

# **Statistical Data Analysis**

---

Dr. Jana de Wiljes

30. November 2021

Universität Potsdam

**Best linear unbiased estimator  
(BLUE)**

---

# Linear estimator

**Def:** A linear estimator has the form

$$\hat{\beta}^L = \mathbf{b} + \mathbf{A}\mathbf{y} \quad (1)$$

where  $\mathbf{b} \in \mathbb{R}^{(p+1) \times 1}$  and  $\mathbf{A} \in \mathbb{R}^{(p+1) \times n}$ .

**Example:** The LS-estimator:

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (2)$$

is a linear estimator with  $\mathbf{b} = \mathbf{0}$  and  $\mathbf{A} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$

**Theorem:** The LS-estimator is BLUE. This means that the LS-estimator has minimal variance among all linear and unbiased estimators  $\hat{\beta}^L$

$$\text{Var}(\hat{\beta}_j) \leq \text{Var}(\hat{\beta}_j^L), \quad j = 0, \dots, p. \quad (3)$$

Furthermore, for an arbitrary linear combination  $\mathbf{c}^\top \hat{\beta}$  it holds that

$$\text{Var}(\mathbf{c}^\top \hat{\beta}) \leq \text{Var}(\mathbf{c}^\top \hat{\beta}^L) \quad (4)$$

LS

# Proof

$$\begin{aligned}
 \text{Proof: } \mathbb{E}[\hat{\beta}^L] &= \mathbb{E}[b + Ax\beta] = \mathbb{E}[b] + \mathbb{E}[Ax\beta + \varepsilon] \\
 &= b + \mathbb{E}[Ax\beta] + A\mathbb{E}[\varepsilon] \\
 &= b + Ax\beta \stackrel{!}{=} \beta \quad \forall \beta \in \mathbb{R}^{p+1} \quad \mathbb{E}[\varepsilon] = 0
 \end{aligned}$$

For the special  $\beta = 0$ :  $b + \underbrace{Ax \cdot 0}_{=0} = 0 \stackrel{!}{=} \beta \Rightarrow b = 0$

$\Rightarrow$  For  $\hat{\beta}^L$  unbiased  $\Rightarrow b = 0$

$$Ax\beta = \beta$$

$\forall \beta \in \mathbb{R}^{p+1}$

$$\Leftrightarrow (Ax - I_{p+1})\beta = 0$$

$$Ax = I_{p+1}$$

$$B \otimes = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$f: V \xrightarrow{B} W$$

$$x \mapsto Bx$$

if  $B$  has full rank  
then  $x = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$  is only solution

- $I_{p+1}$  has full rank:  $\text{rk}(I_{p+1}) = p+1$
  - $X$  has full rank:  $\text{rk}(I_{p+1})$
  - condition:  $\text{rk}(A) = p+1$
- $$\text{rk}(Ax) = \min(\text{rk}(X), \text{rk}(A)) = \text{rk}(I_{p+1}) = p+1$$
- $\Rightarrow$

# Proof

Let the matrix without loss of generality be of the form

$$\rightarrow \boxed{A = (X^T X)^{-1} X^T + B}$$

Inserting into unbiasedness condition  $I_{p+1} = A X$  yields

$$\underline{I_{p+1}} = \underline{AX} = \underbrace{(X^T X)^{-1} X^T X}_{\text{Id}} + BX = \underline{I_{p+1}} + \underbrace{BX}_{\begin{matrix} 1) \\ 0 \end{matrix}}$$

$$\leadsto BX = 0$$

Note:  $\boxed{\text{Cov}(\hat{\beta}^L) = \sigma^2 A A^T}$  (proof too similar to one problem sheet)

$$\text{Cov}(\hat{\beta}^L) = \sigma^2 A A^T$$

$$= \sigma^2 ((X^T X)^{-1} X^T + B) (\underbrace{X(X^T X)^{-1} + B^T}_{\text{Id}})$$

$$= \sigma^2 ((X^T X)^{-1} X^T X (X^T X)^{-1} + (X^T X)^{-1} X^T B^T)$$

$$+ \underbrace{BX(X^T X)^{-1}}_{=0} + B B^T )$$

$$= (BX)^T$$

$$= \boxed{\sigma^2 (X^T X)^{-1} + \sigma^2 B B^T} = \text{Cov}(\hat{\beta}) + \boxed{\sigma^2 B B^T}$$

## Proof

We know BBT is positive semi-definite

$$\text{cov}(\hat{\beta}^L) - \text{cov}(\hat{\beta}) = \sigma^2 B B^T \geq 0$$

$c^T \text{cov}(\hat{\beta}^L)c - c^T \text{cov}(\hat{\beta})c \geq 0 \quad \forall c \in \mathbb{R}^{p+1}$

$$\Rightarrow \text{Var}(c^T \hat{\beta}^L) \geq \text{Var}(c^T \hat{\beta})$$

$$\text{Var}(c^T \hat{\beta}^L) = c^T \text{cov}(\hat{\beta}^L)c \quad \text{and} \quad \text{Var}(c^T \hat{\beta})$$

As  $c$  is arbitrary we can choose for each  $j=0, \dots, p$

$$c = (0, \dots, 1, \dots, 0)^T$$

$$\boxed{\text{Var}(\hat{\beta}_j^L)} = \text{Var}(c^T \hat{\beta}^L) \geq \text{Var}(c^T \hat{\beta}) = \boxed{\text{Var}(\hat{\beta}_j)}$$

□

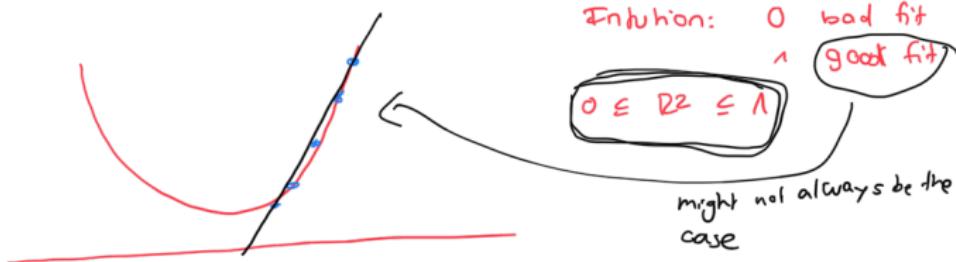
# Coefficient of determination

**Def:** The coefficient of determination is defined by

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (5)$$

The equation is annotated with red circles and arrows. Red circles highlight the terms  $(\hat{y}_i - \bar{y})^2$  and  $(y_i - \bar{y})^2$ . A red arrow points from the term  $\hat{y} = \bar{y}$  to the term  $(\hat{y}_i - \bar{y})^2$ .

and measures the proportion of variability in  $y$  that is accounted for by the statistical model from the overall variation in  $y$ .



# Coefficient of determination

**Lemma:** The method of least squares yields the following geometrical results:

$$\hat{Y} = X\beta + \hat{\epsilon}$$

$$\hat{\epsilon} \sim N(0, \sigma^2)$$

- (i) • The fitted values  $\hat{y}$  are orthogonal to the residuals  $\hat{\epsilon}$ , i.e.,

$$\hat{y}^\top \hat{\epsilon} = 0.$$

$$\hat{Y} = X\hat{\beta} + \hat{\epsilon}$$

- (ii) • The columns of  $X$  are orthogonal to the residuals  $\hat{\epsilon}$ , i.e.,

$$X^\top \hat{\epsilon} = (0)$$

$$X^\top = \begin{pmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} \hat{\epsilon} \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

- (iii) • The residuals are zero on average, i.e.,

$$\sum_{i=1}^n \hat{\epsilon}_i = 0 \quad \text{and} \quad \bar{\hat{\epsilon}} = \frac{1}{n} \underbrace{\sum_{i=1}^n \hat{\epsilon}_i}_{=0} = 0 \quad (6)$$

- (iv) • The mean of the estimated values

$$\bar{\hat{y}} = \frac{1}{n} \sum_{i=1}^n \hat{y}_i \cancel{=} \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \quad (7)$$

# Proof

Proof of (i) Will be using  $H = X \underbrace{(X^T X)^{-1}}_{=H} X^T$ ,  $\hat{y} = X\hat{\beta}$ ,  $\hat{e} = y - \underline{X\hat{\beta}}$

$$\begin{aligned}
 \hat{y}^T \hat{e} &= (X(X^T X)^{-1} \underline{X^T} y)^T \cdot (\underline{y} - \underline{X\hat{\beta}}) \\
 &= y^T \underbrace{X(X^T X)^{-1} X^T}_{=H} \cdot (y - Hy) \\
 &= y^T H (Id - H)y = \cancel{y^T Hy} - y^T \cancel{Hy} \\
 &\quad = H \\
 &\quad = 0 \quad \square \quad \boxed{y} = y
 \end{aligned}$$

$$\begin{aligned}
 \text{Proof of (ii)} \quad X^T \cdot \hat{e} &= X^T(y - \hat{y}) = X^T y - X^T Hy \\
 &= X^T y - \underbrace{X^T X (X^T X)^{-1} X^T y}_{=Id} \\
 &= X^T y - X^T y = 0
 \end{aligned}$$

# Proof

Proof (iii) Remind :  $X = \begin{bmatrix} x_{1,1} & \dots & x_{1,p} \\ \vdots & & \vdots \\ x_{n,1} & \dots & x_{n,p} \end{bmatrix}$

$$\sim \underline{x^0} = (1, \dots, 1)^T$$

$$O = (\underline{x^0})^T \hat{\epsilon} = \underbrace{1^T \hat{\epsilon}}_{\substack{\uparrow \\ \text{true (ii)}}} = \sum_{i=1}^n \hat{\epsilon}_i$$

□

Proof of (iv) : Using (iii) we have

$$\sum_{i=1}^n \underline{y_i} = \sum_{i=1}^n (y_i - \hat{\epsilon}_i) = \sum_{i=1}^n y_i - \underbrace{\sum_{i=1}^n \epsilon_i}_{=0} = \sum_{i=1}^n y_i$$

$$\frac{1}{n} \sum_{i=1}^n \hat{y}_i = \bar{y} \quad \hat{y} = \bar{y}$$

$$y_i = \hat{y}_i + \hat{\epsilon}_i \quad \square$$

# Coefficient of determination

**Lemma:** The following decomposition holds:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n \hat{\epsilon}_i^2 \quad (8)$$

1      R<sup>2</sup>

# Proof

Proof: First we define the  $n \times n$  matrix

$$C = I_n - \frac{1}{n} \underbrace{\mathbf{1}\mathbf{1}^T}_{\substack{n \times 1 \\ 1 \times n}} \quad \begin{matrix} \text{"} \\ \text{(1)} \\ \vdots \\ \text{(1)} \end{matrix} \quad \underbrace{\text{(1 ... 1)}}_{1 \times n}$$

Note that:  $C \cdot C = (I_n - \frac{1}{n} \mathbf{1}\mathbf{1}^T) (I_n - \frac{1}{n} \mathbf{1}\mathbf{1}^T)$

$$\begin{aligned} &= (I_n - 2 \cdot \frac{1}{n} I_n \mathbf{1}\mathbf{1}^T) + \left( \frac{1}{n^2} \mathbf{1} \underbrace{\mathbf{1}^T \mathbf{1}}_{1 \times n \cdot n \times 1} \mathbf{1}^T \right) \\ &= \left( I_n - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right) = C \end{aligned}$$

*GR*  
 $= n$

$C$  is symmetric

Let  $a \in \mathbb{R}^n$  be an arbitrary vector then

$$Ca = \begin{pmatrix} a_1 - \bar{a} \\ a_2 - \bar{a} \\ \vdots \\ a_n - \bar{a} \end{pmatrix}$$

$$\begin{aligned} y &= \hat{y} + \hat{\epsilon} \quad | \cdot C \\ Cy &= C\hat{y} + C\hat{\epsilon} \Rightarrow Cy = C\hat{y} + \hat{\epsilon}^T \quad |^T \\ C\hat{\epsilon} &= \begin{pmatrix} \hat{\epsilon}_1 \\ \hat{\epsilon}_2 \\ \vdots \\ \hat{\epsilon}_n \end{pmatrix} - \bar{\hat{\epsilon}} = \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i = 0 \end{aligned}$$

$$a^T C a = \sum_{i=1}^n (a_i - \bar{a})^2 \quad ①$$

12

# Proof

Using that result we obtain:

$$y^T C C \hat{y} = (\hat{y}^T C + \hat{\epsilon}^T) (C \hat{y} + \hat{\epsilon})$$

$$= \hat{y}^T C C \hat{y} + \hat{y}^T C \hat{\epsilon} + \hat{\epsilon}^T C \hat{y} + \hat{\epsilon}^T \hat{\epsilon}$$

$$C \hat{\epsilon} = \hat{\epsilon}$$

$$C C = C$$

$$\underline{\hat{y}^T C \hat{y}}$$

$$\underbrace{\hat{y}^T \hat{\epsilon}}_{=0}$$

$$\underbrace{\hat{\epsilon}^T \hat{y}}_{=0}$$

$$\underbrace{\hat{\epsilon}^T \hat{\epsilon}}_{=0}$$

- orthogonality follows  
from previous Lemma part (i)

④ and  $C C = C$

$$\rightarrow \underline{y^T C C y} = \underline{y^T C y} = \sum_{i=1}^n (y_i - \bar{y})^2$$

Further we use that  $\underline{\hat{y}_i - \bar{y}}$  and obtain

$$\underline{y^T C \hat{y}} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$\Rightarrow \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n \hat{\epsilon}_i^2$$

□

# Proof

$$Ax = \lambda x$$

$$\boxed{(A - \lambda I)x = 0}$$

↓  
Full rank  
↓

$$\begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\operatorname{rk}(A \cdot B) = \min(\operatorname{rk}(A), \operatorname{rk}(B))$$

$$\boxed{A = \lambda I}$$

$$\begin{pmatrix} 1 & & \\ \vdots & \ddots & \\ 0 & & 1 \end{pmatrix} + \lambda \begin{pmatrix} 0 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 0 \end{pmatrix}$$

# Coefficient of determination

**Lemma:** The coefficient of determination  $R^2$  can be transformed into

$$R^2 = 1 - \frac{\sum_{i=1}^n \hat{\epsilon}_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{\hat{\beta}^\top \mathbf{X}^\top \mathbf{y} - n\bar{y}^2}{\mathbf{y}^\top \mathbf{y} - n\bar{y}^2} \quad (9)$$
$$1 - \frac{\sum_{i=1}^n \hat{\epsilon}_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$
$$\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

# Proof

Proof of second equality : consider numerator and denominator separately

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n \hat{y}_i \cdot \hat{y}_i$$

$$= \hat{\beta} X^T X \hat{\beta} - n \bar{y}^2$$

$$= \hat{\beta} \underbrace{X^T X}_{\text{id}} (X^T X)^{-1} X^T y - n \bar{y}^2$$

$$= \hat{\beta} X^T y - n \bar{y}^2$$

$$\boxed{\sum_{i=1}^n (y_i - \bar{y})^2 = y^T y - n \bar{y}^2}$$

□

# Proof

---

## Coefficient of determination

**Def:** The corrected coefficient of determination  $\bar{R}^2$  is defined by

$$\bar{R}^2 = 1 - \underbrace{\left( \frac{n-1}{n-p-1} \right)}_{\cdot} (1 - R^2) \quad (10)$$

- Remark:
- $0 \leq \bar{R}^2 \leq 1$
  - $\bar{R}^2$  increases with  $p$

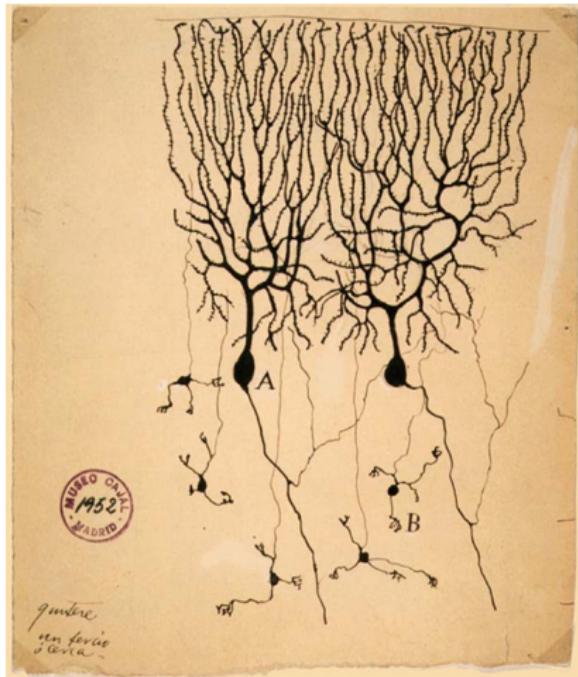
## **Connection to Neural Networks**

---

# Neural Networks

---

# Motivation from biology



By Santiago Ramón y Cajal in 1899 see

[https://de.wikipedia.org/wiki/Santiago\\_Ramón\\_y\\_Cajal](https://de.wikipedia.org/wiki/Santiago_Ramón_y_Cajal) for details

# Neuron

