

Q2 - a

To derive the gradient of cost function. The cost function is the log likelihood

$$l(\theta) = \sum_{i=1}^n \{-\log(1 + e^{-\theta x_i}) + (y_i - 1)\theta x_i\}$$

$$\frac{\partial l(\theta)}{\partial \theta} = \sum_{i=1}^n (y_i x_i - x_i - \left(\frac{e^{-\theta x_i}}{1 + e^{-\theta x_i}}\right) x_i)$$

$$= \sum_{i=1}^n (y_i - \frac{1}{1 + e^{-\theta x_i}}) x_i = F_{dx} \quad \text{this is our gradient decent steps.}$$

for gradient decent, we need define the steps which is the formula above  
And we need the learning rate  $\lambda$ .

the pseudo code will be:

Procedure GD (  $F_{dx}$ ,  $\theta(0)$  )

$\theta \leftarrow \theta(0)$

while not converged do

$\theta \leftarrow \theta + \lambda F_{dx}(\theta)$

return  $\theta$

$F_{dx}$  is a function that returns a matrix of form  $F_{dx}(\theta) = \begin{bmatrix} \frac{\partial l(\theta)}{\partial \theta_1} \\ \frac{\partial l(\theta)}{\partial \theta_2} \\ \vdots \\ \frac{\partial l(\theta)}{\partial \theta_n} \end{bmatrix}$

$\theta \in \mathbb{R}$

$\theta(0)$  is a Random initialization

the returned  $\theta$  is a point in which the log likelihood is the optimum. We run the Algorithm for each and points.