

# Nelineární regrese a tvorba matematických modelů

Jiří Militký, KTM, VŠST Liberec, 461 17 LIBEREC

## 1. Úvod

Nelineární regrese se v technických aplikacích používá velmi často při konstrukci matematických modelů. Pro tyto účaly však již obecné nepostačuje klasické kriterium nejménších čtverců odchylek (MNČ), ale je třeba modelovat i struktury chyb.

V tomto příspěvku je nejdříve stručně pojednáno o jednotlivých fázích matematického modelování s využitím regresních metod. Druhá část je zaměřena na využití některých technik transformace a vážení pro zlepšení rozdělení chyb, resp. snížení vlivu nelinearity modelu na statistickou analýzu odhadů.

## 2. Modelování s využitím regresních metod

Pod pojmem modelování se obecně chápe záměrná činnost přiřazující zkoumanému systému (originál objekt) jiný známý systém (fyzikální, abstraktní) zvaný model. Při konstrukci modelu se využívá vhodných kriterií vyjadřujících míru shody s jistým rysem originálu. V řadě praktických případů se konstruuují modely matematické, které lze relativně snadno použít (např. pro optimalizaci atd.).

Vlastní postup modelování je výrazně ovlivněn typem modelovaného systému.

Relativně nejjednodušší je konstrukce modelů pro nesložité, malé a dobře organizované systémy (laboratorní experimenty v idealizovaných podmínkách). Modely téhoto systému vycházejí z fyzikálně chemických zákonů a závisí na relativně malém počtu parametrů s definovaným fyzikálním významem. Je obyčejně také známo, v jakém rozmezí se mohou parametry pohybovat. Pokud lze závislost mezi výstupní (vy-světlovanou) a vstupními (vysvětlujícími) proměnnými znázornit, jsou patrné nejen trendy, ale také lokální tvarové zvláštnosti. V anglosasské literatuře se často tento

typ úloh označuje jako "curve fitting".

Je zajímavé, že také u velkých, špatně organizovaných systémů (typických pro behavioristická vědy) není postup konstrukce modelů relativně složitý. Vzhledem k existenci řady proměnných a nezachytitelných vlivů se často používá modelů ve tvaru Taylorova rozvoje neznámé funkce obyčejně do maximálně kvadratických členů. Pro zjištění významnosti jednotlivých vysvětlujících proměnných a jejich interakcí tento postup často vyhovuje. Jde o klasický případ regrese (v řadě úloh lineární).

Nejvíce komplikované z hlediska modelování jsou vědy technické (průmyslové experimenty). Zde se obyčejně pracuje s hůře organizovanými "difusními" systémy, kde sice platí fyzikálně chemické zákony, ale jejich vliv je "maskován" méně známými, nepopsanými vlivy. Datově orientované modely typické pro velké i špatně organizované systémy neumožňují stanovení modelových parametrů fyzikální povahy. Na druhé straně jsou však modely vycházející z fyzikálních představ silně zjednodušené a nepopisují obyčejně dobře chování originálu. Typickým jevem je zde multimodelovost a vlastní modelování se stává tématem "umění".

V mnoha případech se postup modelování mění v závislosti na cíli modelování. Je zřejmé, že modely pro účely predikce mohou být tvoreny jako datově závislé (kriteriem je shoda modelu s experimentem). Na druhé straně modely strukturální musí být v souladu s hypotézami a teoriemi o modelovaném systému.

Omezme se v dalším na jednu speciální třídu modelů, kterou lze označit jako deterministické modely s náhodným rušením.

Vysvětlující proměnné  $x$ , nechť jsou deterministické. Tyto proměnné při působení na orginál vyvolají měřitelný důsledek  $y^*$ , který lze vyjádřit deterministickým modelem

$$y^* = \{ (x, \underline{\alpha}) \quad (1)$$

kde  $f(\cdot)$  je modelová funkce a  $\underline{\alpha}$  jsou modelové parametry. Důsledkem způsobu měření veličiny  $y^*$ , kolísání nastavených podmínek a nezachytitelných vlivů je výsledná vysvětlování

veličina  $y$  náhodná

$$y = G(y^*, e) \quad (2)$$

Zde symbol  $e$  označuje souhrnnou chybu (považuje se za náhodnou veličinu charakterizovanou např. hustotou pravděpodobnosti  $p(e)$ ) a  $G(\cdot)$  vyjadřuje pravděpodobnostní model působení poruch.

Obecně je tedy třeba při modelování těchto typů úloh hledat deterministický model  $f(x, a)$  a pravděpodobnostní model  $p(e)$ . Ve velké většině případů se v praxi věnuje pozornost pouze specifikaci modelu  $f(x, a)$  s tím, že se předpokládá aditivní model působení poruch

$$G(f(x, a), e) = f(x, a) + G \cdot e_s \quad (3)$$

Zde je směrodatná odchylka,  $e_s$  je náhodná veličina s nulovou střední hodnotou a jednotkovým rozptylem. Pokud je  $p(e_s)$  hustota pravděpodobnosti normovaného normálního rozdělení lze aplikací metody maximální věrohodnosti dospět ke kritériu MNČ. Neplatí-li však model (3), může být použiti MNČ zdrojem toho, že není jinak korektní model  $f(x, a)$  prakticky akceptován.

Z toho plyne, že podcenění konstrukce pravděpodobnostního modelu je jednou z příčin možného nezdaru celého procesu modelování.

### 3. Postup při modelování

Při modelování deterministických systémů s náhodným rušením se běžně vychází z experimentu provedeného na originálu. Výsledkem je  $N$ -tice hodnot vysvětlované (výstupní) proměnné  $y_i$ , pro zvolené nastavení vysvětlujících proměnných  $x_i$ . Tyto výsledky experimentu lze chápat jako  $N$ -tici bodů v  $E^{m+1}$  rozměrném prostoru.

Vlastní postup modelování se obecně skládá z těchto fází:

1. Definice úlohy
2. Realizace experimentu
3. Konstrukce deterministického modelu
4. Konstrukce pravděpodobnostního modelu

## 5. Adjustace modelu

## 6. Verifikace modelu

V každé fázi záleží zejména na znalostech o modelovaném systému a zkušenostech modelujícího. Platí zde v plné míře výrok "Experimentátor ví vždy o problému více než jeho počítač". Na druhé straně právě použití počítače umožnuje zahrnutí obecnějších typů pravděpodobnostních modelů a diskriminaci mezi nimi.

### 3.1. Definice úlohy

Tato etapa se považuje velmi často za zcela triviální. V mnoha případech je však právě nejasná definice cíle modelování zdrojem řady problémů. Např. úloha "vyhlazení modelu" za účelem následné numerické derivace, resp. integrace modelu (častá třeba v chemii) se dá jednoduše řešit s využitím regresních spline modelů, vyjádřených jednoduše ve tvaru useknutých polynomů [4]

$$E(y/x) = \sum_{j=0}^m a_j x^j + \sum_{L=1}^k b_L (x - t_L)_+^m \quad (4)$$

Zde  $t_L$  jsou tzv. uzlové body (místa styku lokálních polynomů stupně  $m$ ). Symbol  $(x)_+$  je useknutá funkce, pro kterou platí

$$\begin{aligned} (x)_+ &= x && \text{pro } x > 0 \\ (x)_+ &= 0 && \text{pro } x \leq 0 \end{aligned}$$

Je zřejmé, že pro zvolené  $t_L$ ,  $L = 1 \dots k$  je úloha odhadu parametrů  $(a_j, b_L)$ ,  $j = 1 \dots m$ ,  $L = 1 \dots k$  úlohou lineární regrese. Při interaktivním způsobu volby uzlových bodů na displeji lze snadno modelovat i velmi komplikované polymodální závislosti, pro které je nalezení matematických modelů jinak takřka nemožné.

Také v situacích, kdy je třeba hledat optimální podmínky funkce systémů (úloha extremalizace) postačuje místo konstrukce složitých modelů použití polynomů nízkého stupně, případně spline polynomů.

V řadě praktických úloh tedy může správná specifikace cíle značně ulehčit vlastní modelování.

### 3.2. Realizace experimentu

Obecné je třeba při vlastní realizaci experimentu stanovit experimentální oblast (obyčejně jako rozmezí všech vysvětlujících proměnných) a způsob jejího "zmapování", tj. určení poloh jednotlivých složek  $x_i$ .

Plati pravidlo, že dobré naplánovaný experiment značně ulehčí jak numerickou, tak i statistickou stránku vlastní regrese (odhadu parametrů). Na druhé straně řada potíží, zejména při více vysvětlujících proměnných tvoří multikolinearita jako typický případ špatně plánovaného experimentu.

Z celé řady kritérií pro vhodné plánování experimentů, uvedeme pouze tzv. D-optimalitu, kdy se minimalizuje determinant matice  $(J^T J)^{-1}$ . Matice  $J$  má prvky

$$J_{ij} = \frac{\partial f(x_i, \alpha)}{\partial \alpha_j} \quad \begin{matrix} i = 1 \dots m \\ j = 1 \dots m \end{matrix} \quad (5)$$

Obyčejně platí, že pro špatně postavené modely  $f(x, \alpha)$  nevede konstrukce D-optimálního experimentálního plánu k výraznému zlepšení. V řadě praktických úloh se také zpracovávají výsledky již hotových (neoptimálně plánovaných) experimentů.

V těchto situacích nezbývá než vzít v úvahu možné důsledky způsobené např. multikolinearitou [4].

### 3.3. Konstrukce deterministického modelu

I když existují jistá pravidla pro obecné řešení úlohy modelování (uvedená např. v Descartové spisu "Pravidla o řízení rozumu") nelze nalézt universální postup.

Pro případ konstrukce třídy empirických modelů založených na geometrické podobnosti modelu s experimentem existuje celé spektrum technik využívajících různých katalogů nebo atlasů křivek a ploch. S oblibou se používají takové funkce, které lze linearizovat nebo jejichž tvar závisí jednoduše na hodnotách parametrů. Pro konstrukci více-rozměrných empirických modelů lze použít také speciální metody regresní diagnostiky (grafy konstrukční proměnné, přidané proměnné atd.) [4].

Obecně platí, že v praxi se této fázi modelování věnuje maximální pozornost, protože jde o "viditelný" výsledek celé činnosti.

### **3.4. Konstrukce pravděpodobnostního modelu**

Tato etapa modelování je v praxi běžně zanedbávána. Navíc se konstrukce pravděpodobnostního modelu projeví "jen" na odhadech parametrů a jejich statistickém chování.

Jde tedy o "skrytou" část modelu, která však může celý proces modelování rozhodujícím způsobem ovlivnit. Při konstrukci pravděpodobnostního modelu  $G(\cdot)$  lze postupovat třemi způsoby.

A. Sestavit teoretický model působení chyb na základě znalostí a hypotéz o jejich vzniku a statistickém chování.

B. Sestavit empirický, velmi obecný model zahrnující většinu modelů působení chyb a hledat jeho zjednodušení pro dany případ.

C. Na základě ověření předpokladů MNČ (po její aplikaci) provést jejich zpřesnění. To se týká zejména heteroskedasti-city a autokorelačních struktur.

Ve druhé části tohoto příspěvku jsou detailněji rozebrány cesty A a B. Postup ad C je vlastně iterativní aplikace metod regresní diagnostiky, která je popsána např. v práci [4].

### **3.5. Adjustace modelu**

V této etapě se provádí odhad modelových parametrů a jejich statistická analýza (testy hypotéz).

Oblíbený je Bayesovský přístup, který lze často v kontextu modelování chápat jako zpřesňování z priorních informací s využitím informací získaných z experimentu.

Pomocí Bayesovy definice podmíněné pravděpodobnosti lze z posteriorní hustotu pravděpodobnosti  $p(\underline{\alpha}/y)$  vyjádřit pomocí z priorní hustoty pravděpodobnosti modelových parametrů  $p(\underline{\alpha})$  a věrohodnostní funkce  $L(\underline{\alpha})$  ve tvaru

$$p(\underline{\alpha}/y) = \frac{L(\underline{\alpha}) \cdot p(\underline{\alpha})}{p(y)} \quad (6)$$

Bayesovský odhad  $\underline{\alpha}^*$  pak maximalizuje a posteriorní hustotu pravděpodobnosti, tedy

$$\underline{\alpha}^* = \max p(\underline{\alpha}/y) \quad (7)$$

Pokud nejsou k dispozici a priorní informace o parametrech je Bayesovský odhad  $\underline{\alpha}^*$  shodný s odhadem maximalizujícím věrohodnostní funkci

$$\underline{\alpha}^* = \max L(\underline{\alpha}) \equiv \max \ln L(\underline{\alpha}) \quad (8)$$

Za předpokladu, že  $y_i$  jsou nezávislé, stejně rozdělené náhodné veličiny s hustotou pravděpodobnosti  $p(y_i/\underline{\alpha})$  má věrohodnostní funkce tvar

$$L(\underline{\alpha}) = \prod_{i=1}^N p(y_i/\underline{\alpha}) \quad (9)$$

Pokud platí model měření vyjádřený rov. (2) a chyby  $e$  mají hustotu pravděpodobnosti  $p_e(y_i/e)$ , lze snadno určit, že

$$p(y_i/\underline{\alpha}) = p_e[G^{-1}(f(x_i, \underline{\alpha}), y_i)] \left| \frac{\partial G^{-1}(\cdot)}{\partial y_i} \right| \quad (10)$$

kde  $G^{-1}(\cdot)$  je formálně inverzní funkce k funkci  $G(\cdot)$ . Je zřejmé, že pro aditivní model měření

$$G(\cdot) = y^* + e$$

je

$$G'(\cdot) = y - y^*$$

a tedy

$$p(y/\underline{\alpha}) = p_e[y - f(x, \underline{\alpha})] \quad (11)$$

Pro multiplikativní model měření

$$G(\cdot) = y^* \cdot \exp(e)$$

je inverse

$$G'(\cdot) = \ln y - \ln y^*$$

Po dosazení do rov. (10) pak vyjde

$$p(y_i/\underline{\alpha}) = \frac{1}{|y_i|} \cdot p_e[\ln y_i - \ln \{(\underline{x}_i, \underline{\alpha})\}] \quad (12)$$

Pro případ, že  $p_e(\cdot)$  je hustota pravděpodobnosti normálního rozdělení  $N(0, \sigma^2)$ , vyjde po dosazení z rovnice (11) do rovnice (9) a (8), že maximálně věrohodné odhadu minimalizují kriterium MNČ

$$S(\underline{\alpha}) = \sum_{i=1}^n [y_i - \{(\underline{x}_i, \underline{\alpha})\}]^2 \quad (13)$$

Podobně resultuje pro multiplikativní model chyb (rov. (12)) kriterium MNČ v logaritmech.

Pokud nejsou chyby e nezávislé, ale jsou charakterizovány kovarianční maticí  $C_e$  předpokládá se obyčejně, že mají vícerozměrné normální rozdělení. Modeluje se pak pouze kovariační matice.

Z výše uvedeného je patrné, že v řadě případů vede adjustace modelu na problém minimalizace součtu čtverců. Pro tuto úlohu existuje řada speciálních algoritmů (Gauss-Newtonovy a Marquardtovy metody jsou nejznámější), které umožňují vlastní numerickou extremalizaci. Ve složitějších případech je třeba použít obecných optimalizačních metod vhodných pro maximalizaci sestavené věrohodnostní funkce. (Známý je např. program MAXLIK v jazyce GAUSS).

Při statistické analýze se většinou vychází z představy, že lze v okolí odhadu  $\underline{\alpha}^*$  funkci  $f(\underline{x}, \underline{\alpha})$  dostatečně přesně linearizovat. Pak lze vyjádřit kovariační matici odhadů ve tvaru

$$D(\underline{\alpha}) = G^2 (J^T J)^{-1} \quad (14)$$

a provádět statistickou analýzu podle stejných vztahů jako u lineární regrese (při nahradě matice X maticí J). Kvalita linearizace souvisí úzce s nelinearitou regresního modelu. Dá se také vyjádřit pomocí vychýlení odhadů  $\underline{b}$ . Z hlediska jednoduchosti vyjádření je výhodné počítat vychýlení dle Cooka a Tsaiie

$$\underline{b} = (\underline{J}^T \underline{J})^{-1} \underline{J}^T \underline{d} \quad (15)$$

kde  $\underline{d}$  je vektor se složkami

$$d_i = -G^2 \operatorname{tr} [(\underline{J}^T \underline{J})^{-1} \cdot G_i] / 2 \quad (16)$$

V rov. (16) je  $\operatorname{tr}(.)$  stopa matice a matice  $G_i$  má složky

$$G_{j,L} = \frac{\partial f(x_i, \underline{\alpha})}{\partial \alpha_j \partial \alpha_L} \quad j, L = 1, \dots, m \quad (17)$$

Vychýlení odhadů se považuje za akceptovatelné pokud

$$|100 b_j / \alpha_j^*| \leq 1$$

### 3.6. Verifikace modelu

V této etapě se ověřuje, zda model:

- a) vyhovuje datům
- b) vyhovuje k priorním představám o parametrech
- c) má dostatečné predikční schopnosti.

V řadě případů se provádí také diskriminace mezi modely, při které se ještě zohledňuje to, aby byl model co nejjednodušší (s nejmenším počtem parametrů).

Existuje celá řada technik pro vyjádření shody modelu s daty a pro posouzení jeho predikčních schopností.

Jednoduše lze počítat střední kvadratickou chybu predikce dle vztahu

$$MSE = \frac{1}{N} \sum_{i=1}^N [y_i - \{f(x_i, \underline{\alpha}_{(-i)}^*)\}]^2$$

kde  $\underline{\alpha}_{(-i)}^*$  jsou odhadы získané při vyneschání i-tého bodu  $(x_i, y_i)$ .

Zobecněním tohoto postupu jsou metody typu "Cross validation", kdy se počítají odhadы z "trénovacích" dat a chyba predikce z "ověřovacích" dat.

#### 4. Pravděpodobnostní modely měření

Jak již bylo uvedeno lze při konstrukci modelu působení chyb využít buď z priorních představ nebo sestavit obecný empirický model. Zde si ukážme obě tyto techniky.

##### 4.1. Teoretický model chyb

Teoretické modely chyb se tvoří na základě konkretních znalostí o modelovaném systému a měřícím procesu. Většinou se tyto modely týkají typu heteroskedasticity nebo autokorelační struktury chyb. Je např. známo, že u řady přístrojů, které měří v rozsahu několika řádů nelze dodržet požadavek konstantního rozptylu měření  $\sigma^2$ . Je však často docíleno konstantnosti relativní chyby měření definované výrazem

$$s_e = \sigma / f(x, \underline{\alpha})$$

Rozptyl měření je pak modelován jako kvadratická funkce hodnoty regresního modelu

$$D(y_i) = \sigma_0^2 f^2(x_i, \underline{\alpha}) \quad (18)$$

Pokud platí ostatní standardní předpoklady o chybách měření vede použití rovnice (18) k modelu

$$y = f(x, \underline{\alpha}) + \sigma_0 \cdot f(x, \underline{\alpha}) \cdot e_s \quad (19)$$

kde  $e_s$  je opět standardizovaná náhodná veličina s nulovou střední hodnotou a jednotkovým rozptylem. Rovnice (19) je případem nelineárního heteroskedastického modelu, pro který vede použití metody maximální věrohodnosti k úloze minimalizace kriteria

$$S_W = \sum_{i=1}^n [y_i - f(x_i, \underline{\alpha})]^2 / f^2(x_i, \underline{\alpha}) \quad (20)$$

Pokud se dá předpokládat, že chyby měření jsou malé, lze místo kriteria  $S_W$  použít jednodušší kriterium čtverců relativních odchylek

$$S_R = \sum_{i=1}^N \left[ 1 - \{(\bar{x}_i, \omega) / y_i\} \right]^2 \quad (22)$$

Typickým případem vzniku speciální autokorelační struktury chyb je experimentování na jednom vzorku. Příkladem je např. sledování průběhu chemických reakcí, kdy se z reaktoru odebírá v určitých intervalech vzorky, které se analyzuji na měřicím přístroji [4].

Celková chyba  $e_i$  se v těchto situacích skládá z chyb způsobených fluktuací podmínek procesu  $u_j$  (ve všech časech až do i-tého) a chyby měření  $v_i$ . To lze vyjádřit vztahem

$$e_i = \sum_{j=1}^i u_j + v_i \quad (23)$$

Rov. (23) vyjadřuje fakt, že se chyby způsobené neidealitami procesu kumulují.

Pro jednoduchost předpokládejme, že

a) chyby procesu a chyby měření jsou vzájemně nezávislé

$$E(u_j, v_i) = \emptyset \quad j = 1 \dots i$$

b) chyby procesu jsou vzájemně nezávislé  $E(u_i, u_j) = \emptyset$

Mají konstantní rozptyl  $D(u_i) = G_u^2$  a nulovou střední hodnotu  $E(u_i) = \emptyset$

c) chyby měření  $v_i$  jsou také vzájemně nezávislé  $E(v_i, v_j) = \emptyset$

Mají konstantní rozptyl  $D(v_i) = G_v^2$  a nulovou střední hodnotu  $E(v_i) = \emptyset$ .

Rov. (23) se dá vyjádřit také ve tvaru

$$e_i = e_{i-1} + w_i \quad (24)$$

kde  $e_{i-1}$  je chyba v (i-1)ním čase a  $w_i$  je souhrnná chyba odpovídající přechodu systému z (i - 1)ního do i-tého stavu. Na základě výše uvedeného platí, že

$$E(w_i) = \emptyset \quad D(w_i) = G_u^2 + G_v^2 \approx \tau^2$$

Chyby  $e_i$  však již nejsou nezávislé, a platí pro ně

$$e_i = \sum_{j=1}^i w_j \quad e_1 = w_1$$

Maticově lze pak vektor chyb  $\underline{e}$  vyjádřit ve tvaru

$$\underline{e} = \mathbf{A} \underline{w} \quad (25)$$

kde  $\mathbf{A}$  je dolní trojúhelníková matici s jedničkami na a pod hlavní diagonálou.

Pro kovarianční matici chyb  $\underline{e}$  pak platí

$$\mathbf{C}_e = E(\underline{e} \cdot \underline{e}^T) = \sigma^2 \mathbf{A} \cdot \mathbf{A}^T \quad (26)$$

a její inverze je rovna

$$\mathbf{C}_e^{-1} = \sigma^{-2} (\mathbf{A}^{-1})^T (\mathbf{A}^{-1}) \quad (27)$$

Vzhledem ke speciální struktuře matice  $\mathbf{A}$  je matice  $\mathbf{A}^{-1}$  bidiagonální.

Diagonální prvky jsou rovny 1 a prvky poddiagonálního pásu jsou rovny -1. Ostatní prvky jsou nulové.

Pokud mají chyby  $\underline{e}$  normální rozdělení je sdružená hustota pravděpodobnosti  $p(y/\underline{\alpha})$  rovna

$$p(y/\underline{\alpha}) = (2\pi\sigma^2)^{-\frac{N}{2}} \cdot \exp[-0.5(\underline{y}-\underline{f})^T \mathbf{C}_e^{-1} (\underline{y}-\underline{f})] \quad (28)$$

kde vektor  $\underline{f}$  má prvky  $f(x_i, \underline{\alpha})$ .

Vzhledem k tomu, že rovnice (28) definuje přímo věrohodnostní funkci, je zřejmé, že maximálné věrohodné odhady lze získat minimalizací výrazu v exponenciále. Ten je možno vyjádřit ve tvaru

$$S_L = (\underline{y}-\underline{f})^T \mathbf{C}_e^{-1} (\underline{y}-\underline{f}) = \sum_{i=1}^N [L_i - K_i(\underline{\alpha})]^2 \quad (29)$$

Dále platí

$$L_i = y_i - y_{i-1} \quad K_i(\underline{\alpha}) = f(x_i, \underline{\alpha}) - f(x_{i-1}, \underline{\alpha})$$

$$y_0 = f(x_0, \underline{\alpha}) = \emptyset$$

Kriterium  $S_L$  tedy odpovídá metodě nejménších čtverců pro první diference.

Opuštěním předpokladů normality, resp. aditivity chyb lze

dospět ke komplikovanějším modelům  $G(\cdot)$ , kterým odpovídají složitější kriteria regrese.

Tak pro případ modelování neizotermních kinetických procesů se ukázalo jako vhodné kriterium čtverců logaritmů prvních diferencí (kumulativní procesní chyby a multiplikativní model měření).

Složitější modely chyb zahrnující i chyby ve vysvětlujících proměnných jsou uvedeny v [1].

#### 4.2. Empirický model chyb

Při konstrukci empirického modelu chyb je snahou nalézt takový výraz, který by zahrnoval pokud možno všechny typické struktury chyb.

S výhodou se zde využívá vhodné mocninné transformace, která může často zajistit zkonstantnění rozptylu a zlepšení symetrie rozdělení chyb (přiblížení k normalitě). Klasickým příkladem je Box-Coxova rodina transformací závisející na jednom parametru  $c$ . Tato rodina transformací je definována vztahy

$$z^{(c)} = (z^c - 1)/c \quad c \neq 0$$

$$z^{(c)} = \ln z \quad c = 0$$

S využitím této transformace lze vyjádřit celou třídu modelů měření ve tvaru

$$y^{(c)} = f^{(c)}(\underline{x}, \underline{\alpha}) + G \cdot e_s \quad (30)$$

Je zřejmé, že pro  $c = 1$  rezultuje z rovnice (30) aditivní model měření a pro  $c = 0$  multiplikativní model měření.

Model (30) má tu zajímavou vlastnost, že  $f(\underline{x}, \underline{\alpha})$  představuje podmíněný medián veličiny  $y$  nezávisle na velikosti  $c$  [3]. Toto platí pro případ, že  $e_s$  v modelu (30) má symetrické rozdělení.

V práci [3] bylo ukázáno pomocí Taylorova rozvoje modelu (30), že tento model je možno vyjádřit také ve tvaru

$$y \approx f(\underline{x}, \underline{\alpha}) + G \int^{1-c} f'(\underline{x}, \underline{\alpha}) \cdot e_s + O(G^2) \quad (31)$$

Pro malé  $G$  je tedy model (30) ekvivalentní nelineárnímu heteroskedastickému regresnímu modelu.

Platí tedy např., že multiplikativní model měření ( $c = 0$ ) je přibližně ekvivalentní modelu heteroskedasticity vyjádřenému rov (18) tj. případu konstantních relativních chyb měření.

V některých případech model (30) neumožňuje simultánní přibližení k normalitě (resp. symetrickému rozdělení) a odstranění heteroskedasticity.

Pak lze provést další rozšíření o funkci heteroskedasticity a dospět k modelu

$$y^{(c)} = f^{(c)}(\underline{x}, \underline{\alpha}) + G \cdot g(x, d) e_s \quad (32)$$

Funkce  $g(x, d)$  se často volí ve tvaru [2]

$$g(x, d) = x^d$$

nebo ve tvaru [5]

$$g(x, d) = f^d(\underline{x}, \underline{\alpha})$$

Komplikovanější typy  $g(x, d)$  jsou vhodné pro případ, kdy se uvažují také chyby ve vysvětlujících proměnných [1].

Lze dokázat, že model (32) odpovídá přibližně nelineárnímu heteroskedastickému regresnímu modelu

$$E(y) = f(\underline{x}, \underline{\alpha})$$

$$D(y) = G^2 g^2(x, d) \cdot f^{2(1-c)}(\underline{x}, \underline{\alpha})$$

Za předpokladu, že v modelu (32) mají již chyby  $e_s$  přibližně normované normální rozdělení je možné vyjádřit logaritmus věrohodnostní funkce ve tvaru

$$\begin{aligned} \ln L &= \sum_{i=1}^N \left\{ (c-1) \cdot \ln(y_i) - \ln[G g(x_i, d)] \right\} - \\ &\quad - \frac{1}{2G^2} \sum_{i=1}^N \left\{ [y_i^{(c)} - f^{(c)}(\underline{x}_i, \underline{\alpha})] / g(x_i, d) \right\}^2 \end{aligned} \quad (33)$$

Pro pevné hodnoty parametrů  $c, d$ , a lze z rovnice (33) nalézt maximálně věrohodný odhad  $G^2$  analyticky

$$\hat{G}^2(c, d, \underline{\alpha}) = N^{-1} \sum_{i=1}^N \left\{ [y_i^{(c)} - f^{(c)}(\underline{x}_i, \underline{\alpha})] / g(x_i, d) \right\}^2$$

Po dosazení do rov. (33) resultuje koncentrovaná věrohodnostní funkce

$$\ln L^* = \sum_{i=1}^N \left\{ (c-1) \ln y_i - \ln [\hat{G}(c, d, \underline{\alpha}) \cdot g(x_i, d)] \right\} - N/2 \quad (34)$$

Maximalizace rovnice (34) se dá provádět pomocí obecných extremizačních programů nebo kombinací metody vážených nejmenších čtverců se systematickými změnami parametrů  $c$ ,  $d$ . Po určení odhadů  $c^*$ ,  $d^*$ ,  $\underline{\alpha}^*$  je vhodné ověřit, zda nelze model vyjádřený rov. (32) zjednodušit.

Pro tyto účely je možné např. využít 100  $(1-\alpha)$  %ní oblasti spolehlivosti pro parametry  $(c, d)$  [5].

Označíme-li  $(c^*, d^*)$  odhady odpovídající maximu věrohodnostní funkce (34) můžeme zkonstruovat oblast spolehlivosti pro  $(c, d)$  hledáním všech dvojic  $(c_0, d_0)$ , pro které platí, že

$$2 [\ln L^*(\hat{c}, \hat{d}) - \ln L^*(c_0, d_0)] \leq \chi^2(2)$$

kde  $\chi^2(2)$  je 100  $(1-\alpha)$  %ní kvantil chí kvadrát rozdělení. V práci [2] je doporučeno použít pro testaci významnosti parametrů  $c$ ,  $d$  speciálního F-testu.

Pro případ  $d = 0$  existuje možnost jednoduššího přibližného odhadu parametru  $c$  (z linearizace obou stran modelu (30)) a testace jeho významnosti. Postačují přitom pouze výsledky nelineární regrese klasickou MNČ.

Transformace typu (30) obyčejně vede také k snížení vnitřní nelinearity regresního modelu.

Pro odstranění (nebo snížení) nelinearity způsobené parametry je buď možné provést reparametrisaci modelu nebo mocninnou transformaci odhadu  $\underline{\alpha}^*$ . V obou případech lze s výhodou kontrolovat snížení nelinearity modelu pomocí vychýlení  $\hat{p}$  definovaného rovnicí (15).