

G H RAISONI COLLEGE OF ENG., Nagpur
END SEMESTER EXAMINATION: SUMMER - 2021
SESSION 2020 - 2021 (ONLINE)

DEPARTMENT: ARTIFICIAL INTELLIGENCE

SEM/SEC : 4th A DATE : 11/03/2021

SUBJECT : MLA COURSE CODE : FAIZL 201

ROLL NO: A - 58

NAME: SHIVAM TAWARI

REG NO: 2019AAIE1117028

Q.1. QOT

A.

1. Clustering :

Clustering is a machine learning technique in which it groups the unlabelled dataset.

Example: Netflix uses clustering to show the recommendations to the users. It seeks to substantially improve the accuracy of predictions.

Pg. no. 1 ~~Other~~

MLA - 11AM - 1 PM BAZL

about how much someone is going to love a movie based on their movie preferences.

2. Classification :

Classification algorithm is a supervised learning technique that is used to identify the category of new observations on the basis of training data. Then it learns from the given dataset or observations and then classifies new observations into a number of classes or groups.

Example: When we want to classify between birds and reptile. Our program classify them into two different classes

3. Supervised Learning :

Supervised Learning is trained well using 'labelled' training data

B.G.N.O.2 ~~After~~ ^{After} ~~Answers~~

and basis of that data, machine predict the output.

Example :

Spam filtering can be done using supervised machine learning algorithm. It will classify it into spam or not spam.

4. Unsupervised Learning :

Unsupervised learning is a ML technique in which models are not supervised using training dataset.

Example :

'Unsupervised' can be used in customers segments in marketing data. Being able to determine segments of customers helps in marketing teams approach these customers segments in unique ways.

5. Hypothesis Space:

The hypothesis space utilized by an ML system is the arrangement of all hypothesis that may be returned by it.

Example:

The input space is in the given example 2^4 , i.e. the number of possible inputs. The hypothesis space is $2^{2^4} = 65536$ because for each set of features of the input space two outcomes (0 and 1) are possible.

Q.2.

Overfitting

- fitting the data too well.

- Features are noisy and uncorrelated to concept

Underfitting

- Learning too little of the true concept.

- Features don't capture concept

B.M.A. 4th

- Modeling process very powerful.
- There is too much search.
- Good performance on training data.
- Poor generalization to other data.

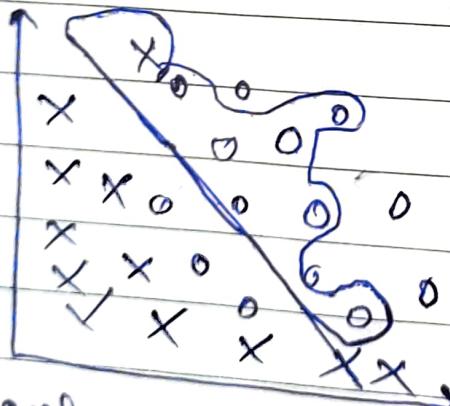
There is too much bias in model

Too ~~much~~ little search to fit model.

Poor performance on training data.

Poor generalization to other data.

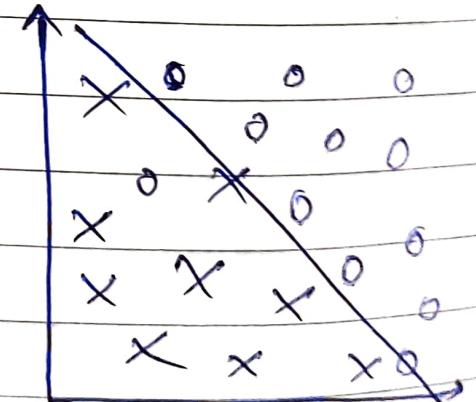
Example :



Model performing too well!

Overfitting

Accuracy : 100 %



Bad performance!

Underfitting

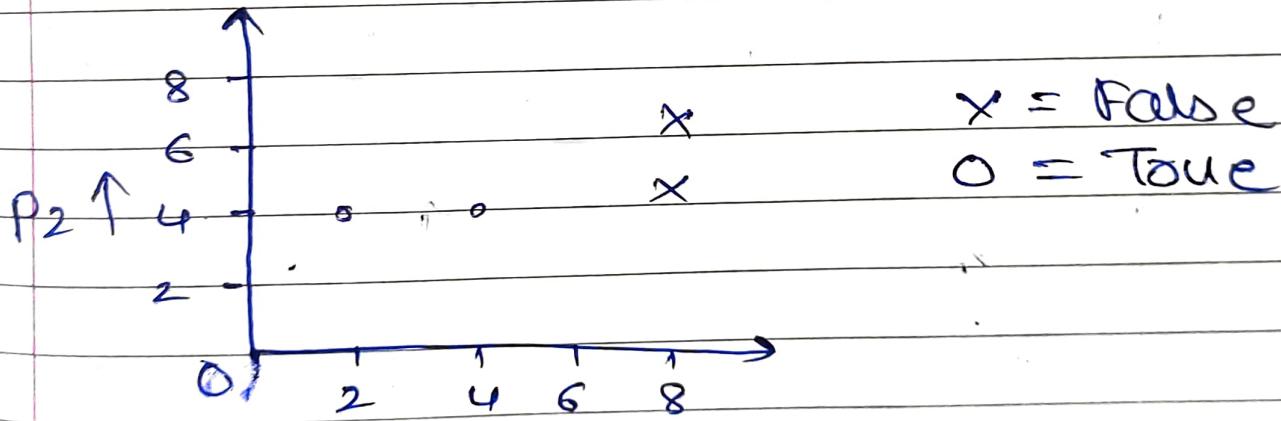
Accuracy : 50 %

2.
CO2

a. Given $\Rightarrow K = 3$

P1	P2	Class
7	7	False
7	4	False
3	4	True
1	4	True

Graph:



P1 \rightarrow

New point: $P_1 = 3$ and $P_2 = 7$

Comparing new observation with every other observation in dataset.

Euclidean distance measure is used to select first $k = 3$

$$\text{dist} = \frac{\text{observations}}{\sqrt{(P_1 - P_1')^2 + (P_2 - P_2')^2}}$$

$P_1 \cdot P_2$	$\text{dist}(3, 7)$	Class
7 7	$\sqrt{(7-3)^2 + (7-7)^2} = 4$	False
7 4	$\sqrt{(7-3)^2 + (4-7)^2} = 5$	False
3 4	$\sqrt{(3-3)^2 + (4-7)^2} = 3$	True
1 4	$\sqrt{(1-3)^2 + (4-7)^2} = 3.6$	True

Closest $k = 3$ neighbours are True (3, 4), True (1, 4), False (7, 7)

\therefore By majority voting, new observation will be true (3; 7)

Q.3.

CQ3.

- a. The goal of SVM algorithm is to create the best fit line or decision boundary

that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future. This best decision boundary is called a hyperplane.

SVM chooses the extreme vectors that help in creating the hyperplane.

Linear SVM is used for linearly separable data, which means the dataset can be classified into two classes by using a single straight line, then such data is termed as linearly separable data, and

Q 3

COS

b. Let R be rain:

$$\textcircled{a} P_R(R) = \frac{5}{365} = \frac{1}{73}$$

\textcircled{b} When 90% rains:

$$\textcircled{c} P_R(F|R) = \frac{9}{10}$$

\textcircled{d} When rain is 10%:

$$P_R(F|\bar{R}) = \frac{1}{10}$$

$$P_R(R) + P_R(\bar{R}) = 1$$

$$P_R(\bar{R}) = \frac{72}{73}$$

365 69

ENLA - IIAM - IPM - BAZL-203

Pg no 9

So by Bayes theorem:

$$P_B(A|B) = \frac{P_B(A) P_B(B|A)}{\sum P_B(A_i) P_B(B|A_i)}$$

Substituting into Bayes theorem:

$$P_B(R|F) = \frac{P_B(R) P_B(F|R)}{P_B(R) P_B(F|R) + P_B(\bar{R}) P_B(F|\bar{R})}$$

$$= \frac{1/73 \times 9/10}{(1/73 \times 9/10) + (72/73 \times 1/10)}$$

$$= \frac{9/730}{9/730 + 72/730}$$

$$= \frac{9}{9+72}$$

$$= \frac{1}{9} \approx 11\%$$

So, there is 11% chance that it will rain on the Friday of Madhusi's wedding.

Ans.

Ans.

Q. 4.

CO 4.

A.

$A = x_1$	$\times 2$
1.4	1.65
1.8 -	1.975
-1.4	-1.775
-2.0	-2.525
-3.0	-3.95
2.4	3.075
1.5	2.015
2.3	2.75
-3.2	-4.05
-4.1	-4.85
<hr/>	
$M = \text{mean} : -0.45$	$\therefore -0.5625$

centered matrix $C = A - M$

$x_1 + 0.4S$	$x_2 + 0.567S$
1.8S	2.217S
2.0S	2.502S
-0.9S	-1.207S
-1.5S	-1.957S
-2.5S	-3.282S
2.8S	3.642S
1.9S	2.592S
2.7S	2.378
-3.6S	-4.282S

① Covariance :

$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x}) \times (y_i - \bar{y})}{N}$$

$$= 7.88$$

$$\text{Cov}(x, x) = \frac{\sum (x_i - \bar{x})^2}{N}$$

$$= 6.42$$

$$\text{Cov}(x, y) = \frac{\sum (y_i - \bar{y})^2}{N}$$

$$= 9.95$$

$$\text{Cov}(xy, x) = \text{Cov}(x, y)$$

$$= 7.98$$

$$\therefore \text{Covariance matrix} = \begin{bmatrix} 6.42 & 7.98 \\ 7.98 & 9.95 \end{bmatrix}$$

Eigen values & vectors:

$$(A - \lambda I) v = 0$$

$$|A - \lambda I| = 0$$

$$\begin{vmatrix} 6.42 - \lambda & 7.98 \\ 7.98 & 9.95 - \lambda \end{vmatrix} = 0$$

$$(6.42 - \lambda)(9.95 - \lambda) - (7.98)(7.98) = 0$$

Pg no. 13

Solving eqⁿ. we get

Eigen values: $\lambda_1 = 0.00746$
 $\lambda_2 = 16.36$

Putting λ in eigen decomposition.

$$\begin{bmatrix} 6.42 - \lambda & 7.98 \\ 7.98 & 9.98 - \lambda \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 0$$

Vectors $\Rightarrow \begin{bmatrix} -0.779 \\ 0.626 \end{bmatrix}$ for $\lambda_1 = 0.007$

$$\begin{bmatrix} -0.626 \\ -0.779 \end{bmatrix} \text{ for } \lambda_2 = 16.36$$

Q.5
Q.5

4. Kmeans algorithm is an iterative algorithm that divides a group of n datasets into k subgroups / clusters based on the similarity

and their mean distance from
the centroid of that
particular subgroup formed.

Algorithm of k-means:

Step 1: Select the value of k ,
to decide the number
of clusters to be
formed.

Step 2: Select random k points
which will act as
centroids.

Step 3: Assign data point, based
on their distance from
the randomly selected
points, to the nearest
centroid which will
form the predefined
clusters.

Step 4: Place a new centroid of each cluster.

Step 5: Repeat step no. 3, which reassign each datapoint to the new closest centroid of each cluster.

Step 6: FINISH.

Gaussian Mixture Models assume that there are a certain number of gaussian distributions and each represents a cluster, hence, a GMM tends to group the data points belonging to a single distribution together.

We have taken unlabelled input data, which means its not categorized and corresponding outputs are not given.

Now it is fed to the machine learning model in to train it. Firstly, it will interpret the raw data to find the hidden patterns from data and then apply clustering.