) a směrnice  $b_1$ , (má být  $\beta_1 = 1$ ).

Data: Obsah formaldehydu ve vodě [mg/l] polarograficky  $x$ , redox-titrací  $y$ :

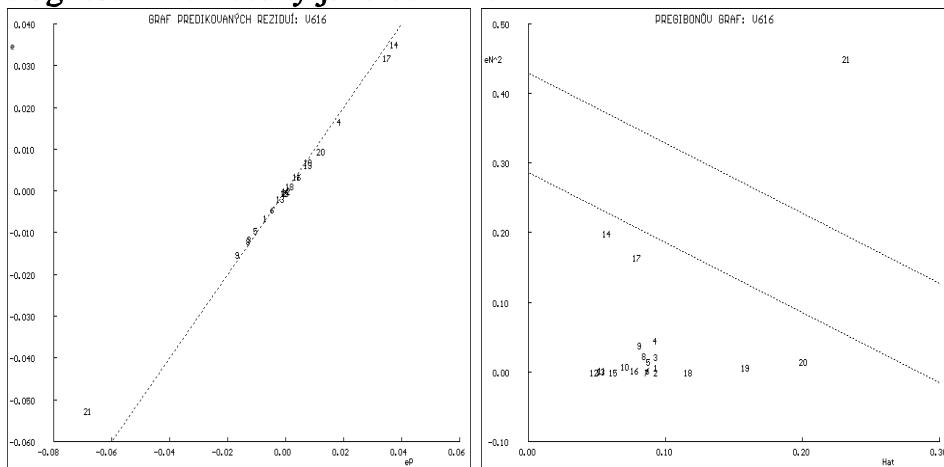
x:	76.8	117	129.1	160.1	236.4	258	284.2	303.2	386.2
	474.3	532.4	937.6	2654.3					
y:	80.5	112.6	128	152.2	239.4	250	287	307.8	391.7
	480.2	530.8	934.2	2647.2					

Řešení:

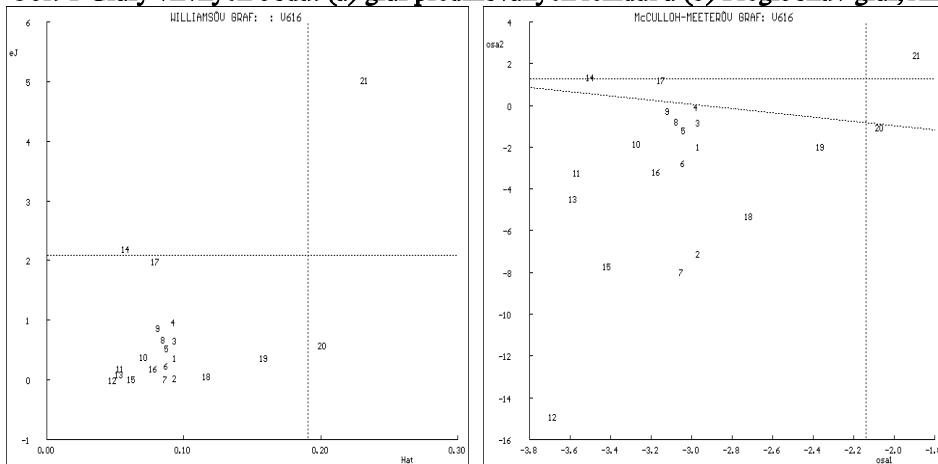
**1. Odhadování parametrů:** klasickou metodou nejmenších čtverců (MNČ) byly nalezeny odhady parametrů, úseku  $\beta_0$  a směrnice  $\beta_1$ . Studentův  $t$ -test ukázal, že úsek (absolutní člen)  $\beta_0$  je statisticky nevýznamný, zatímco směrnice  $\beta_1$  je statisticky významná.

Odhad	Směrodatná odchylka	Test H0: $B[j] = 0$ vs. HA: $B[j] \neq 0$	t-kriterium	hypoteza H0 je	Hlad. význam.
B[ 0]	1.0432E-02	5.6588E-03	1.8436E+00	Akceptována	0.081
B[ 1]	9.2170E-01	1.9331E-02	4.7681E+01	Zamítnuta	0.000

**2. Regresní diagnostika:** párový korelační koeficient  $R = 0.99585$  ukazuje, že navržený lineární regresní model je statisticky významný. Vysoká hodnota koeficientu determinace  $D = R^2 (99.17\%)$ , představující procento bodů, vyhovujících regresnímu modelu ukazuje, že všechny body výtečně korespondují s modelem přímky. Střední kvadratická chyba predikce  $MEP = 4.14E-04$  a Akaikeovo informační kritérium  $AIC = -166.78$  se užívají k rozlišení mezi několika navrženými modely. Za optimální se považuje model, pro který dosahuje  $MEP$  a  $AIC$  minimální hodnotu. Regresní diagnostika dále obsahuje pomůcky a postupy pro interaktivní analýzu (a) dat, (b) modelu, (c) metody, což jsou složky tzv. regresního tripletu. Věrohodnost nalezených odhadů parametrů  $\beta_0, \beta_1$  lze posoudit na základě grafu regresního modelu. Grafy vlivných bodů jsou schopny indikovat přítomnost odlehlých hodnot a extrémů. Graf predikovaných reziduí ukazuje na odlehlé body č. 21, 14, 17. Pregibonův graf ukazuje na silně vlivný bod č. 21. Williamsův graf indikuje č. 14 a 21 jako odlehlé body a extrémy č. 20, 21. McCulloh-Meeterův graf dokazuje odlehlé body č. 14, 17, 21, extrémy č. 20, 21. Konečně L-R graf dokazuje odlehlé body č. 14, 17, 21 a současně extrém č. 20. Lze uzavřít, že body č. 14, 21 jsou většinou diagnostik indikovány jako odlehlé.



Obr. 1 Grafy vlivných bodů: (a) graf predikovaných reziduí a (b) Pregibonův graf, ADSTAT



Obr. 2 Grafy vlivných bodů: (a) Williamsův graf a (b) McCulloh-Meeterův graf, ADSTAT 1.25

### 3. Konstrukce zpřesněného modelu:

(a) Po odstranění bodů č. 14, 17, 21 byly nalezeny odhady parametrů zpřesněného modelu.

Odhad	Směrodatná odchylka	Test H0: $B[j] = 0$ vs. HA: $B[j] \neq 0$	t-kriterium	hypoteza H0 je	Hlad. význam.
B[ 0 ]	6.3907E-03	2.4125E-03	2.6490E+00	Zamítnuta	0.018
B[ 1 ]	9.4034E-01	9.3557E-03	1.0051E+02	Zamítnuta	0.000

Zpřesněný model (v závorce je uveden odhad směrodatné odchylky parametru)

$$y = 0.00639 (0.00241) + 0.9403 (0.0094) x$$

je doložen statistickými charakteristikami: střední kvadratická chyba predikce  $MEP = 6.15E-05$  a Akaikeho informační kritérium  $AIC = -174.03$  dosáhly nižších hodnot, čímž dokazují kvalitnější model než předešlý. Rezidua nyní vykazují normální rozdělení a nevykazují trend, stále však vykazují heteroskedasticitu, a proto lze doporučit užití metody vážených nejmenších čtverců.

(b) Užitím statistické váhy ( $w_i = 1/y_i^2$ ) kompenzujeme heteroskedasticitu v datech. Obdržíme nové správnější odhady parametrů. Opravený model má tvar, (v závorce je vždy uveden odhad směrodatné odchylky parametru)  $y = 0.00213 (0.00197) + 0.9462 (0.0731) x$ . Jelikož došlo ke snížení rozhodujících kritérií, střední kvadratické chyby predikce  $MEP = 5.12E-05$  a Akaikeho informačního kritéria  $AIC = -181.63$  lze považovat tyto odhady za lepší než předešlé.

**4. Zhodnocení kvality modelu:** nalezený model má tvar, (v závorce je vždy uveden odhad směrodatné odchylky parametru)  $y = 0.00213 (0.00197) + 0.9462 (0.0731) x$  a intervalový odhad parametrů úseku  $\beta_0$  a směrnice  $\beta_1$  bude

$$b_0 - t_{1-\alpha/2}(18) \sqrt{D(b_0)} \leq \beta_0 \leq b_0 + t_{1-\alpha/2}(18) \sqrt{D(b_0)}$$

a po dosazení  $0.00213 - 2.12 \times 0.00197 \leq \beta_0 \leq 0.00213 + 2.12 \times 0.00197$  vyjde  $-0.00205 \leq \beta_0 \leq 0.00630$ . Tento interval spolehlivosti úseku regresní přímky zahrnuje nulu, takže lze úsek  $\beta_0$  považovat za nulový.

Analogicky dosazením do intervalu spolehlivosti směrnice obdržíme nerovnost

$$0.9462 - 2.12 \times 0.0731 \leq \beta_1 \leq 0.9462 + 2.12 \times 0.0731$$

a po výpočtu  $0.7912 \leq \beta_1 \leq 1.1012$ . Jelikož tento interval obsahuje jedničku, lze považovat směrnici  $\beta_1$  za jednotkovou.

Lze uzavřít, že úsek regresní přímky lze považovat za nulový  $\beta_2 = 0$  a směrnice  $\beta_1$  není významně odlišná od jedničky. Výsledky nové metody se proto statisticky významně neliší od metody standardní.

### **3. Vzorová úloha: Vliv parametrov na obsah kadmia v potravinárské pšenici (M6.19)**

Ukážeme postup analýzy vícerozměrného lineárního regresního modelu.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_m x_m,$$

kde  $\beta_0, \beta_1, \beta_2, \dots, \beta_m$  jsou odhadované parametry: u vzorků potravinářské pšenice byl zjištován obsah kadmia v zrnu  $y$  v závislosti na obsahu kadmia v otrubách  $x_1$ , ve stonku s listy  $x_2$  a v kořenovém systému  $x_3$ . Vyšetřením regresního tripletu nalezněte nejlepší vícerozměrný regresní model. Využijte k tomu regresní diagnostiku a pomocí parciálních regresních a parciálních reziduálních grafů diskutujte významnost jednotlivých parametrů v modelu stejně jako i jejich fyzikální smysl.

*Data:* Obsah v otrubách  $x_1$  [mg/l], ve stonku s listy  $x_2$  [mg/l] a v kořenovém systému  $x_3$  [mg/l], obsah kadmia v zrnu  $y$  [mg/l]:

**Řešení:**

**1. Návrh modelu:** začíná se vždy od nejjednoduššího modelu, u kterého vystupují  $x_1, x_2, x_3$  v prvních mocninách a nevyskytují se žádné interakční členy. Na začátku analýzy vždy zařadíme i absolutní člen  $\beta_0$ , takže pro daná data bude navržený regresní model tvaru

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

**2. Předběžná analýza dat:** polohu a proměnlivost proměnných  $y, x_1, x_2, x_3$  charakterizuje *průměr a směrodatná odchylka* hodnot každé proměnné. *Párový korelační koeficient*  $y$  vs.  $x_1$ ,  $y$  vs.  $x_2$ ,  $y$  vs.  $x_3$  ukazuje na vysokou korelaci, všechny tři nezávislé proměnné  $x_1, x_2, x_3$  jsou se závisle proměnnou  $y$  spjaty silnou lineární závislostí. *Párové korelační koeficienty mezi dvojicemi vysvětlujících proměnných* ukazují na silnou korelaci i mezi nezávisle proměnnými. Nejsilnější lineární vztah existuje mezi  $x_1$  vs.  $x_2$ , a u  $x_1$  vs.  $x_3$ , a  $x_2$  vs.  $x_3$  je rovněž silná korelace. Blok INDIKACE MULTIKOLINEARITY vykazuje ve všech kritériích multikolinearitu, protože je značná korelace mezi  $x_i$  a  $x_j$ .

Proměnná	Průměr	Směrodatná odchylka	Párový korelační koeficient	Spočtená hladina významnosti
y	6.0125E+00	4.8734E+00	1.0000	-----
x1	4.8937E+00	3.5692E+00	0.9837	0.000
x2	5.7813E+00	4.5296E+00	0.9935	0.000
x3	5.0813E+00	3.8782E+00	0.9948	0.000
<b>Párové korelační koeficienty mezi dvojicemi vysvětlujících proměnných</b>			<b>Spočtená hladina významnosti</b>	
x1 versus x2:	9.9344E-01		0.000	
x1 versus x3:	9.8693E-01		0.000	
x2 versus x3:	9.8847E-01		0.000	
<b>INDIKACE MULTIKOLINEARITY:</b>				
Č [j] Vlastní čísla korel. maticy I[j]	Čísla podmínenosti K[j]	Variance inflation faktor VIF[j]	Vícenás. korel. koef. pro X[j]	
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	
<b>Maximální číslo podmíněnosti K:</b>		4.6141E+02		

*Ná pověda: K[j], K > 1000 indikuje silnou multikolinearitu, VIF[j] > 10 indikuje silnou multikolinearitu.*

**3. Odhadování parametrů:** klasickou metodou nejmenších čtverců (MNČ) byly nalezeny nejlepší odhady čtyř parametrů  $\beta_0, \beta_1, \beta_2, \beta_3$ . Studentův  $t$ -test ukázal, že absolutní člen  $\beta_0$  je statisticky nevýznamný, zatímco ostatní parametry statisticky významné jsou. To je v souladu i s biologickou interpretací:  $\beta_0$  se týká zbytkového obsahu kadmia v zrnu, když je obsah kadmia v otrubách ( $x_1 = 0$ ), ve stonku ( $x_2 = 0$ ) a v kořenovém systému ( $x_3 = 0$ ) nulový. Je zřejmé, že když v celé rostlince bude obsah kadmia nulový, musí být nulový obsah i v zrnu a zbytkový obsah proto nemá ani biologický smysl. Absolutní člen  $\beta_0$  je proto nutno ve zpřesněném modelu vyněchat.

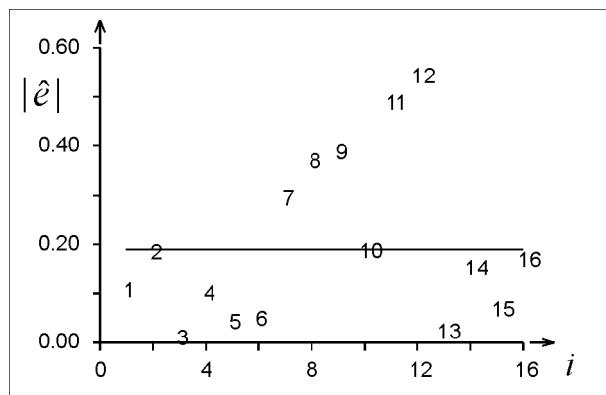
Odhad	Směrodatná odchylka	Test H0: B[j] = 0 vs. HA: B[j] $\neq$ 0		
		t-kriterium	hypoteza H0 je	Hlad. význam.
B[0] -7.2666E-02	1.3791E-01	-5.2692E-01	Akceptována	0.608
B[1] -6.8505E-01	1.9165E-01	-3.5746E+00	Zamítnuta	0.004
B[2] 8.9619E-01	1.6072E-01	5.5761E+00	Zamítnuta	0.000
B[3] 8.3769E-01	1.3322E-01	6.2879E+00	Zamítnuta	0.000

**4. Základní statistické charakteristiky:** vícenásobný korelační koeficient  $R = 0.99858$  ukazuje, že navržený lineární regresní model je statisticky významný. Vysoká hodnota koeficientu determinace  $D = R^2 = 99.716\%$  ukazuje, že všechny body výtečně korespondují s modelem. Predikovaný koeficient determinace  $R_p^2 = 99.527\%$  má podobný význam jako koeficient determinace, je však vyčíslen jinak, místo  $RSC$  se ve vztahu užije  $MEP$ . Střední kvadratická chyba predikce  $MEP = 0.2101$  a Akaikovo informační kritérium  $AIC = -36.18$  se užívají k rozlišení mezi několika navrženými modely. Za optimální se považuje model, pro který dosahuje  $MEP$  a  $AIC$  minimální hodnotu.

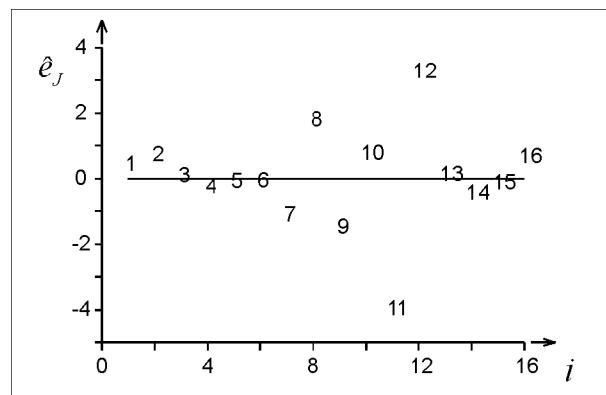
**5. Regresní diagnostika:** obsahuje pomůcky a postupy pro interaktivní analýzu (a) dat, (b) modelu, (c) metody, což jsou složky tzv. *regresního tripletu*.

**5.1 Data:** skládá se z analýzy několika druhů grafických diagnostik a tabulek různých druhů reziduí.

(a) *Analýza klasických reziduí* není příliš spolehlivá, protože klasická rezidua jsou korelovaná, s nekonstantním rozptylem, jeví se normálnější než náhodné chyby (*efekt supernormality*) a nemusí indikovat silně odlehlé hodnoty. Grafická analýza  $\hat{e}$  vs.  $\hat{y}_P$  je však schopna indikovat podezřelé body, trend, a nekonstantnost podmíněného rozptylu tj. heteroskedasticitu. Míry polohy a rozptýlení klasických reziduí by měly dosahovat hodnot blízkých experimentálnímu šumu. *Odhad směrodatné odchylky*  $s(e)$  by se měl blížit svou velikostí experimentální chybě, kterou je zatížena závisle proměnná. *Odhad šíkmosti* a *odhad špičatosti* by měly dokazovat normální rozdělení reziduí, normalitu.



Obr. 3 (a) Indexový graf absolutních hodnot klasických reziduí, **ADSTAT 3.0**



Obr. 3 (b) Indexový graf Jackknife reziduí, **ADSTAT 3.0**

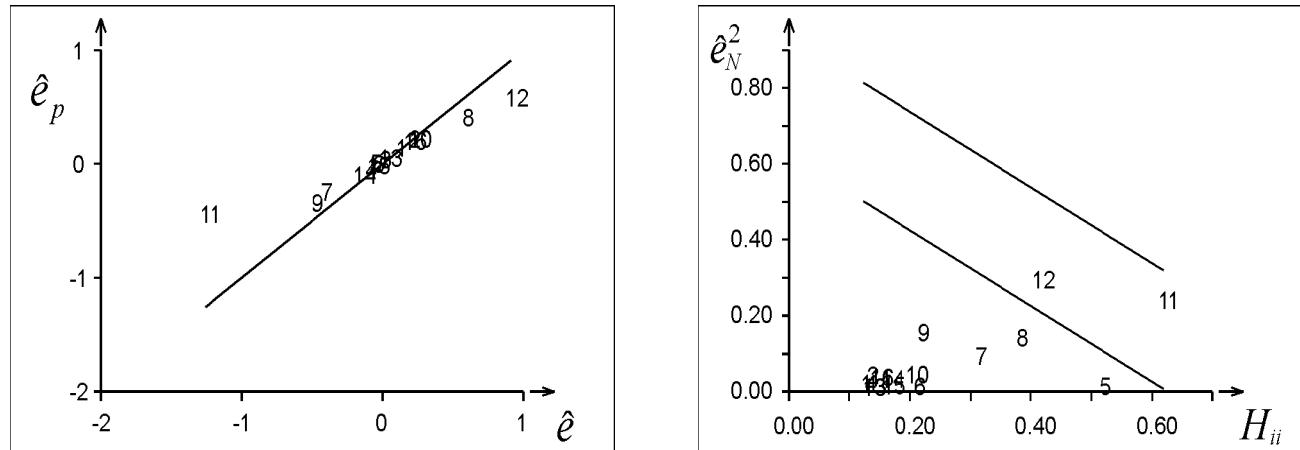
Rezidualní součet čtverců	: 1.0114E+00
Průměr absolutních hodnot reziduí	: 1.9048E-01
Průměr relativních reziduí	: 4.3750E+00
Odhad reziduálního rozptylu	: 8.4282E-02
Odhad směrodatné odchylky reziduí	: 2.9031E-01
Odhad šíkmosti reziduí	: 9.2354E-02
Odhad špičatosti reziduí	: 2.8755E+00

(b) *Analýza ostatních reziduí:* Jackknife rezidua indikují odlehlé body, diagonální prvky  $H_{ii}$  od projekční matice  $H$  a diagonální prvky  $H_{mii}$  od zobecněné projekční matice  $H_m$  pouze extrémy. Ostatní druhy reziduí a kritéria v tabulce pak obojí. Jackknife rezidua  $e_J[i]$  ukazují, že bod č. 11 a 12 je odlehlý, stejně tak i Cookova vzdálenost  $D[i]$ , Atkinsova vzdálenost  $A[i]$  ukazují na č. 8, 11, 12, kritérium  $DF[i]$  na č. 8, 11, 12 a věrohodnostní vzdálenosti  $LD(b)[i]$  a  $LD(s^2)[i]$  na č. 11 a  $LD(b, s^2)[i]$  na č. 11 a 12. Diagonální prvky  $H[i, i]$  projekční matice  $H$

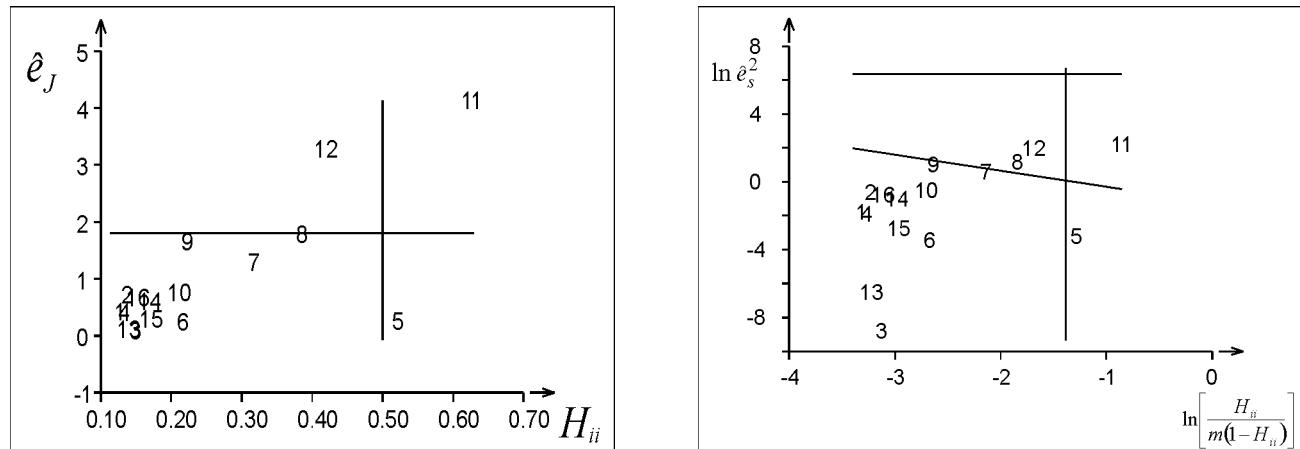
ukazují na extrémy č. 5 a 11, a diagonální prvky  $H_m[i, i]$  zobecněné projekční matice  $\mathbf{H}_m$  pak na extrémy č. 11 a 12.

(c) *Grafy vlivných bodů* jsou schopny indikovat a současně i testovat, dokazovat přítomnost odlehlých hodnot a extrémů. *Graf predikovaných reziduí* ukazuje na odlehlé body č. 8, 11, 12 a částečně na extrémy č. 11 a 12. *Pregibonův graf* ukazuje na středně vlivné body č. 11 a 12. *Williamsův graf* indikuje č. 11 a 12 jako odlehlé body a extrémy č. 11 a 12. *McCulloh-Meeterův graf* dokazuje odlehlé body č. 11 a 12, extrémy č. 8 a 11. Konečně *L-R graf* dokazuje odlehlé body č. 8, 11 a 12 a extrémy č. 11, 12 a 9.

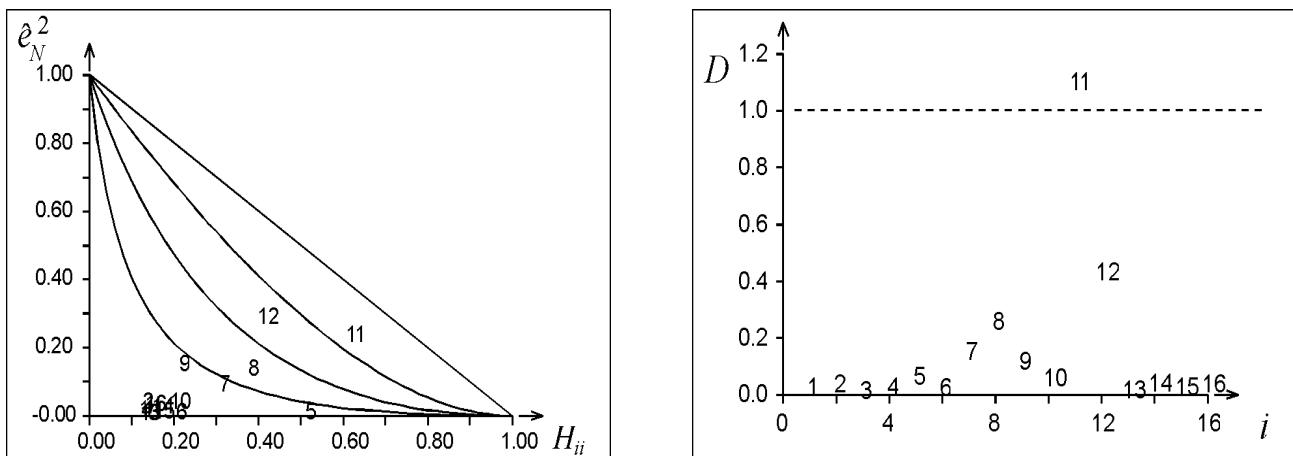
Lze uzavřít, že body č. 11 a 12 jsou většinou diagnostik prokázány za odlehlé, a proto je třeba je dále prověřit nebo z výběru vyloučit.



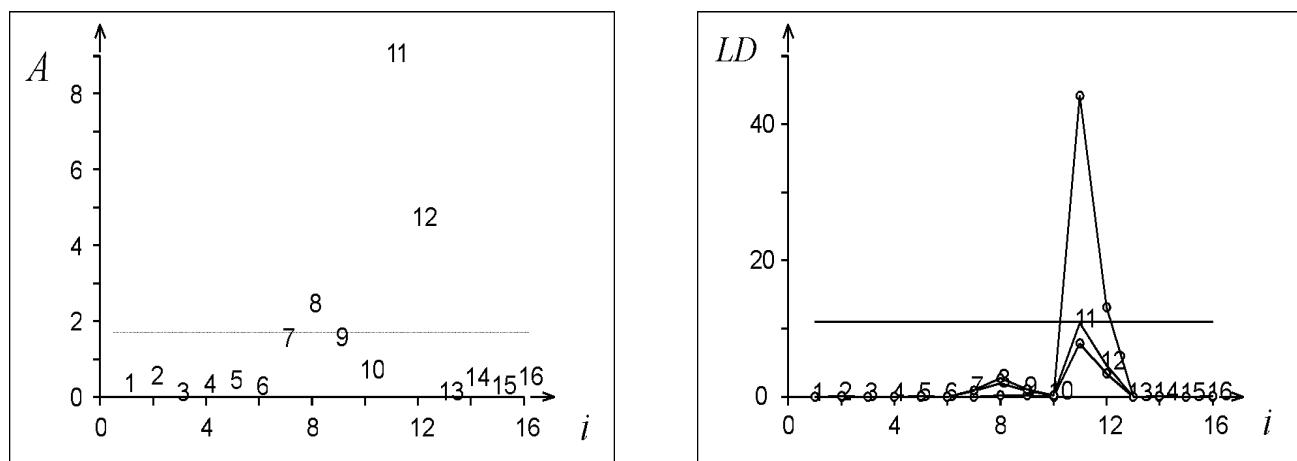
Obr. 4 Grafy vlivných bodů: (a) Graf predikovaných reziduí a (b) Pregibonův graf, **ADSTAT 3.0**



Obr. 5 Grafy vlivných bodů: (a) Williamsův graf a (b) McCulloh-Meeterův graf, **ADSTAT 3.0**

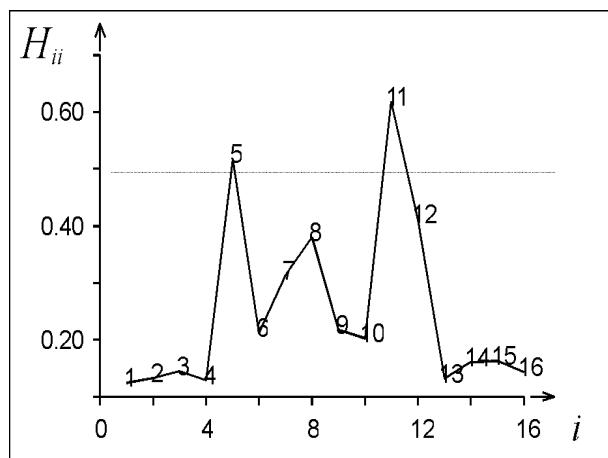


Obr. 6 Grafy vlivných bodů: (a) L-R graf, (b) Indexový graf Cookovy vzdálenosti, **ADSTAT 3.0**



Obr. 7 Grafy vlivných bodů: (a) Indexový graf Atkinsonovy vzdálenosti, (b) Indexový graf věrohodnostní vzdálenosti, **ADSTAT 3.0**

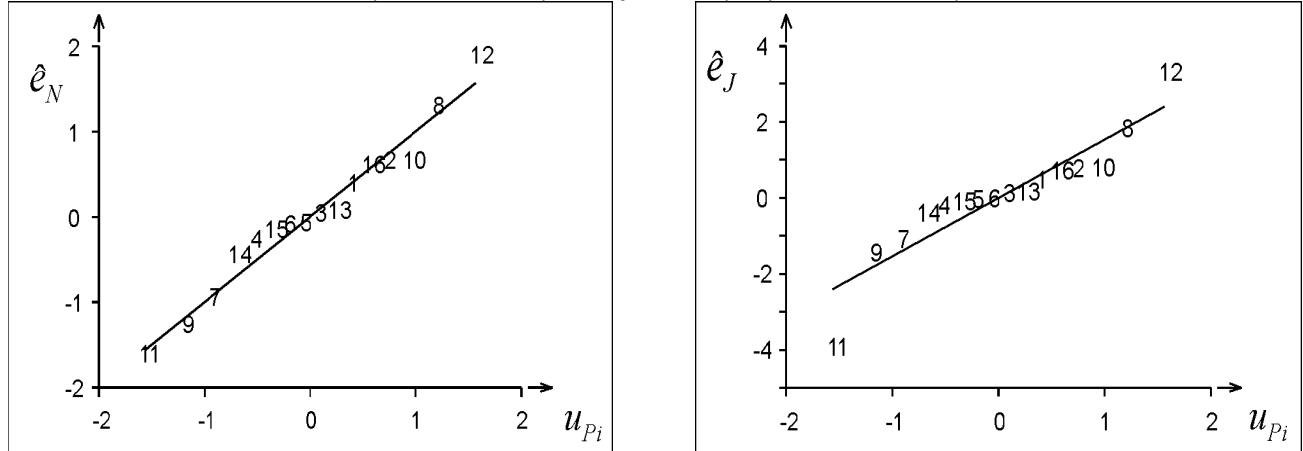
(d) **Indexové grafy** upozorňují na podezřelé body. Andrewsův indexový graf a graf normovaných reziduí ukazují na podezřelé body č. 5, 8, 11, 12. Indexový graf prvků projekční matice  $H$  pak na podezřelé extrémy č. 5, 8, 11.



Obr. 8 Grafy vlivných bodů: (a) Indexový graf prvků H projekční matice, **ADSTAT 3.0**

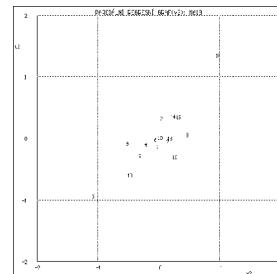
(e) **Rankitové grafy** ukazují vedle normality rozdělení dotyčných reziduí i na vlivné (zde odlehlé) body. Graf normovaných reziduí ukazuje na začátku č. 11 a na konci č. 12 jako odlehle

body. Graf Jackknife reziduí č. 11 a 12 jako odlehlé. Po odstranění dvou odlehlých bodů č. 11 a 12 lze konstatovat, že zbytek dat nevykazuje odchylky od normality.

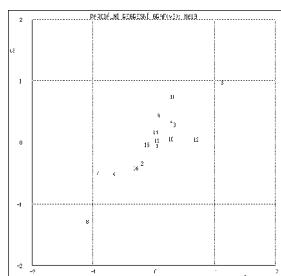


Obr. 9 Rankitové Q-Q grafy: (a) Graf normovaných reziduí, (b) Graf Jackknife reziduí, **ADSTAT 3.0**

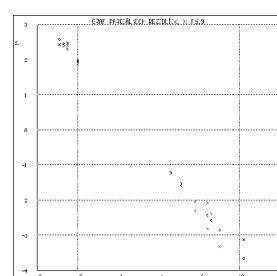
**5.2 Model:** *Parciální regresní grafy*, ale především *parciální reziduální grafy* ukazují na čisté lineární závislosti jednotlivých nezávisle proměnných. Vedle posouzení závislosti navrženého regresního modelu umožňují také indikovat vlivné body, a to č. 5, 8, 11 a 12. Navržený model se jeví stran členů  $\beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$  správný, pouze  $\beta_0$  je nadbytečné.



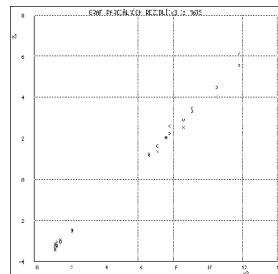
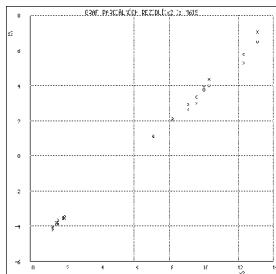
Obr. 10 Parciální regresní grafy: (a) pro proměnnou  $x_1$ , (b) pro proměnnou  $x_2$ , **ADSTAT 1.25**



Obr. 10 (c) pro proměnnou  $x_3$ , **ADSTAT 1.25**



Obr. 11 (a) Parciální reziduální graf pro proměnnou  $x_1$ , **ADSTAT 1.25**



Obr. 11 Parciální reziduální grafy (b) pro proměnnou  $x_2$ , (c) pro proměnnou  $x_3$ , ADSTAT 1.25

**5.3 Metoda:** do této části patří vyšetření splnění základních předpokladů metody nejmenších čtverců (MNČ), za kterých by měla vést k nejlepším nestranným lineárním odhadům regresních parametrů:

*Fisher-Snedecorův test významnosti regrese* potvrdil, že navržený model je přijat jako významný, jinými slovy: závisle proměnná  $y$  a nezávisle proměnné  $x_1, x_2, x_3$  jsou v lineární závislosti.

*Scottovo kritérium multikolinearity* ukazuje, že navržený model není korektní s ohledem na vazby mezi proměnnými.

*Cook-Weisbergův test heteroskedasticity* dokazuje, že rezidua vykazují heteroskedasticitu (nekonstantnost rozptylu).

*Jarque-Berraův test normality reziduí* ukazuje, že klasická rezidua vykazují Gaussovo rozdělení.

*Waldův test autokorelace* ukazuje, že klasická rezidua jsou autokorelována.

*Znaménkový test* prokazuje, že znaménko klasických reziduí se dostatečně střídá, a proto rezidua nevykazují žádný trend.

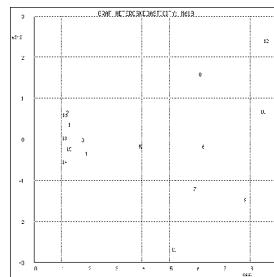
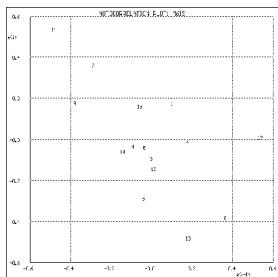
#### TESTOVÁNÍ REGRESNÍHO TRIPLETU (DATA + MODEL + METODA):

Fisher-Snedecorův test významnosti regrese, F	: 1.4050E+03
Tabulkový kvantil, F(1-alpha,m-1,n-m)	: 3.4903E+00
Závěr: Navržený model je přijat jako významný.	
Spočtená hladina významnosti	: 0.000
Scottovo kriterium multikolinearity, M	: 9.6119E-01
Závěr: Navržený model není korektní.	
Cook-Weisbergův test heteroskedasticity, Sf	: 1.9926E+01
Tabulkový kvantil, Chi^2(1-alpha,1)	: 3.8415E+00
Závěr: Rezidua vykazují heteroskedasticitu.	
Spočtená hladina významnosti	: 0.000
Jarque-Berraův test normality reziduí, L(e)	: 3.3085E-02
Tabulkový kvantil, Chi^2(1-alpha,2)	: 5.9915E+00
Závěr: Normalita je přijata.	
Spočtená hladina významnosti	: 0.984
Waldův test autokorelace, Wa	: 1.0320E+01
Tabulkový kvantil, Chi^2(1-alpha,1)	: 3.8415E+00
Závěr: Rezidua jsou autokorelována.	
Spočtená hladina významnosti	: 0.001

Znamékový test, Dt	: -1.9739E-01
Tabulkový kvantil, N(1-alpha/2)	: 1.6449E+00
Závěr: Rezidua nevykazují trend.	
Spočtená hladina významnosti	: 0.422

Graf autokorelace vykazuje přibližně mrak bodů reziduů.

Graf heteroskedasticity vykazuje klín, a proto rezidua vykazují heteroskedasticitu, nekonstantnost rozptylu.



Obr. 12 (a) Graf autokorelace, a (b) Graf heteroskedasticity, ADSTAT 1.25

**6. Konstrukce zpřesněného modelu:** (a) Po odstranění bodů č. 11 a 12 a absolutního členu  $\beta_0$  byly nalezeny nové odhady parametrů zpřesněného modelu.

Odhad	Směrodatná odchylka	Test H0: $B[j] = 0$ vs. HA: $B[j] \neq 0$ t-kriterium	Hypoteza H0 je	Hlad. význam.
B[0]	0.0000E+00	-----	-----	-----
B[1]	-1.1808E+00	3.8271E-01	-3.0854E+00	Zamítnuta
B[2]	1.2454E+00	2.3610E-01	5.2751E+00	Zamítnuta
B[3]	9.1666E-01	1.5049E-01	6.0910E+00	Zamítnuta

Zpřesněný model (v závorce je uveden odhad směrodatné odchylky parametru)

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

je doložen statistickými charakteristikami: vícenásobný korelační koeficient  $R = 0.9990$ , koeficient determinace  $D = 99.80\%$  a predikovaný koeficient determinace  $R^2_p = 99.79\%$  dosáhly vesměs vysokých hodnot. Střední kvadratická chyba predikce  $MEP = 6.078E-02$  a Akaikeho informační kritérium  $AIC = -43.472$  dosáhly nižších hodnot, což dokazuje lepší model než předešlý.

(b) Užitím statistické váhy ( $w_i = 1/y_i^2$ ) kompenzujeme heteroskedasticitu v datech. Obdržíme nové odhady parametrů, v nichž však parametr  $\beta_1$  vychází jako statisticky nevýznamný. Opravený model má tvar (v závorce je vždy uveden odhad směrodatné odchylky parametru)  $y = 0.62 (0.17) x_2 + 0.36 (0.15) x_3$ . Jelikož došlo k významnému snížení rozhodujících kritérií, střední kvadratické chyby predikce  $MEP = 1.23E-02$  a Akaikeho informačního kritéria  $AIC = -62.09$ , lze považovat tyto odhady za lepší než předešlé.

**7. Zhodnocení kvality modelu:** porovnáním hodnot regresní diagnostiky lze snadno provést zhodnocení regresního tripletu dosaženého lineárního regresního modelu pro upravená data, zbavená odlehlých hodnot a upravený regresní model bez absolutního členu a metodou vážených nejmenších čtverců. Nalezený model má tvar (v závorce je vždy uveden odhad směrodatné odchylky parametru)

$$y = 0.62 (0.17) x_2 + 0.36 (0.15) x_3$$

čili obsah kadmia v zrnu potravinářské pšenice je funkcí pouze obsahu kadmia ve stonku a v kořenovém systému a není funkcí obsahu kadmia v otrubách a dále nemá smysl uvádět ani zbytkový obsah kadmia v zrnu při nulovém obsahu kadmia ve zbytku rostlinky.

### **Doporučená literatura a software:**

- (1) Milan Meloun a Jiří Militký: *Statistické zpracování experimentálních dat*, Plus Praha 1994, resp. East Publishing Praha 1998.
- (2) Milan Meloun a Jiří Militký: *Statistické zpracování experimentálních dat - Sbírka úloh*, Univerzita pardubice 1996.
- (3) ADSTAT 1.25, 2.0 a verze 3.0, TriloByte Statistical Software Pardubice, 1992, 1993, 1999.