

to je data vzhľadom na fit / prediction dilemma

(93)

- najmä efekt nastáva takto pri

$$\hat{\beta}_i - \hat{\beta}_{(-i)} = (X^T X)^{-1} x_i \hat{e}_{(-i)}$$

rozhodnúť musíme či je „malý“ počet je „býť“ dobrý, ale musíme
byť „veľký“ počet je tiež veľký.

Myšlienka

- je pre porovnanie modelov možno dva rôzne vzťahy (x, y) a (x', y')
vzhľadom k rôznym závislostiam
- vzhľadom na to, že disponujeme originálnymi dátami
- budú mať rôzne vplyvy na i -tú hodnotu x na model
- ak máme veľké k_{ii} indikujú, že i -tá pos. má veľký vplyv
a veľké reziduálne hodnoty môžu naznačovať model
- myšlienka, že vedieť, keď kombinujeme tieto dva faktory
- používame pri výpočte a PRESS reziduál, kde budeme sledovať
je veľký vplyv na i -tú pozorovanú na $\hat{\beta} - \hat{\beta}_i$

DFBETAS

$\hat{\beta} - \hat{\beta}_{(-i)}$ má vplyv na i -tú pos. na odhad β
(budeme sledovať pri rôznych analýzách)

príj.
$$\hat{\beta} - \hat{\beta}_{(-i)} = (X^T X)^{-1} x_i \frac{\hat{e}_i}{1 - k_{ii}}$$

a) vplyv na i -tú pozorovanú na β_j

$$\hat{\beta}_j - \hat{\beta}_{(-i)j} = \frac{r_{ji} \hat{e}_i}{1 - k_{ii}}, \text{ kde } r_{ji} \text{ je } (j, i) \text{ prvok matice}$$
$$R = (X^T X)^{-1} X^T$$

$\hat{\beta}_j$ - le' poz. bodene porovnaní se influenční, ^{na β_j} pokud

(99)

$\hat{\beta}_j - \hat{\beta}_{(-j)}$ bude velké

protože $\hat{\beta}_j$ je m. vel., "velké" bychom měli měřit relativně

vzhledem k s.d. ($\hat{\beta}_j$) což je $\sigma \sqrt{w_j}$, $w_j = (X^T X)^{-1}_{jj}$

pokud je odhadneme pomocí $\hat{\sigma}_{(-j)} \sqrt{w_j}$, dostaneme definici

$$\underline{\underline{DFBETAS_{j,i} = \frac{\hat{\beta}_j - \hat{\beta}_{(-j)}}{\hat{\sigma}_{(-j)} \sqrt{w_j}} = \frac{r_{ji} \hat{e}_i}{\sqrt{w_j} \hat{\sigma}_{(-j)} (1 - h_{ii})} = \frac{r_{ji}}{\sqrt{w_j}} \frac{\hat{e}_i}{\sqrt{1 - h_{ii}}}}}$$

kde \hat{e}_i je est. studentizované residuum

- kombinuje obě velké residua \hat{e}_i a velké h_{ii}

- ~~možno~~ je to měřítko pro limitní hodnoty:

$\hat{\beta}_j$ - le' porovnání je porovnáním se influenční na odhad β_j ,

pokud $|DFBETAS_{j,i}| > \frac{2}{\sqrt{n}}$

- máme $(m+1) \times m$ hodnot pro rovnání, zjednodušíme

b) vliv $\hat{\beta}_j$ - le' porovnání na celý vektor $\hat{\beta}$

posvítit náhled pomocí na vektor $\hat{\beta} - \hat{\beta}_{(-j)}$

Což rovnice

$$D_{ji} = \frac{(\hat{\beta} - \hat{\beta}_{(-j)})^T M (\hat{\beta} - \hat{\beta}_{(-j)})}{(m+1) c}$$

kde M je PD matice
a c je norm. konst.

nepřesáhneji vektor je $M = X^T X$ a $c = \sigma_m^2$

Cožom vztahemost:

$$\underline{\underline{D_{ji} = \frac{(\hat{\beta} - \hat{\beta}_{(-j)})^T X^T X (\hat{\beta} - \hat{\beta}_{(-j)})}{(m+1) \sigma_m^2}}}$$

$$D_i = \frac{1}{(m+1)\sigma_m^2} \left(\frac{\hat{e}_i}{1-h_{ii}} \right)^2 \underbrace{x_i^T (x^T x)^{-1} x_i}_{=h_{ii}} = \frac{1}{m+1} \frac{h_{ii}}{1-h_{ii}} \underbrace{\frac{\hat{e}_i^2}{\sigma_m^2 (1-h_{ii})}}_{(\hat{r}_i)^2}$$

výpočetná formule

$$\underline{\underline{D_i = \frac{\hat{r}_i^2}{m+1} \left(\frac{h_{ii}}{1-h_{ii}} \right)}}$$

Pozn: 100(1- α)% simultánná IS pre β je

$$C(\alpha) = \left\{ \beta \mid \frac{(\hat{\beta} - \beta)^T x^T x (\hat{\beta} - \beta)}{(m+1)\sigma_m^2} \leq F_{1-\alpha}(m+1, m-m-1) \right\}$$

$$\text{km. } \beta_{(i)} \in C(\alpha) \Leftrightarrow D_i \leq F_{1-\alpha}(m+1, m-m-1)$$

to je metrika pre RULE OF THUMB:

i -le' poz. je influenčný, keďže $D_i > F_{\frac{\alpha}{2}}(m+1, m-m-1)$

(pre veľkým m, n $F_{\frac{\alpha}{2}} \approx 1$, jednoduchšie povedať $D_i > 1$)

Pozn: Takže platí

$$D_i = \frac{(\hat{\beta} - \hat{\beta}_{(-i)})^T (\hat{\beta} - \hat{\beta}_{(-i)})}{(m+1)\sigma_m^2}$$

km. čo' se dá' získať jeho miera
influeny na celkovú ~~predikciu~~ predikciu
~~chýbu~~

DIFFITS: aká i -le' prispieva na $\hat{\sigma}_i$

$$\underline{\underline{DIFFITS_i = \frac{\hat{\sigma}_i - \hat{\sigma}_{(-i)}}{\hat{\sigma}_{(-i)} \sqrt{h_{ii}}} = \dots = \hat{r}_i \sqrt{\frac{h_{ii}}{1-h_{ii}}}}}$$

odhad s.d. \hat{r}_i

RULE OF THUMB: i -le' poz. je influenčný, keďže $|DIFFITS_i| > 3 \sqrt{\frac{m+1}{n-m-1}}$

POZV: min influence ~ R:

(36)

DFBETAS - `dfbetas()` DFFITS - `dfbifs()`

~~cooks~~ Cookova vzdálenost D_i - `cooks.distance()`

Leverage h_{ii} - `hatvalues()`

všechny funkce influence.measures()

ma' navíc (conversion ratio)

pravidla:

je-li pro je influenční pokud:

$$|DFBETAS| > 1, \quad |DFFITS| > 3 \sqrt{\frac{m+1}{n-m-1}}$$

$$D_i > F_{0,5}(m+1, n-m-1), \quad h_{ii} > 3 \frac{m+1}{n}$$

TRANSFORMACE

- pokud není splněn některý z předpokladů modelu
linearity, normality chyb, homoskedasticita
ještě lze zkusit je pokud se transformací nepřipraví,
~~ale~~ aby transformovaný model byl předpoklady alespoň
"přibližně" splněn

Transformace vyřešení problémů

hledáme funkci $h(\cdot)$ tak, aby model $y_i^* = h(y_i) = \beta_0 + \sum_{j=1}^m x_{ij} \beta_j + \epsilon_i$
splňoval předpoklady

3 hlavní důvody pro transformaci Y :

- 1) transformace šlých měření tak, aby ~~odstranili~~ pokrývali celý R ,
což může odstranit problémy s podmínkami na β

např. studie dopravy plic (FEV data, $FEV > 0$)

- čtěl: logem, aby model reprodukoval rozpočet hodnoty
(\Rightarrow restriktivní podmínky na par. β)

na stejné modelování $y^* = \log FEV$

pohled na jiné počty a 0 je možná hodnota, která se používá

$$y^* = \log(y+1) \quad \text{nebo} \quad \text{obecně} \quad y^* = \log(y+c)$$

2) transformace y , aby její rozdělení bylo "více" normální

- typický to momentálně pohled se ustálil rozdělení hodnot y více symetrické

často se setkáváme s rozděleními vychýlenými spravo

(obvykle se to stane, pokud měříme nějakou fyz. veličinu, která může nabývat pouze kladných hodnot)

transformace $y^* = \log y$ nebo $y^* = y^2$, $2 < 1$ budou redukovat toto vychýlení

typický postup: začít s hodnotou 2 blízkou 1, pak postupně hodnotu 2 dohodnout menší dosáhnout "přibližně" symetrické rozdělení

3) může nejzákladnější motivace je pohled na dosahovat konstantní rozptyl přes všechnu pozorování

- např. pro fyzikální veličinu s kladnými hodnotami se často stane, že rozptyl bude malý pro $y \approx 0$ a větší pro y větší (ne jen se zvyšuje, se stane hodnot y je menší na blízké hodnoty)

říká se tomu positive mean - variance relationship

- nepřímá měření kladných veličin se také často vyskytují

forma' boljicima variace

$$CV(Y) = \frac{p.d. Y}{E(Y)}$$

(98)

čisto logična in harkomhu' mezi p'ugab' no' p.d.

variabilnih vs. p'ogibajene relativno sp'ite no' absolutno

matematično se razume, če $\text{Var } Y = \phi E(Y)^2 = \phi y^2$ pri n'ajbolj ϕ

- pri odločanju' o'lohu mezi $E(Y)$ in $\text{Var } Y$ se čisto p'ouživajo

merilne' transformace $Y^* = Y^2$ (pri $Y > 0$)

Transformace: $\leftarrow \dots Y^3 \quad Y^2 \quad Y \quad \sqrt{Y} \quad \log Y \quad \frac{1}{\sqrt{Y}} \quad \frac{1}{Y} \quad \frac{1}{Y^2} \dots \rightarrow$

Priloga 2: $\begin{matrix} & 3 & 2 & 1 & \frac{1}{2} & 0 & -\frac{1}{2} & -1 & -2 \end{matrix}$

• p'obud $\text{Var } Y$ k'len'
in p'obud $E Y$

• p'obud $\text{Var } Y$ kot a p'obud $E Y$

OBEČKE: p'edp'ohladajene v'el' $\text{Var } Y = \phi V(Y)$

a n'usajene transformace $Y^* = h(Y)$

Tajlogično razvoj 1. n'adu funkcije $h(Y)$ in k'otično y

$$Y^* = h(Y) \approx h(y) + h'(y)(Y - y)$$

in čisto p'ogibajene, če $\text{Var } Y^* \approx (h'(y))^2 \cdot \text{Var } Y$

transformace $Y^* = h(Y)$ k'otično k'otično p'obudajene stabilizirajene razp'ede,

p'obud $h'(y)$ in p'obudajene imen' $(\text{Var } Y)^{-\frac{1}{2}} = V(y)$

• p'obud $V(y) = y^2 \Rightarrow$ transformace stabilizirajene razp'ede in $\log(Y) = h(Y)$
p'obud $h'(y) = \frac{1}{y}$

• p'obud $V(y) = y \Rightarrow$ stab. transformace in $h(Y) = \sqrt{Y}$,
p'obud $h'(y) = \frac{1}{2\sqrt{y}}$

$$\left(h(y) = \int \frac{dy}{V(y)} \right)$$

- ani nejvíc užitečnou transformací je $y^* = \log(y)$
jedním z důvodů je i dobrá interpretovatelnost par. β

pozn (interpretace parametrů LM):

a) klasický LM: $Y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m$

jednotková změna proměnné $x_j \Rightarrow$ změna EY o β_j jednotek
(pro ostatních proměnných stejných)

$$\left(\begin{array}{ccc} X = (1, x_1, \dots, x_m) & , & X_{\text{new}} = (1, x_1, \dots, x_j+1, \dots, x_m) \\ \downarrow & & \downarrow \\ EY & & EY_{\text{new}} \\ EY_{\text{new}} - EY = \beta_j \end{array} \right)$$

b) LM pro $\log Y$: $\log Y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m + \varepsilon$
 $\varepsilon \sim N(0, \sigma^2)$

předpoklad je to stejný model, proměnná ε , $\bar{\varepsilon}$

$\log Y \sim N(\mu, \sigma^2) \Rightarrow Y \sim LN(\mu, \sigma^2)$ a tedy
 $EY = e^{\mu + \frac{\sigma^2}{2}}$

předikce pro $E \log Y$ je $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_m x_m$
předikce pro EY bude $e^{\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_m x_m + \frac{\hat{\sigma}^2}{2}}$

uvážejme opět jednotkovou změnu x_j ($x_j \rightarrow x_j + 1$)

$$\frac{EY_{\text{new}}}{EY} = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_j x_j + \hat{\beta}_j + \dots + \hat{\beta}_m x_m + \frac{\hat{\sigma}^2}{2}}}{e^{\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_m x_m + \frac{\hat{\sigma}^2}{2}}} = e^{\hat{\beta}_j}$$

jednotková změna proměnné $x_j \Rightarrow$ multiplikativní změna EY o $e^{\hat{\beta}_j}$ -krát

jinak řečeno: $100(e^{\hat{\beta}_j} - 1)$ je procentuální změna EY způsobená jednotkovou změnou x_j