**1. Prove Bayes' Theorem. Briefly explain why it is useful for machine learning problems, i.e., by converting posterior probability to likelihood and prior probability.**

**Bayes' theorem** is a formula that describes how to update the probabilities of hypotheses when given evidence. If A and B denote two events, P(A|B) denotes the conditional probability of A occurring, given that B occurs. The two conditional probabilities P(A|B) and P(B|A) are in general different. Bayes theorem gives a relation between P(A|B) and P(B|A).

$$P (A \mid B) = P (A \cap B) / P(B)$$

Similarly,

$$P (B \mid A) = P (A \cap B) / P(A)$$

Thus,

$$P (A \cap B) = P (A \mid B) * P (A) = P (B \mid A) * P(B)$$

Equating the two yields:
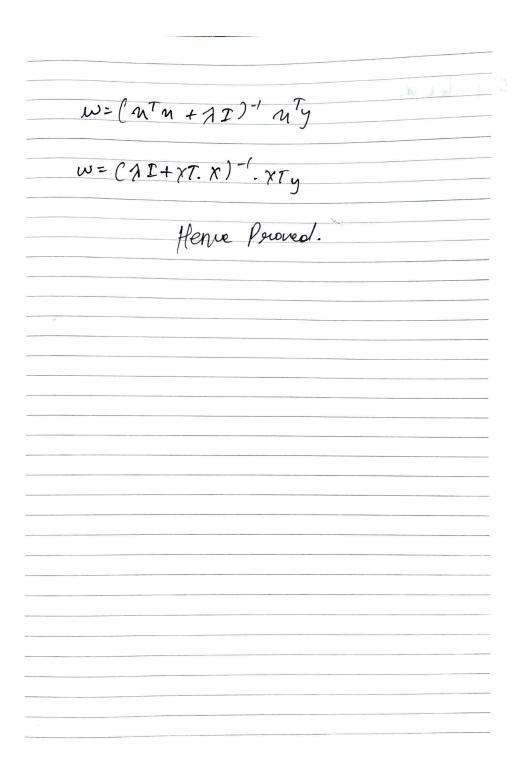
$$P(A|B) = P(A) P(B|A) P(B)$$

This equation, is known as Bayes Theorem is the basis of statistical inference

P(A) and P(B) are probabilities of A and B which are mutually exclusive. Hence, P(B|A) and P(A|B) are conditional probabilities. In the same lines, P(B|A) is Likelihood and P(A|B) is Posterior probability.

In machine learning, we have attributes, response variables and predictions or classifications. Using this algorithm, we will deal with the probability distributions of the variables in the dataset and predict the probability of the response variable belonging to a particular value, given the attributes of a new instance. We are interested in defining the best hypothesis from some space H, given the observed training data D. One way to specify what we mean by the best hypothesis is to say that we demand the most probable hypothesis, given the data D plus any initial knowledge about the prior probabilities of the various hypotheses in H. Bayes theorem provides a direct method for calculating such probabilities.

**2. In Lecture 3-1, we gave the normal equation (i.e., closed-form solution) for linear regression using MSE as the cost function. Prove that the closed-form solution for Ridge Regression is $w = (\lambda I + X^T \cdot X)^{-1} \cdot X^T \cdot y$, where $I$ is the identity matrix, $X = (x(1), x(2), \dots, x(m))^T$ is the input data matrix, $x(i) = (1, x1, x2, \dots, xn)$ is the $i^{th}$ data sample, and $y = (y(1), y(2), \dots, ym)$. Assume the hypothesis function $hw(x) = w0 + w1x1 + w2x2 + \cdots + wnxn$, and $y(j)$ is the measurement of $hw(x)$ for the $j$th training sample. The cost function of the Ridge Regression is $E(w) = MSE(w) + \lambda \frac{}{2}\sum^m_{i=1} wi2$. [Hint: please refer to the proof of the normal equation of linear regression].**

Cost function of Ridge Regression:-

$$E(\omega) = MSE(\omega) + \frac{\lambda}{2} \sum^m_{i=1} \omega_i^2$$

$$= \frac{1}{m} \sum^m_{i=1} \left(h\omega(x_i) - y_i\right)^2 + \frac{\lambda}{2} \sum^m_{i=1} \omega_i^2$$

$$E(\omega) = \frac{1}{m} \left(X\omega - y\right)^T \left(X\omega - y\right) + \frac{\lambda}{2} \omega^T\omega$$

neglecting $\frac{1}{m}$ & $\frac{1}{2}$ we get

$$= \left((X\omega)^T - y^T\right)\left(X\omega - y\right) + \lambda \omega^T\omega$$

$$= (X\omega)^T(X\omega) - (X\omega)^T(y) - (X\omega)^T(y) + y^Ty + \lambda\omega^T\omega$$

$$= \omega^TX^TX\omega - 2(X\omega)^Ty + y^Ty + \lambda\omega^T\omega$$

Now, $\frac{dE}{d\omega} = 0$

∴ On differentiating we get

$$\left(X^TX + \lambda I\right)\omega = X^Ty$$

Multiplying $\left(X^TX + \lambda I\right)^{-1}$ on both sides

$$w = (u^T u + \lambda I)^{-1} u^T y$$

$$w = (\lambda I + X^T \cdot X)^{-1} \cdot X^T y$$

Hence Proved.

3. [10 points] Recall the multi-class Softmax Regression model on page 16 of Lecture 3-3. Assume we have K different classes. The posterior probability is $\hat{p}_k = \delta(s_k(x))_k = \frac{\exp(s_k(x))}{\sum_{j=1}^{K} \exp(s_j(x))}$ for $k = 1, 2, \dots, K$, where $s_k(x) = \theta_k^T \cdot x$, and input $x$ is an $n$-dimension vector.

1) To learn this Softmax Regression model, how many parameters we need to estimate? What are these parameters?

2) Consider the cross-entropy cost function $J(\Theta)$ (see page 16 of Lecture 3-3) of $m$ training samples $\{(x_i, y_i)\}_{i=1,2,\dots,m}$. Derive the gradient of $J(\Theta)$ regarding to $\theta_k$ as shown in page 17 of Lecture 3-3

**To learn Softmax Regression Model, we need to estimate (n+1) parameters and the parameters are $\Theta_0, \Theta_1, \Theta_2, \Theta_3 \dots \Theta_n$**

3)

$$\hat{p}_k = \delta(s_k(m))_k = \frac{\exp(s_k(m))}{\sum_{j=1}^{k} \exp(s_j(m))}$$

where $s_k(m) = \theta_k^T \cdot m$

Training to minimize the cost function of cross entropy :-

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

$$= -\frac{1}{m} \sum_{i=1}^{m} \sum_{k=1}^{k} y_{(k)}^{(i)} \log\left(\frac{\exp(s_k(m^{(i)}))}{\sum_{j=1}^{k} \exp(s_j(m^{(i)}))}\right)$$

$$= -\frac{1}{m} \sum_{i=1}^{m} \left( \sum_{k=1}^{k} y_k^{(i)} \log(\exp(s_k(m^{(i)}))) - \right.$$

$$\left. \sum_{k=1}^{k} y_{(k)}^{(i)} \log\left(\sum_{j=1}^{k} \exp(s_j(m^{(i)}))\right) \right)$$

$y_k^{(i)} = 1$ if the ith instance belongs to class $k$ ; 0 otherwise.

Softmax Regression gradient for cross-entropy
cost function:

$$\nabla_{\theta_k} J(\theta) = \frac{1}{M} \sum_{i=1}^{M} \left( \hat{p}_k^{(i)} - y_{(k)}^{(i)} \right) u^{(i)}$$

$$= -\frac{1}{M} \sum_{i=1}^{M} u^{(i)} - \frac{1}{\sum_{j=1}^{K} \exp(\theta_j^T u^{(i)})} \exp(\theta_k^T u^{(i)}) u^i$$

$$= -\frac{1}{M} \sum_{i=1}^{M} \left( 1 - \hat{p}_k^{(i)} \right) u^{(i)}$$

$$= \frac{1}{M} \sum_{i=1}^{M} \left( \hat{p}_k^{(i)} - 1 \right) u^{(i)}$$

$$= \frac{1}{M} \sum_{i=1}^{M} \left( \hat{p}_k^{(i)} - y_k^{(i)} \right) u^{(i)} = \nabla_{\theta_k} J(\theta)$$

Hence Proved 4

**Parth Parab**                   **10444835**                   **CPE 695 Assignment 2**