

### 3 STATICKÉ MODELY

#### 3.1. Korelačná a regresná analýza

- ⇒ Vhodná na **vyšetrovanie statických**, ale aj **dynamických systémov**.  
⇒ Skúmame funkčný vzťah, podľa ktorého sa mení **závisle premenná y pri zmenách nezávislých veličín**  $u = (u_1, u_2, \dots, u_n)^T$ .

Predpokladajme, že skutočný funkčný vzťah má nasledovný tvar:

$$y = F(u, \theta) + v$$

t.j. predpokladáme, že **výsledok experimentu**  $y \in \mathbb{R}^1$ , bude funkciou **systematicky pôsobiacich nezávislých veličín**  $u \in \mathbb{R}^n$  (n-faktorový problém), a tiež od **náhodne pôsobiacich faktorov**  $v \in \mathbb{R}^1$ .  $\theta$  je **vektor (neznámych) parametrov** a  $F$  je **funkcia štruktúry** identifikovaného systému.

- ⇒ **Regresná analýza** sa zaoberá **určením štruktúry**, t.j. stanovením, ktoré nezávislé premenné a v akom charaktere väzieb do funkčnej závislosti vstupujú.

**Korelačná analýza** kvantifikuje **intenzitu vplyvov nezávislých premenných** na závislú premennú (**koefficienty**).

- ⇒ **Štruktúru** reálneho systému **presne poznať nebudeme**, preto budeme predpokladať, že je vyjadrená v tvare **mocninového radu**, tvoriaceho tzv. **teoretickú regresnú rovnicu**

$$y = \theta_0 + \sum_{i=1}^n \theta_i u_i + \sum_{i=1}^n \sum_{j=1}^i \theta_{ij} u_i u_j + \dots + v$$

čo je funkcia **nelineárna voči nezávislým premenným** tvoriacim vektor  $u = (u_1, u_2, \dots, u_n)^T$  nazývaný **regresor**, ale **lineárna voči regresným koefficientom**  $\theta_0, \theta_i, \theta_{ij}$ .

- ⇒ **Model** tohto systému nech popisuje funkčná závislosť opäť **v tvare mocninového radu s konečným stupňom s bez náhodnej zložky**

$$\hat{y} = \varphi(u, \hat{\theta})$$

kde  $\hat{\theta}$  je **odhad vektora parametrov**

- ⇒ **n-faktorový polynomiálny model stupňa s**

**Počet neznámych parametrov vrátane absolútneho člena** pre model, ktorý má **n faktorov (= vstupov)** a je vyjadrený **polynómom stupňa s** je

$$1 + k = \binom{n+s}{s} \quad \text{kde} \quad \binom{x}{y} = \left( \frac{x!}{y! (x-y)!} \right) \text{ je kombinačné číslo}$$

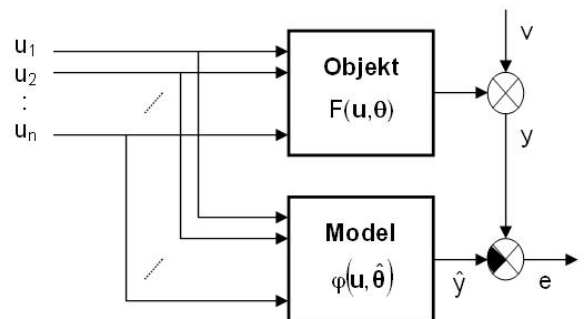
Pri väčšom počte faktorov  $n$  a so vzrastom stupňa polynómu  $s$  počet neznámych parametrov veľmi rýchle narastá, z praktického hľadiska úplne stačí uvažovať  $1 \leq s \leq 3$ , pričom:

- ak  $s = 1$  – **lineárna regresia**

$$\hat{y} = \hat{\theta}_0 + \hat{\theta}_1 u_1 + \dots + \hat{\theta}_n u_n,$$

- ak  $s = 2$  – **kvadratická regresia**

$$\hat{y} = \hat{\theta}_0 + \hat{\theta}_1 u_1 + \dots + \hat{\theta}_n u_n + \hat{\theta}_{n+1} u_1^2 + \dots + \hat{\theta}_{2n} u_n^2 + \hat{\theta}_{2n+1} u_1 u_2 + \dots + \hat{\theta}_k u_{n-1} u_n,$$



- ak  $s = 3$  – **kubická regresia**

$$\hat{y} = \hat{\theta}_0 + \sum_{i=1}^n \hat{\theta}_i u_i + \sum_{i=1}^n \sum_{j=1}^i \hat{\theta}_{ij} u_i u_j + \sum_{i=1}^n \sum_{j=1}^i \sum_{k=1}^j \hat{\theta}_{ijk} u_i u_j u_k,$$

⇒ Pri jednofaktorových modeloch môže byť aj **polynomiálna závislosť**

$$\hat{y} = \hat{\theta}_0 + \hat{\theta}_1 u + \hat{\theta}_2 u^2 \dots + \hat{\theta}_k u^k$$

Na model **lineárny** voči regresným koeficientom môžeme **pretransformovať** aj niektoré **nelineárne závislosti**, napr.

$$z = c \cdot x_1^{\hat{\theta}_1} \cdot x_2^{\hat{\theta}_2} \dots x_n^{\hat{\theta}_n} \rightarrow \hat{y} = \hat{\theta}_0 + \hat{\theta}_1 u_1 + \dots + \hat{\theta}_n u_n$$

$$\text{kde } \hat{y} = \ln z \quad \hat{\theta}_0 = \ln c \quad u_i = \ln x_i \quad i = 1, 2, \dots, n$$

alebo

$$z = c \cdot e^{(\hat{\theta}_1 u_1 + \dots + \hat{\theta}_n u_n)} \rightarrow \hat{y} = \hat{\theta}_0 + \hat{\theta}_1 u_1 + \dots + \hat{\theta}_n u_n$$

$$\text{kde } \hat{y} = \ln z \quad \hat{\theta}_0 = \ln c$$

⇒ Vo **všeobecnosti** možno **všetky** uvedené **modely vyjadriť** v tvare **lineárnej regresie**

$$\hat{y} = \hat{\theta}_0 f_0(u) + \hat{\theta}_1 f_1(u) + \hat{\theta}_2 f_2(u) \dots + \hat{\theta}_k f_k(u)$$

$$\hat{y} = \hat{\theta}^T \mathbf{f}(u)$$

$$\hat{\theta} = (\hat{\theta}_0, \hat{\theta}_1, \dots, \hat{\theta}_k)^T \text{ je vektor odhadovaných parametrov}$$

$$\mathbf{f}(u) = (f_0(u), f_1(u), \dots, f_k(u))^T \text{ je regresný vektor}$$

**Poznámka:** všetky vektory sú stĺpcové, preto sú transponované, keď sú zapísané ako riadkové.

⇒ **Predpokladajme, že uskutočníme celkom**  $N \gg 1+k$  **meraní**, pri každom z nich budeme určovať hodnotu výstupu zo sústavy  $y_i$  aj hodnoty zložiek vektora  $\mathbf{f}(u_i)$

$$\mathbf{f}(u_i) = (f_0(u), f_1(u), \dots, f_k(u))^T = (1, h_{1i}, h_{2i}, \dots, h_{ki})^T = \mathbf{h}_i \quad i = 1, \dots, N$$

**Potom v i-tom experimente vyjadríme hodnotu odchýlky** medzi nameraným výstupom a výstupom z modelu

$$e_i = y_i - \hat{y}_i = y_i - \hat{\theta} \mathbf{h}_i$$

čím sa dostaneme k **preurčenému systému rovníc (t.j. máme viac rovníc ako neznámych)**

$$\begin{aligned} e_1 &= y_1 - \hat{\theta}_0 - \hat{\theta}_1 h_{11} - \hat{\theta}_2 h_{21} - \dots - \hat{\theta}_k h_{k1} \\ e_2 &= y_2 - \hat{\theta}_0 - \hat{\theta}_1 h_{12} - \hat{\theta}_2 h_{22} - \dots - \hat{\theta}_k h_{k2} \\ &\vdots \\ e_N &= y_N - \hat{\theta}_0 - \hat{\theta}_1 h_{1N} - \hat{\theta}_2 h_{2N} - \dots - \hat{\theta}_k h_{kN} \end{aligned}$$

alebo v maticovom tvare

$$\mathbf{e} = \mathbf{y} - \mathbf{H}\hat{\theta} = \mathbf{e}(\hat{\theta}) \quad \text{kde} \quad \mathbf{e} = \mathbf{y} - \hat{\mathbf{y}} = (e_1, e_2, \dots, e_N)^T$$

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

$$\mathbf{H} = \begin{pmatrix} 1 & h_{11} & h_{21} & \dots & h_{k1} \\ 1 & h_{12} & h_{22} & \dots & h_{k2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & h_{1N} & h_{2N} & \dots & h_{kN} \end{pmatrix} = \begin{pmatrix} \mathbf{h}_1^T \\ \mathbf{h}_2^T \\ \vdots \\ \mathbf{h}_N^T \end{pmatrix}$$

⇒ **Vektor neznámych parametrov  $\hat{\theta}$**  z tohto preurčeného systému rovníc budeme určovať **minimalizáciou súčtu štvorcov odchýlok (metóda najmenších štvorcov, MNŠ – ako prvý ju použil Gauss, publikoval ju v r. 1795)**

$$Q(\hat{\theta}) = \frac{1}{2} \sum_{i=1}^N e_i^2 = \frac{1}{2} \mathbf{e}(\hat{\theta})^T \mathbf{e}(\hat{\theta}) = \frac{1}{2} (\mathbf{y} - \mathbf{H}\hat{\theta})^T (\mathbf{y} - \mathbf{H}\hat{\theta}) \quad \text{účelová funkcia}$$

**Optimálnu hodnotu  $\hat{\theta}^*$** , ktorú budeme považovať za **odhad hľadaných parametrov**, určíme z **nulovej hodnoty gradientu účelovej funkcie**

$$\nabla_{\hat{\theta}} Q(\hat{\theta}) \Big|_{\hat{\theta}^*} = \mathbf{H}^T (\mathbf{H}\hat{\theta}^* - \mathbf{y}) = 0$$

$$\mathbf{H}^T \mathbf{H} \hat{\theta}^* = \mathbf{H}^T \mathbf{y} \quad \Leftrightarrow \quad \boxed{\hat{\theta}^* = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{y}} \quad \text{Gaussov vzorec (1809)}$$

Aby tento vektor zabezpečil **minimum účelovej funkcie**, musí byť matica  **$\mathbf{R} = \mathbf{H}^T \mathbf{H}$  pozitívne definitná**, čo je splnené, pretože súčin nenulovej matice s jej transpozíciou je vždy symetrická a pozitívne definitná matica

$$\nabla_{\hat{\theta}} \nabla_{\hat{\theta}} Q(\hat{\theta}) \Big|_{\hat{\theta}^*} = \mathbf{H}^T \mathbf{H} = \mathbf{R} \quad \text{informačná matica}$$

$$\mathbf{P} = \mathbf{R}^{-1} = (\mathbf{H}^T \mathbf{H})^{-1} \quad \text{disperzná matica}$$

$$\mathbf{H}^T \mathbf{H} \hat{\theta}^* = \mathbf{H}^T \mathbf{y} \quad \text{normálna sústava rovníc}$$

 **Príklad** → *priklady\_regresia.pdf + cvičenia*

⇒ **MATLAB**

### 1. Curve Fitting Toolbox

Umožňuje **spracovanie údajov** (vyhladenie, odstránenie niektorých vzoriek), **aproximáciu funkciou** a **porovnanie výsledkov** graficky aj numericky

**Možnosti použitia:**

- GUI – spúšťa sa príkazom *cftool*
- Matlab – príkazy  
`fit(x,y,FT)`

### 2. Basic Fitting GUI

v okne **Figure** vybrať **Tools – Basic Fitting**

⇒ **Vlastnosti odhadu:**

Predpokladajme, že:

1.  **$\text{hod}(\mathbf{H}) = 1 + k \leq N$**  (počet pozorovaní je väčší ako počet neznámych parametrov modelu)
2. **parazitný šum spíňa:**

**stredná hodnota** šumu sa **rovná nule**, t.j. vylučujeme výskyt jednosmernej zložky aj nestacionárnych šumov

$$E\{v_i\} = 0 \quad i = 1, 2, \dots, N$$

jednotlivé **zložky** šumov sú navzájom **nekorelované**

$$E\{v_i v_j\} = 0 \quad i \neq j$$

**disperzia šumov je konštantná**

$$E\{v_i^2\} = \sigma_i^2 = \sigma^2 \quad i = 1, 2, \dots, N$$

**potom kovariančná matica má tvar**

$$\text{cov}(\mathbf{v}) = E\{\mathbf{v}\mathbf{v}^T\} = \sigma^2 \mathbf{I}$$

**Gaussova – Markovova veta:** Ak platia predpoklady 1. a 2., potom  $\hat{\boldsymbol{\theta}}^* = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{y}$  je výdatný nevychýlený odhad (BLUE) .

### BLUE (Best Linear Unbiased Estimator)

- **lineárny** vzhľadom k dátam
- **nevychýlený** (unbiased)
- **výdatný** (best, minimum variance)

#### a) Linearita

$$\hat{\boldsymbol{\theta}}^* = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{y}$$

odhad je **lineárny vzhľadom k nameraným výstupom**

#### b) Nevychýlenosť odhadu parametrov

Pre výstup reálneho systému platí  $\mathbf{y} = \mathbf{H}\boldsymbol{\theta} + \mathbf{v}$

$$\hat{\boldsymbol{\theta}}^* = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{y} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T (\mathbf{H}\boldsymbol{\theta} + \mathbf{v}) = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{H}\boldsymbol{\theta} + (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{v} = \boldsymbol{\theta} + (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{v}$$

$$E\{\hat{\boldsymbol{\theta}}^*\} = \boldsymbol{\theta} + E\{(\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{v}\}$$

$$\begin{aligned} \hat{\boldsymbol{\theta}}^* - \boldsymbol{\theta} &= (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{v} = \\ &= \mathbf{P} \mathbf{H}^T \mathbf{v} \end{aligned}$$

Keďže  $E\{\mathbf{v}\} = \mathbf{0}$ , potom platí  $E\{\hat{\boldsymbol{\theta}}^*\} = \boldsymbol{\theta}$ .

Ak sa parazitný šum  $\mathbf{v}$  riadi rozdelením  $N(0, \sigma^2)$ , potom vektor  $\hat{\boldsymbol{\theta}}^*$  je **nevychýleným odhadom** vektora  $\boldsymbol{\theta}$ .

#### c) Výdatnosť odhadu parametrov

Kovariančná matica odhadu je

$$\text{cov}(\hat{\boldsymbol{\theta}}^*) = E\{(\hat{\boldsymbol{\theta}}^* - \boldsymbol{\theta})(\hat{\boldsymbol{\theta}}^* - \boldsymbol{\theta})^T\} = E\{\mathbf{P} \mathbf{H}^T \mathbf{v} \mathbf{v}^T \mathbf{H} \mathbf{P}\} = \mathbf{P} \mathbf{H}^T E\{\mathbf{v} \mathbf{v}^T\} \mathbf{H} \mathbf{P} = \sigma^2 \mathbf{P} \mathbf{H}^T \mathbf{H} \mathbf{P} = \sigma^2 \mathbf{P} \mathbf{P}^{-1} \mathbf{P} = \sigma^2 \mathbf{P}$$

$$\boxed{\text{cov}(\hat{\boldsymbol{\theta}}^*) = \sigma^2 \mathbf{P}}$$

čo znamená, že **výdatnosť odhadu sa zvyšuje so vzrastom hladiny užitočných signálov  $\mathbf{u}$**  a naopak **klesá so zväčšovaním disperzie šumu  $\sigma^2$** .

Tiež  $\hat{\mathbf{y}}$  je **nevychýleným** odhadom  $\mathbf{y}$ , pričom jeho kovariancia je

$$\boxed{\text{cov}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}^T \mathbf{P} \mathbf{H}}$$

Problém **zabezpečenia výdatnosti hľadaných odhadov** je teda spojený s otázkou **vhodného modelu** pre aproximáciu určovanej funkčnej závislosti a tiež so starostlivou **prípravou experimentu**.

**Experiment** by mal byť zostavený tak, aby

- prinášal **maximum užitočnej informácie**
- bol realizovaný s **minimálnymi nákladmi**
- bol **vyhodnocovaný** metódami **matematickej štatistiky**.

**Objektivita a presnosť** modelu sa zvyšuje **priamo úmerne s počtom experimentov**, avšak **ekonomická stránka** identifikácie smeruje k **menšiemu počtu experimentov**.