

• na stránke - skúškové otázky

dotaz: extrakcia premenných \rightarrow skonstruovať q nových Y z pôvodných q X-ov

(3) Metódy selekcie premenných (feature selection)

- už nie konštrukcia nových, ale: vyber podmnožiny $\{1, \dots, q\}$

\rightarrow pre potreby najlepšího modelu

tu: lineárna regresia: $y = X\beta + \epsilon$
 $n \times 1$ $n \times q$ $q \times 1$ $n \times 1$

i-ty riadok: $y_i = x_i^T \beta + \epsilon_i$

$y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$, $\epsilon = \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix}$

$X = \begin{pmatrix} 1 & \dots & \\ \vdots & & \\ 1 & \dots & \end{pmatrix} = \begin{pmatrix} - & x_1^T & - \\ & \vdots & \\ - & x_n^T & - \end{pmatrix}$
 $n \times q$ $n \times q$

$x_i^T = (1 \ x_{i1} \dots \ x_{iq})$: riadky
: porovnaniam

X_j : stĺpce, premenné \rightarrow vyberáme
 $j = 1, \dots, q-1$

$= \begin{pmatrix} 1_n, X_1, \dots, X_{q-1} \end{pmatrix}$
 $n \times q$

Kombinatorická selekcia:

- vyber $q \leq q$ premenných, kt. najlepšie popisujú dáta

\rightarrow intercept a $(q-1)$ X-ov

prečo nie všetky?

kvalita predikcie: $x_{new} \in \mathbb{R}^q$: nezáhodne!

y_{new} : náhodná hodnota z rozdelenia

$y = f(x, \beta) + \epsilon$

netušíme; dúfame, že je to $x^T \beta$

expected prediction error:

$\Rightarrow EPE = E \left[(y_{new} - x_{new}^T \hat{\beta})^2 \right] = \text{var} + (\text{Bias})^2$

$\text{var}(y_{new} - x_{new}^T \hat{\beta}) = E \left[(y_{new} - x_{new}^T \hat{\beta})^2 \right] - \left(E \left[y_{new} - x_{new}^T \hat{\beta} \right] \right)^2$
 rozptyl: var EPE Bias: výdychos

• ak máme veľa X_j v modeli $\Rightarrow x^T \beta$ dobre popisuje $f(x, \beta)$; málo Bias

$\hookrightarrow \Rightarrow$ typický je ale utedy veľká variácia predikcií:

↳ ⇒ typický je ale vždy vyšší variace predikci: vyšší Var

→ celková chyba predikce $\left\{ \begin{array}{l} \text{výchylka: } (\text{Bias}(x^T \hat{\beta}))^2 \\ \text{rozptylnost: } \text{var}(y_{\text{new}} - x_{\text{new}}^T \hat{\beta}) \end{array} \right.$

- bias-variance trade-off

Prk: umělý příklad: $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ (M1)

$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \varepsilon_i$ (M2)

$y_i = \beta_0 + \dots + \beta_k x_i^k + \varepsilon_i$ (Mk)

- ako merať, či sú premenné dobré?

$RSS = \|y - X\beta\|^2 = \sum (y_i - x_i^T \beta)^2 \rightarrow \min_{\beta}$: "chyba" regresie : súčet štvorcov rezidui

• pridávanie premenných \Rightarrow znižovanie RSS : čím viac premenných data X

→ pre modely s rovnakým # prem.: môžeme porovnať pomocou RSS - ten s nižším je lepší

→ pre rôzne # prem. (q) → iné kritérium

↳ s penalizáciou za # premenných

↳ informačné kritéria (information criteria) (IC)

IC → min

• $AIC = \ln\left(\frac{RSS}{n}\right) + \frac{2q}{n}$: Akaike inf. crit. (An inf. crit.)

• $BIC = \ln\left(\frac{RSS}{n}\right) + \frac{q}{n} \ln(n)$: Bayes inf. crit. (Schwarz inf. crit.)

• $C_p = \frac{RSS}{\hat{\sigma}^2} - n + 2q$: Mallows's C_p ↳ tiež: SIC

↑

↳ v orig. článku $P =$ monotónne indexové $\{ \xi_0, \dots, \xi_{17} \}$

št. sme vybrali do modelu

$$\hat{\sigma}^2 = S_{Full}^2 = \frac{RSS_{Full}}{n-k}$$

↳ asymptoticky je to ekvivalentné s AIC

↳ z plného modelu

všeobecne: $AIC = -2 \log L_k + 2q$

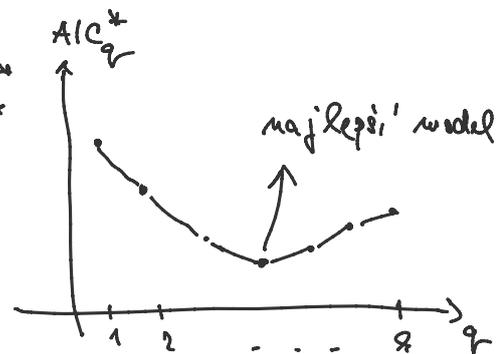
všeobecné: $AIC = -2 \log L(\hat{\theta}) + 2q$
 $BIC = -2 \log L(\hat{\theta}) + q \log(n)$

Āko najst model s min IC?

i) všetky podmnožiny: všetky komb. perm. z $\{1, \dots, p-1\} \Rightarrow 2^{p-1}$ modelov

- M_1^* : najlepší model pre $q=1$ (len intercept)
 - M_2^* : " " " " $q=2$ (intercept + jeden člen)
 - \vdots
 - M_p^* : " " " " $q=p$ (plný model)
- } Stačí RSS

- na základe IC: vyberiem najlepší z M_1^*, \dots, M_p^*



ii) krokové metody (stepwise selection)

• 2^p je príliš veľa pre veľké "p"

→ vyskúšať iba niektoré

- dopredná selekcia:

- začnem s úbohým modelom ($y_i = \beta_0 + \epsilon_i$)
- v každom kroku pridám prem., kt. spôsobí najväčší pokles IC
- opakujem až kým sa nestane: prídanie ľub. prem. spôsobí nárast IC

(- najviac 2^p modelov) → posledný model = výsledok selekcie

- spätná selekcia: začnem s plným

(backward): v každom kroku odoberiem tu, ktorá spôsobí najväčší pokles IC

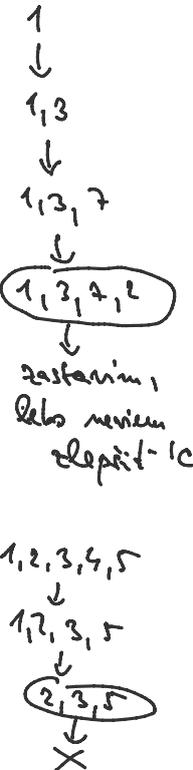
: ak kým sa nestane: odobratie ľub. prem. zvýši IC

- objektívna selekcia: môžem pridať aj odobrať

Poznámky:

- krokové metody - rýchlejšie
- nemusia najst min IC

- klasická inferencia (testy, int. spoľahlivosti) nie je relevantná po použití selekcie prem.
 → výberom premenných sa de facto ...



- klasická inferencia (testy, int. spolehlivost) nie je relevantná po použití selekcie prem.
 - výberom premenných sa deformuje rozdelenie
- ak máme skupinu prem., ktoré nedáme oddeliť (napr. kategorická ⇒ dummy), tak ich pridáme / odhadujeme naset (a formu prispôbime výpočty)
- 3 triky pre rýchle počítanie selekcie - netreba vždy znovu počítať β , RSS, ...

višobne: veľa modelov môže byť dost dobrých

⇒ tu vybraný by sa nemal používať na príliš silnú interpretáciu (na predikciu je ale vhodný)

Kvalita modelu (pre predikciu)

• aký je dobrý model na predikciu: EPE = ?

→ vieme ju odhadnúť: pomocou testovacej vzorky

- na zač.: rozdělíme dáta na 2 časti:

- trénovacia (Train): odhadujeme modely, vyberáme najlepší
- testovacia (Test): posúdime kvalitu vybraného modelu

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

\downarrow \downarrow
 $\begin{pmatrix} y_{\text{train}} \\ y_{\text{test}} \end{pmatrix}$ $\begin{pmatrix} X_{\text{train}} \\ X_{\text{test}} \end{pmatrix}$

$y_{\text{train}} | X_{\text{train}} \rightarrow$ finálny model $\Rightarrow \hat{\beta}$

testovacia chyba

$$MSE_{\text{test}} = \frac{1}{n_{\text{test}}} \sum_{i \in \text{Test}} (y_i - x_i^T \hat{\beta})^2 = \frac{1}{n_{\text{test}}} \|y_{\text{test}} - X_{\text{test}} \hat{\beta}\|^2$$

mean squared error veľkosť test. vzorky

: priemerná chyba predikcia

→ môžeme porovnať s inými modelmi, určiť relatívnu chybu predikcie

- napr.: trénovacie : 90% dát
- : testovacie : 10% dát

Rko: Online news - články + mashable.com

- dáta o nich

- modelujem # zdieľaných článkov (shares)

a) iba niektoré premenné → výber i vyskúšanie všetkých podmnožín (manuálne)

→ R-ová funkcia: regsubsets

~~...função~~ função : `regsubsets`