

①  $(x_1, y_1), (x_2, y_2), (x_3, y_3) \dots (x_m, y_m)$

~~1~~  $h_\theta(x) = \theta_0 + \theta_1 x$

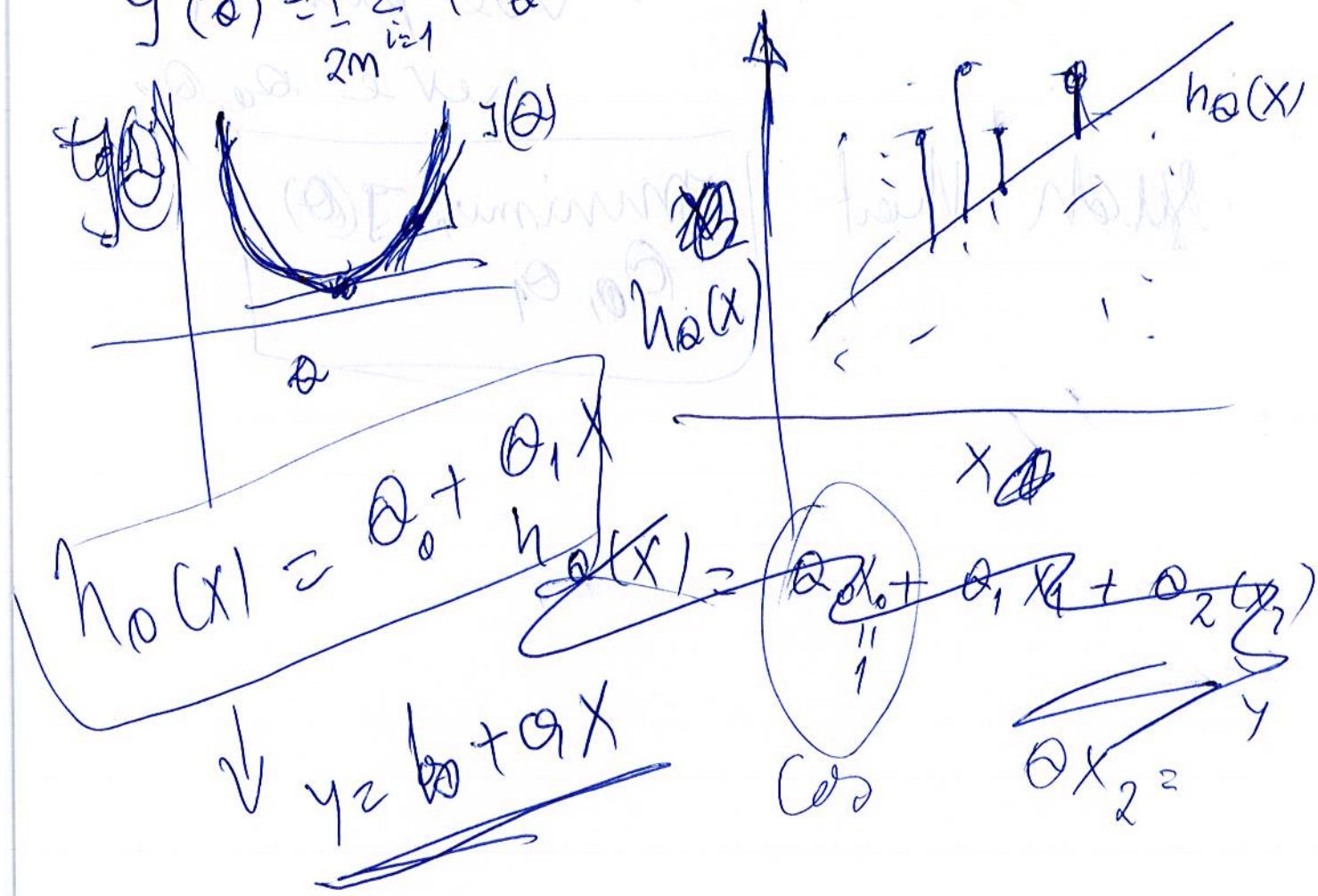
~~2~~  $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 \quad x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$

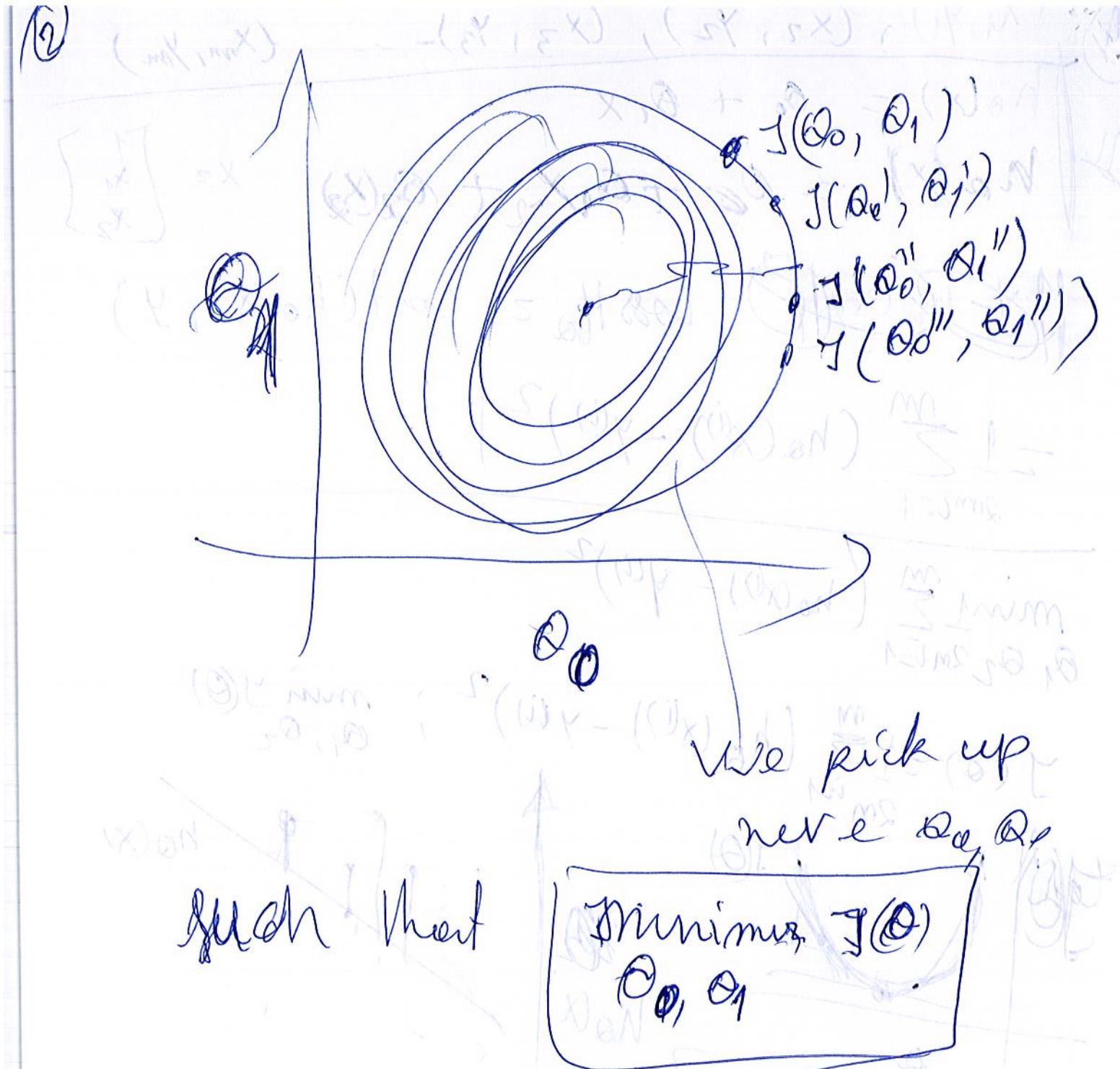
~~3~~  $\text{Loss } H_\theta = \text{Cost}(h_\theta(x), y)$

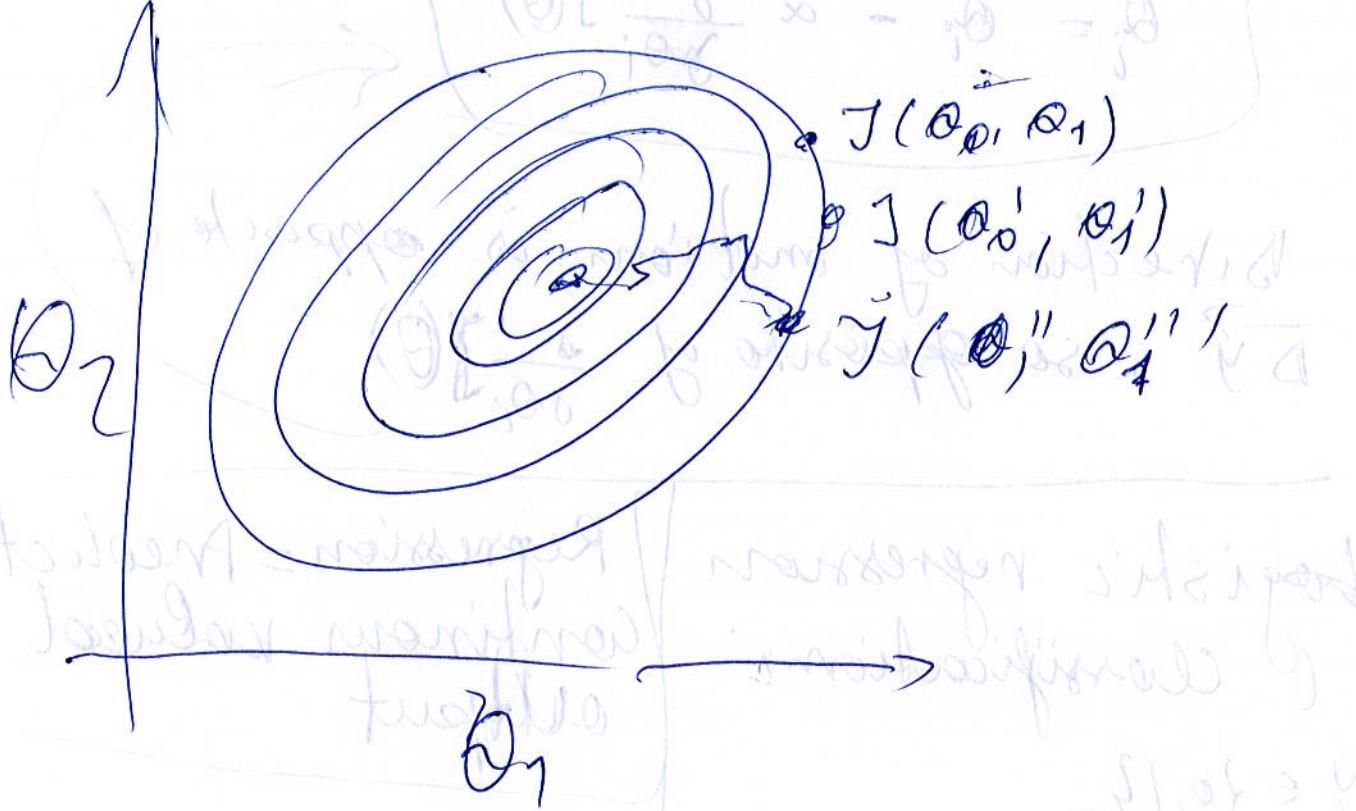
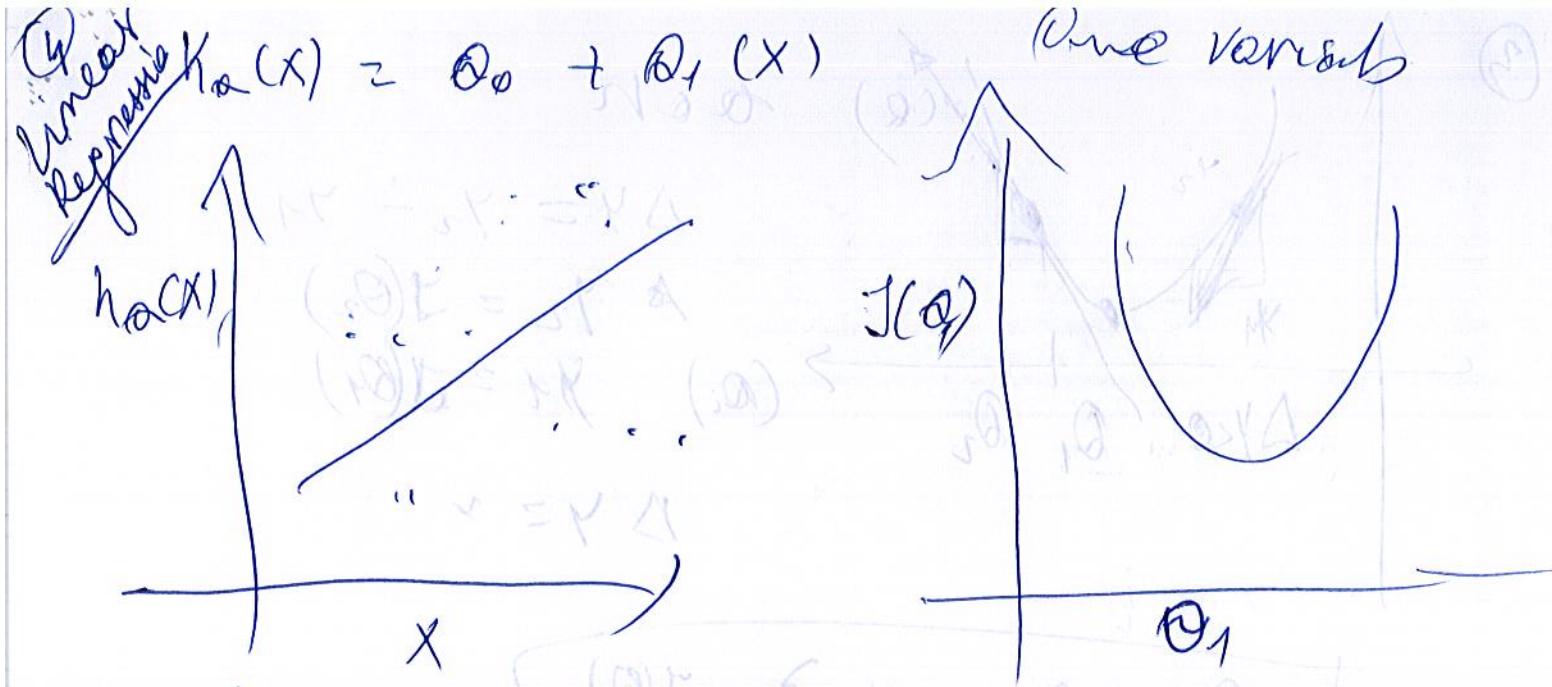
$$= \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\min_{\theta_0, \theta_1} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 ; \min_{\theta_0, \theta_1} J(\theta)$$





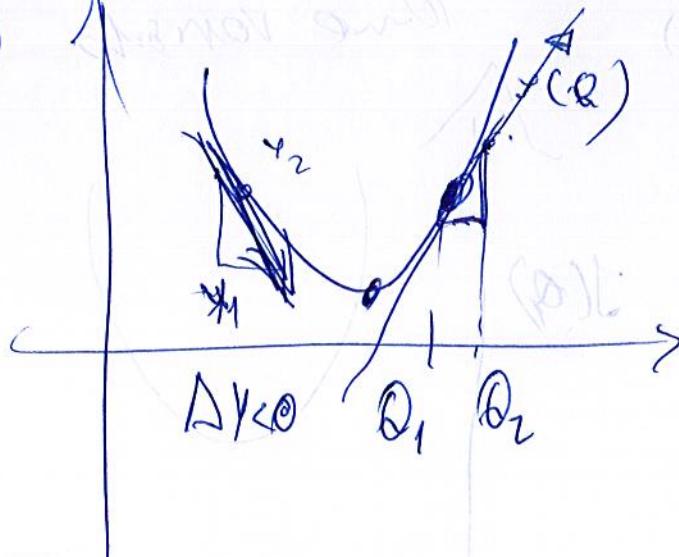


$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

~~linear~~ regression  
Predict continuous valued output

③

 $\Delta \theta$ 

$$\Delta y = y_2 - y_1$$

$$y_2 = J(\theta_2)$$

$$y_1 = J(\theta_1)$$

$$\Delta y = \dots$$

$$\theta_i = \theta_i - \alpha \frac{\partial}{\partial \theta_i} J(\theta)$$

Direction of motion is opposite of  $\nabla J$ , so opposite of  $\frac{\partial}{\partial \theta_i} J(\theta)$

Logistic regression  
classification,

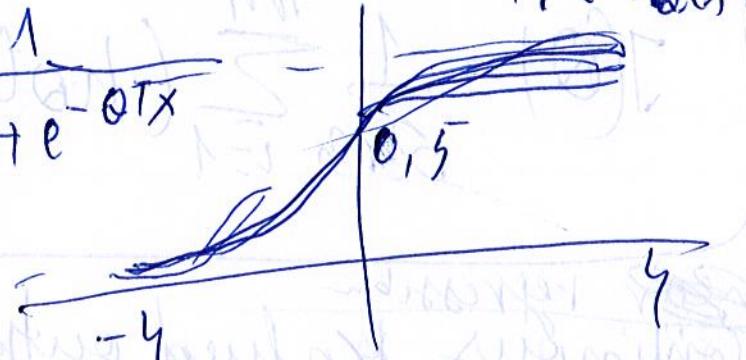
$$y \in \{0, 1\}$$

Regression = Predict  
continuous valued  
output

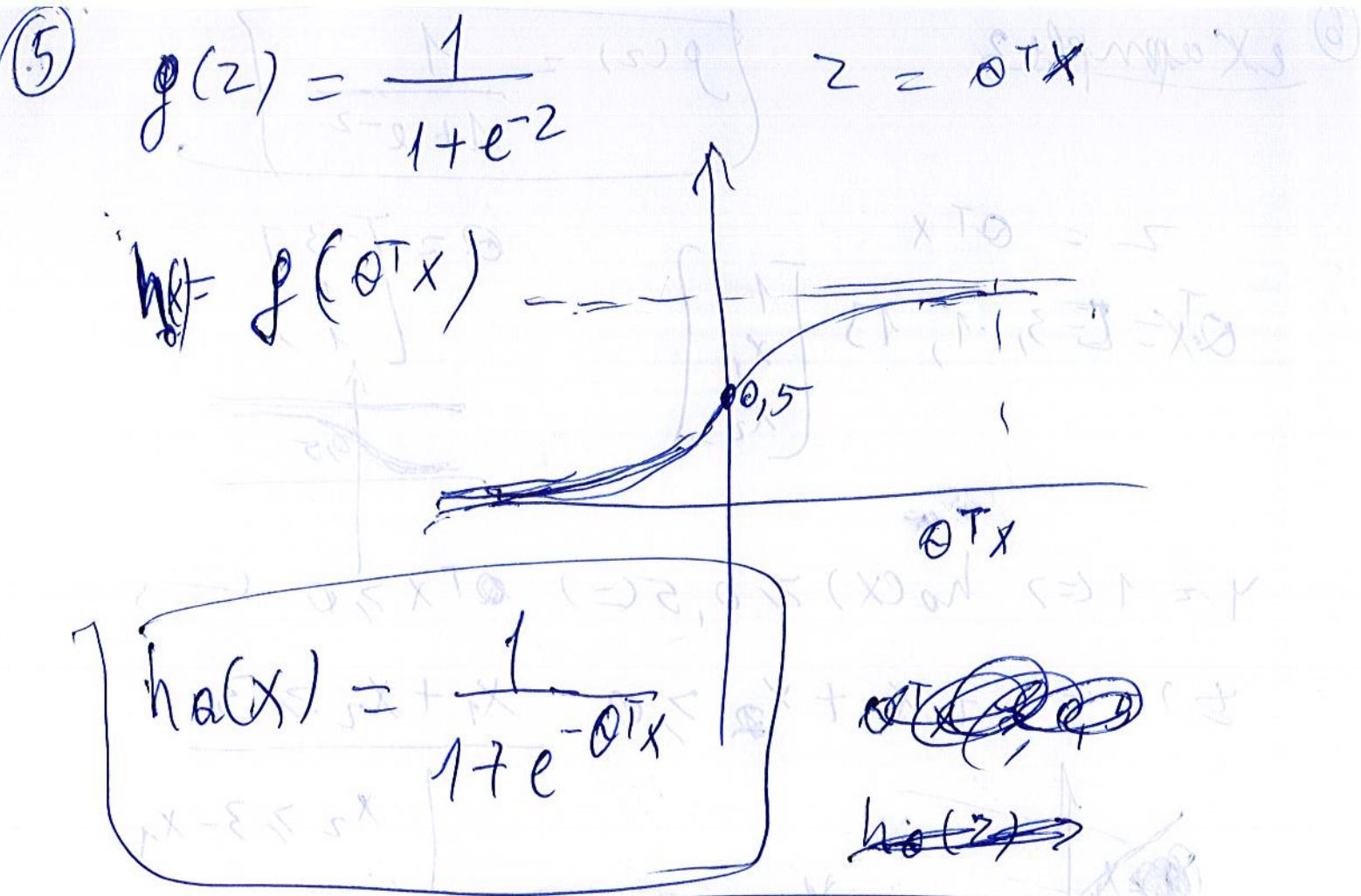
$$0 \leq h_\theta(x) \leq 1$$

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$



$$z = \theta^T x$$



$$y=1 \text{ if } h_\theta(x) \geq 0,5 \quad \theta^T x \geq 0$$

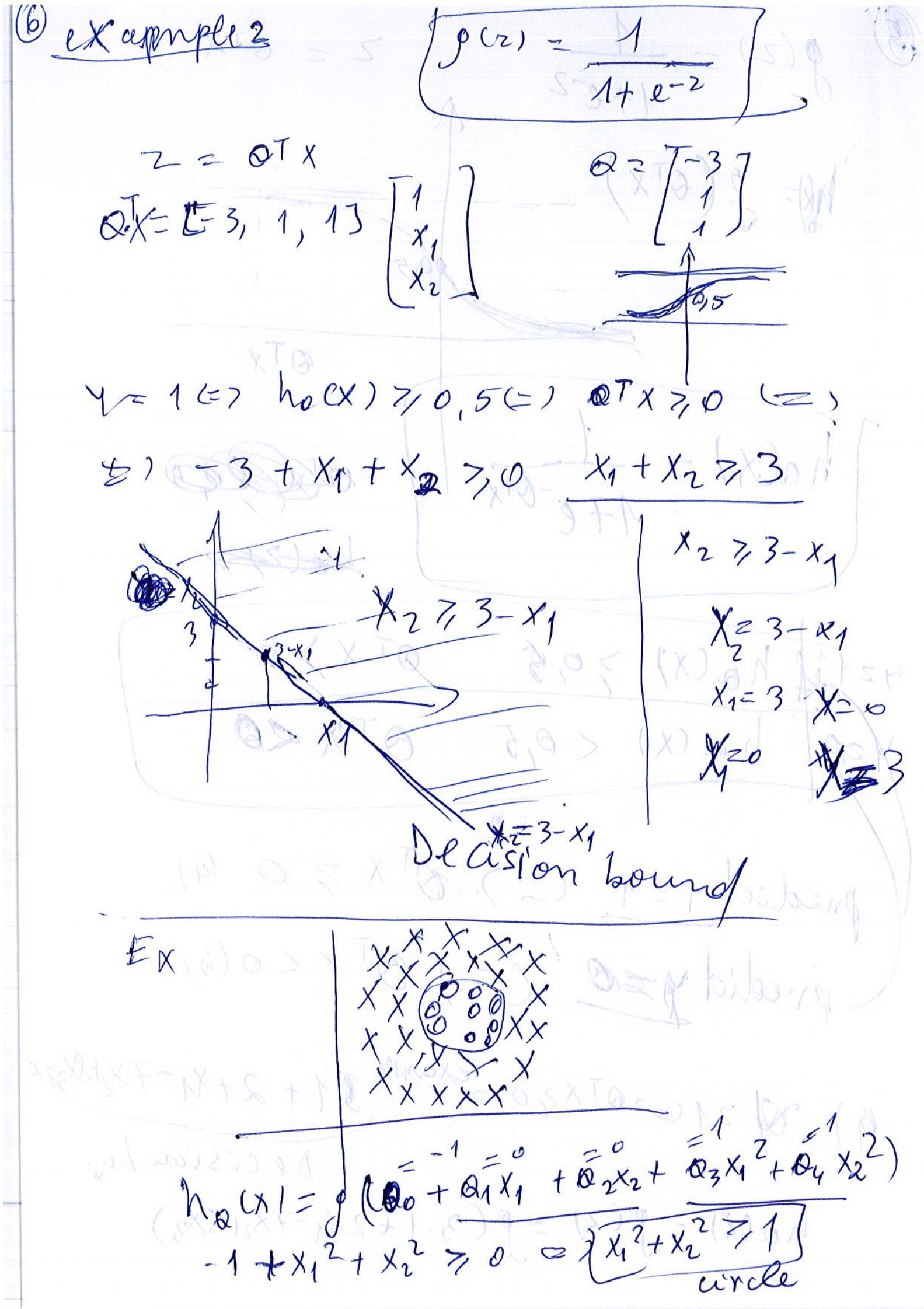
$$y=0 \quad \text{if } h_\theta(x) < 0,5 \quad \theta^T x < 0$$

predict  $y=1$   $\Leftrightarrow \theta^T x \geq 0$  (a)

predict  $y=0$   $\Leftrightarrow \theta^T x < 0$  (b)

(a)  ~~$z_1 \Rightarrow \theta^T x \geq 0$~~   $\stackrel{\text{example}}{=} 3 \cdot 1 + 2 + x_1 - 7x_2 + 8x_3 \geq 0$   
 Decision by

$$h_\theta(x) = g(z) = g(3 \cdot 1 + 2 + x_1 - 7x_2 + 8x_3)$$



Linear regression:  $h_{\theta}(x) = [0_0, \dots, 0_m] \cdot \begin{bmatrix} t_0 \\ t_1 \\ \vdots \\ t_m \end{bmatrix}$

$$t_i \in \mathbb{R}$$

$$t \in \{x_i, x_i^2, x_2 x_1, x_3^2 x_2^2, \dots\}$$

t is monomials

$$J(\theta) = \frac{1}{m} \left[ \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right] \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

$$\text{Cost}(h_{\theta}(x^{(i)}), y^{(i)}) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Logistic Regression i.e. classification

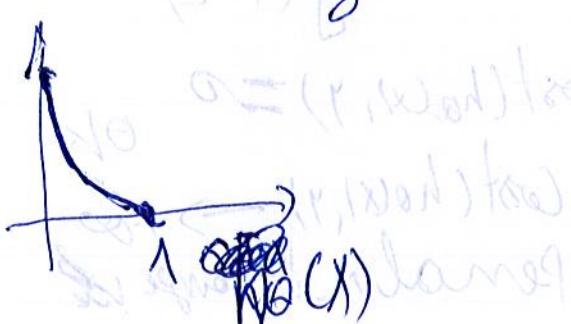
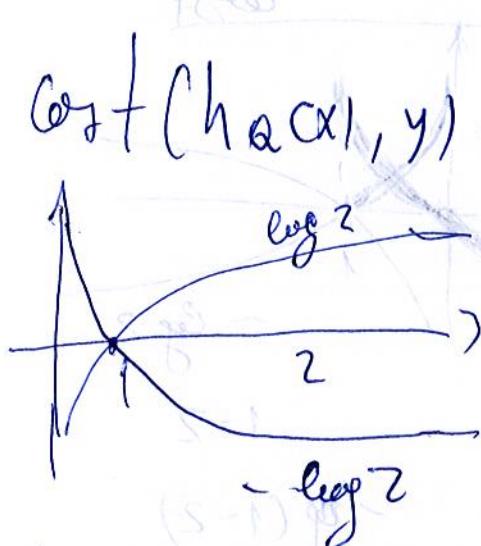
$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & y=1 \\ -\log(1-h_{\theta}(x)) & y=0 \end{cases}$$

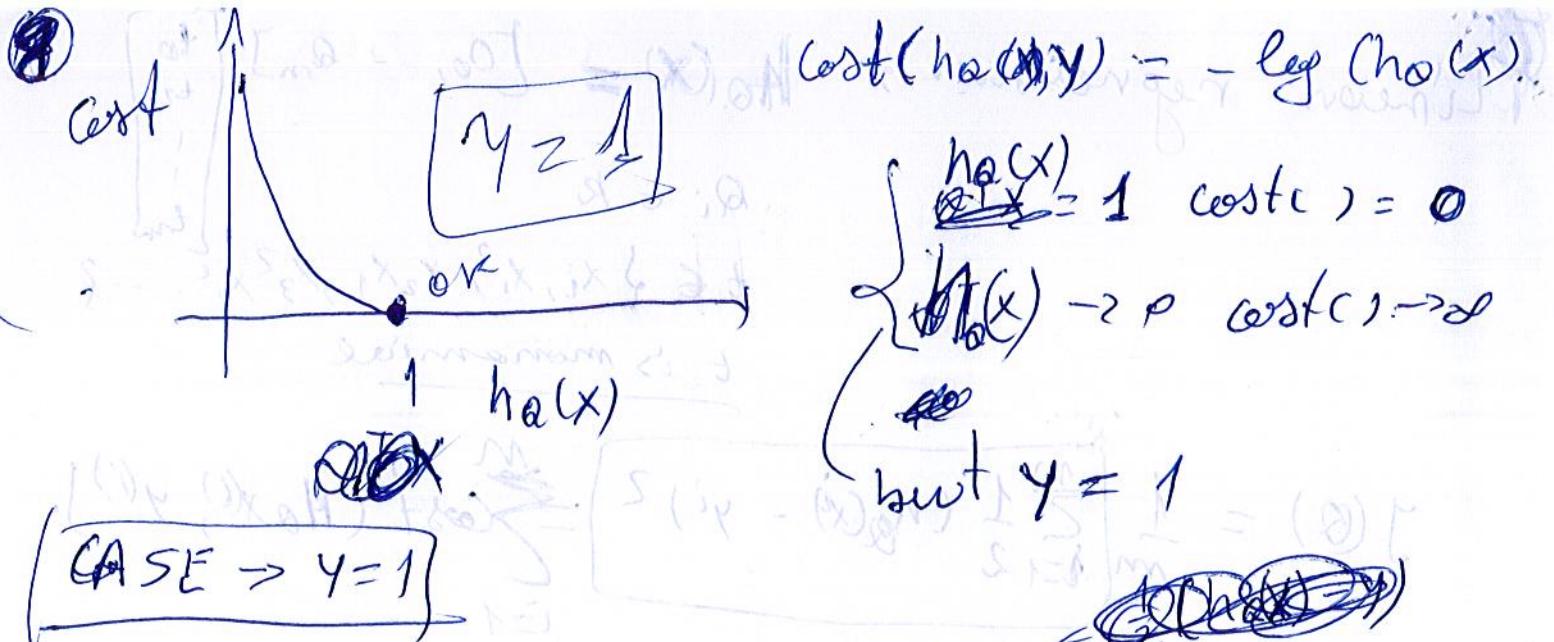
$$\text{Cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

$$0 \leq h_{\theta}(x) \leq 1, \text{ because}$$

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

SO  $h_{\theta}(x)$  still  $\approx g(\theta^T x)$   
where  $g$  is sigmoid function  
or ReLU



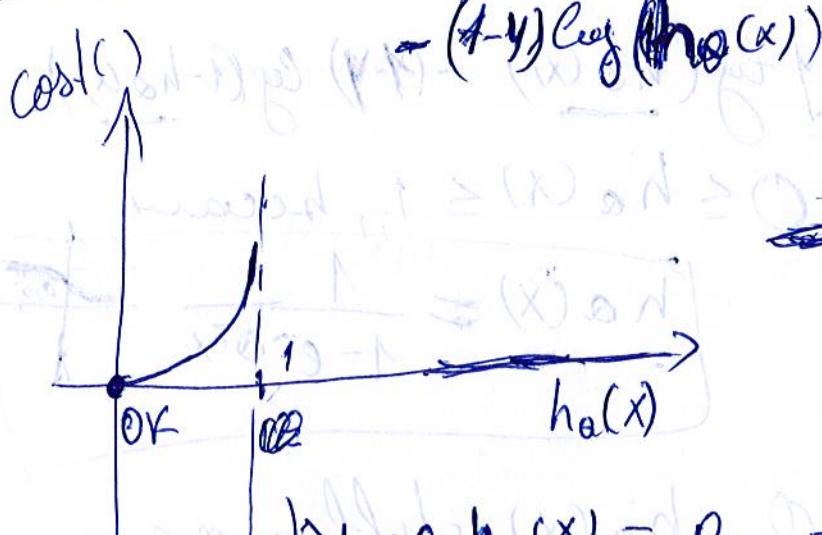
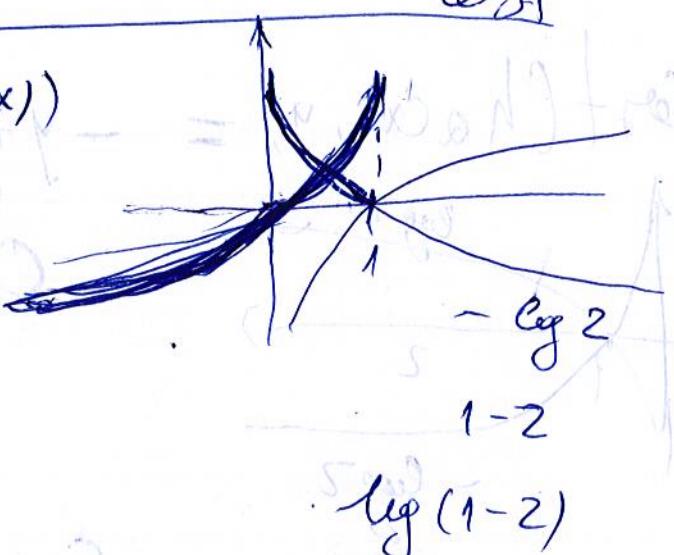


a)  $y=1, h_0(x)=1 \Rightarrow \text{cost}(h_0(x), y)=0$  is OK

b)  $y=1, h_0(x) \rightarrow \infty \Rightarrow \text{cost}(h_0(x), y) \rightarrow \infty$

a)  $P(y=1 | x; \theta) = 1$  OK

b)  $P(y=1 | x; \theta) = 0$  Penalize be larger



$y=0, h_0(x)=0 \Rightarrow \text{cost}(h_0(x), y)=0$

$y=0, h_0(x) \neq 1 \rightarrow \text{cost}(h_0(x), y) \rightarrow +\infty$  OK  
Penalize by large value

## ⑨ Logistic Regression Cost Function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

→ this is more general than logistic regression

Let  $x = x^{(i)}$ ,  $y = y^{(i)}$ ,  $y \in \{0, 1\}$  (i.e. logistic regression)  
then

$$\text{Cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$$

$$= -[y \log(h_\theta(x)) + (1-y) \log(1-h_\theta(x))]$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))]$$

## Logistic Regression Cost Function

### Gradient Descent (general notes)

6. G.D. is used to find those value of  $\theta \in \mathbb{R}^{n+1}$  for which  $J(\theta)$  achieves minimal value.

$$\text{Let } \theta \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} = \begin{bmatrix} 1 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

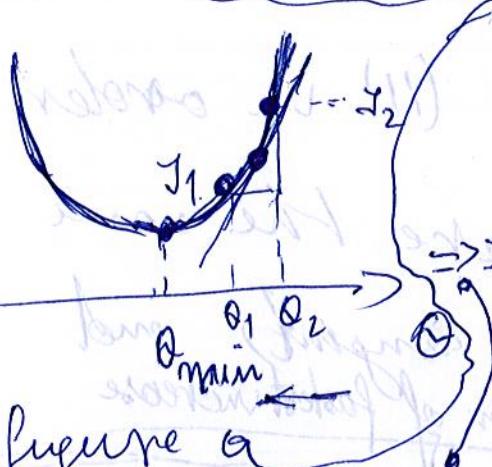


figure a

(case b)

$$\begin{aligned} \Delta J &= J_2 - J_1 \geq 0 \\ \Delta \theta &= \theta_2 - \theta_1 \geq 0 \end{aligned} \quad \frac{\Delta J}{\Delta \theta} \geq 0 \Rightarrow$$

$$\frac{dJ}{d\theta} \geq 0$$

- is positive
- decreases as  $\theta$  value when we move from right toward  $\theta_{\min}$

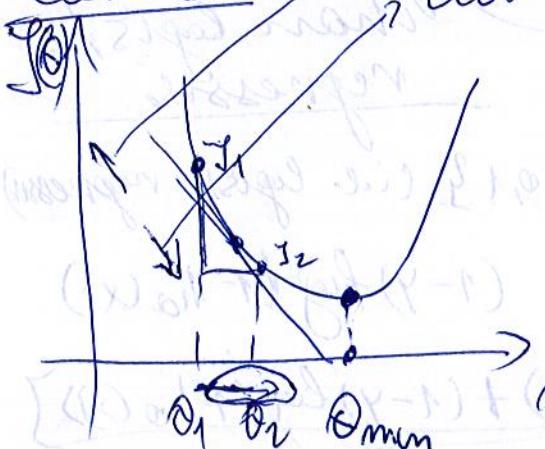
R.S.  $J'(\theta_1) = \lim_{\theta_2 \rightarrow \theta_1} \frac{J(\theta_2) - J(\theta_1)}{\theta_2 - \theta_1}$   
(i.e.  $\frac{dJ}{d\theta_1} = \lim_{\theta_2 \rightarrow \theta_1} \frac{dJ}{d\theta}$ )

V. i.e. in order to move toward  $\theta_{\min}$  we should remove from actual value of  $\theta$  a positive value

Thus we can decide  $\theta = \theta - \alpha \frac{\partial J}{\partial \theta}$  (1)

(i.e. we decrease actual  $\theta$ )

Case 5 away from min  $J(\theta)$ , closer to min  $J(\theta)$



$$\Delta J = J_2 - J_1 < 0 \quad \left. \begin{array}{l} \\ \end{array} \right\} \Delta \theta = \theta_2 - \theta_1 > 0 \Rightarrow$$

$$\Rightarrow \frac{\Delta J}{\Delta \theta} < 0 \Rightarrow \frac{\partial J}{\partial \theta} < 0$$

$$\frac{\partial J}{\partial \theta} = 0 \Leftrightarrow \theta = \theta_{\min}$$

$\frac{\partial J}{\partial \theta}$  } is negative  
decreases as absolute value  
When we move toward  $\theta_{\min}$  from left.

for this function  
we know that this local  
min,  $\theta_{\min}$ , is general  
minimum as function  
is convex.

Hence, to go toward  $\theta_{\min}$ , we add a positive value to actual value of  $\theta$  to achieve  $\theta_{\min}$ . (so we increase actual  $\theta$ )

$$\frac{\partial J}{\partial \theta} > 0 = \frac{\partial J}{\partial \theta} > 0 \Rightarrow \theta > \theta$$

thus  $\theta := \theta - \alpha \frac{\partial J}{\partial \theta}$  (2) in order to come closer to  $\theta_{\min}$ .

We see that (1) and (2) are the same

~~decrease~~  
derivative  $\frac{\partial J}{\partial \theta} \rightarrow$  shows the slope (tangent) and direction of fastest increase

More generally  $\theta \in \mathbb{R}^{n+1}$

$$\nabla J(\theta) = \left( \frac{\partial J}{\partial \theta_0}, \frac{\partial J}{\partial \theta_1}, \frac{\partial J}{\partial \theta_2}, \dots, \frac{\partial J}{\partial \theta_n} \right) \rightarrow$$

Slope gradient of the direction of fastest increase

# (11) HOW TO CALCULATE $\frac{\partial J}{\partial \theta_j}$ ?

We need  $\frac{\partial J}{\partial \theta_j}$  while tuning  $\theta_j$  ( $\forall j = 0 - n$ )  
in order to  $J(\theta)$  converge to  $\theta_{min} = [\theta_{0min}, \dots, \theta_{nmin}]$

I.e. We use  $\frac{\partial J}{\partial \theta_j}$ , for example, during G. Descent.  
or any other algorithm with similar purpose

I way → calculate directly  $\frac{\partial J}{\partial \theta_j}$   
as you do in calculus

Course or ↳ Linear regression, where

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

Let denote  $U = h_\theta(x^{(i)}) - y^{(i)}$  | i.e.  $U$  is function of  $\theta$  and  $x$  ( $y$  is const.)

$$\bullet (U^2)' = 2U \cdot U' \quad | \text{ if we derivate by } \theta$$

$$\frac{\partial (U^2)}{\partial \theta} = 2U \cdot \frac{\partial U}{\partial \theta}$$

in fact we are interested to derivate by  $\theta$

the let denote  $U(\theta) := h_\theta(x^{(i)}) - y^{(i)}$

$$\bullet (\sum U^2)' = 2 \sum U^2 \cdot U'$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m U^2(\theta, x^{(i)}) \quad ; \quad \text{cost}^{(i)} = \frac{1}{2} U^2(\theta, x^{(i)})$$

$$x^{(i)} = (x_0^{(i)}, x_1^{(i)}, x_2^{(i)}, \dots, x_m^{(i)}) \quad ; \quad U(\theta) = U(\theta, x^{(i)}) = h_\theta(x^{(i)}) - y^{(i)}$$

(12) Let's derivate initially only  $\text{cost}^{(i)} = \frac{1}{2} V^2(\theta, x^{(i)})$   
 $(\forall i=1..m)$   $x^{(i)} \in \mathbb{R}^{m+1}$   $x_0^{(i)} = 1 + i$

Let's omit superscript  $(i)$  for a moment w.r.t.

let  $x = x^{(i)} = (x_0, x_1, x_2, \dots, x_m)$   $y = y^{(i)}$

$$\frac{\partial \text{cost}}{\partial \theta_j} = \frac{\partial V^2(\theta, x)}{\partial \theta_j} = \frac{1}{2} 2V \cdot \frac{\partial V}{\partial \theta_j} = V \cdot \frac{\partial (h_\theta(x) - y)}{\partial \theta_j} =$$

$$= V \cdot \frac{\partial h_\theta(x)}{\partial \theta_j} \quad (\text{because } y \text{ is constant as we derivative by } \theta_j)$$

$$= V \cdot \frac{\partial (\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_m x_m)}{\partial \theta_j} = V \cdot (\theta_0 + \theta_1 + \dots + x_j + \theta_m)$$

$$= V \cdot x_j = (h_\theta(x) - y) \cdot x_j$$

i.e.  $\frac{\partial \text{cost}}{\partial \theta_j} = (h_\theta(x) - y) x_j \rightarrow \text{this is true for any } x^{(i)} \Rightarrow x^{(i)} = x$

$$\Rightarrow \frac{\partial \text{cost}^{(i)}}{\partial \theta_j} = h_\theta(x^{(i)}) - y^{(i)} x_j^{(i)}$$

that how we calculate in this scenario

$$\frac{\partial J}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$\Rightarrow \frac{\partial J}{\partial \theta_j}$  for Linear Regress.

where  $h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_m x_m$

(P.S. We use the formula above during gradient descent)

$$\theta_j = \theta_j - \alpha \frac{\partial J}{\partial \theta_j} \quad \theta \in \mathbb{R}^{m+1}$$

$$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$$

$$= \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix} - \alpha \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_m} \end{bmatrix}$$

$$\begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_m} \end{bmatrix} \text{ def } = \frac{\partial J}{\partial \theta} \rightarrow \text{gradient of cost } J(\theta).$$

$$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix} - \alpha \cdot \nabla$$

$\parallel \quad \parallel$   
 $\theta \in \mathbb{R}^{m+1}$

$$\nabla = \left[ \begin{array}{c} \sum_{i=1}^m (\hat{h}_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \sum_{i=1}^m (\hat{h}_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)} \\ \vdots \\ \sum_{i=1}^m (\hat{h}_\theta(x^{(i)}) - y^{(i)}) x_m^{(i)} \end{array} \right]$$

Note: For parallel update, keep  $\theta$  unchanged  
 when  $\theta = \theta_{\text{new}}$ , because we use  $\theta$  in  $\hat{h}_\theta(x)$ .

$$\begin{aligned} \nabla &= \sum_{i=1}^m [(\hat{h}_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}] = (\theta_0 x_0^{(1)} + \dots + \theta_m x_m^{(1)} - y^{(1)}) x_0^{(1)} \\ &\quad + (\theta_0 x_0^{(2)} + \dots + \theta_m x_m^{(2)} - y^{(2)}) x_0^{(2)} \\ &\quad + \dots \\ &\quad + (\theta_0 x_0^{(m)} + \dots + \theta_m x_m^{(m)} - y^{(m)}) x_0^{(m)} \\ &= (\theta^T x^{(1)} - y^{(1)}) x_0^{(1)} + \\ &\quad + (\theta^T x^{(2)} - y^{(2)}) x_0^{(2)} \\ &\quad + \dots \\ &\quad + (\theta^T x^{(m)} - y^{(m)}) x_0^{(m)} \end{aligned}$$

$$[\theta_0 \dots \dots \theta_m] \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_m^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ x_0^{(m)} & x_1^{(m)} & \dots & x_m^{(m)} \end{bmatrix} = \theta^T x$$

$$x = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}, \quad y^T = [y^{(1)} \dots y^{(m)}]$$

$$\theta^T x - y^T = [\theta^T x^{(1)} - y^{(1)}, \theta^T x^{(2)} - y^{(2)}, \dots, \theta^T x^{(m)} - y^{(m)}]$$

$$j=0: (\theta^T x - y^T) \cdot x_0^{(i)} \quad \circledcirc (X^T)^{(i)}$$

$$\boxed{\frac{1}{m} (\theta^T x - y^T) x^T = \nabla}$$

$$\theta = \theta - \alpha \nabla$$

j20

(4)

$$\theta = \theta - \alpha \cdot \nabla J$$

$$[\nabla J] = \frac{1}{m} (\theta^T x - y^T) x^T$$

where  $\nabla J = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_m} \end{bmatrix}$

While to calculate  
we accomplish  
 $\frac{1}{m} (\theta^T x - y^T) x^T$

We use this  
formula to  
calculate  $\nabla J$

When  $\theta^T = [\theta_0, \dots, \theta_n]; y^T = [y^{(1)}, \dots, y^{(m)}]$

$$X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(m)} \end{bmatrix} \in \mathbb{R}^{(m+1) \times m}$$

$$X^T = \begin{bmatrix} x^{(1)} & x^{(0)} \\ x^{(2)} & x^{(1)} \\ x^{(3)} & x^{(2)} \\ \vdots & \vdots \\ x^{(m)} & x^{(m-1)} \end{bmatrix} \in \mathbb{R}^{m \times (m+1)}$$

$$x^T x = x^T x = x^T x$$

$$x^T x = x^T x = x^T x$$

$$x^T \in \mathbb{R}^{m \times (m+1)}$$

Cali:  $\theta = \theta - \alpha \cdot \nabla J$  where  $\nabla J = \frac{1}{m} \cdot (\theta^T x - y^T) x^T$

$$\nabla J - \theta = \theta$$

$$\nabla J = \frac{1}{m} \cdot x^T (y - \theta^T x)$$

Note: If we define: ~~Denote~~  $\Omega = \begin{bmatrix} \alpha_0 \\ 1 \\ \vdots \\ \alpha_m \end{bmatrix}$

but as  $\Omega = \begin{bmatrix} \alpha_0 \\ 1 \\ \vdots \\ \alpha_m \end{bmatrix}$  and  $y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^{m \times 1}$

and  $x^{(i)} \in \mathbb{R}^{m+1}$

$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_m^{(i)} \end{bmatrix}$

This is transpose version of what I used on previous pages

$X = \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \Rightarrow X^T = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | & | \end{bmatrix}$

$X \in \mathbb{R}^{m \times (m+1)}$

$X^T \in \mathbb{R}^{(m+1) \times m}$

Then the previous equality  $\nabla J = \frac{1}{m} (\Omega^T X - Y^T) X^T$   
 will look different.  $X$  here transpose version of the one in previous pages.

$$\begin{aligned} i=0 & \sum_{j=1}^m (\Omega^T x^{(i)} - y^{(i)}) x_0^{(i)} \\ i=1 & \sum_{j=1}^m (\Omega^T x^{(i)} - y^{(i)}) x_1^{(i)} \\ i=2 & \sum_{j=1}^m (\Omega^T x^{(i)} - y^{(i)}) x_2^{(i)} \end{aligned}$$

$$X = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & x_2^{(1)} & \dots & x_m^{(1)} \\ x_0^{(2)} & x_1^{(2)} & x_2^{(2)} & \dots & x_m^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_0^{(m)} & x_1^{(m)} & x_2^{(m)} & \dots & x_m^{(m)} \end{bmatrix}$$

$$x_0 \cdot \Omega = \sum_{j=1}^m x_j^{(i)} \cdot \alpha_j$$

$$= X \cdot \Omega \in \mathbb{R}^m$$

i-th column of  $X \cdot \Omega$   $= \Omega^T x^{(i)} \mid x_0^{(i)} \text{ where } i: 1 \dots m \text{ is } x_0^{(i)} \text{ to } x_m^{(i)}$   
 The  $i$ -th column of  $X$  or  
 $i$ -th row of  $X^T$

1) Let  $\beta$  be fixed then  $x_j^{(0 \dots m)}$  is the  $j$ -th column of  $X$ .  
 or the  $j$ -th row of  $X^T$ .

$$\sum (\beta^T x^{(i)} - y^{(i)}) x_{(j)}^{(i)} = ?$$

$$\cancel{x \cdot \beta - y} - \begin{bmatrix} x_0^{(1)} \beta \\ x_1^{(2)} \beta \\ \vdots \\ x_m^{(m)} \beta \end{bmatrix} - \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} = \begin{bmatrix} x_0^{(1)} \beta - y^{(1)} \\ x_1^{(2)} \beta - y^{(2)} \\ \vdots \\ x_m^{(m)} \beta - y^{(m)} \end{bmatrix} = B$$

2)  $\beta^T x^{(i)} - y^{(i)}$  is the  $i$ -th row of  $B$

$$y = 0 \dots m$$

$$\sum_{i=1}^m (\beta^T x^{(i)} - y^{(i)}) x_{(i)}^{(i)} = \begin{bmatrix} x_0^{(1)} \beta - y^{(1)} \\ x_1^{(2)} \beta - y^{(2)} \\ \vdots \\ x_m^{(m)} \beta - y^{(m)} \end{bmatrix} = B - \cancel{x \cdot \beta - y}$$

$$X^T = \begin{bmatrix} x_0^{(1)} & x_0^{(2)} & x_0^{(3)} & \cdots & x_0^{(m)} \\ x_1^{(1)} & x_1^{(2)} & x_1^{(3)} & \cdots & x_1^{(m)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_n^{(1)} & x_n^{(2)} & x_n^{(3)} & \cdots & x_n^{(m)} \end{bmatrix}$$

$$\sum_{i=0}^m (x^{(i)} \beta - y^{(i)}) x_{(i)}^{(i)} = \begin{bmatrix} (x_0^{(1)} \beta - y^{(1)}) x_{(1)}^{(1)} \\ (x_0^{(2)} \beta - y^{(2)}) x_{(2)}^{(2)} \\ \vdots \\ (x_0^{(m)} \beta - y^{(m)}) x_{(m)}^{(m)} \end{bmatrix} = \sum_{i=0}^m (x^{(i)} \beta - y^{(i)}) = \cancel{x \cdot \beta - y} = B$$

$$\frac{1}{m} \left[ \sum_{i=1}^m (\beta^T x^{(i)} - y^{(i)}) x_{(i)}^{(i)} + \sum_{j=1}^m (\beta^T x^{(j)} - y^{(j)}) x_j^{(i)} + \sum_{i=1}^m (\beta^T x^{(i)} - y^{(i)}) x_m^{(i)} \right] = \cancel{\frac{1}{m} X^T B} = \cancel{\frac{1}{m} X^T (x \cdot \beta - y)} = \cancel{\frac{1}{m} X^T (x \cdot \beta - y)} = \nabla J$$

$$\frac{1}{m} X^T \cdot (\cancel{\frac{1}{m} X^T (x \cdot \beta - y)}) = \cancel{\frac{1}{m}}$$

$$\nabla J = \cancel{\frac{1}{m} X^T (x \cdot \beta - y)}$$

$$(1) \nabla J = \frac{1}{m} X^T (X \cdot \theta - y)$$

if  $x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_m^{(i)} \end{bmatrix}$  and  $X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}$

$$(2) \nabla J = \frac{1}{m} (\theta^T X' - y^T) X^T$$

if  $X' = X^T \Rightarrow X = (X')^T$

Check if our calculation was ok.

$$\frac{1}{m} (X^T (X \cdot \theta - y))^T = \frac{1}{m} (X \cdot \theta - y)^T \cdot (X^T)^T = \frac{1}{m} (X \cdot \theta - y)^T \cdot X = \nabla$$

$$\frac{1}{m} (\theta^T X^T - y^T) \cdot X = \frac{1}{m} (\theta^T X' - y^T) \cdot X^T \stackrel{(1)}{=} (2)$$

P.S. If we write  $h_{\theta}(x) = \begin{bmatrix} h_{\theta}(x^{(1)}) \\ \vdots \\ h_{\theta}(x^{(m)}) \end{bmatrix}$ , then

$$\nabla J = \boxed{\frac{1}{m} X^T (X \cdot \theta - y)} = \boxed{\frac{1}{m} X^T (h_{\theta}(x) - y)}$$

Note:  $\nabla J = \frac{\partial J}{\partial \theta} = \frac{\partial J(\theta)}{\partial \theta} = \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_m} \end{bmatrix} = \left( \frac{\partial J(\theta)}{\partial \theta_1}, \dots, \frac{\partial J(\theta)}{\partial \theta_m} \right)$

$$\boxed{\theta := \theta - \alpha \nabla J} \quad \equiv \quad \boxed{\theta := \theta - \frac{\alpha}{m} X^T (h_{\theta}(x) - y)}$$

to your favorite in standard form given (3) as  
 $(P - \theta \cdot X)^T X = \nabla$

(18)

# SUMMARY

How to calculate  $\frac{\partial J}{\partial \theta_j}$  for linear

Regression of type  $h_{\theta}(x) = \sum_{j=0}^m \theta_j x_j$ ?  $x_0 = 1$

I way

Like in calculus exercises, we derivat and find out

$$\frac{\partial J}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

(1)

$$\nabla J = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_m} \end{bmatrix}$$

Note:

$\nabla J$  → is gradient of cost function

Gradient shows the direction of fastest increase of function.  
 (This is proven in calculus)  
 at a certain point.

so we have point (current point)  $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$

and we derivate at this point and we get formula (1) which we apply to concrete values of  $\theta, x, y$  and we get a concrete value which we use to update  $\theta$  with the purpose of achieving (getting to)  $\theta_{\min}$ .

The (2) using (1) and matrix calculus, we can write and calculate in compact way as following:

$$\nabla J = \frac{1}{m} x^T (x \cdot \theta - y)$$

$\theta \in \mathbb{R}^{m+1}, y \in \mathbb{R}^m$   
 $x \in \mathbb{R}^{m \times (m+1)}$  Then we calculate  $\theta = \theta - \alpha \nabla J$

⑨ Case b Logistic Regression  $\rightarrow \frac{\partial J(\theta)}{\partial \theta_j} = ?$

$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}} \quad (1)$$

$$g(u) = \frac{1}{1+e^{-u}} \quad (2)$$

sigmoid function  
 $g'(u) = g(u)(1-g(u))$

$$\frac{e^u}{e^u + 1} \quad (3)$$

$$g(\theta^T x) = \frac{e^{\theta^T x}}{e^{\theta^T x} + 1}$$

$$\frac{e^{\theta^T x}}{e^{\theta^T x} + 1} = \frac{y}{(1-y)} \cdot g(\theta^T x) \cdot g'(u) \cdot u'(x) \quad (4)$$

$$\theta \in \mathbb{R}^{n+1}, x \in \mathbb{R}^{(n+1)}, x^{(i)} \in \mathbb{R}^{(n+1)}$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))] \quad (4)$$

$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}} = g(\theta^T x) \quad h_{\theta}(x) = g(\theta^T x) \quad (5)$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(g(\theta^T x^{(i)})) + (1-y^{(i)}) \log(1-g(\theta^T x^{(i)}))] \quad (6)$$

$$\log(z) = \ln(z) = \log_e(z)$$

$$\text{Cost}^{(i)} = -[y^{(i)} \log(g(\theta^T x^{(i)})) + (1-y^{(i)}) \log(1-g(\theta^T x^{(i)}))]$$

Simplex

$$\text{Cost}^{(i)} = -[y \log(g(\theta^T x)) + (1-y) \log(1-g(\theta^T x))]$$

$$\text{Cost}_a^{(i)} = y \log g(\theta^T x) ; \quad \text{Cost}_b^{(i)} = (1-y) \log(1-g(\theta^T x))$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \text{Cost}^{(i)} = \frac{1}{m} \sum_{i=1}^m (\text{Cost}_a^{(i)} + \text{Cost}_b^{(i)})$$

$$\frac{\partial \text{Cost}_a^{(i)}}{\partial \theta_j} = \frac{\partial [y \log(g(\theta^T x))]}{\partial \theta_j} = \frac{y}{g(\theta^T x)} \cdot \frac{\partial g(\theta^T x)}{\partial \theta_j} = \frac{y}{g(\theta^T x)} \cdot (1-g(\theta^T x)) \cdot g'(x)$$

$$\frac{\partial g(\theta^T x)}{\partial \theta_j} = y(1-g(\theta^T x)) \cdot x_j = yx_j - yg(\theta^T x)x_j$$

$$= yx_j - yh_{\theta}(x)x_j \quad (9)$$

$$G(z) = h_{\theta}(x) = g(z) \quad z = \theta^T x$$

$$\begin{aligned}
 \frac{\partial \text{Cost}_b^{(i)}}{\partial \theta_j} &= \frac{\partial (1-y) \log(1-g(\theta^T x))}{\partial \theta_j} = \frac{1-y}{1-g(\theta^T x)} \cdot \frac{(1-g(\theta^T x))}{\partial \theta_j} = \\
 &= \frac{1-y}{1-g(\theta^T x)} \times -\frac{\partial g(\theta^T x)}{\partial \theta_j} = \frac{y-1}{1-g(\theta^T x)} \cdot (1-g(\theta^T x))g(\theta^T x) \cdot \frac{\partial (\theta^T x)}{\partial \theta_j} = \\
 &= (y-1)g(\theta^T x) \cdot x_j = (y-1)h_\theta(x) \cdot x_j = \\
 &= y h_\theta(x) \cdot x_j - h_\theta(x) x_j \quad (10)
 \end{aligned}$$

$$(9) + (10) = \frac{\partial (\text{Cost}_a^{(i)} + \text{Cost}_b^{(i)})}{\partial \theta_j} = \frac{y x_j - h_\theta(x) x_j}{m} = \text{Cost}^{(i)}$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{m} \left[ \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)} \right] \quad \Rightarrow \quad \frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m [(h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}]$$

see on previous pages  
how we go to  
matrix form.

$$\nabla J(\theta) = \frac{\partial J(\theta)}{\partial \theta} = \frac{1}{m} X^T [h_\theta(X) - Y]$$

because  $\nabla J(\theta) = \frac{\partial J(\theta)}{\partial \theta_j}$

$$\begin{aligned}
 \nabla J(\theta) &= \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_m} \end{bmatrix} = \begin{bmatrix} \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)} \\ \vdots \\ \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_m^{(i)} \end{bmatrix}
 \end{aligned}$$

R.S

$$\begin{aligned}
 f(z) &= \frac{1}{1+e^{-z}} \stackrel{f' = 1+e^{-z}}{\Rightarrow} f'(z) = -\frac{f'}{f^2} = -\frac{(1+e^{-z})'}{(1+e^{-z})^2} = \frac{e^{-z}}{(1+e^{-z})^2} \\
 &= \frac{e^{-z}}{(1+e^{-z})} \cdot \frac{1}{(1+e^{-z})} = \frac{e^{-z}}{1+e^{-z}} \cdot f(z) = \frac{1-1+e^{-z}}{1+e^{-z}} f(z) = e^{-z}
 \end{aligned}$$

$$\begin{aligned}
 &= 1 - \frac{e^{-z}}{1+e^{-z}} \cdot f(z) = (1-f(z))f(z); \text{ if } z = z(x), \text{ then} \\
 &\quad \frac{\partial f}{\partial x} = f'(z(x)) = (1-f(z))f'(z) \cdot z'(x)
 \end{aligned}$$

$$\left(\frac{1}{1+e^{-z}}\right)^l = \frac{e^{-z}}{(1+e^{-z})^2} = \frac{e^{-z}}{1+e^{-z}} \cdot \frac{1}{1+e^{-z}} = \left(\frac{e^{-z}}{1+e^{-z}}\right)^2 \stackrel{1+e^{-z}-1}{1+e^{-z}} \circ g(z)$$

$$1 - g(z) = \frac{1+e^{-z}}{1+e^{-z}} - \frac{1}{1+e^{-z}} = \frac{e^{-z}}{1+e^{-z}} \quad (1-g(z))$$

$$1 - \frac{1}{1+e^{-z}} = \frac{e^{-z}}{1+e^{-z}} = e^{-z} \cdot \frac{1}{1+e^{-z}}$$

$$1 - g(z) = e^{-z} \circ g(z)$$

$$\frac{e^{-z}}{1+e^{-z}} = \frac{1+e^{-z} - 1}{1+e^{-z}} = \frac{1+e^{-z}}{1+e^{-z}} - \frac{1}{1+e^{-z}} = 1 - \frac{1}{1+e^{-z}}$$

$$\frac{1}{g(\theta^T x)} \frac{\partial f(\theta^T x)}{\partial \theta_j} = \frac{1}{(1-g(\theta^T x))} \frac{\partial g(\theta^T x)}{\partial \theta_j} \frac{\partial \theta^T x}{\partial \theta_j}$$

$$(1 - g(\theta^T x)) \cdot x_j$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = 0$$

$$\frac{1}{m} \sum_{i=1}^m x^T(x\cdot\theta - y) = 0 \Leftrightarrow x^T x \cdot \theta - x^T y = 0$$

$$x^T x \cdot \theta = x^T y \quad \Rightarrow \boxed{\theta = (x^T x)^{-1} x^T y}$$

Normal equations

$$\frac{\partial J}{\partial \theta_3} = (\hat{y}^{(i)} - y^{(i)}) \theta^{(3)}$$

$$\frac{\partial J}{\partial \theta^{(2)}} = (\hat{y}^{(i)} - y^{(i)}) \theta^{(3)} \cdot g'(z^{(2)}) \circ \theta^{(2)}$$

$$\frac{\partial J}{\partial \theta^{(1)}} = (\hat{y}^{(i)} - y^{(i)}) \theta^{(3)}$$

$$f(x) = ((x^T \theta) \hat{y} - b)$$

## Summary

$\boxed{\frac{\partial J}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x) - y) x_j}$  → For both linear regression and logistic regression  $\frac{\partial J}{\partial \theta_j}$  looks the same, but the difference is on  $h_\theta(x)$ .

Linear regression  $\rightarrow h_\theta(x) = \theta^T x$  (or  $x\theta$ )

Logistic regression  $\rightarrow h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$  (or  $\frac{1}{1 + e^{x\theta}}$ )

$$h_\theta(x) = \sigma(\theta^T x)$$

$$h_\theta(x) = \sigma(z), \text{ where } z = \theta^T x \quad (z = \theta^T x + b)$$

$$h_\theta(x) = f(z)$$

$$\text{where } f(z) = \sigma(z)$$

$$b = \theta_0$$

