

# Statistical Data Analysis

---

Jana de Wiljes

`wiljes@uni-potsdam.de`

`www.dewiljes-lab.com`

31. Oktober 2022

Universität Potsdam

# Linear regression

---

# Model for simple linear regression

**Model:**

$$Y_i = f(X_i, \beta) + \epsilon_i, \quad i = 1, \dots, n \quad (1)$$

where  $\epsilon_i$  are iid with  $\mathbb{E}[\epsilon_i] = 0$  and  $\text{Var}(\epsilon_i) = \sigma^2$

**Data:** it is possible to observe realisations

$$(y_i, x_i) \quad i = 1, \dots, n \quad (2)$$

**Goal:** estimate parameters  $\beta$  of the function to obtain approximative  $f(x, \hat{\beta})$

**Note:** note that  $f$  approximates  $\mathbb{E}[Y_i|X_i]$

# Model for simple linear regression

**Model:**

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i, \quad i = 1, \dots, n \quad (3)$$

where  $\epsilon_i$  are iid with  $\mathbb{E}[\epsilon_i] = 0$  and  $\text{Var}(\epsilon_i) = \sigma^2$

**Data:**

$$(y_i, x_i) \quad i = 1, \dots, n \quad (4)$$

**Goal:** estimate  $f(x, \hat{\beta}) = \hat{\beta}_0 + \hat{\beta}_1 x$

# The Ordinary Multiple Linear Regression Model

## Model:

$$Y_i = \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \beta_3 X_{i,3} + \cdots + \beta_p X_{i,p} + \epsilon_i, \quad i = 1, \dots, n \quad (5)$$

where  $\epsilon_i$  are iid with  $\mathbb{E}[\epsilon_i] = 0$  and  $\text{Var}(\epsilon_i) = \sigma^2$

## Data:

$$(y_i, x_i) \quad i = 1, \dots, n \quad (6)$$

**Goal:** estimate  $\hat{f}(x_1, \dots, x_p, \hat{\beta}_1, \dots, \hat{\beta}_p) = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots, \hat{\beta}_p x_p$

# Multivariate Random Variables

**Def:** Let  $\mathbf{X}$  be a vector of (univariate) random variables, i.e.,  $\mathbf{X} = (X_1, \dots, X_p)^\top$  with  $\mathbb{E}[X_i] = \mu_i$ .  $\mathbf{X}$  is called a multivariate random variable and we denote  $\mathbb{E}[\mathbf{X}] = \mu$

## Note:

- Variance  $\text{Var}(X_i) = \mathbb{E}[(X_i - \mathbb{E}(X_i))^2] = \mathbb{E}[(X_i - \mathbb{E}(X_i))(X_i - \mathbb{E}(X_i))]$
- Covariance  $\text{Cov}(X_i, X_j) = \mathbb{E}[(X_i - \mathbb{E}(X_i))(X_j - \mathbb{E}(X_j))]$

# Least squares Estimation: minimize the sum of squared errors

Least squares estimation: minimize the sum of squared errors

$$L(\beta) = \sum_{i=1}^N (y_i - x_i^\top \beta_i)^2 = \sum_{i=1}^N \epsilon_i^2 = \epsilon^\top \epsilon \quad (7)$$

with respect to  $\beta \in \mathbb{R}^{p+1}$

$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^N \epsilon_i^2 = \epsilon^\top \epsilon = (Y - X\beta)^\top (Y - X\beta) \quad (8)$$

$$Y = \underbrace{\begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}}_{\in \mathbb{R}^{N \times 1}}, \quad X = \underbrace{\begin{bmatrix} 1 & x_{11} & \cdots & x_{1p} \\ \vdots & & & \\ 1 & x_{N1} & \cdots & x_{Np} \end{bmatrix}}_{\in \mathbb{R}^{N \times p+1}}, \quad \beta = \underbrace{\begin{bmatrix} \beta_1 \\ \vdots \\ \beta_N \end{bmatrix}}_{\in \mathbb{R}^{p+1 \times 1}}, \quad \epsilon = \underbrace{\begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_N \end{bmatrix}}_{\in \mathbb{R}^{N \times 1}},$$

**Def:** The covariance of the multivariate random variable  $\mathbf{X}$  is defined by

$$\Sigma := \text{Cov}(\mathbf{X}) = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top] \quad (9)$$

**Example:**

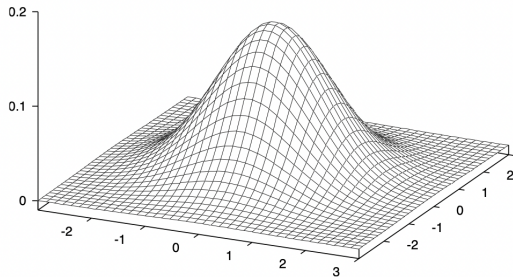
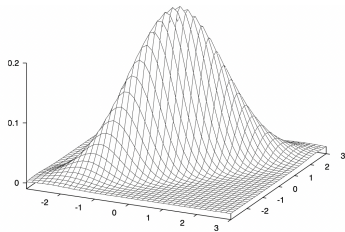
$$\Sigma = \begin{pmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) \end{pmatrix} \quad (10)$$

Properties of  $\Sigma$ :

- quadratic
- symmetric
- positive-semidefinite



# Multivariate Normal Distribution



$$\mathbf{X} \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (11)$$

**Lemma:** Let  $\mathbf{B}$  be an  $n \times (p+1)$  matrix. Then the matrix  $\mathbf{B}^\top \mathbf{B}$  is symmetric and positive semi-definite. It is positive definite, if  $\mathbf{B}$  has full column rank. Then, besides  $\mathbf{B}^\top \mathbf{B}$  also  $\mathbf{B}\mathbf{B}^\top$  is positive semi-definite.

**Theorem:** The LS-estimator of the unknown parameters  $\beta$  is

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (12)$$

if  $\mathbf{X}$  has full column rank  $p+1$ .

**Proof:**



**Proposition:** The hat-matrix  $\mathbf{H} = (h_{ij})_{1 \leq i, j \leq b}$  has the following properties:

1.  $\mathbf{H}$  is symmetric
2.  $\mathbf{H}$  is idempotent, i.e.,  $\mathbf{H}\mathbf{H} = \mathbf{H}$
3.  $rk(\mathbf{H}) = tr(\mathbf{H}) = p + 1$
4.  $0 \leq h_{ii} \leq 1, \quad \forall i = 1, \dots, n$
5. the matrix  $\mathbf{I}_n - \mathbf{H}$  is also symmetric and idempotent with  $rk(\mathbf{I}_n - \mathbf{H}) = n - p - 1$

**Theorem:** The ML-estimator of the unknown parameters  $\sigma^2$  is  $\hat{\sigma}_{ML}^2 = \frac{\hat{\epsilon}\hat{\epsilon}}{n}$  with  $\hat{\epsilon} = \mathbf{y} - \mathbf{X}\hat{\beta}$ .



**Proposition:** For the ML-estimator  $\hat{\sigma}_{ML}^2$  of  $\sigma^2$  the following property holds:

$$\mathbb{E}[\sigma_{ML}^2] = \frac{n - p - 1}{n} \sigma^2 \quad (13)$$





**Proposition:** The adjusted estimator

$$\hat{\sigma}_{ad}^2 = \frac{\hat{\epsilon}\hat{\epsilon}}{n - p - 1} \quad (14)$$

of the unknown parameter  $\sigma^2$  can be written as

$$\hat{\sigma}_{ad}^2 = \frac{\mathbf{y}^\top \mathbf{y} - \hat{\beta}^\top \mathbf{X}^\top \mathbf{y}}{n - p - 1} \quad (15)$$



**Proposition:** The LS-estimator  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$  is equivalent to the ML-estimator based on maximization of the log-likelihood

$$l(\beta, \sigma^2) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) \quad (16)$$



**Proposition:** The LS-estimator  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$  and the REML-estimator  $\hat{\sigma}^2 = \frac{1}{n-p-1} \hat{\epsilon}^\top \hat{\epsilon}$  the following properties hold:

1.  $\mathbb{E}[\hat{\beta}] = \beta$ ,  $\text{Cov}(\hat{\beta}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$
2.  $\mathbb{E}[\hat{\sigma}^2] = \sigma^2$

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

$$\hat{\beta}^L = \mathbf{b} + \mathbf{A}\mathbf{y} \quad (17)$$

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

**Example:** The LS-estimator:

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (18)$$

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