

# **MEM-264 Applied Statistics**

**Department of Mathematics and Applied Mathematics, University of Crete**

Costas Smaragdakis (kesmarag@uoc.gr)

21th lecture - 23-04-2021

# Linear Model

Sample  $\{(X_1, Y_1), \dots, (X_n, Y_n)\} \stackrel{k=1}{\text{iid.}}$

$X_i$  - independent random variable

$Y_i$  - dependent random variable

## Linear regression

In the general case we have the following random sample

$$Y_i = a + \sum_{j=1}^k b^{(j)} X_i^{(j)} + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2), \quad i = 1, \dots, n$$

$\{(X_1^{(1)}, X_1^{(2)}, X_1^{(3)}, \dots, X_1^{(k)}, Y_1), \dots, (X_n^{(1)}, \dots, X_n^{(k)}, Y_n)\}$   
 $\epsilon_1, \dots, \epsilon_n \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$

$$\tilde{p} = [a, b^{(1)}, b^{(2)}, \dots, b^{(k)}]^T \in \mathbb{R}^{k+1}$$

↑  
model's parameters

## The least squares estimator

Dataset:  $\{(X_1^{(1)}, \dots, X_1^{(k)}, Y_1), \dots, (X_n^{(1)}, \dots, X_n^{(k)}, Y_n)\}$  we wish to find an estimator

- ▶ Minimize the following function of  $p$ :

$$(\tilde{y} - \tilde{X}\tilde{p})^T (\tilde{y} - \tilde{X}\tilde{p})$$

$$\hat{p} = [\hat{a}, \hat{b}^{(1)}, \dots, \hat{b}^{(k)}]$$

$$\hat{y} = \hat{a} + \sum_{j=1}^k \hat{b}^{(j)} X^{(j)}$$

↑  
the estimator

Given  $(X^{(1)}, \dots, X^{(k)}) \xrightarrow{\text{estimator}} \hat{y} \approx y$

$$\underline{\tilde{y}} = [y_1, \dots, y_n]^T, \quad \underline{\tilde{X}} = \begin{bmatrix} 1 & X_1^{(1)} & \dots & X_1^{(k)} \\ 1 & X_2^{(1)} & \dots & X_2^{(k)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_n^{(1)} & \dots & X_n^{(k)} \end{bmatrix}$$

By multiplying the  $i$ -th row of the matrix with  $b$

$$(\underline{\tilde{y}} - \underline{\tilde{X}} \underline{\tilde{p}})^T (\underline{\tilde{y}} - \underline{\tilde{X}} \underline{\tilde{p}}) = \|\underline{\tilde{y}} - \underline{\tilde{X}} \underline{\tilde{p}}\|_2^2$$

$$[1, X_i^{(1)}, \dots, X_i^{(k)}] \begin{bmatrix} \hat{\alpha} \\ \hat{b}^{(1)} \\ \vdots \\ \hat{b}^{(k)} \end{bmatrix} = \hat{\alpha} + \sum_{j=1}^k \hat{b}^{(j)} X_i^{(j)}$$



the minimum value of the norm is taken when  $\underline{\tilde{p}} = (\underline{\tilde{X}}^T \underline{\tilde{X}})^{-1} \underline{\tilde{X}}^T \underline{\tilde{y}}$

or

$$(\underline{\tilde{X}}^T \underline{\tilde{X}}) \underline{\tilde{p}} = \underline{\tilde{X}}^T \underline{\tilde{y}}$$

$k+1$  equations and  $k+1$  unknowns

exercise/derivation

$$\frac{\partial}{\partial \underline{\tilde{p}}} (\underline{\tilde{y}} - \underline{\tilde{X}} \underline{\tilde{p}})^T (\underline{\tilde{y}} - \underline{\tilde{X}} \underline{\tilde{p}}) = 2 (\underline{\tilde{X}}^T \underline{\tilde{X}} \underline{\tilde{p}} - \underline{\tilde{X}}^T \underline{\tilde{y}}) = 0$$

$$\frac{\partial^2}{\partial \underline{\tilde{p}} \partial \underline{\tilde{p}}^T} (\underline{\tilde{y}} - \underline{\tilde{X}} \underline{\tilde{p}})^T (\underline{\tilde{y}} - \underline{\tilde{X}} \underline{\tilde{p}}) = 2 \underline{\tilde{X}}^T \underline{\tilde{X}} > 0 \quad \text{if } \underline{\tilde{X}} \text{ a invertible matrix. (non singular)}$$

Therefore

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

example

$\{(1, 2, 3), (2, 5, 9)\}$   
independent      dependent

$$\underline{y} = \begin{bmatrix} 3 \\ 9 \end{bmatrix}$$

$$\underline{X} = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 5 \end{bmatrix}$$

$$\underline{\beta} = \begin{bmatrix} \alpha \\ \beta^{(1)} \\ \beta^{(2)} \end{bmatrix}$$

$$\begin{matrix} (X^T X) & \overset{3 \times 1}{P} = & X^T \underline{y} \\ \underset{\sim}{\sim} & & \underset{\sim}{\sim} \\ \begin{matrix} 3 \times 2 & 2 \times 3 \\ 3 \times 3 \end{matrix} & & \begin{matrix} 3 \times 2 & 2 \times 1 \end{matrix} \end{matrix}$$

the solution of this problem defines the estimator.

$$\hat{y} = \hat{\alpha} + \hat{\beta}^{(1)} X^{(1)} + \hat{\beta}^{(2)} X^{(2)}$$

(1.5, 1.5, ?)

(end of the first part of this lecture!) —



**Sigmoid function** (logistic function)

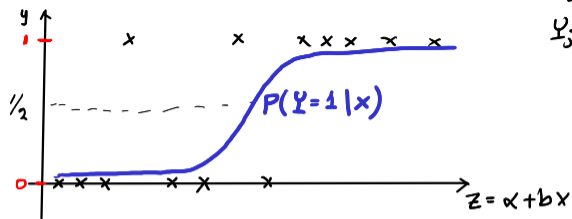
$$0 < \text{sigmoid}(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} < 1$$

**Logistic regression**

$$\log \frac{p(Y_j = 1 | X_j^{(1)}, \dots, X_j^{(k)})}{1 - p(Y_j = 1 | X_j^{(1)}, \dots, X_j^{(k)})} = a + \sum_{j=1}^k b^{(j)} X_j + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2), \quad i = 1, \dots, n$$

$$\pi_j = P(Y_j = 1 | X_j^{(1)}, \dots, X_j^{(k)})$$

$\{(x_1, y_1), \dots, (x_n, y_n)\}$  iid  $(x_j, y_j)$ ,  $x_j$  Continuous Random Variable  
 $y_j \sim \text{Be}(\theta)$  Bernoulli



$$P(Y=1|x) = \text{Sigmoid}(z(x)) = \frac{e^{z(x)}}{1 + e^{z(x)}}$$

$$z(x) = \log \frac{P(Y=1|x)}{\underbrace{1 - P(Y=1|x)}_{P(Y=0|x)}}$$

$$\log \frac{P(Y=1|x)}{1 - P(Y=1|x)} = \alpha + bx + \varepsilon$$

$$\log \frac{\pi_j}{1 - \pi_j} = \alpha + bx_j + \varepsilon_j$$

$$P(\underline{y} | \underline{x}) = \prod_{j=1}^n P(y_j | x_j) = \prod_{j=1}^n P(y_j = 1 | z_j)^{y_j} P(y_j = 0 | z_j)^{1 - y_j} =$$

$$\underline{y} = (y_1, \dots, y_n)^T$$

$$\underline{x} = (x_1, \dots, x_n)^T$$

$$= \prod_{j=1}^n \pi_j^{y_j} (1 - \pi_j)^{1 - y_j}$$

$$\log P(\underline{y} | \underline{x}) = \log \prod_{j=1}^n \pi_j^{y_j} (1 - \pi_j)^{1 - y_j} = \sum_{j=1}^n y_j \log \pi_j + (1 - y_j) \log (1 - \pi_j)$$

Categorical cross entropy BCE =  $-\frac{1}{N} \sum_{j=1}^n y_j \log \pi_j + (1 - y_j) \log (1 - \pi_j)$