

6

LINEÁRNÍ REGRESNÍ MODELY

Při budování 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 výsledky metodou nejmenších čtverců své vlastnosti. Užití lineární regresní analýzy se týká následujících možností:

1. *Popis dat*: hledáme vztah, lineární regresní model, který sumarizuje soubor dat. 2. *Určení parametrů*: nejběžnějším cílem regresní analýzy je vyčíslení nejlepších odhadů neznámých parametrů regresního modelu. Uživatel navrhne regresní model a regresní analýzou se snaží model prokázat. Často tento cíl překrývá i ostatní záměry regresní analýzy.

3. *Predikce*: nejdůležitějším cílem regresní analýzy je predikce, vyčíslení hodnot závisle proměnných. Bývá to často cena, dodací lhůta, účinnost, obložnost v nemocnici, výtěžek reakce, síla kovu, atd. Predikce jsou důležité i v plánování, monitoringu, vyhodnocování chemických procesů atd. Je však řada předpokladů a kvalifikací, které se musí respektovat v regresním modelu a datech. Často se například nesmí extrapolovat mimo rozsah dat. Intervalové odhady vyžadují dodržení předpokladu normality. Metoda nejmenších čtverců (MNČ) má svých sedm důležitých předpokladů, které je třeba respektovat a dodržet.

4. *Řízení*: regresní modely lze využít také k monitoringu a řízení systémů, například ke kalibraci měřicího systému. Když využijeme regresní model k řídicím účelům, nezávisle proměnné musí být vztaheny k závisle proměnné kauzálním způsobem.

5. *Výběr (volba) proměnných*: volba proměnných sleduje ty nezávisle proměnné, které vysvětlují významný objem proměnlivosti závisle proměnné. V řadě aplikací nejde o jednorázový proces, ale o spojitý proces výstavby modelu.

Základní předpoklady metody nejmenších čtverců: statistické vlastnosti odhadů $\hat{y}_p, \hat{\epsilon}$, \mathbf{b} závisí na splnění jistých základních předpokladů.

Předpoklady metody nejmenších čtverců:

I. *Regresní parametry β mohou nabývat libovolných hodnot.* V praxi však existují často 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 hodnotu 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 g mají nulovou střední hodnotu $E(g) = 0$.* To musí u korelačních modelů platit vždy. U regresních modelů se může stát, že $E(g) = K, i = 1, \dots, n$, což znamená, že model neobsahuje absolutní člen. Po jeho zavedení bude $E(\mathbf{g}) = 0$, kde $\mathbf{g} = y_i - \hat{y}_{P,i} - K$.

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

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

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

Pokud platí předpoklady I až VI, jsou odhady \mathbf{b} parametrů $\boldsymbol{\beta}$ 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.

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* (kritika dat, kritika modelu a kritika metody 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 MNC (složka *metoda*).

Základní rozdíl mezi regresní diagnostikou a klasickými testy spočívá v tom, že u regresní diagnostiky není třeba přesně formulovat alternativní hypotézu H_A . 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 zde slouží pouze 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 regresní 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 z kap. 2. Přes svoji

jednoduchost umožňuje regresní 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 náhodnými veličinami.

Kvalita dat úzce souvisí s užitým regresním modelem. Při posuzování se sleduje především výskyt *vlivných bodů*, které mohou být hlavním zdrojem řady problémů, jako je zkreslení odhadů a růst rozptylů až k naprosté nepoužitelnosti regresních modelů. Vlivné body lze rozdělit do tří skupin:

a) *Hrubé chyby*, které jsou způsobeny měřenou veličinou (*vybočující pozorování*) nebo nevhodným nastavením vysvětlujících proměnných (*extrémy*). Hrubé chyby jsou obvykle důsledkem chyb při manipulaci s daty.

b) *Body s vysokým vlivem* (tzv. golden points) jsou speciálně vybrané body, které byly přesně změřeny, a které obvykle rozšiřují predikční schopnosti modelu.

c) *Zdánlivě vlivné body* vznikají jako důsledek nesprávně navrženého regresního modelu.

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 vysvětlujících (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. O jejich výsledném vlivu však především rozhoduje to, že jsou 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 \mathbf{H} (detaily v učebnici⁷²).

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 *reziduí v* regrese y na sloupcích matice $\mathbf{X}_{(i)}$ a reziduí *u* regrese x_i na sloupcích matice $\mathbf{X}_{(i)}$. Přitom matice $\mathbf{X}_{(i)}$ vznikne z matice \mathbf{X} vynecháním i -tého sloupce \mathbf{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 $\hat{R}_{y,x_j(x)}$.

c) Rezidua v parciálním regresním grafu jsou shodná s klasickými rezidui $\hat{\epsilon}_i$ pro nedělený model.

d) V grafu jsou indikovány vlivné body a i některá porušení předpokladů metody nejmenších čtverců (heteroskedasticita).

Parciální reziduální grafy se označují také jako grafy "*komponenta + reziduum*". Parciální reziduální grafy však poskytují poněkud odlišné informace než parciální regresní grafy:

a) Směrnice lineární závislosti je rovna b_j a úsek je nulový. Lineární závislost pak ukazuje na vhodnost navržené proměnné x_j v modelu.

b) Rezidua regresní přímky jsou přímo rezidua $\hat{\epsilon}_i$ pro nedělený model.

c) Pokud je úhel mezi x_j a některými sloupci matice $X_{(j)}$ malý (*multikolinearita*), ukazuje parciální reziduální graf nesprávně malý rozptyl kolem regresní přímky $b_j x_j$ a dochází navíc i k potlačení efektu vlivných bodů.

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), popř. 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-odhady*. 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 zatíženy 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 F_y^2 (vysvětlovaná proměnná) a F_x^2 (vysvětlující proměnné), $K = F_y^2/F_x^2$. Pro odhad směrnice regresní přímky $y = a x + b$ pak platí

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

kde

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

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

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

$$S_{yx} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}).$$

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 uvedených vztahů k odhadům minimalizujícím kolmé vzdálenosti (*ortogonální 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 hodnotí*, (**General Principal Component Regression**) vedoucí k systému vychýlených odhadů, kde vychýlení je řízeno jedním parametrem.

Postup výstavby lineárního regresního modelu (ADSTAT)

(definice regresních diagnostik a ostatních statistických pojmů jsou v učebnici⁷²):

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 $x_j x_k$. Pouze v případech, u kterých je předem známo, že model má obsahovat funkce vysvětlujících proměnných, může být výchozí model dle těchto požadavků upraven.

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 i vlivné body, které způsobují multikolinearitu.

Podle volby uživatele se provedou požadované transformace původních proměnných. Zadává se, zda model obsahuje absolutní člen. Uživatel může volit polynomičtí transformaci zadáním stupně polynomu, Taylorův rozvoj do 2. stupně a lineární model s interakcemi. Uživatel může zadat libovolnou mocninu původních proměnných včetně logaritmu. Ostatní typy transformací se provádějí při přípravě dat k výpočtu v datovém editoru. K odstranění případné heteroskedasticity, vzniklé nelineární transformací proměnné y , je možné zadat nestatistické váhy, jež odpovídají kvazilinearizaci.

Provádí se sestavení korelační matice R a její rozklad na vlastní čísla a vlastní vektory. Jsou vypočteny faktory *VIF* (variation inflation factor) k indikaci multikolinearity a dále jsou vyčíslena seříděná vlastní čísla. K určení inverzní matice R^{-1} se užívá metoda racionálních hodnotí GPCR 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 intervalu $10^{-5} \leq P \leq 10^{-3}$.

3. Odhadování parametrů: odhadování parametrů modelu se provádí metodou racionálních hodnot GPCR s volbou $P = 10^5$, což je vlastně MNČ. Ze zobecněné inverzní matice R^+ jsou určovány odhady parametrů \mathbf{b} , jejich směrodatné odchylky $\sqrt{D(\mathbf{b}_j)}$ 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 , popř. posoudit linearitu modelu.

4. Regresní diagnostika: identifikace vlivných bodů je prováděna využitím pěti rozličných grafů, a to *grafů Wiliamsova, Pregibonova, McCullohova-Meeterova, L-R, a grafu predikovaných reziduí*. Dále musíme ověřit 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 využitím parciálních regresních grafů nebo grafů "komponenta + reziduum". Obvykle jsou využívány následující tabulky:

Tabulka výsledků obsahuje hodnoty predikce \hat{y}_j , rozptylů predikce $D(\hat{y}_j)$ a relativní odchylky predikce od experimentálních dat. Je uvedena i průměrná absolutní, resp. relativní odchylka a reziduální suma čtverců RSC . Následuje statistická analýza klasických reziduí.

Tabulka reziduí obsahuje klasická rezidua \hat{e}_j , normovaná rezidua \hat{e}_{Ni} , standardizovaná rezidua \hat{e}_{Si} a Jackknife rezidua \hat{e}_{Ji} . Je uveden odhad autokorelačního koeficientu reziduí prvního řádu k_1 .

Tabulka vlivných bodů obsahuje veličiny H_{ii} , $H_{ii}^{\{}$, D_i , A_i , DF_i , $LD_i(\mathbf{b})$, $LD_i(F^2)$ a $LD_i(\mathbf{b}, F^2)$. Hvězdičkou bývají označeny hodnoty silně vlivných bodů.

5. Konstrukce zpřesněného modelu: při využití

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

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

7. Kalibrační modely: u kalibračních modelů se pro daný signál y^* vypočte hodnota x^* spolu se svým konfidenčním intervalem. Před vlastním užitím kalibračního modelu je vhodné určit limitu detekce a limitu stanovení, které určují použitelnou dolní hranici kalibračního modelu nebo odpovídající metody. Postup obsahuje

- (a) Návrh modelu.
- (b) Statistickou analýzu reziduí.
- (c) Výpočet derivací a integrálů.
- (d) Určení kalibračních mezí.
- (e) Sestavení kalibrační tabulky.

8. Testování různých hypotéz: ve zvláštních případech, jako je porovnání několika přímk atd., se provádí testování pomocí dalších testů k ověřování rozličných typů hypotéz.

Uživatel může při interaktivní práci s počítačem některé tabulky nebo grafy vynechat. Na základě analýzy vlivných bodů a reziduí lze provést i vypuštění některých bodů a výpočty pak zopakovat. Podrobný popis, vzorce a statistické testy najde čtenář v doporučené učebnici⁷².