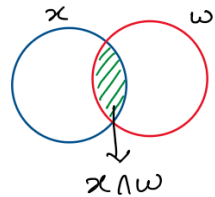


1. Prove Bayes' Theorem. Briefly explain why it is useful for machine learning problems, i.e., by converting posterior probability to likelihood and prior probability.

let "x" be the feature

"w" be the class we want to classify



Probability of having feature x provided class w have already classified

$$P(x|w) = \frac{P(x \cap w)}{P(w)} \quad - (1)$$

similarly,

$$P(w|x) = \frac{P(x \cap w)}{P(x)} \quad - (2)$$

from (1) and (2)

$$P(x \cap w) = P(x|w)P(w) = P(w|x)P(x)$$

$$\therefore P(w|x) = \frac{P(x|w)P(w)}{P(x)}$$

→ this is Bayes Theorem

where  $P(x|w)$  = likelihood probability as we are calculating likelihood of having x given w.  
[we have all the information provided to calculate  $P(x|w)$ ]

$P(w)$  = prior probability, it is probability of occurring(w)

$$P(x) = \text{evidence probability} \left( \sum_{i=1}^I P(x|w_i)P(w_i) \right)$$

= It is just the multiplication of likelihood x prior for all classes.

This Bayes theorem is very useful in ML application as we have to predict the posterior probability (predicting class given feature x already happened)

where we could calculate,  $P(x|w)$ ,  $P(w)$  and  $P(x)$  using the data available.

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  $\mathbf{w} = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ , where  $\mathbf{I}$  is the identity matrix,  $\mathbf{X} = (x(1), x(2), \dots, x(m))^T$  is the input data matrix,  $x(i) = (1, x_1, x_2, \dots, x_n)$  is the  $i$ th data sample, and  $\mathbf{y} = (y(1), y(2), \dots, y(m))$ . Assume the hypothesis function  $h_{\mathbf{w}}(x) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$ , and  $y(j)$  is the measurement of  $h_{\mathbf{w}}(x)$  for the  $j$ th training sample. The cost function of the Ridge Regression is  $E(\mathbf{w}) = \text{MSE}(\mathbf{w}) + \lambda \sum_{i=1}^m w_i^2$ . [Hint: please refer to the proof of the normal equation of linear regression. [ Note: Please use the following rectified definition of MSE when you prove:  $\text{MSE}(\mathbf{w}) = \sum (\mathbf{w}^T \mathbf{x}(i) - y(i))_{i=1}^m$ . ]

Given:

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

$$\text{where, } \text{MSE}(\mathbf{w}) = \sum_{i=1}^m (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2$$

$\therefore E(\mathbf{w})$  becomes

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

Representing above equation in matrix equivalent form.

$$\begin{aligned} E(\mathbf{w}) &= \frac{1}{2} (\mathbf{w}^T \mathbf{x} - y)^T (\mathbf{w}^T \mathbf{x} - y) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \\ &= \frac{1}{2} \left[ (\mathbf{w}^T \mathbf{x})^T - y^T \right] (\mathbf{w}^T \mathbf{x} - y) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \\ &= \frac{1}{2} \left[ (\mathbf{w}^T \mathbf{x})^T \cdot \mathbf{w}^T \mathbf{x} - y^T \cdot \mathbf{w}^T \mathbf{x} - (\mathbf{w}^T \mathbf{x})^T y \right] - y^T \cdot y + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \end{aligned}$$

$$= \frac{1}{2} [\omega \cdot x^T \cdot \omega^T \cdot x - y^T \omega^T x - \omega x^T \cdot y - y^T \cdot y] + \frac{\lambda}{2} \omega^T \omega$$

$$\frac{\partial E(\omega)}{\partial \omega} = \frac{1}{2} [2 x^T \omega \cdot x - y^T x - x^T \cdot y + 2\lambda \omega] \quad \text{--- ①}$$

putting  $y^T x = x^T y$  as

$$y = \begin{bmatrix} \quad \end{bmatrix}_{m \times 1} \quad \text{and} \quad x = \begin{bmatrix} \quad \end{bmatrix}_{m \times 1}$$

$$y^T = \begin{bmatrix} \quad \end{bmatrix}_{1 \times m} \quad x^T = \begin{bmatrix} \quad \end{bmatrix}_{1 \times m}$$

$$y^T x = \begin{bmatrix} \quad \end{bmatrix}_{1 \times m} \begin{bmatrix} \quad \end{bmatrix}_{m \times 1} \quad x^T y = \begin{bmatrix} \quad \end{bmatrix}_{1 \times m} \begin{bmatrix} \quad \end{bmatrix}_{m \times 1}$$

$$y^T x = \begin{bmatrix} \quad \end{bmatrix}_{1 \times 1} \quad x^T y = \begin{bmatrix} \quad \end{bmatrix}_{1 \times 1}$$

$\therefore$  Equation ① becomes,

$$= \frac{1}{2} [2 x^T \omega \cdot x - 2 x^T \cdot y + 2\lambda \omega]$$

$$= x^T \omega \cdot x - x^T \cdot y + \lambda \omega$$

$$\frac{\partial E(\omega)}{\partial \omega} = 0$$

$$x^T \omega \cdot x - x^T y + \lambda \omega = 0$$

$$x^T \omega \cdot x - \lambda \omega = x^T y$$

$$\omega (x^T \cdot x - \lambda I) = x^T y \quad \dots \dots \text{Since } AI = A \text{ (matrix multiplication property)}$$

$$\omega = \frac{x^T y}{(x^T \cdot x - \lambda I)}$$

$$\omega = (x^T \cdot x - \lambda I)^{-1} \cdot x^T y$$

3. Recall the multi-class Softmax Regression model on page 16 of Lecture 3-3. Assume we have  $K$  different classes. The posterior probability is  $\hat{p}_k = \delta(s_k(x))_{k=1,2,\dots,K}$  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?

We need to estimate  $\theta$  parameters i.e.,  $\theta_0, \theta_1, \theta_2, \dots, \theta_{n+1}$  where  $\theta_0$  is bias parameter. These parameters are the numbers we use to adjust the features in order to predict the output.

$$\text{for eq } y = \theta_0 + \theta_1 \omega_1 + \theta_2 \omega_2 + \dots + \theta_n \omega_n$$

↑  
output

where  $n$  = number of classes

∴ therefore there will be  $n+1$  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

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

We know,

$$p_k = \delta(s_k(x))_k = \frac{\exp(s_k(x))}{\sum_{j=1}^K \exp(s_j(x))}$$

$$\text{where } s_k(x) = \theta_k^T \cdot x$$

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

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

putting  $y_k^{(i)} = 1$  as  $y_k = 1$  ... if  $i$ th instance belongs to  $k$  } lecture slide 3-3 pg 16  
0 otherwise

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[ \underbrace{\sum_{k=1}^K \log [\exp(\theta_k^T \cdot x)]}_{= \theta_k^T x} - \sum_{k=1}^K \log \left[ \sum_{j=1}^K \exp(\theta_j^T \cdot x) \right] \right]$$

=  $\theta_k^T x$  ... since  $\ln(e^x) = x$

$$J(\theta) = -\frac{1}{M} \sum_{i=1}^M \left[ \sum_{k=1}^K \theta_k^T x^i - \sum_{k=1}^K \log \left[ \sum_{j=1}^K \exp(\theta_j^T x^i) \right] \right]$$

$$\frac{\partial J(\theta)}{\partial \theta_k} = -\frac{1}{M} \sum_{i=1}^M \left[ x^{(i)} - \frac{1}{\sum_{j=1}^K \exp(\theta_j^T x)} \cdot \exp(\theta_k^T x^i) x^{(i)} \right]$$

$$= -\frac{1}{M} \sum_{i=1}^M \left[ 1 - \frac{\exp(\theta_k^T x)}{\sum_{j=1}^K \exp(\theta_j^T x)} \right] x^{(i)}$$

$$= -\frac{1}{M} \sum_{i=1}^M \left[ 1 - \frac{\exp(s_k(x))}{\sum_{j=1}^K \exp(s_j(x))} \right] x^{(i)} \quad \dots \quad \text{as } \theta_k^T x = s_k(x) \\ \theta_j^T x = s_j(x)$$

$$= -\frac{1}{M} \sum_{i=1}^M \left[ 1 - p_k^{(i)} \right] x^{(i)} \quad \dots \quad \text{as } p_k = \frac{\exp(s_k(x))}{\sum_{j=1}^K \exp(s_j(x))}$$

$$= \frac{1}{M} \sum_{i=1}^M \left[ p_k^{(i)} - 1 \right] x^{(i)}$$