# Solutions:

**1.** Let $A$ and $B$ be two events,

if the conditional probability of $A$ given $B$ has already occurred is $P(A \mid B)$

Baye's Theorem states

$$P(B \mid A) = \frac{P(A \mid B) P(B)}{P(A)}$$

$$\boxed{= \frac{P(A \mid B) P(B)}{P(A \mid B) P(B) + P(A \mid B^c) P(B^c)}} \to eq^n \, (1)$$

**Proof :**

we know $\quad P(A \mid B) = \dfrac{P(A \cap B)}{P(B)} \quad \to \boxed{P(A \cap B) = P(A \mid B) P(B)} \to eq^n \, (2)$.

Now, $\quad P(A) = P((A \cap B) \cup (A \cap B^c))$

$$= P(A \cap B) + P(A \cap B^c) \to (3)$$

from $eq^n$ (2) and (3), we can say that

$$P(A) = P(A \mid B) P(B) + P(A \mid B^c) \cdot P(B^c) \quad \left[\substack{\text{substituting} \\ \text{results from (2)}}\right]$$

Now, $\because (A \cap B) \, \& \, (A \cap B^c)$ are ~~and base~~ mutually exclusive, we have,

$$P(B \mid A) = \frac{P(B \cap A)}{P(A)} = \frac{P(A \cap B)}{P(A)} \quad \left[\substack{\because A \cap B \\ = B \cap A}\right]$$

from (2), we have $\quad P(A \cap B)$

$$P(A \cap B) = \frac{P(A \mid B) P(B)}{P(A)} \qquad \frac{P(A \mid B) P(B)}{P(A)}$$

**Hence,** $\quad \boxed{P(B \mid A) = \dfrac{P(A \mid B) P(B)}{P(A)}} \qquad$ QED.

Bayes Theorem provides a method to visualize relationship between data and a model. A machine Learning algorithm is aimed at decyphering the structured relationship between the data and a model.

Thus, Bayes Theorem of cos for computing conditional probability is very useful in ML, in determing the methods which will make more mathematically accurate prediction.

2.

Using MSE as the cost function, and minimizing it, we found the normal equation,

Similarly, to find the Ridge Regression closed form solution, we must

$\longrightarrow$ Determine the value of $w$ for which Corresponding cost is minimum.

Given: $X = \begin{bmatrix} x^1 \\ x^2 \\ \vdots \\ x^m \end{bmatrix}$ is the i/p data matrix,

$x^i = (1, x_1, x_2, ..., x_n) \longrightarrow i^{th}$ sample

$y = (y^1, y^2, ..., y^m) \longrightarrow y^j =$ measurement of hyp, $h_w(x)$ for $j^{th}$ sample

$h_w(x) = w_0 + w_1 x_1 + ... w_n x_n$

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

Finding partial derivate wrt $w$.

$$\frac{\partial E(w)}{\partial w} = \frac{\partial}{\partial w} \left[ \sum_i (w^T x^i - y^i)^2 + \lambda \sum_i w_i^2 \right]$$

$$= 2 \sum_i (w^T x^i - y^i) + 2\lambda \sum_i w_i$$

Premultiplication with $X^T$ $= 2(X^T X)w - X^T y) + 2\lambda I w \begin{bmatrix} \because \lambda \text{ is a} \\ \text{scalar} \end{bmatrix}$

Now, taking second derivative, we find.

$$\frac{\partial^2 (E(\omega))}{\partial \omega^2} = 2X^T x + 2\lambda I \rightarrow \cancel{X^T}$$

$$= 2|X|^2 + 2\lambda I \rightarrow [\because A^T A = |A|^2]$$

so for $X^T x$,

$$\therefore \frac{\partial^2 (E(\omega))}{\partial \omega^2} > 0 \text{ for } \lambda > 0 . \quad [\because I \text{ is} > 0]$$

Hence, equating first derivative wrt '$\omega$' will give us $\omega$ for which $E(\omega)$ or the cost is minimum.

$$\therefore \frac{\partial E(\omega)}{\partial \omega} = 0 = 2X^T x\omega + 2\lambda I\omega - 2X^T y$$

By the commutative property. $\Rightarrow$ $X^T y = X^T x\omega + 2\lambda I\omega$.

Premultiplication with
$(X^T x + \lambda I)^{-1}$ on both sides. $X^T y \text{ or } X^T x$

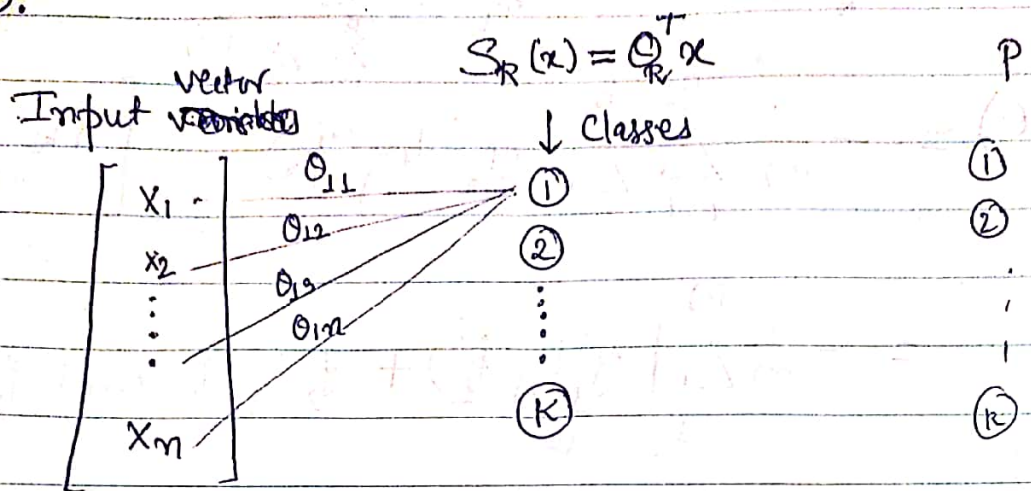we get $(X^T x + \lambda I)^{-1} X^T y = (\lambda I + X^T x)^{-1}.(\lambda I + X^T x)\omega$

$$= I \cdot \omega$$

$$= \omega .$$

$$\therefore \omega_{(for min)(cost)} = (\lambda I + X^T x)^{-1}. X^T y .$$

$\therefore$ for Ridge reg$^n$, $(\lambda I + X^T x) X^T y$ is the closed form solution.

## 3.

$$S_k(x) = \theta_k^T x$$

Input vector ~~variable~~

$$
\begin{bmatrix}
X_1 \\
X_2 \\
\vdots \\
X_n
\end{bmatrix}
$$

$\theta_{11}$
$\theta_{12}$
$\theta_{13}$
$\theta_{1n}$

↓ Classes

① 
② 
⋮ 
Ⓚ

P

① 
② 
⋮ 
Ⓚ

$$P_k = \frac{e^{S_k(x)}}{\sum\limits_{j=1} e^{S_j(x)}}$$

1) So for each class, in order to learn this softmax Regression model, we need to estimate 'm' parameters; for the m corresponding input.

∴ In total, we need $(m \times k)$ parameters (for all k classes).

2) $\quad S_R(x) = \Theta_R^T \cdot x \quad ; \quad \hat{P}_R = \dfrac{e^{(S_R(x))}}{\sum\limits_{j=1}^{K} e^{(S_k(x))}}$

$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^{m} \sum_{R=1}^{K} y_R^i \log(\hat{P}_R) \qquad \left[ \begin{array}{l} y_{.k}^i = 1, \text{ if } i \text{ bot} \\ \quad = 0, \text{ otherwise} \end{array} \right.$$

expanding inner Summation,

$$= -\frac{1}{m} \sum_{i=1}^{m} \left[ y_1^i \log(\hat{P}_1^i) + y_2^i \log(\hat{P}_2^i) + \ldots + y_R^i \log(\hat{P}_R^i) + \ldots \right]$$

$$= -\frac{1}{m} \sum_{i=1}^{m} y_R^i \log(\hat{P}_R^{(i)})$$

or, $\quad J(\Theta) = -\dfrac{1}{m} \sum\limits_{i=1}^{m} \sum\limits_{j=1}^{k} y_j^i \log(\hat{P}_i^j)$

$$\Rightarrow \frac{\partial J(\Theta)}{\partial \Theta} = -\frac{1}{m} \sum_{i=1}^{m} \sum_{R} y_R^i \frac{\partial[\log P_R]}{\partial \Theta}$$

$$= -\frac{1}{m} \sum_{i=1}^{m} \sum_{R} y_R^i \cdot \frac{1}{P_R^i} \cdot \frac{\partial P_R^i}{\partial \Theta} \cdot \qquad \left[ \text{Chain rule} \right.$$

$$= -\frac{1}{m} \sum_{i=1}^{m} \left\{ -y_R^i(1-\hat{P}_R^i)x^i - \sum_{R \neq i} x^i \left[ y_R^i \frac{1}{P_R}(-P_R P_i) \right] \right.$$

Now,

$$\frac{\partial J(\theta)}{\partial \theta} = -\frac{1}{m} \sum_{i=1}^{m} \left\{ -y_k^i (1-p_k^i) + \sum y_R (p_i^i) \right\} x^i$$

$$= \frac{1}{m} \sum_{i=1}^{m} \left[ -y_k^i + y_k^i p_R^i + \sum y_R (p_k^i) \right] x^i$$

$$= \frac{1}{m} \sum_{i=1}^{m} \left\{ p_k^i \left( \sum_R y_R \right) - y_k^i + y_k^i p_k^i \right\} x^i$$

$$= \frac{1}{m} \sum_{i=1}^{m} \left( \hat{p}_R^i - y_R^i \right) x^i$$