



Narenjan B. Rokaya

TRIBHUVAN UNIVERSITY  
SCHOOL OF MATHEMATICAL SCIENCES  
Balkhu, Kathmandu

Internal Assessment Answer Sheet

Level : Masters.....

Semester : Second.....

Name : Arpan Saphota.....

Roll No. : 07.....

Subject : Applied Machine Learning.....

Date : 20.80...../06...../20.....

Marks Obtained

Q. No. 1	
Q. No. 2	
Q. No. 3	
Q. No. 4	
Q. No. 5	
Q. No. 6	
Q. No. 7	
Q. No. 8	
Q. No. 9	
Q. No. 10	
<b>TOTAL</b>	

## Group 'A'

### 2. Batch Gradient Descent

- ① In Batch gradient descent, samples of training data are taken.
- ② It gives the average effect.
- ③ It converges slowly.
- ④ Gives Optimal solution.

### Stochastic Gradient Descent

- ① In Stochastic Gradient Descent individual data are taken.
- ② It gives the individual effect.
- ③ Convergence is faster as compared to Batch Gradient Descent.
- ④ Give solution but not optimal.

⑤ If uses,

$$w_j^* = w_j^* + \alpha(y - h(x))x_j$$

⑤ If uses;

$$B_j = B_j + \alpha(y - h_B(x))x_j$$

3. Precision is one of the evaluation metric for the classification task, where it gives the ratio of positive prediction over the positive values.

Let us consider a confusion matrix.

		Actual	
		T	F
Predicted	T	TP      FP	
	F	FN      TN	

from the above confusion matrix,

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

Hence, Precision gives the true positive value to sum of True positive and false positive ratio. It give more insight to the algorithm since, accuracy will not be efficient alone to evaluate the classification models.

Similarly,

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Recall is also another evaluation metric for the classification algorithm, where it gives the ratio of True positive value to the sum of true positive and false negative value. Also Recall gives the sensitivity of the classification algorithm.

4. K-means clustering is the Unsupervised Machine learning Algorithm, in which the k-number of clusters will be formed from the dataset which undergoes the k-means clustering.

for k-means clustering, we have the following procedure :-

1. Initialize the initial clusters  $c_1, c_2, \dots, c_k$ , where  $k$  is the number of clusters.
2. Get the distance of each data point with the cluster  $c_1, c_2, \dots, c_k$ . Here, we can use any distance measure like:
  - Euclidean distance  $\Rightarrow d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$  (mostly useful)
  - Manhattan distance
  - Cosine Similarity Measure
  - Jaccard distance, etc.
3. Compare the distance with the clusters  $c_1, c_2, \dots, c_k$  and get the minimum distance from them in which the cluster lie between  $c_1, c_2, \dots, c_k$ .
4. Repeat the step 3 for all the datapoints.
5. From the new assigned clusters, get the centers by taking the averages, which will be the next cluster for next iteration.  
i.e.  $c'_1, c'_2, \dots, c'_k$ .
6. Repeat the process until we get the datapoints in the same clusters repeatedly.
7. Finally, we get the clusters :
 
$$C^A = \{d_1, d_2, \dots, d_m\}$$

$$C^B = \{d'_1, d'_2, \dots, d'_{m'}\}$$

$$\vdots$$

$$C^k = \{d''_1, d''_2, \dots, d''_{m''}\}$$
 Compute the center of the obtained clusters.

1. Supervised learning is a costly approach compared to unsupervised learning. In supervised learning, there are labelled data which undergoes the model training process. In unsupervised learning, there are unlabelled data which undergoes the model building process. During this process, supervision is required for the labelled data, whereas in unsupervised learning, supervision is not required. Due to this, the computational cost for training the unsupervised machine learning algorithm will be lesser as compared to the supervised machine learning algorithm. Extra resources like memory, processors, etc. are required in order to train the supervised machine learning models. In unsupervised machine learning, there the model itself learns the pattern and performs the prediction. No extra input should be given except the training data. However in supervised machine learning algorithm, we have to give the labelled data and supervise the model to learn the pattern as per the given data. This, supervised machine learning is costly approach compared to the unsupervised machine learning.

5. Logistic Regression and Softmax regression both are the regression algorithm. where in logistic regression we perform the regression of binary class i.e. (two classes). whereas in softmax regression we perform the regression of multi-class i.e. greater than two.

In logistic regression,

$$h(x) = \frac{1}{1+e^{-x}} ; \text{Sigmoid function}$$

is used.

But in softmax regression.

$$h(x) = \frac{e^x}{\sum e^x} ; \text{Softmax function}$$

is used.

In logistic regression, we get the values in between 0 & 1 where as in softmax regression, we get the value between positive infinity and negative infinity.

Example:

If we want to do the ~~one~~-decision in two class like Yes or No, True or False then we use logistic regression.

If we want to get the decision in multi-class value like between 1 to 100 or more for any kind of prediction, then we use softmax regression.

6.

Logistic Regression for classification tasks.

In the linear regression, we use,  $h(x) = \hat{y} = w_1x_1 + w_2x_2 + \dots + w_nx_n = w^T x$ . But in logistic regression we use probability  $P$ .

where  $x_1, x_2, \dots, x_n$  are the inputs and  $w_1, w_2, \dots, w_n$  are the weights.

While using the probability values, our range is limited to 0 and 1., so, we use odds,

i.e.

$$\text{odds} = \frac{P}{1-P}$$

$\Rightarrow$  It is nothing but the ratio of probability to the probability value which is complementing value to the  $P$ . Here, the range will be 0 to positive infinite, so we are still restricted to the range,

To extend the range we take log,

$\log\left(\frac{P}{1-P}\right)$ , So, our range becomes positive infinite to negative infinite.

Now,

$$\log\left(\frac{P}{1-P}\right) = w^T x$$

Taking exponential on both sides, we get.

$$\exp\left(\log\left(\frac{P}{1-P}\right)\right) = \exp(w^T x)$$

$$\Rightarrow \frac{P}{1-P} = \exp(w^T x)$$

$$\Rightarrow P = \exp(w^T x) - p \cdot \exp(w^T x).$$

$$\Rightarrow P + p \exp(w^T x) = \exp(w^T x)$$

$$\Rightarrow P(1 + \exp(w^T x)) = \exp(w^T x)$$

$$\Rightarrow p = \frac{\exp(w^T x)}{1 + \exp(w^T x)}$$

Dividing by  $\exp(w^T x)$

$$\Rightarrow p = \frac{\exp(w^T x) / \exp(w^T x)}{(1 + \exp(w^T x)) / \exp(w^T x)}$$

$$= \frac{1}{\frac{1}{\exp(w^T x)} + \frac{\exp(w^T x)}{\exp(w^T x)}}$$

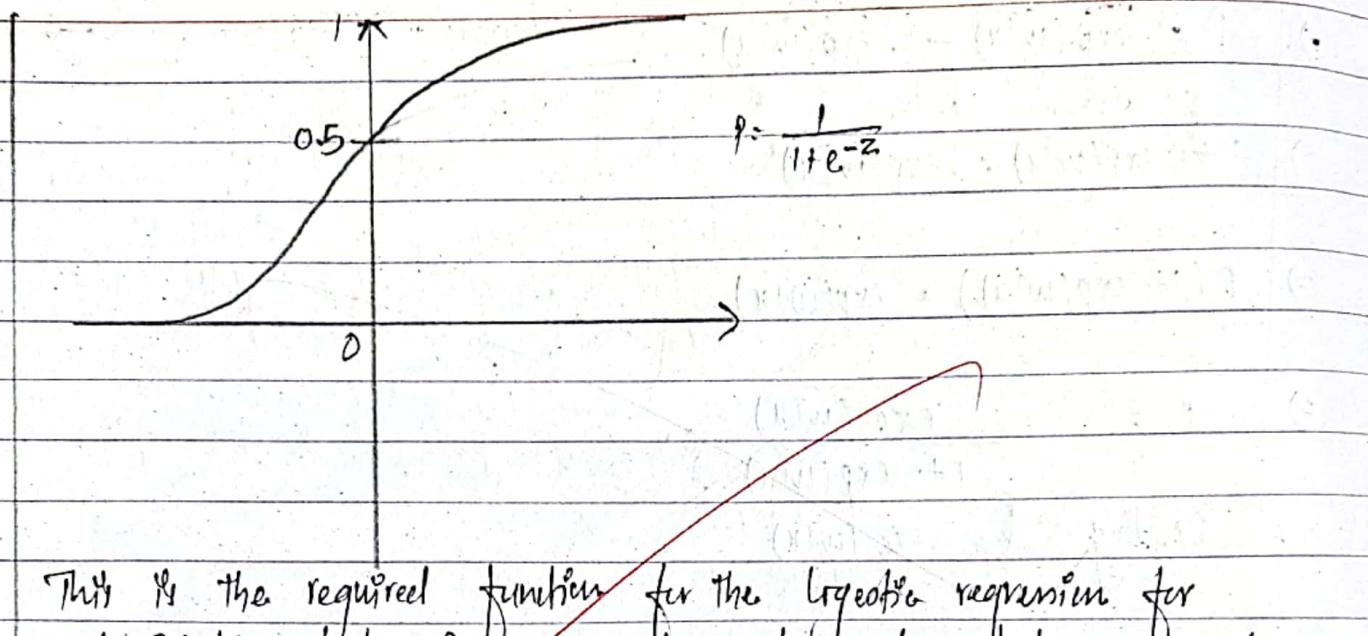
$$= \frac{1}{\exp(-w^T x) + 1}$$

$$= \frac{1}{1 + \exp(-w^T x)}$$

$$\therefore p = \frac{1}{1 + \exp(-z)}$$

where  $z = w^T x$

Here,  $p = \frac{1}{1 + \exp(-z)}$  is the sigmoid function.



This is the required function for the logistic regression for classification task. It will takes the values between  $-\infty$  to  $+\infty$  and classify between the 0 and 1 as shown in the graph above.

7. The various form of evaluation metrics used for regression are:

- ① Mean Square Error (MSE)
- ② Mean Absolute Error (MAE)
- ③ Root Mean Square Error (RMSE)
- ④ R-square Error (R<sup>2</sup>-error)

①. Mean Square Error (MSE)

In the regression algorithm, we have,  
 $\text{error} = Y_{\text{actual}} - Y_{\text{predicted}}$ .

The sum of total error,  $e = e_1 + e_2 + \dots + e_n$ .

Mean Square Error is given by,

$$MSE = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$$

where  $m$  is the number of data.

$\hat{y}_i$  = Predicted value.

$y_i$  = actual value.

It is the sum of squared of the errors.

### (ii) Mean Absolute Error (MAE)

In MSE, we take square, but taking square is not only an approach, we can take the absolute value for predicted and actual value i.e.  $|\hat{y}_i - y_i|$

This will also help us to avoid the error nullification problem.

So,

$$\text{Mean Absolute Error (MAE)} = \frac{1}{m} \sum_{i=1}^m |\hat{y}_i - y_i|$$

### (iii) Root Mean Square Error (RMSE)

In this root mean square error, we take the root value of Mean square error (MSE). This is done to scale out the error,

i.e.

$$RMSE = \sqrt{\frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2}$$

In order to scale out the scaling factor  $1/2$ , root is taken.

### (ii) R-Square Error ( $R^2$ -Error)

In this error, variance of model and the variance of the average is taken.  
i.e.

$$R^2\text{-Error} = 1 - \frac{\text{Variance(model)}}{\text{Variance(average)}}$$

where,  $\text{Variance(model)} = \sum_{i=1}^m (\hat{y}_i - y_i)^2$

and,  $\text{Variance(average)} = \sum_{i=1}^m (\bar{y} - y_i)^2$

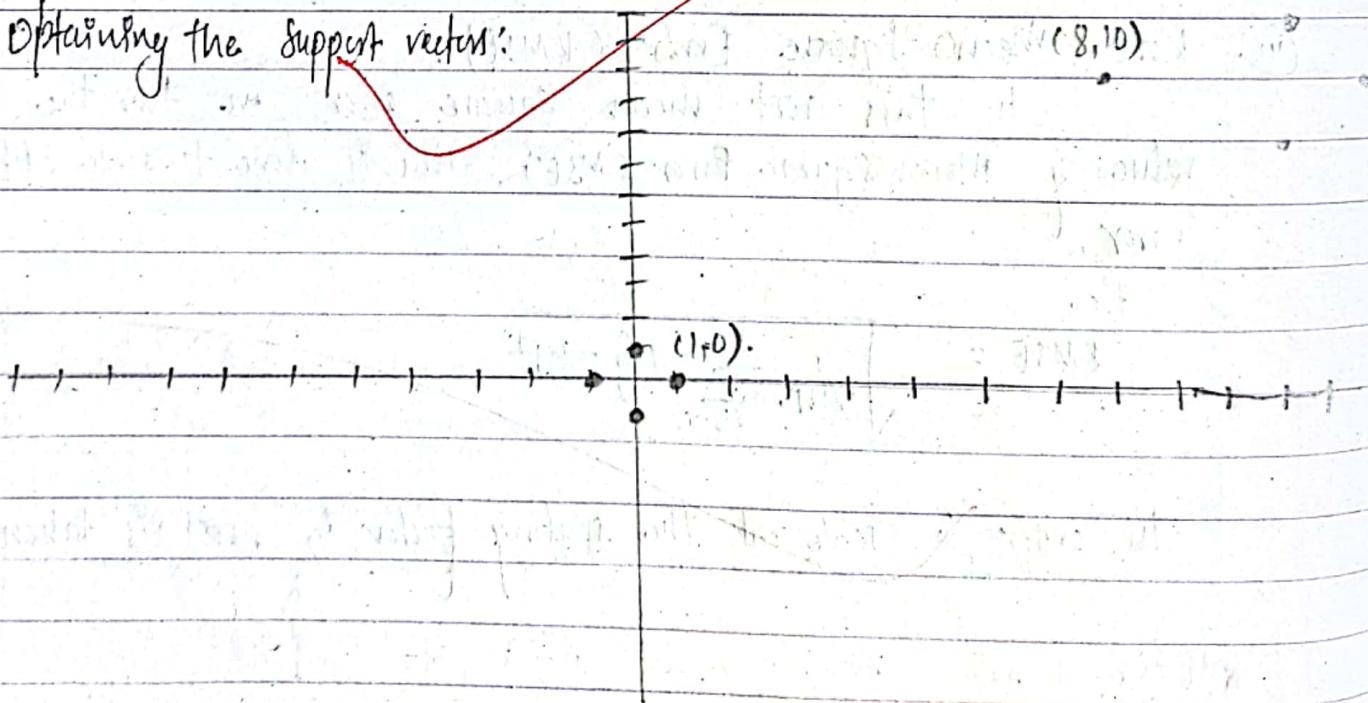
8:

Given,

Negatively labelled data:  $(1, 0), (0, 1), (-1, 0), (0, -1)$

Positively labelled data:  $(8, 10), (10, 8), (12, 10), (10, 12)$

Obtaining the support vector:



From the graph, we can see, the two support vectors.

$$S_1 : (1, 0) \Rightarrow \text{Negatively labelled.}$$

$$S_2 : (8, 10) \Rightarrow \text{Positively labelled.}$$

We take two variables i.e.  $\alpha, \beta$  then, the equation becomes,

$$\alpha S_1 S_1 + \beta S_1 S_2 = -1$$

$$\alpha S_2 S_1 + \beta S_2 S_2 = 1$$

Adding the bias  $1$  to the support vectors,

$$\text{i.e. } S_1 = (1, 0, 1)$$

$$S_2 = (8, 10, 1).$$

Performing the dot product for the linear kernel SVM,

$$S_1 S_2 = (1, 0, 1) \cdot (8, 10, 1) = 8 + 0 + 1 = 9$$

$$S_2 S_1 = (8, 10, 1) \cdot (1, 0, 1) = 8 + 0 + 1 = 9. \quad 1 + 1 = 2.$$

$$S_2 S_2 = (8, 10, 1) \cdot (8, 10, 1) = 64 + 100 + 1 = 165$$

Now the equation becomes,

$$\cancel{\alpha - \beta - 1} \quad 2\alpha + 9\beta = -1$$

$$\cancel{\alpha + 165\beta = 1}$$

On solving, we get

$$\alpha = -0.69$$

$$\beta = 0.044$$

Now, the weight becomes.

$$w = \alpha S_1 + \beta S_2$$

$$= -0.69(1, 0, 1) + 0.044(8, 10, 1)$$

$$w = (-0.338, 0.44, -0.646)$$

Here,

$$w_1 = -0.338, w_2 = 0.44 \text{ & } w_0 = -0.646.$$

Now,

$$y = w_1x_1 + w_2x_2 + w_0$$

$y = -0.338x_1 + 0.44x_2 + (-0.646)$ . required hyperplane equation.

In order to predict the class (8, 6).

$$y = -0.338 \times 8 + 0.44 \times 6 - 0.646$$

$$y = -0.71$$

Thus, it falls under negatively labelled data.

10.	Points.	1	2	3	4	5.
1	0					
2	9	0				
3	3	7	0			
4	6	5	9	0		
5	11	10	(2)	5.	0	

Point 3 & 5 ; i.e. the minimum so, we combine it as cluster

Points.	1	2	{3,5}	4.
1	0			
2	9	0		
{3,5}	(3)	7	0	
4	6	5	9	0

Point {3,5} and 1 is the minimum, we combine it as cluster

So

Point.  $\{1, 3, 5\}$ . 2 4

$\{1, 3, 5\}$ . 0

2 9 0

4. 6 (5) 0

Point 2 4 6 are minimum so, we combine it as cluster.

Point.  $\{1, 3, 5\}$ .  $\{2, 4\}$ .

$\{1, 3, 5\}$  0

$\{2, 4\}$ . (6) 0

Finally,

Point  $\{1, 3, 5, 2, 4\}$ .

$\{1, 3, 5, 2, 4\}$ . 0



Hence the cluster formed. cut at second Cut at third.

We draw the dendrogram as:

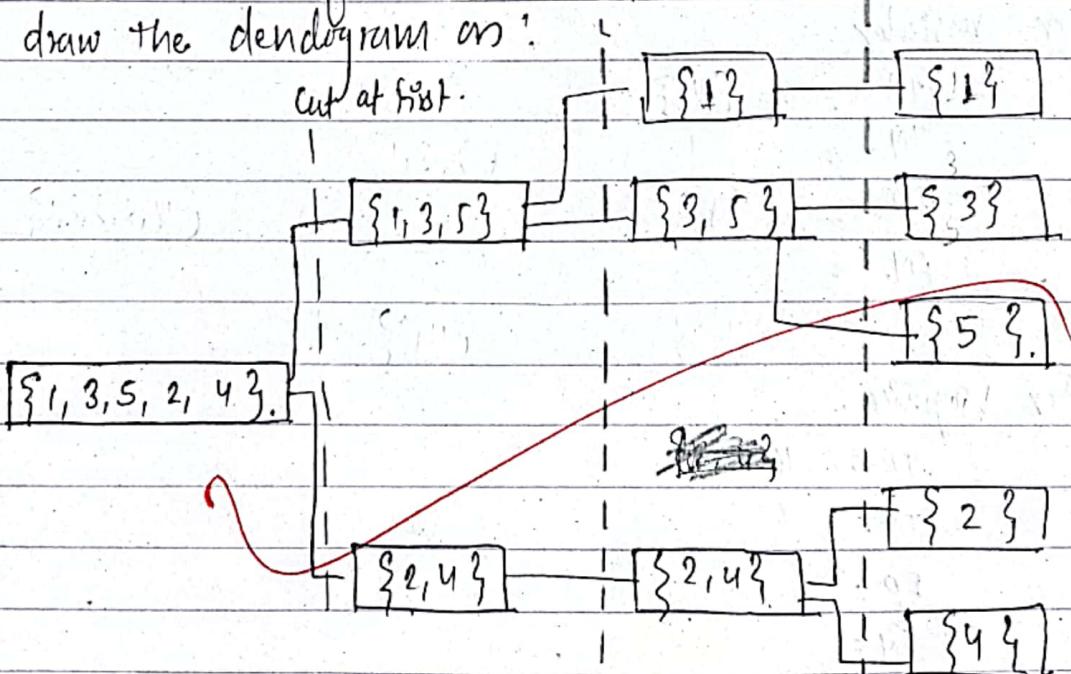


Fig.: Dendrogram.

In this, cut we determine the value of k for k-means clustering.

9. Given, Confusion Matrix.

Predicted.

Setosa Versicolor Virginica

Actual	Setosa	21	2	0
Setosa				
Versicolor		12	4	12
Virginica	5	3	16	

For, Setosa,

$$TP = 21$$

$$TN = 2 + 0 = 2.$$

$$FP = 2 + 0 = 2.$$

$$FN =$$

A	T	F
T	TP	FP
F	FN	TN

For Versicolor:

$$TP = 4$$

$$TN =$$

$$FP =$$

$$FN =$$

For Virginica.

$$TP = 16.$$

$$TN =$$

$$FP =$$

$$FN =$$

In order to compute the F<sub>1</sub> Score,

$$F_1\text{-Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

where Precision is the total actual and predicted positive i.e. True Positive to the ratio of True and False Positive.

Recall is the 1 - sensitivity.

Hence, F<sub>1</sub>-Score gives the Harmonic Mean of the confusion matrix.

