ISMET OKAN CELIK    CWID: 11472265
HOMEWORK-2

## Question-1:

$$E(w) = MSE(w) + \frac{\partial}{2} \sum_{i=1}^{m} w_i^2$$

$$= \frac{1}{m} \sum_{i=1}^{m} (h_w(x) - y)^2 + \frac{\partial}{2} \sum_{i=1}^{m} w_i^2$$

$$= \frac{1}{m} (w^T x_1 - y_1 \quad w^T x_2 - y_2 \cdots w^T x_m - y_m) \begin{pmatrix} w^T x_1 - y_1 \\ w^T x_2 - y_2 \\ | \\ w^T x_m - y_m \end{pmatrix} + \frac{\partial}{2} w^T w$$

we can ignore $\frac{1}{m}$ and $\frac{1}{2}$ because they are constant values.

$$= \left( w^T (x_1 \quad x_2 \cdots x_m) - (y_1 \quad y_2 \cdots y_m) \right) \begin{pmatrix} w^T x_1 - y_1 \\ w^T x_2 - y_2 \\ | \\ w^T x_m - y_3 \end{pmatrix} + \partial w^T w$$

$$= \left( w^T X^T - Y^T \right) \cdot \left( w^T X^T - Y^T \right)^T + \partial w^T w$$

$$= \left( w^T X^T - Y^T \right) \cdot (Xw - Y) + \partial w^T w$$

$$= w^T X^T X w - w^T X^T Y - Y^T X w + Y^T Y + \partial w^T w$$

$$E(w) = w^T w (X^T X + \partial) - 2 w^T X^T Y + Y^T Y$$

If we take the partial derivative of the equation in order to optimize it

$$\frac{\partial E(w)}{\partial w} = 2w (X^T X + \partial) - 2 X^T Y = 0$$

$$= 2 \left( w (X^T X + \partial) - X^T Y \right) = 0$$

$$= w (X^T X + \partial) = X^T Y$$

$$w = \frac{X^T Y}{(X^T X + \partial)} \quad \rightarrow \quad w = (X^T X + \partial)^{-1} X^T Y$$

Property of Identity Matrix $I$: $\boxed{I \cdot \partial = \partial}$

Proof: $\boxed{w = (\partial I + X^T X)^{-1} \cdot X^T Y}$

# Question 2:

## 1)



$$\theta_1^T \cdot X_0 = z_1 \qquad \theta_1^T \cdot X_1 = z_1$$
$$\theta_2^T \cdot X_0 = z_2 \qquad \theta_2^T \cdot X_1 = z_2$$
$$\vdots$$
$$\theta_k^T \cdot X_0 = z_k \qquad \theta_k^T \cdot X_1 = z_k$$

As we can see on above for $X_0$ value we need to estimate k number of value, also it is same for $X_1$ value. We have k+1 number of x.

In order to learn Softmax Regression model we need to estimate $k \cdot (n+1)$ paramater

## 2)

$$\hat{p}_k = \delta(\delta_k(x))_k = \frac{\exp(\delta_k(x))}{\sum\limits_{j=1}^{k} \exp(S_j(x))} \quad , \quad S_k(x) = \theta_k^T \cdot X$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^{m} \sum_{k=1}^{k} y_k^{(i)} \cdot \log(\hat{p}_k^{(i)})$$

$$= -\frac{1}{m} \sum_{i=1}^{m} \sum_{k=1}^{k} y_k^{(i)} \cdot \log\left[\frac{\exp(\theta_k^T X^{(i)})}{\sum\limits_{j=1}^{k} \exp(\theta_j^T \cdot X^{(i)})}\right]$$

$$= -\frac{1}{m} \sum_{i=1}^{m} \sum_{k=1}^{k} y_k^{(i)}\left(\log(\exp(\theta_k^T X^{(i)})) - \log\left(\sum_{j=1}^{k} \exp(\theta_j^T \cdot X^{(i)})\right)\right)$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^{m} \sum_{k=1}^{k} y_k^{(i)} \cdot \theta_k^T \cdot X^{(i)} + \frac{1}{m}\sum_{i=1}^{m} \sum_{k=1}^{k} y_k^{(i)} \cdot \log\left(\sum_{j=1}^{k}\exp(\theta_j^T \cdot X^{(i)})\right)$$

$$\frac{\partial J(\theta)}{\partial \theta} = -\frac{1}{m}\sum_{i=1}^{m} y_k^{(i)} \cdot X^{(i)} + \frac{1}{m}\sum_{i=1}^{m} 1 \cdot \frac{1}{\sum\limits_{j=1}^{k} \exp(\theta_j^T \cdot X^{(i)})} \cdot \exp(\theta_k^T X^{(i)}) \cdot X^{(i)}$$

$$\nabla J(\theta) = \frac{1}{m}\sum_{i=1}^{m}\left(\underbrace{\frac{\exp(\theta_k^T \cdot X^{(i)})}{\sum\limits_{j=1}^{k}\exp(\theta_j^T \cdot X^{(i)})}} - y_k^{(i)}\right) \cdot X^{(i)}$$

$$\nabla J(\theta) = \frac{1}{m}\sum_{i=1}^{m}\left(\hat{p}_k - y_k^{(i)}\right) \cdot X^{(i)}$$