

## Potlačení kolinearity

- získání dalších dat
- změna formulace modelu
  - někdy pomůže centrování proměnných
  - vynechání proměnných (může vést k vychýleným odhadům)
- kolinearita  $\Rightarrow$  velké rozptyly odhadů
  - Gauss-Markov říká, že OLS je BLUE parametru  $\beta$
  - tzn. zmenšené rozptyly  $\hat{\beta}$  je možné jen použitím **nelineárních** nebo **vychýlených** odhadů

## Hřebenová regrese (Ridge regression)

- zůstaneme u standardizovaného modelu, tzn. uvažujeme  $\mathbf{X}^*$ ,  $\mathbf{Y}^*$
- pro jednoduchost značení budeme používat  $\mathbf{X}$ ,  $\mathbf{Y}$  a předpokládat  $h(\mathbf{X}) = m$
- hřebenová regrese - uměle se zvedne diagonála  $\mathbf{X}^T \mathbf{X}$  a tím se potlačí kolinearita

- uvažujme  $\mathbf{X}^T \mathbf{X} + \delta \mathbf{I}_m$ , kde  $\delta > 0$  je malé
  - $\mathbf{X}^T \mathbf{X}$  je symetrická, tzn. ex. OG matice  $\mathbf{Q}$  taková, že  $\mathbf{X}^T \mathbf{X} = \mathbf{Q} \mathbf{S}^2 \mathbf{Q}^T$ ,  
kde  $\mathbf{S}^2$  je diagonální matice, mající na diagonále vlastní čísla,  $\lambda_i$ , matice  $\mathbf{X}^T \mathbf{X}$
  - protože  $h(\mathbf{X}^T \mathbf{X}) = m$ , je  $\mathbf{X}^T \mathbf{X}$  PD a  $\lambda_i > 0, \forall i \in \hat{m}$
  - $\mathbf{X}^T \mathbf{X} + \delta \mathbf{I}_m = \mathbf{Q}(\mathbf{S}^2 + \delta \mathbf{I}_m) \mathbf{Q}^T$  to je ale regulární matice
  - odtud plyne

$$(\mathbf{X}^T \mathbf{X} + \delta \mathbf{I}_m) \mathbf{Q} = \mathbf{Q}(\mathbf{S}^2 + \delta \mathbf{I}_m)$$

tzn. čísla na diagonále  $\mathbf{S}^2 + \delta \mathbf{I}_m$  jsou vlastní čísla matice  $\mathbf{X}^T \mathbf{X} + \delta \mathbf{I}_m$

- neboli  $\mathbf{X}^T \mathbf{X} + \delta \mathbf{I}_m$  je regulární a má vlastní čísla  $\lambda_i + \delta$
- jejich velikost i poměry lze měnit volbou  $\delta$  - hlavní myšlenka HR
- místo odhadu  $\hat{\beta}^{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$  budeme studovat odhad

$$\hat{\beta}^{R,\delta} = (\mathbf{X}^T \mathbf{X} + \delta \mathbf{I}_m)^{-1} \mathbf{X}^T \mathbf{Y}$$

### VĚTA 6.3

Vychýlení odhadu  $\hat{\beta}^{R,\delta}$  je

$$\text{bias}(\hat{\beta}^{R,\delta}) = -\delta(\mathbf{X}^T \mathbf{X} + \delta \mathbf{I}_m)^{-1} \beta.$$

Pro střední kvadratickou chybu odhadu  $\hat{\beta}^{R,\delta}$  platí

$$MSE(\hat{\beta}^{R,\delta}) = (\mathbf{X}^T \mathbf{X} + \delta \mathbf{I}_m)^{-1} \left[ \sigma^2 \mathbf{X}^T \mathbf{X} + \delta^2 \beta \beta^T \right] (\mathbf{X}^T \mathbf{X} + \delta \mathbf{I}_m)^{-1}.$$

Důkaz.

## VĚTA 6.4

Ve standardizovaném modelu s plnou hodnotí pro  $0 < \delta < \frac{2\sigma^2}{\|\beta\|^2}$  platí, že

$$\text{Cov}\hat{\beta}^{LS} - \text{MSE}(\hat{\beta}^{R,\delta}) \text{ je PD matice.}$$

## POZNÁMKA 6.7

- ukázali jsme  $E\|\hat{\beta}\|^2 > \|\beta\|^2$   $(+\sigma^2\text{tr}(\mathbf{X}^T\mathbf{X})^{-1})$
- můžeme se snažit  $\hat{\beta}$  nějak zkrátit (za cenu ztráty nestrannosti)
- jiná interpretace hřebenového odhadu:
  - hledejme  $\beta$ , které pro dané  $\delta$  minimalizuje

$$\varphi(\beta) = \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \delta\|\beta\|^2$$

- dá se ukázat, že řešením je  $\hat{\beta}^{R,\delta}$

## VĚTA 6.5

Nechť  $\mathbf{PSQ}^T$  je singulární rozklad matice  $\mathbf{X}$ . Potom platí

$$\text{Cov}(\hat{\boldsymbol{\beta}}^{R,\delta}) = \sigma^2 \sum_{i=1}^m \left( \frac{s_i}{s_i^2 + \delta} \right)^2 \mathbf{q}_i \mathbf{q}_i^T = \sigma^2 \sum_{i=1}^m \frac{\lambda_i}{(\lambda_i + \delta)^2} \mathbf{q}_i \mathbf{q}_i^T.$$

Důkaz.

## POZNÁMKA 6.8

- pokud je  $\mathbf{X}$  špatně podmíněná, bude min. jedno  $s_i^2$  malé, tedy  $\frac{1}{s_i^2}$  vystupující ve  $\text{Var}(\hat{\beta}_j^{LS})$  bude velké
- ale  $\frac{s_i^2}{(s_i^2 + \delta)^2}$  může být zase malé!

## Volba parametru $\delta$

### 1) Hřebenová stopa (ridge trace):

- vypočítají se  $\hat{\beta}^{R,\delta}$  pro různé hodnoty  $\delta$
- vykreslí se graf jednotlivých složek  $\hat{\beta}^{R,\delta}$  v závislosti na  $\delta$
- doporučuje se volit takové  $\delta$ , pro které se grafy „stabilizují“ (subjektivní)

### 2) Harmonic mean estimator: $\hat{\delta} = \frac{ms_n^2}{\|\hat{\beta}\|^2}$ ,

kde  $\hat{\beta}$  a  $s_n^2$  jsou klasické (OLS) odhady parametrů  $\beta, \sigma^2$

## POZNÁMKA 6.9

- existuje spousta dalších metod pro odhad  $\delta$ , vlastnosti většinou zkoumány jen pomocí simulačních studií
- pokud nemáme nějakou apriorní informaci o  $\beta$ , použití hřebenové regrese nezaručí zlepšení OLS

## PŘÍKLAD 6.1 (Data Cement)

```
XX <- cbind(scale(x1),scale(x2),scale(x3),scale(x4))
kappa(XX,exact = TRUE)

## 37.10634

ridge <- lm.ridge(scale(y) ~ XX, lambda = seq(0,0.3,0.001))

plot(ridge)

select(ridge)

## modified HKB estimator is 0.08499604
## modified L-W estimator is 0.05830686
## mallest value of GCV at 0.3
```

