

# Logistic regression (binary classification)

feature/design matrix

$\bar{X}$

target

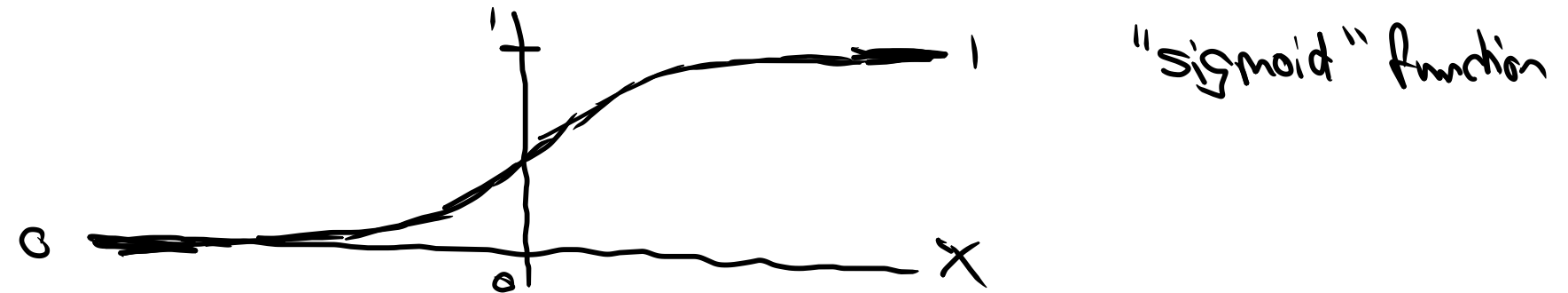
$\bar{y}$

|   |          |  |  |  |          |
|---|----------|--|--|--|----------|
|   | 1        |  |  |  | n        |
|   | $x_{11}$ |  |  |  |          |
|   |          |  |  |  |          |
|   |          |  |  |  |          |
|   |          |  |  |  |          |
|   |          |  |  |  |          |
| m |          |  |  |  | $x_{mn}$ |

|   |      |
|---|------|
| 0 | 0,00 |
| 1 | 1,00 |
| 1 | 1,00 |
|   |      |
| ⋮ |      |
| 0 | 0,00 |

model predicts the probability of class 1 :

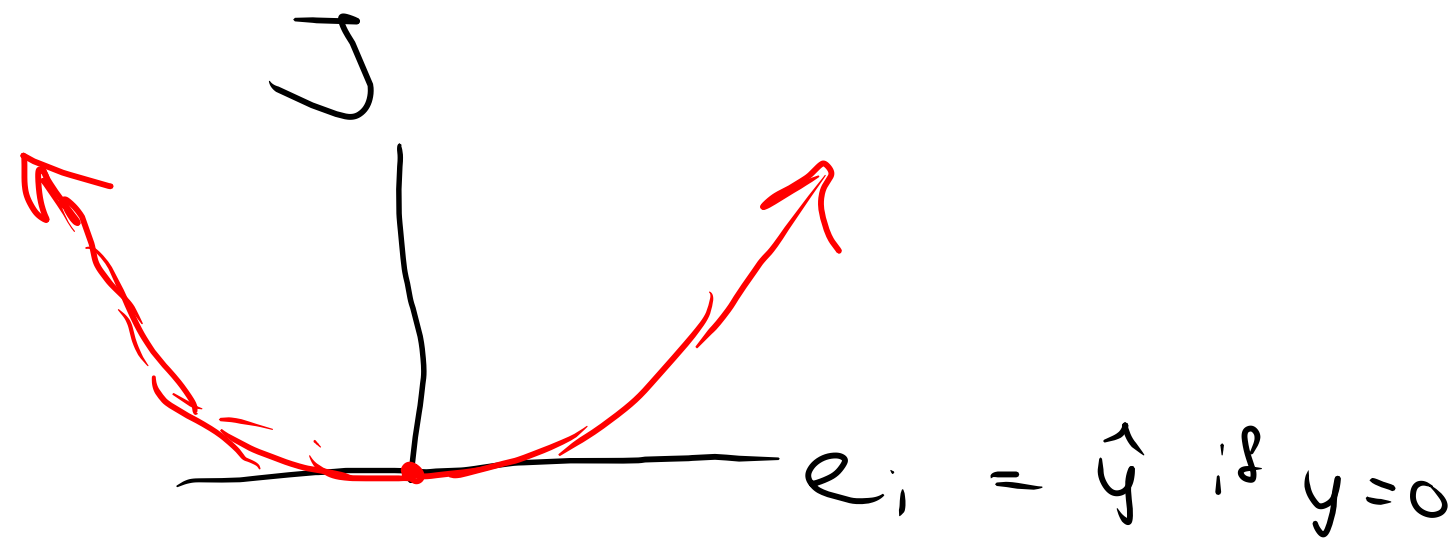
$$\hat{y}_i = \sigma(x_{i1}\theta_1 + x_{i2}\theta_2 + \dots + x_{in}\theta_n) = \sigma(\bar{X} \cdot \bar{\theta})$$



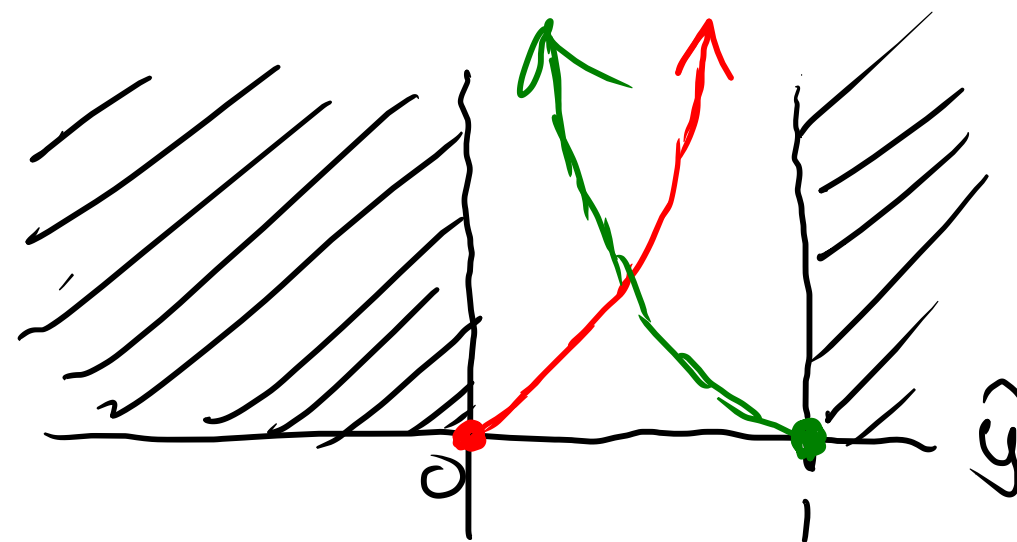
logistic function:  $\sigma(x) = \frac{1}{1 + e^{-x}}$

Cost Function:

regression:  $J = \frac{1}{2m} \sum_i \overbrace{(\hat{y}_i - y_i)^2}^{e_i}$



classification: 
$$\begin{cases} -\ln(1-\hat{y}) & \text{if } y=0 \\ -\ln(\hat{y}) & \text{if } y=1 \end{cases}$$



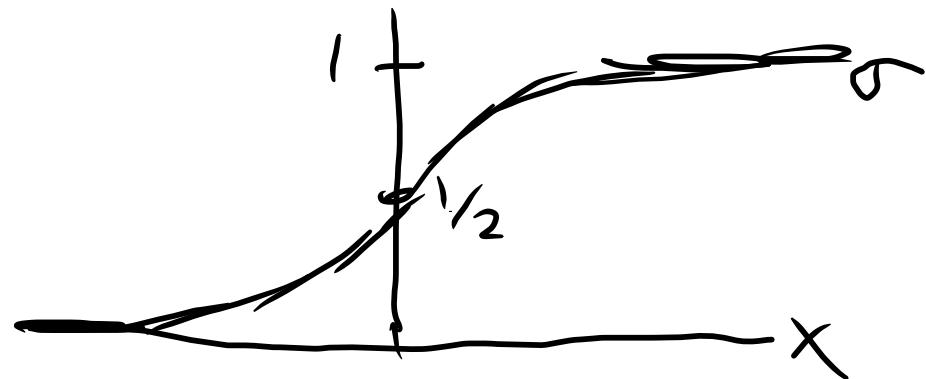
if  $y=0$   
 $y=1$

"cross-entropy"

alternative: 
$$J = \frac{1}{m} \sum_i \left( -y_i \ln(\hat{y}_i) - (1-y_i) \ln(1-\hat{y}_i) \right)$$

Intermezzo:

logistic function  $\sigma(x)$

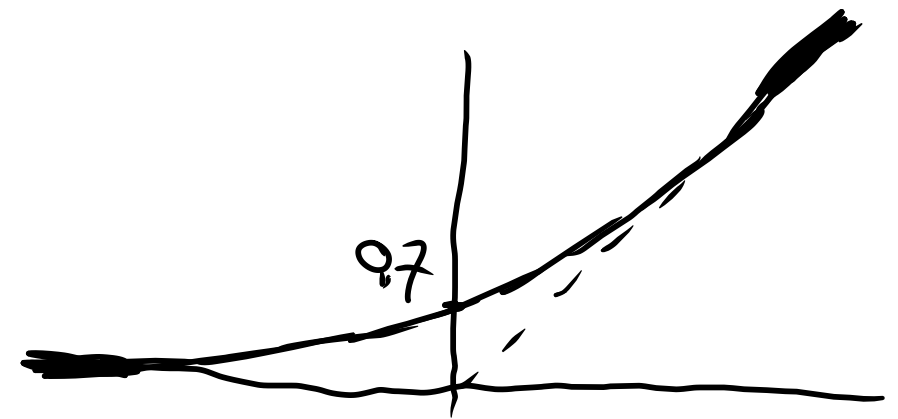


$$\sigma(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{1+e^x}$$

(A)  $\sigma(-x) = 1 - \sigma(x) \Leftrightarrow \sigma(x) + \sigma(-x) = 1$

(B) 
$$\begin{aligned} h(\sigma(x)) &= \ln\left(\frac{1}{1+e^{-x}}\right) = \ln(1) - \ln(1+e^{-x}) \\ &= -\ln(1+e^{-x}) = -S(-x) \end{aligned}$$

Softplus function  $S(x)$



$$S(x) = \ln(1+e^x)$$

(C)  $S'(x) = \frac{d}{dx} \ln(1+e^x) = \frac{e^x}{1+e^x} = \sigma(x)$

Optimize the model  $\hat{y} = \sigma(\bar{X}\bar{\theta})$  by means of gradient descent:  $\theta_k \leftarrow \theta_k - \alpha \cdot \frac{\partial J}{\partial \theta_k}$

$$\text{If } y=0: \theta_k \leftarrow \theta_k - \alpha \cdot \frac{\partial}{\partial \theta_k} \left( -\frac{1}{n} \ln(1 - \hat{y}_i) \right) = \theta_k + \frac{\alpha}{n} \frac{\partial}{\partial \theta_k} \ln(1 - \sigma(\sum_j x_{ij} \theta_j)) \stackrel{A}{=}$$

$$\stackrel{A}{=} \theta_k + \frac{\alpha}{n} \frac{\partial}{\partial \theta_k} \ln(\sigma(-\sum_j x_{ij} \theta_j)) \stackrel{B}{=} \theta_k - \frac{\alpha}{n} \frac{\partial}{\partial \theta_k} S(\sum_j x_{ij} \theta_j) \stackrel{C}{=} \theta_k - \frac{\alpha}{n} \sigma(\sum_j x_{ij} \theta_j) \cdot x_{ik}$$

$$= \theta_k - \frac{\alpha}{n} (\hat{y}_i - 0) \cdot x_{ik}$$

$$\text{If } y=1: \theta_k \leftarrow \theta_k - \alpha \cdot \frac{\partial}{\partial \theta_k} \left( -\frac{1}{n} \ln(\hat{y}_i) \right) = \theta_k + \frac{\alpha}{n} \frac{\partial}{\partial \theta_k} \ln(\sigma(\sum_j x_{ij} \theta_j)) \stackrel{B}{=} \theta_k - \frac{\alpha}{n} \frac{\partial}{\partial \theta_k} S(-\sum_j x_{ij} \theta_j) \stackrel{C}{=}$$

$$\stackrel{C}{=} \theta_k + \frac{\alpha}{n} \sigma(-\sum_j x_{ij} \theta_j) x_{ik} \stackrel{A}{=} \theta_k + \frac{\alpha}{n} (1 - \sigma(\sum_j x_{ij} \theta_j)) x_{ik} = \theta_k + \frac{\alpha}{n} (1 - \hat{y}_i) x_{ik} =$$

$$= \theta_k - \frac{\alpha}{n} (\hat{y}_i - 1) x_{ik}$$

$$\text{In general: } \theta_k \leftarrow \theta_k - \frac{\alpha}{n} (\hat{y}_i - y_i) x_{ik} = \theta_k - \frac{\alpha}{n} x_{ik} e_i \Rightarrow \underline{\underline{\bar{\theta} \leftarrow \bar{\theta} - \frac{\alpha}{n} \bar{X}^T \bar{e}}}$$

This is identical to the update-rule for linear regression