

# Statistical Data Analysis

---

Dr. Jana de Wiljes

1. Dezember 2021

Universität Potsdam

## **Best linear unbiased estimator (BLUE)**

---

**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

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







## Coefficient of determination

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

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

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:

- The fitted values  $\hat{\mathbf{y}}$  are orthogonal to the residuals  $\hat{\mathbf{e}}$ , i.e.,  $\hat{\mathbf{y}}^\top \hat{\mathbf{e}} = 0$ .
- The columns of  $\mathbf{X}$  are orthogonal to the residuals  $\hat{\mathbf{e}}$ , i.e.,  $\mathbf{X}^\top \hat{\mathbf{e}} = 0$
- The residuals are zero on average, i.e.,

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

- The mean of the estimated values

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





**Lemma:** The following decomposition holds:

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







## Coefficient of determination

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

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







## Coefficient of determination

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

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

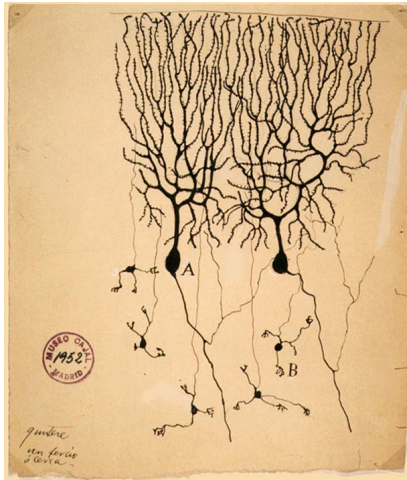
## Connection to Neural Networks

---

# Neural Networks

---

# Motivation from biology



By Santiago Ramn y Cajal in 1899 see

[https://de.wikipedia.org/wiki/Santiago\\_Ramn\\_y\\_Cajalfordetails](https://de.wikipedia.org/wiki/Santiago_Ramn_y_Cajalfordetails)

# Neuron

