

MULTIPLE LR: ESTIMATION

parameters to estimate : β_1, \dots, β_p } \Rightarrow denote with $\theta = (\beta_1, \dots, \beta_p, \sigma^2) \Rightarrow$ parameter space $\Theta = \mathbb{R}^p \times (0, +\infty)$

data : random sample (y_1, \dots, y_n) ; covariates (x_{i1}, \dots, x_{ip}) for $i = 1, \dots, n$

MODEL : $Y_i \sim N(\mu_i, \sigma^2)$ independent for $i = 1, \dots, n$

with $\mu_i = \beta_1 x_{i1} + \dots + \beta_p x_{ip}$

density $f(y_1, \dots, y_n) = \prod_{i=1}^n f(y_i) = \prod_{i=1}^n \phi(y_i; \mu_i, \sigma^2)$ with $\mu_i = \underline{x}_i^T \underline{\beta}$

likelihood

$$L(\theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2} (y_i - \mu_i)^2\right\}$$

$$= (2\pi)^{-\frac{n}{2}} (\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu_i)^2\right\} \quad \mu_i = \underline{x}_i^T \underline{\beta}$$

$$= (2\pi)^{-\frac{n}{2}} (\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \underline{x}_i^T \underline{\beta})^2\right\}$$

log-likelihood

$$\ell(\theta) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \underline{x}_i^T \underline{\beta})^2$$

$$= -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \underline{x}_i^T \underline{\beta})^2$$

$$\sum_{i=1}^n (y_i - \underline{x}_i^T \underline{\beta})^2 = (\underline{y} - \underline{X}\underline{\beta})^T (\underline{y} - \underline{X}\underline{\beta}) = S(\underline{\beta}) \quad \text{sum of squared residuals}$$

for fixed σ^2 maximizing the likelihood is equivalent to minimizing $S(\underline{\beta})$, independently of the value of σ^2

$$\Rightarrow \hat{\underline{\beta}} = \underset{\underline{\beta} \in \mathbb{R}^p}{\operatorname{argmin}} S(\underline{\beta})$$

Similar to the simple linear model, the ML estimators are the same that we obtain minimizing the sum of squared residuals (OLS estimation).

To find $\hat{\underline{\beta}}$ we need to solve $\frac{\partial}{\partial \underline{\beta}} S(\underline{\beta}) = 0$

$$\text{where } S(\underline{\beta}) = (\underline{y} - \underline{X}\underline{\beta})^T (\underline{y} - \underline{X}\underline{\beta}) =$$

$$= \underline{y}^T \underline{y} - \underline{y}^T \underline{X}\underline{\beta} - \underline{\beta}^T \underline{X}^T \underline{y} + \underline{\beta}^T \underline{X}^T \underline{X} \underline{\beta}$$

$$= \underline{y}^T \underline{y} - 2\underline{y}^T \underline{X}\underline{\beta} + \underline{\beta}^T \underline{X}^T \underline{X} \underline{\beta}$$

useful properties of derivatives :

consider : \underline{a} ($p \times 1$) vector of constants

A ($p \times p$) matrix of constants

$$\cdot \frac{\partial}{\partial \underline{\beta}} \underline{a}^T \underline{\beta} = \underline{a} \quad (p \times 1)$$

$$\cdot \frac{\partial}{\partial \underline{\beta}} \underline{\beta}^T A \underline{\beta} = 2A \underline{\beta} \quad (p \times 1)$$

Hence,

$$\frac{\partial}{\partial \underline{\beta}} S(\underline{\beta}) = \frac{\partial}{\partial \underline{\beta}} (\underline{y}^T \underline{y} - 2\underline{y}^T \underline{X}\underline{\beta} + \underline{\beta}^T \underline{X}^T \underline{X} \underline{\beta})$$

$$= 0 - 2\underline{X}^T \underline{y} + 2\underline{X}^T \underline{X} \underline{\beta}$$

$$\Rightarrow \frac{\partial}{\partial \underline{\beta}} S(\underline{\beta}) = 0$$

$$\Rightarrow \underline{X}^T \underline{X} \underline{\beta} = \underline{X}^T \underline{y}$$

$$\Rightarrow \hat{\underline{\beta}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

notice that

$$-2\underline{X}^T (\underline{y} - \underline{X}\underline{\beta}) = 0 \Rightarrow \begin{cases} \underline{x}_1^T (\underline{y} - \underline{X}\underline{\beta}) = 0 \\ \vdots \\ \underline{x}_p^T (\underline{y} - \underline{X}\underline{\beta}) = 0 \end{cases}$$

"normal equations"

Remark

To solve the equation, $\underline{X}^T \underline{X}$ has to be nonsingular (invertible).

This is ensured by assumption ③ of ABSENCE of MULTICOLLINEARITY (i.e. $\operatorname{rank}(\underline{X}) = p$).

We have found a critical point. Is it a minimum?

$$\text{Hessian: } \frac{\partial^2}{\partial \underline{\beta} \partial \underline{\beta}^T} S(\underline{\beta}) = \frac{\partial^2}{\partial \underline{\beta} \partial \underline{\beta}^T} (-2\underline{X}^T \underline{y} + 2\underline{X}^T \underline{X} \underline{\beta}) = 2\underline{X}^T \underline{X} \Big|_{\underline{\beta} = \hat{\underline{\beta}}} = 2\underline{X}^T \underline{X} \quad \text{has to be positive definite}$$

Recall: \underline{A} is positive definite if $\forall \underline{a} \neq \underline{0}, \underline{a}^T \underline{A} \underline{a} > 0$

does it hold for $\underline{X}^T \underline{X}$?

$$\underline{a}^T \underline{X}^T \underline{X} \underline{a} = (\underline{X} \underline{a})^T (\underline{X} \underline{a}) \geq 0 \quad \text{and it is } = 0 \Leftrightarrow \underline{X} \underline{a} = \underline{0}$$

since we required \underline{X} to have full rank $\Rightarrow \underline{X} \underline{a} = \underline{0} \Leftrightarrow \underline{a} = \underline{0}$

$$\Rightarrow \underline{a}^T \underline{X}^T \underline{X} \underline{a} > 0 \Rightarrow 2\underline{X}^T \underline{X} \text{ is positive definite}$$

$$\Rightarrow \hat{\underline{\beta}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y} \text{ is the minimum of } S(\underline{\beta})$$

and the MAXIMUM LIKELIHOOD ESTIMATE

The maximum likelihood estimator is $\hat{\underline{\beta}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$

• ESTIMATE of σ^2

$$\ell(\theta) = \ell(\underline{\beta}, \sigma^2) = -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\underline{y} - \underline{X}\underline{\beta})^T (\underline{y} - \underline{X}\underline{\beta})$$

$$\ell(\hat{\underline{\beta}}, \sigma^2) = -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\underline{y} - \underline{X}\hat{\underline{\beta}})^T (\underline{y} - \underline{X}\hat{\underline{\beta}})$$

$$\frac{\partial}{\partial \sigma^2} \ell(\hat{\underline{\beta}}, \sigma^2) = -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} (\underline{y} - \underline{X}\hat{\underline{\beta}})^T (\underline{y} - \underline{X}\hat{\underline{\beta}})$$

$$\frac{\partial}{\partial \sigma^2} \ell(\hat{\underline{\beta}}, \sigma^2) = 0$$

$$\Rightarrow -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} (\underline{y} - \underline{X}\hat{\underline{\beta}})^T (\underline{y} - \underline{X}\hat{\underline{\beta}}) = 0$$

$$\Rightarrow -\frac{1}{2(\sigma^2)^2} [n\sigma^2 - (\underline{y} - \underline{X}\hat{\underline{\beta}})^T (\underline{y} - \underline{X}\hat{\underline{\beta}})] = 0 \Rightarrow \hat{\sigma}^2 = \frac{(\underline{y} - \underline{X}\hat{\underline{\beta}})^T (\underline{y} - \underline{X}\hat{\underline{\beta}})}{n} = \frac{\underline{e}^T \underline{e}}{n}$$

$$\frac{\partial^2}{\partial (\sigma^2)^2} \ell(\hat{\underline{\beta}}, \hat{\sigma}^2) = \frac{n}{2(\hat{\sigma}^2)^2} - \frac{1}{2(\hat{\sigma}^2)^3} (\underline{y} - \underline{X}\hat{\underline{\beta}})^T (\underline{y} - \underline{X}\hat{\underline{\beta}})$$

$$= \frac{n}{2(\hat{\sigma}^2)^2} - \frac{1}{(\hat{\sigma}^2)^3} \cdot n\hat{\sigma}^2$$

$$\frac{\partial^2}{\partial (\sigma^2)^2} \ell(\hat{\underline{\beta}}, \hat{\sigma}^2) \Big|_{\sigma^2 = \hat{\sigma}^2} \Rightarrow \frac{n}{2(\hat{\sigma}^2)^2} - \frac{n\hat{\sigma}^2}{(\hat{\sigma}^2)^3} = \frac{n}{2(\hat{\sigma}^2)^2} - \frac{n}{(\hat{\sigma}^2)^2} = -\frac{n}{2(\hat{\sigma}^2)^2} < 0 \quad \text{it's a max}$$

$$\text{The maximum likelihood estimator is } \hat{\underline{\Sigma}}^2 = \frac{(\underline{Y} - \underline{X}\hat{\underline{\beta}})^T (\underline{Y} - \underline{X}\hat{\underline{\beta}})}{n} = \frac{\underline{E}^T \underline{E}}{n}$$

Similarly to the case of the simple linear model, one can show that $\hat{\underline{\Sigma}}^2$ is biased:

$$\mathbb{E}[\hat{\underline{\Sigma}}^2] = \frac{n-p}{n} \sigma^2$$

We can define an UNBIASED estimator of the variance:

$$s^2 = \frac{(\underline{Y} - \underline{X}\hat{\underline{\beta}})^T (\underline{Y} - \underline{X}\hat{\underline{\beta}})}{n-p} = \frac{\underline{E}^T \underline{E}}{n-p} = \frac{n}{n-p} \hat{\underline{\Sigma}}^2$$

the denominator is
 $n - \# \text{ columns of } \underline{X}$

Remarks

$$\cdot \text{the normal equations imply } \begin{cases} (\underline{y} - \underline{X}\hat{\underline{\beta}})^T \underline{x}_1 = 0 \rightarrow \underline{e}^T \underline{x}_1 = 0 \\ \vdots \\ (\underline{y} - \underline{X}\hat{\underline{\beta}})^T \underline{x}_p = 0 \rightarrow \underline{e}^T \underline{x}_p = 0 \end{cases} \Rightarrow \underline{e}^T \underline{X} = 0$$

\Rightarrow orthogonality between the residuals and the columns of \underline{X}

• if we include the intercept $\underline{x}_1 = \underline{1}$

$$\underline{e}^T \underline{x}_1 = 0 \Rightarrow \underline{e}^T \underline{1} = 0 \Rightarrow \sum_{i=1}^n e_i = 0 \Rightarrow \bar{e} = 0 \Rightarrow \text{the residuals have mean } = 0$$

• the vector of the predicted values is

$$\hat{\underline{y}} = \hat{\underline{\mu}} = \underline{X}\hat{\underline{\beta}}$$

$$= \underline{X} \cdot (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y} = \underline{P} \cdot \underline{y}$$

The matrix \underline{P} is:

$$\cdot \text{symmetric: } \underline{P}^T = (\underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T)^T = \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T = \underline{P}$$

$$\cdot \text{idempotent: } \underline{P} \cdot \underline{P} = \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T = \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T = \underline{P}$$

\underline{P} is called the PROJECTION MATRIX (details in the next class...)

• the vector of the residuals is

$$\underline{e} = \underline{y} - \hat{\underline{y}} = \underline{y} - \underline{X}\hat{\underline{\beta}} = \underline{y} - \underline{P}\underline{y} = (\underline{I}_n - \underline{P})\underline{y}$$

The matrix $(\underline{I}_n - \underline{P})$ is:

$$\cdot \text{symmetric: } (\underline{I}_n - \underline{P})^T = \underline{I}_n^T - \underline{P}^T = \underline{I}_n - (\underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T)^T = \underline{I}_n - \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T = \underline{I}_n - \underline{P}$$

$$\cdot \text{idempotent: } (\underline{I}_n - \underline{P})(\underline{I}_n - \underline{P}) = \underline{I}_n - \underline{P} - \underline{P} + \underline{P}^2 = \underline{I}_n - 2\underline{P} + \underline{P}^2 = \underline{I}_n - 2\underline{P} + \underline{P} = \underline{I}_n - \underline{P}$$