

1. According to the Product Rule $p(X, Y) = p(Y|X)p(X)$,
 $p(X, Y) = p(X|Y)p(Y)$
 thus, $p(X|Y)p(Y) = p(Y|X)p(X)$

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$

And according to the Sum Rule: $p(X) = \sum p(X, Y)$

$$p(X) = \sum p(X|Y)p(Y)$$

Thus, the Bayes' Theorem is proved.

Most of decisions in life are made with incomplete information, and we have only limited information to training our model. Bayes' Theorem can predict probabilities based on past data by converting posterior probability to likelihood and prior probability.

2. $MSE(w) = \frac{1}{m} \sum_{i=1}^m (w^T \cdot x^{(i)} - y^{(i)})^2$
 $E(w) = MSE(w) + \frac{\lambda}{2} \sum_{i=1}^m w_i^2$

$$= \frac{1}{m} \sum_{i=1}^m (w^T \cdot x^{(i)} - y^{(i)})^2 + \frac{\lambda}{2} \sum_{i=1}^m w_i^2$$

in the linear regression, we let

$$\frac{\partial [(w^T \cdot x - y)^2]}{\partial w} = 0$$

$$\frac{\partial [(w^T \cdot x - y)(w^T \cdot x - y)]}{\partial w} = 0 \quad \text{so, } w^T x^T x - 2y^T x = 0$$

$$\text{thus, } \hat{w} = (x^T \cdot x)^{-1} \cdot x^T \cdot y$$

thus, in the ridge regression we let

$$\frac{\partial [(w^T \cdot x - y)(w^T \cdot x - y) + \frac{\lambda}{2} w^2]}{\partial w} = 0$$

$$\text{so, } w^T x^T x - 2y^T x + \lambda I w = 0$$

$$\text{thus, } w = (x^T \cdot x + \lambda I)^{-1} \cdot x^T \cdot y$$

3. 1) We need to ^{estimate} ~~parameter~~ one parameter, it is θ .

2)
$$P_k = \frac{\exp(S_k(x))}{\sum_{j=1}^K \exp(S_j(x))} \quad S_k(x) = \theta^T \cdot x$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(p_k^{(i)})$$

$$\Delta_{\theta_k} J(\theta) = \frac{\partial J(\theta)}{\partial p_k} \cdot \frac{\partial p_k}{\partial S_k} \cdot \frac{\partial S_k}{\partial \theta_k}$$

$$\frac{\partial J(\theta)}{\partial p_k} = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \frac{1}{p_k^{(i)}}$$

$$\frac{\partial p_k}{\partial S_k} = \frac{\partial \left[\frac{\exp(S_k(x))}{\sum_{j=1}^K \exp(S_j(x))} \right]}{\partial S_k} = \frac{\exp(S_k(x)) \cdot \sum_{j=1}^K \exp(S_j(x)) - \exp(S_k(x))^2}{\left[\sum_{j=1}^K \exp(S_j(x)) \right]^2}$$

$$= \frac{\exp(S_k(x))}{\sum_{j=1}^K \exp(S_j(x))} \cdot \left(1 - \frac{\exp(S_k(x))}{\sum_{j=1}^K \exp(S_j(x))} \right)$$

$$= p_k \cdot (1 - p_k)$$

$$\frac{\partial S_k}{\partial \theta_k} = x^{(i)}$$

$$\begin{aligned} \therefore \Delta_{\theta_k} J(\theta) &= -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \frac{1}{p_k^{(i)}} [p_k^{(i)} \cdot (1 - p_k^{(i)})] \cdot x^{(i)} \\ &= -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} (1 - p_k^{(i)}) \cdot x^{(i)} \end{aligned}$$

$\because y_k^{(i)} = 1$ if the i th instance belongs to class k ,

$$\therefore \Delta_{\theta_k} J(\theta) = \frac{1}{m} \sum_{i=1}^m (p_k^{(i)} - y_k^{(i)}) \cdot x^{(i)}$$