

- dotazník - čo by ste zmenili, čo sa vám nepáči
- ANKETA
- projekt - deň pred stretnutím
- stretnutie - online
- asi dnes nestihneme  
⇒ náhradná video-prednáška na stránke

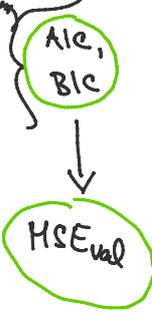
Titulok: kombinatorická selekcia

$$y = X\beta + \varepsilon$$

$$X_{n \times k} = (1_n, X_1, \dots, X_{k-1})$$

→ výber prem. z  $X_1, \dots, X_{k-1}$

- všetky podmnožiny:  $2^{k-1}$
- dopredná selekcia
- spätná sel.
- "dojstá ma" sel.



Online News:

- všetky podmnožiny
  - manuálne (pre iba niektoré  $X_j$ )
  - regsubsets

Namiesto IC: validácia, crossvalidácia

validácia: dáta do 3 častí

- 1.) trénovacia vzorka: na nej odhadujeme  $\beta$  (fitujeme model) →  $\hat{\beta}$  (train)
- 2.) validácia vzorka: výber modelu (val)
- 3.) testovacia (test): posúdenie kvality finálneho modelu

→  $\hat{y}_{i \in val}$ : modelom predikujeme  $\hat{y}_i = x_i^T \hat{\beta}$  ↔ porovnáme s realitou:  $y_i$

$$MSE_{val} = \frac{1}{n_{val}} \sum_{i \in val} (y_i - x_i^T \hat{\beta})^2 = \frac{1}{n_{val}} \|y_{val} - X_{val} \hat{\beta}\|^2 \rightarrow \min \uparrow \text{výberom modelu}$$

↑  
veľkosť val

↳ validačná chyba

- namiesto modelu s najnižšou AIC/BIC skúsime model s najnižou validačnou chybou

- všetky podmnožiny
- akčné met.
  - dopr.
  - spät.
  - objektívna selekcia

Formule:

•  $MSE_{train} = \|y_{train} - X_{train} \beta\|^2 \rightarrow \min_{\beta} \Rightarrow \hat{\beta}$  : trénovacia chyba

•  $MSE_{val} = \|y_{val} - X_{val} \hat{\beta}\|^2 \rightarrow \min_M \Rightarrow \text{Model}$  : val. chyba

•  $MSE_{test} = \|y_{test} - X_{test} \hat{\beta}\|^2 \Rightarrow \text{kvalita modelu}$  : testovacia chyba

rozdeľenie dát: napr. 80% - 10% - 10%  
 50% - 25% - 25%  
 (train - val - test)

Krosvalidácia:

maľo dát => nechceme odčleniť val. vzorku

=> validácia tržkou pomocou tréningovej vzorky: krosvalidácia

- rozdelíme train na k častí (napr. k=5) rovnakej veľkosti:  $y^{(1)}, \dots, y^{(k)}$   
 $X^{(1)}, \dots, X^{(k)}$

$$\begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(k)} \end{pmatrix} = \begin{pmatrix} X^{(1)} \\ \vdots \\ X^{(k)} \end{pmatrix} \beta + \begin{pmatrix} \epsilon^{(1)} \\ \vdots \\ \epsilon^{(k)} \end{pmatrix}$$

$\forall i = 1, \dots, k$ : vezmeme všetky časti okrem i-tej (4/5 dát)

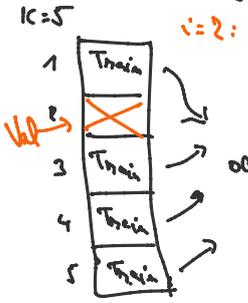
-> na nich odhadneme model

-> predikujeme  $y^{(i)}$  pomocou toho modelu => MSE = chyba:

$$\text{chyba}_i = \frac{1}{n_i} \| y^{(i)} - X^{(i)} \cdot \hat{\beta}^{(i)} \|^2, \quad \hat{\beta}^{(i)} = (X_{-i}^T X_{-i})^{-1} X_{-i}^T y_{-i}$$

$\uparrow$  veľkosť i-tej časti                       $\uparrow$   $\uparrow$   $\uparrow$   $\uparrow$   
 všetky okrem i-tej

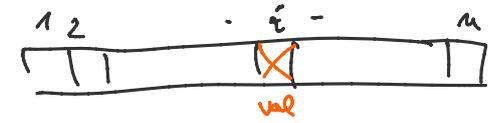
$$\Rightarrow \text{val. chyba} = \frac{1}{k} \sum_{i=1}^k \text{chyba}_i = \frac{1}{k} \sum_{i=1}^k \frac{1}{n_i} \| y^{(i)} - X^{(i)} \hat{\beta}^{(i)} \|^2$$



odhadneme =>  $\hat{\beta}$  -> potom  $y^{(2)}$  vs.  $X^{(2)} \hat{\beta}$

vyber model

• volba k: napr. k=5, k=10, k=n



↓ leave-one-out crossvalidation

4 RKO - Forest fíns:

|       |   |   |   |     |
|-------|---|---|---|-----|
| X \ y | 1 | 2 | 3 | ... |
| 1     |   |   |   |     |
| 2     |   |   |   |     |
| 3     |   |   |   |     |

validácia, krosvalidácia - manuálne

↳ n RKO: balice & casot - ...

variance, covariance - matice

↳ n řek: balice & caset - ista forma spaitaj seledca

- funguje trochu imat - na řadej z k  
sice nejako hodnoti premene, a potom  
odhaduje premene's najbolsim celkem  
skore

### Penalitzovane' regresie

$$y = X\beta + \varepsilon, \quad X = (1 \ x_{11} \dots \ x_{n1}) = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix}$$

- Lamb. seledca - pre řadku prem.  $\beta$ : v modeli  $\left\{ \begin{array}{l} \text{ano} \\ \text{nie} \end{array} \right.$

- veliki skokov: nestabilni voči perturbacii dat

inji pristup: model so vřetkyimi prem., ale zmenšime hodnoty  $\hat{\beta}$  → spojito, nie iba bud  
nula, alebo nenula

→ & LS pridáme penaltu:

$$\hat{\beta} \in \operatorname{argmin} \left( \|y - X\beta\|^2 + \lambda \cdot p(\beta) \right)$$

penalitzacny' parameter  $\lambda \geq 0$        $\uparrow$        $\uparrow$        $\uparrow$    
 rodena' penalta (evyčajne norma)

•  $\lambda = 0 \Rightarrow$  OLS

•  $\lambda \rightarrow \infty \Rightarrow \hat{\beta} = \vec{0}$

•  $\beta_0$  meliceme skencovat: najprv vřetko vycentrujeme:  $y \mapsto y - \bar{y}$   
 $x_j \mapsto x_j - \bar{x}_j \quad \forall j$

: potom už nepotrebujem  $\beta_0$

### 1. Ridge regression - Hehenova' regresia

$$p(\beta) = \|\beta\|_2^2 \Rightarrow \hat{\beta}_{\text{ridge}} \in \operatorname{argmin}_{\beta} \|y - X\beta\|^2 + \lambda \|\beta\|^2 = \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \cdot \sum_{j=1}^p \beta_j^2$$

$$\begin{aligned} &= (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta \\ &= y^T y - 2y^T X\beta + \underbrace{\beta^T X^T X \beta + \lambda \beta^T \beta}_{q(\beta)} \rightarrow \min_{\beta} \end{aligned}$$

$$\frac{\partial}{\partial \beta} y^T y = 0$$

$$\frac{\partial}{\partial \beta} \beta^T A \beta = 2A\beta$$

$$\frac{\partial q}{\partial \beta} = 0 - 2(y^T X)^T + 2X^T X \beta + 2\lambda \beta = 0 \quad /: 2$$

... - vT.

$$\frac{\partial \mathcal{L}}{\partial \beta} = 0 - 2(y^T X) + 2X^T X \beta + 2\lambda \beta = 0 \quad /: 2$$

$$(X^T X + \lambda I) \beta = X^T y$$

$$\hat{\beta}_{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y$$

$$\begin{pmatrix} \lambda & & 0 \\ & \ddots & \\ 0 & & \lambda \end{pmatrix}$$

- pôvodná motivácia:  $X^T X + \lambda I$  je lepšie podmienená ako  $X^T X$  (lepšie sa invertuje)  
(Hoerl & Kennard, 1970)

$\hat{\beta}_{\text{ridge}}$  ako riešenie úlohy s obmedzením:

$$\begin{cases} \min \|y - X\beta\|_2^2 \\ \|\beta\|_2^2 \leq c \end{cases} \quad (U)$$

Lagrange:  $\mathcal{L} = \|y - X\beta\|_2^2 + \mu(\beta^T \beta - c)$

riešenie (U):  $\frac{\partial \mathcal{L}}{\partial \beta} = 0$

to je presne  $f(\beta)$ , akurát namiesto  $\mu$  máme  $\lambda$  (a je tu navyše  $-\mu c$ , čo neplyne na riešenie)

$\Rightarrow \hat{\beta}_{\text{ridge}}$  je riešením (U)

$\Rightarrow \hat{\beta}_{\text{ridge}}$ : minimalizácia RSS, ale máme obmedzenie na to, ako veľký môže  $\hat{\beta}$  byť

