

- pokud budeme chtít model použít nejen k vysvětlení vztahu mezi proměnnými, ale také pro predikci, hodila by se míra vyjadřující, jak dobře model predikuje
- šlo by použít IS nebo IP, museli bychom předem znát body, ve kterých chceme predikovat
- nejjednodušší přístup, jak měřit prediktivní přesnost modelu, by byl analýza reziduí pro predikce hodnot v nových  $x$   
obecně ale nemáme data  $y$  v těchto bodech
- jedna možnost je použít data, která máme k dispozici  
**postup:** vynecháme jedno pozorování, naladíme model bez tohoto pozorování a porovnáme predikovanou a pozorovanou hodnotu pro vynechané pozorování
- předp., že vynecháme  $i$ -té pozorování a označme

$\hat{\beta}_{(-i)}$  - odhad  $\beta$  v modelu s vynechaným  $i$ -tým pozorováním ( $M_{(-i)}$ )

$\hat{y}_{(-i)}$  - predikovanou hodnotu modelem  $M_{(-i)}$  v bodě  $\mathbf{x}_i^T$ , tzn.  $\hat{y}_{(-i)} = \mathbf{x}_i^T \hat{\beta}_{(-i)}$   
potom

$$\hat{e}_{(-i)} = y_i - \hat{y}_{(-i)}, \quad i = 1, \dots, n,$$

nazýváme  $i$ -té PRESS reziduum

## Poznámka 4.3

Otázka je, jak počítat  $\hat{e}_{(-i)}$ ,  $i = 1, \dots, n$

- pro  $n$  velké, se to zdá náročný problém, pro každé  $i$  je třeba naladit nový model
- naštěstí to nebude nutné, ukážeme totiž  $\hat{e}_{(-i)} = \frac{\hat{e}_i}{1 - h_{ii}}$

Označme  $\mathbf{x}_i^T$   $i$ -tý řádek matice  $\mathbf{X}$ ,  $\mathbf{X}_{(-i)}$  matici  $\mathbf{X}$  bez  $i$ -tého řádku a  $h_{ii} = \mathbf{H}_{ii}$ .

## VĚTA 4.1

Jestliže  $h_{ii} \neq 1$ , potom  $(\mathbf{X}_{(-i)}^T \mathbf{X}_{(-i)})^{-1} = (\mathbf{X}^T \mathbf{X})^{-1} + \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1}}{1 - h_{ii}}$ .

Důkaz.

---

### Věta z LA: (Sherman-Morrison-Woodbury)

Necht  $\mathbf{A}$  je  $n \times n$  invertibilní matice a nechť  $\mathbf{z}$  je  $n \times 1$  sloupcový vektor. Jestliže  $\mathbf{z}^T \mathbf{A}^{-1} \mathbf{z} \neq 1$ , potom matice  $\mathbf{B} = \mathbf{A} - \mathbf{z}\mathbf{z}^T$  je invertibilní a platí

$$\mathbf{B}^{-1} = \mathbf{A}^{-1} + \frac{\mathbf{A}^{-1} \mathbf{z} \mathbf{z}^T \mathbf{A}^{-1}}{1 - \mathbf{z}^T \mathbf{A}^{-1} \mathbf{z}}.$$

---

## VĚTA 4.2

Nechť  $\widehat{e}_{(-i)}$  je  $i$ -té PRESS reziduum. Potom  $\widehat{e}_{(-i)} = \frac{\widehat{e}_i}{1 - h_{ii}}, \quad i = 1, \dots, n.$

Důkaz.

### VĚTA 4.3

1) Nechť  $\widehat{\beta}_{(-i)}$  značí LSE parametru  $\beta$  v modelu bez  $i$ -tého pozorování. Potom platí

$$\widehat{\beta} - \widehat{\beta}_{(-i)} = \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \widehat{e}_i}{1 - h_{ii}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \widehat{e}_{(-i)}.$$

2) Pro součet residiálních čtverců  $SSE_{(-i)}$  v modelu bez  $i$ -tého pozorování platí

$$SSE_{(-i)} = \sum_{j=1}^n \widehat{e}_j^2 - \frac{\widehat{e}_i^2}{1 - h_{ii}}.$$

---

Důkaz.

## Důsledek 4.1

V modelu  $(**)$  s  $m + 1$  parametry  $\beta$  a bez  $i$ -tého pozorování platí

$$E[SSE_{(-i)}] = (n - m - 2)\sigma^2, \quad \text{to znamená} \quad \hat{\sigma}_{(-i)}^2 = \frac{SSE_{(-i)}}{n - m - 2} \quad \text{je nestranný odhad } \sigma^2.$$

Dále pak

$$\hat{\sigma}_{(-i)}^2 = \frac{(1 - h_{ii})(n - m - 1)s_n^2 - \hat{e}_i^2}{(1 - h_{ii})(n - m - 2)} = \frac{1}{n - m - 2} \left( SSE - \frac{\hat{e}_i^2}{1 - h_{ii}} \right),$$

kde  $s_n^2 = \frac{1}{n-m-1} SSE$  (pro plný model).

---

Důkaz.

## Poznámka 4.4

- dá se ukázat, že  $SSE_{(-i)}$  a  $\hat{e}_i$  jsou nezávislé náhodné veličiny
- protože  $\frac{SSE_{(-i)}}{\sigma^2} \sim \chi^2(n-m-2)$  a  $\frac{\hat{e}_i}{\sigma\sqrt{1-h_{ii}}} \sim N(0,1)$
- dostaneme  $\frac{\hat{e}_i}{\hat{\sigma}_{(-i)}\sqrt{1-h_{ii}}} \sim t(n-m-2)$

## Tvrzení 4.1

Uvažujme model (\*\*), kde  $h(\mathbf{X}) = m+1$  a  $\mathbf{e} \sim N_n(0, \sigma^2 \mathbf{I}_n)$ . Nechť pro  $i \in \hat{n}$  platí, že  $h_{ii} \neq 1$ .

Potom  $i$ -té (externě) studentizované reziduum

$$\hat{t}_i \sim t(n-m-2).$$

## Poznámka 4.5

- $\hat{t}_i$  lze použít pro test hypotézy, zda je  $i$ -té pozorování odlehlé (outlier), tedy

$$H_0 : i\text{-té pozorování není odlehlé v modelu } M \quad \times \quad H_1 : i\text{-té pozorování je odlehlé v } M,$$

- kde **odlehlé** značí odlehlé vzhledem k  $M$ :  $\mathbf{Y} \sim N_n(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n)$ :
  - a) střední hodnota  $i$ -tého pozorování se nerovná té dané modelem,
  - b) pozorovaná hodnota  $Y_i$  je neobvyklá za platnosti  $M$ .
- $H_0$  zamítneme, pokud  $|\hat{t}_i| > t_{1-\frac{\alpha}{2}}(n-m-2)$
- pokud test použijeme na všechna pozorování, je potřeba aplikovat nějakou korekci na vícenásobné testování, např. Bonferroni

#### Poznámka 4.6 (Vztah mezi $\hat{e}_{(-i)}$ a $\hat{t}_i$ )

- $\hat{e}_{(-i)} = \frac{\hat{e}_i}{1 - h_{ii}} \Rightarrow E \hat{e}_{(-i)} = 0, \quad \text{Var } \hat{e}_{(-i)} = \frac{\sigma^2}{1 - h_{ii}}$
- standardizované PRESS reziduum  $\frac{\hat{e}_{(-i)}}{\sqrt{\text{Var } \hat{e}_{(-i)}}} = \frac{\hat{e}_i}{\frac{\sigma}{\sqrt{1 - h_{ii}}}} = \frac{\hat{e}_i}{\sigma \sqrt{1 - h_{ii}}} = r_i$
- pokud použijeme  $\hat{\sigma}_{(-i)}^2$  jako odhad  $\sigma^2$ , pak **studentizovaná PRESS rezidua**

$$\frac{\hat{e}_i}{\hat{\sigma}_{(-i)} \sqrt{1 - h_{ii}}} = \hat{t}_i$$

## Poznámka 4.7

- $\widehat{e}_{(-i)} = \frac{\widehat{e}_i}{1 - h_{ii}} \Rightarrow$  pokud  $i$ -té pozorování má velké  $h_{ii}$ , bude  $\widehat{e}_{(-i)}$  mnohem větší, než  $\widehat{e}_i$ ,
  - pozorování s velkým  $h_{ii}$  jsou dobře modelována, ale měřeno  $\widehat{e}_{(-i)}$  mohou špatně predikovat
  - to je další ukázka **fit/prediction dilema**
- stejný efekt nastává také pro  $\widehat{\beta} - \widehat{\beta}_{(-i)} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \widehat{e}_{(-i)}$   
rozdíl může být „malý“, pokud je „fit“ dobrý, ale může být také „velký“, pokud je  $h_{ii}$  velké.

# Míry influence

- i pro perfektní model mohou dva různé vzorky  $(\mathbf{x}, \mathbf{y})$  a  $(\mathbf{x}', \mathbf{y}')$  vést k různým závěrům
- většinou máme k dispozici jen originální data
- bude nás zajímat vliv  $i$ -tého řádku  $\mathbf{X}$  na model
- už víme, že velké  $h_{ii}$  indikuje, že  $i$ -té pozorování má velký vliv a velká rezidua naznačují možnou neadekvátnost modelu
- míry, které zavedeme, budou kombinovat tyto dva faktory
- přístup z PRESS reziduů, tzn. jak velký vliv má vynechání  $i$ -tého pozorování na  $\hat{\beta}$  a  $\hat{\mathbf{y}}$

## DFBETAS:

vliv vynechání  $i$ -tého pozorování na odhad  $\hat{\beta}$  měří rozdíl

$$\hat{\beta} - \hat{\beta}_{(-i)} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \frac{\hat{e}_i}{1 - h_{ii}},$$

bude tedy základem pro naši analýzu

a) vliv  $i$ -tého pozorování na  $\hat{\beta}_j$ :

- $\hat{\beta}_j - \hat{\beta}_{(-i)j} = \frac{r_{ji} \hat{e}_i}{1 - h_{ii}}$ , kde  $r_{ji}$  je  $(j, i)$ -tý prvek matice  $\mathbf{R} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$
- $i$ -té pozorování budeme považovat za **influenční** na  $\beta_j$ , pokud  $\hat{\beta}_j - \hat{\beta}_{(-i)j}$  bude velké
- protože  $\hat{\beta}_j$  je náhodná veličina, „velké“ bychom měli měřit relativně vzhledem k  $s.d.(\hat{\beta}_j) = \sigma \sqrt{v_j}$ , kde  $v_j = (\mathbf{X}^T \mathbf{X})_{jj}^{-1}$
- pokud ji odhadneme pomocí  $\hat{\sigma}_{(-i)} \sqrt{v_j}$ , dostaneme definici

$$\text{DFBETAS}_{j,i} = \frac{\hat{\beta}_j - \hat{\beta}_{(-i)j}}{\hat{\sigma}_{(-i)} \sqrt{v_j}} = \frac{r_{ji} \hat{e}_i}{\sqrt{v_j} \hat{\sigma}_{(-i)} (1 - h_{ii})} = \frac{r_{ji}}{\sqrt{v_j}} \frac{\hat{t}_i}{\sqrt{1 - h_{ii}}},$$

kde  $\hat{t}_i$  je (ext.) studentizované reziduum

- kombinuje efekt velkého rezidua  $\hat{t}_i$  a velkého  $h_{ii}$
- jedna možnost pro limitní hodnoty:  $i$ -té pozorování je považováno za **influenční** na odhad  $\beta_j$ , pokud

$$|\text{DFBETAS}_{j,i}| > \frac{2}{\sqrt{n}}$$

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

b) vliv  $i$ -tého pozorování na celý vektor  $\hat{\beta}$ :

- použití nějaké normy na vektor  $\hat{\beta} - \hat{\beta}_{(-i)}$
- Cook navrhl

$$D_i = \frac{(\hat{\beta} - \hat{\beta}_{(-i)})^T \mathbf{M} (\hat{\beta} - \hat{\beta}_{(-i)})}{(m+1)c},$$

kde  $\mathbf{M}$  je PD matici a  $c$  normalizační konstanta

- nejužívanější volba je  $\mathbf{M} = \mathbf{X}^T \mathbf{X}$  a  $c = s_n^2$
- Cookova vzdálenost:

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

- dosazením dostaneme

$$D_i = \frac{1}{(m+1)s_n^2} \left( \frac{\hat{e}_i}{1-h_{ii}} \right)^2 \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i = \frac{1}{m+1} \frac{h_{ii}}{1-h_{ii}} \frac{\hat{e}_i^2}{s_n^2(1-h_{ii})}$$

- výpočetní tvar potom je

$$D_i = \frac{\hat{r}_i^2}{m+1} \left( \frac{h_{ii}}{1-h_{ii}} \right) \quad (\hat{r}_i \text{ jsou interně studentizovaná (standardizovaná) rezidua})$$

## POZNÁMKA 4.8

- $100(1 - \alpha)\%$  simultánní IS pro  $\beta$  je:

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

- tzn.  $\hat{\beta}_{(-i)} \in C(\alpha) \Leftrightarrow D_i \leq F_{1-\alpha}(m+1, n-m-1)$
- to je motivace pro **RULE OF THUMB**:

*i*-té pozorování je influenční, jestliže  $D_i > F_{\frac{1}{2}}(m+1, n-m-1)$

(pro většinu  $m, n$  je  $F_{\frac{1}{2}} \approx 1$ , zjednodušení pravidla  $D_i > 1$ )

---

## POZNÁMKA 4.9

Také platí, že

$$D_i = \frac{(\hat{\mathbf{y}} - \hat{\mathbf{y}}_{(-i)})^T (\hat{\mathbf{y}} - \hat{\mathbf{y}}_{(-i)})}{(m+1)s_n^2},$$

tzn.  $D_i$  se dá chápat jako míra influence na celkovou predikci

**DFFITS:** vliv  $i$ -tého pozorování na  $\hat{y}_i$

$$\text{DFFITS}_i = \frac{\hat{y}_i - \hat{y}_{(-i)}}{\hat{\sigma}_{(-i)} \sqrt{h_{ii}}} = \dots = \hat{t}_i \sqrt{\frac{h_{ii}}{1 - h_{ii}}}$$

**RULE OF THUMB:**  $i$ -té pozorování je influenční, pokud  $|\text{DFFITS}_i| > 3\sqrt{\frac{m+1}{n-m-1}}$

Poznámka 4.10 (Míry influence v 

- DFBETAS - `dfbetas()`, DFFITS - `dffits()`,
- Cookova vzdálenost  $D_i$  - `cooks.distance()`, potenciál  $h_{ii}$  - `hatvalues()`
- vše shrnuje funkce `influence.measures()`
- používaná pravidla:  $i$ -té pozorování je označeno za influenční, pokud

$$|\text{DFBETAS}_{j,i}| > 1, \quad |\text{DFFITS}_i| > 3\sqrt{\frac{m+1}{n-m-1}}, \quad D_i > F_{\frac{1}{2}}(m+1, n-m-1), \quad h_{ii} > 3\frac{m+1}{n}.$$