

## 11 Assignment No 4."

(P1) What is PAC model?

Ans: - PAC stands for probably Approximately Correct (PAC) Learning.

- A good learner will learn with high probability and close approximation to the target concept.
- with high probability the selected hypothesis will have lower the error ("Approximately Correct") with the parameter  $\epsilon, \delta, \gamma$ .
- PAC Learning requires,
  - small parameters  $\epsilon, \delta, \gamma$
  - with probability at least  $(1-\delta)$ , a system learn the concept with error at most  $\epsilon$ .
- $\epsilon$  is upper-bound on the error in accuracy i.e. the hypothesis with error less than  $\epsilon$ .  
Accuracy  $\geq 1 - \epsilon$

-  $\epsilon$  give the probability of failure in achieving the accuracy  $\gamma$  ( $0 < \gamma < 1$ ), the hypothesis generated is approx. correct at least  $\gamma - \delta$  of the time  
confidence  $\gamma - \delta$

(P2) Explain Sample Complexity of finite hypothesis Space

→ - We seek to generalize Occam's razor to

infinite hypothesis space. To do so, we look at the set of behavior  $\mathcal{T}_{\mathcal{H}}(s)$  of hypotheses from  $\mathcal{H}$  on a sample  $s$ .

$$\mathcal{T}_{\mathcal{H}}(s) = \{ \{h(x_1), \dots, h(x_m)\} : h \in \mathcal{H}\}$$

for  $s = \{x_1, \dots, x_m\}$

$\mathcal{T}_{\mathcal{H}}(m) = \max_{S: |S|=m} |\mathcal{T}_{\mathcal{H}}(S)|$  defines the growth function  $H$ .

our goal is to modify Occam's razor to get a bound of the form

$$\text{err}(h) \leq 0 \left( \ln \mathcal{T}_{\mathcal{H}}(2m) + \ln \frac{1}{\delta} \right)$$

first, some definitions, for our proof of this bound, we fix  $\delta$  and let  $D$  denote our target distribution.

$$\text{Step 1: } \Pr[B' \cap B] \geq 1/2$$

$$\text{Step 2: } \Pr[B] \leq 2 \Pr[B']$$

$$\text{Step 3: } \Pr[B''] = \Pr[B']$$

$$\text{Step 4: } \Pr[b(h) \cdot S, S'] \leq e^{-m\delta/2}$$

Define the event  $B' = [\exists h \in \mathcal{H} : h \text{ consist with } S \text{ and } M(h, S') \geq \frac{m\delta}{2}]$ .

Q3

Define & explain "shattering" a set of finite hypothesis spaces

- Shattering is a key notion in ML that refers to a classifier's capacity to accurately distinguish any arbitrary labeling of group of points.
- Strictly speaking, a classifier breaks a collection of points if it can divide them into all viable binary categories.
- The greatest no of points that a classifier is capable of shattering is specified by the VC dimension.
- When classifier is able to properly distinguish any potential labeling of the points, it is said to be "shattering" a collection of points.

Q4

Explain ensemble learning model in detail.

- Ensemble learning refers to a ML approach in which the predictions from multiple models are merged to enhance the accuracy of the ultimate forecast.
- Simple Ensemble Techniques
  - 1) Max Voting
  - 2) Averaging
  - 3) Weighted Averaging.

1) Max Voting :- This method generally

used for classification problems. In this technique, multiple models are used to make predictions for each data point.

### 2) Averaging :-

- This method is similar to max voting technique, multiple prediction are made for each data point in averaging.
- In this method; we take an average of prediction from all models and use it to make the final prediction.

### 3) Weighted Average :-

- This is an extension of averaging method, All methods are assigned different weight for defining the important of each model for prediction.
- Advanced Ensemble techniques
  - 1) Stacking :-
    - It uses predictions from multiple method to build a new model. It will make prediction test set.

### 2) Blending :-

- It follows same approach that stacking following but uses only a holdout set from the train set to make predictions.

3) Bagging :-

— It combines the result of multiple model to get a generalized result.

4) Boosting :-

— Boosting is a sequential process, where each subsequent model attempts to correct the error of the previous model.

Q5)

Explain VC dimension.

→ — The Vapnik-Chervonenkis dimension, more commonly known as VC dimension, is model capacity measurement used in statistics and ML.

— It is termed informally as a measure of model's capacity, It is used frequently to guide the model selection process while developing ML applications.

— How to find VC dimension.

— Let us consider binary classification model, which state that for all points  $(a, b)$ , such that  $a < x < b$ , label them as 1, otherwise label them as 0.

$$h(x) = 1 \text{ if } a < x < b$$

$$h(x) = 0, \text{ otherwise}$$

$$(a, b) \in \mathbb{R}^2$$

we take two points  $m$  &  $n$  for these two points, there can be  $2^2$  distinct labels in binary classification. list cases follow:

$$h(m) = 0; h(n) = 0$$

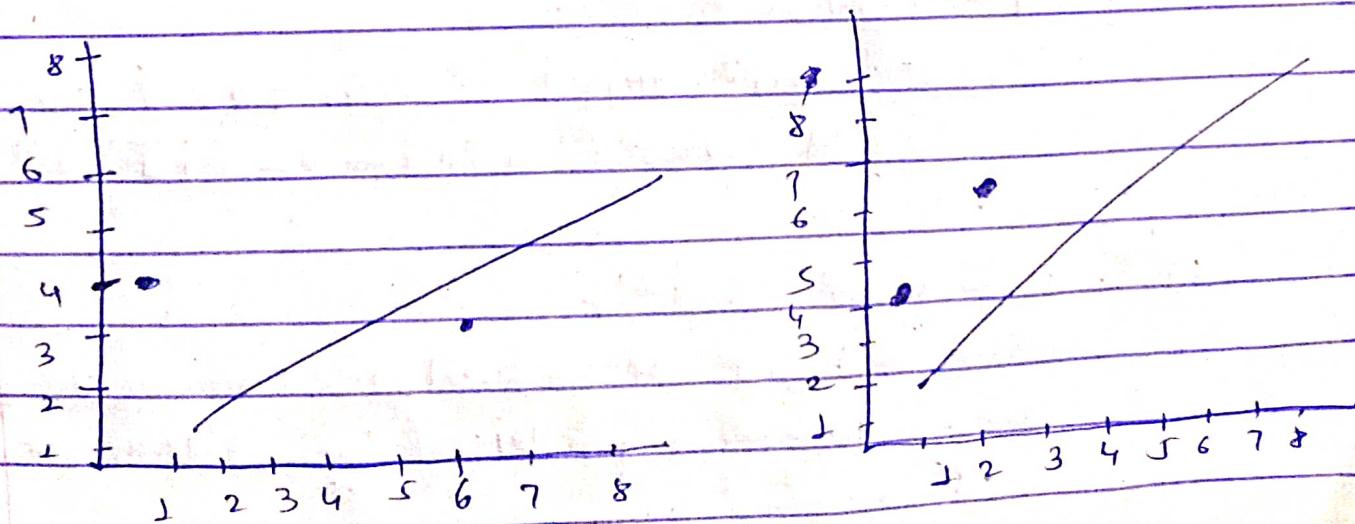
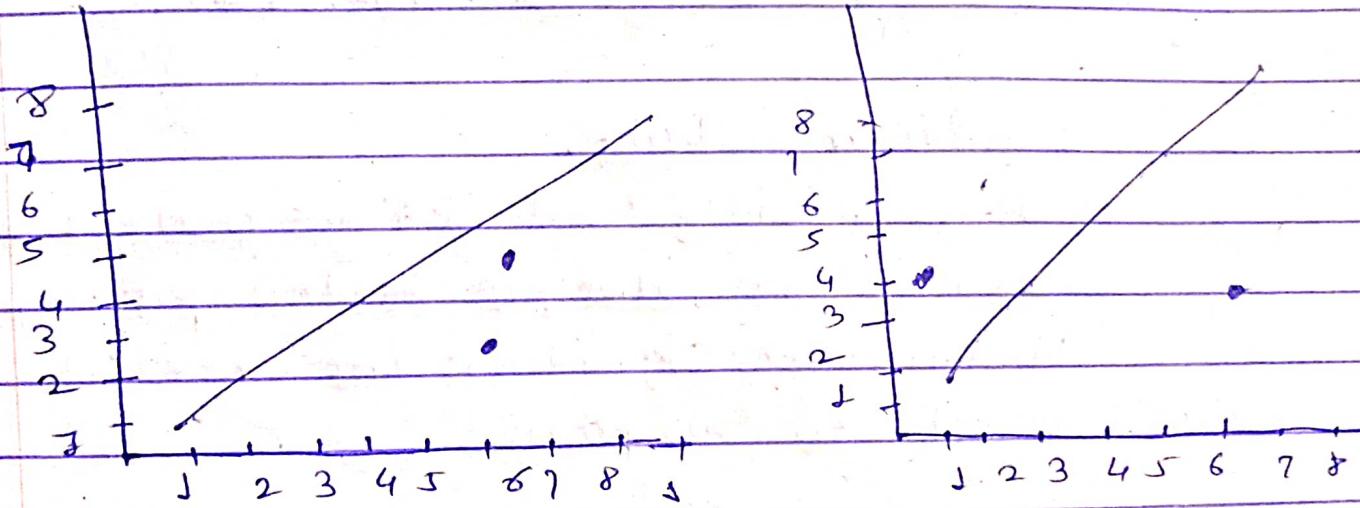
$$h(m) = 0; h(n) = 1$$

$$h(m) = 1; h(n) = 0$$

$$h(m) = 1; h(n) = 1$$

— we can observe that for all the possible labelling variations of  $m$  and  $n$ .

— The model can divide the points into two segments.



# Assignment No 5

Q1)

Write note on :-

a) Manhattan distance

- Manhattan distance is distance measured that is calculated by taking the sum of distance between x and y co-ordinate.

- This distance also known as Manhattan length. In other word measured along axes at right angles.

formula :

Manhattan distance =

$$|x_1 - x_2| + |y_1 - y_2|.$$

b) Euclidean distance.

- It is a widely used distance metric. It works on the principle of the pythagoras theorem and signifies the shortest distance between two points.

formula :

Euclidean Distance = (sum for i to

$$N \cdot (\text{abs}(\sqrt{x_i} - \sqrt{y_i}))^{1/p}$$

c) Elbow method :-

- The elbow method plots the value of the cost function produced by different values of k.

- Q2) - Explain Gaussian Mixture model in detail.
- - Suppose there are  $K$  clusters so,  $\mu$  &  $\Sigma$  and  $\pi$  are also estimated for each  $K$ .  
 Had it been only one distribution, they would have been estimated by the maximum-likelihood method.
- But since there are  $K$  such clusters and the probability density density is defined as a linear function of densities of all these  $K$  distributions i.e.

$$p(x) = \sum_{k=1}^K \pi_k G_k(x | \mu_k, \Sigma_k).$$

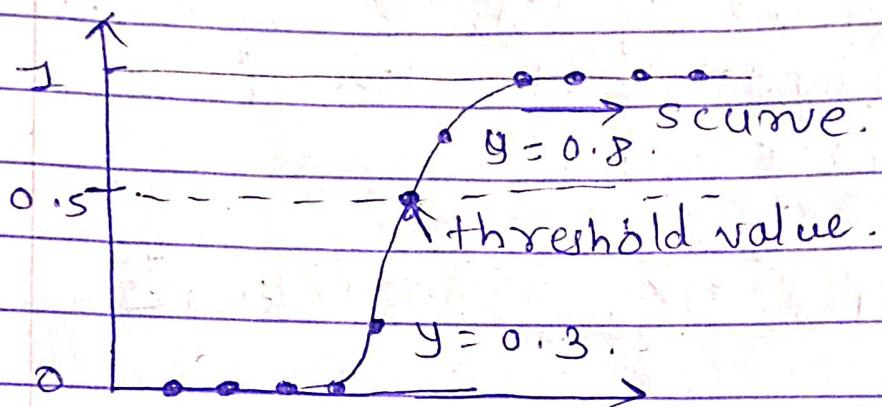
where  $\pi_k$  is mixing coefficient for  $k^{\text{th}}$  distribution, for estimating the parameters by the maximum log-likelihood method, compute  $p(x | \mu, \Sigma, \pi)$ .

- Q3) Explain Maximum likelihood & least square error hypothesis

- - Maximum Likelihood Estimation (MLE) is a probabilistic based approach to determine values for the parameters of the model.
- Parameters could be defined as blueprint for the model because on the algorithm works. MLE is widely used technique in ML.

- Least square is commonly used method in regression analysis for estimating the unknown parameters by creating a model will minimize the sum of squared errors between the observed data & the predicted data.

e.g. of Maximum Likelihood Estimation.



Q4) Predict the class of new point  $x=2$  &  $y=1$  using kNN algorithm assume  $k=3$  ?

→

$x$	$y$	class
-1	1	-
0	1	+
0	2	-
1	-1	-
1	0	+
1	2	+
2	2	-
2	3	+

$$(x_1, y_1) = (1, 1) \text{ } \& \text{ } (x_1, y) = (-1, 1)$$

$$\begin{aligned}
 d_1 &= \sqrt{(x_1 - x)^2 + (y_1 - y)^2} \\
 &= \sqrt{(-1 - 1)^2 + (1 - 1)^2} \\
 &= \sqrt{4} \\
 &= 2
 \end{aligned}$$

$$\begin{aligned}
 d_2 &= (x_1, y_1) = (1, 1) \text{ } \& \text{ } (x_1, y) = (0, 1) \\
 &= \sqrt{(x_1 - x)^2 + (y_1 - y)^2} \\
 &= \sqrt{(1 - 0)^2 + (1 - 1)^2} \\
 &= 1
 \end{aligned}$$

$$\begin{aligned}
 d_3 &= (x_1, y_1) = (1, 1) \text{ } \& \text{ } (x_1, y) = (0, -1) \\
 &= \sqrt{(x_1 - x)^2 + (y_1 - y)^2} \\
 &= \sqrt{1 + 1} \\
 &= \sqrt{2}
 \end{aligned}$$

$$\begin{aligned}
 d_4 &= (x_1, y_1) = (1, 1) \text{ } \& \text{ } (x_1, y) = (1, -1) \\
 &= \sqrt{(x_1 - x)^2 + (y_1 - y)^2} \\
 &= \sqrt{(1 - 1)^2 + (1 - (-1))^2} \\
 &= \sqrt{4} \\
 &= 2
 \end{aligned}$$

$$\begin{aligned}
 d_5 &= (x_1, y_1) = (1, 1) \text{ } \& \text{ } (x_1, y) = (1, 0) \\
 &= \sqrt{(x_1 - x)^2 + (y_1 - y)^2} \\
 &= \sqrt{(1 - 1)^2 + (1 - 0)^2} \\
 &= 1
 \end{aligned}$$

$$\begin{aligned}
 d_6 &= (x_1, y_1) = (1, 1) \text{ & } (x_1, y) = (1, 2) \\
 &= \sqrt{(x_1 - x)^2 + (y_1 - y)^2} \\
 &= \sqrt{(1-1)^2 + (1-2)^2} \\
 &= 1
 \end{aligned}$$

$$\begin{aligned}
 d_7 &= (x_1, y_1) = (1, 1) \text{ & } (x_1, y) = (2, 2) \\
 &= \sqrt{(x_1 - x)^2 + (y_1 - y)^2} \\
 &= \sqrt{(1-2)^2 + (1-2)^2} \\
 &= \sqrt{2}
 \end{aligned}$$

$$\begin{aligned}
 d_8 &= (x_1, y_1) = (1, 1) \text{ & } (x_1, y) = (2, 2) \\
 &= \sqrt{(x_1 - x)^2 + (y_1 - y)^2} \\
 &= \sqrt{(1-2)^2 + (1-3)^2} \\
 &= 2.4
 \end{aligned}$$

$d_1$	2	
$\rightarrow d_2$	1	
$d_3$	1.414	
$d_4$	2	
$\leftarrow d_5$	1	
$\rightarrow d_6$	<del>1.414</del>	
$d_7$	1.414	
$d_8$	2.414	

predicted value class is positive  
 all  $d_2, d_5$  &  $d_6$  are positive.

Q5) What is hierarchical clustering? Consider following distance matrix and apply hierarchical clustering to cluster u, v, w, x, y

	u	v	w	x	y
u	0	1	2	2	3
v	1	0	2	4	3
w	2	2	0	1	5
x	2	4	1	0	3
y	3	3	5	3	0

	uv	w	x	y
uv	0	2	2	3
w	2	0	1	5
x	2	1	0	3
y	3	5	3	0

$$d(w, (uv))$$

$$\min(d[wu], d[wv])$$

$$\min(2, 2)$$

$$d(x, uv)$$

$$\min(d[xu], d[xv])$$

$$\min(2, 4)$$

$$d(y, uv)$$

$$\min(d[yu], d[yv])$$

$$(3, 3)$$

uvw	4	vw	2	y	1
x	0	1	0	3	3
y	3	3	3	0	0

$d(x, (uvw))$

$\min d \{ d(xu), d(xv), d(xw) \}$

$\min d(2, 4, 1)$

$d(y, (uvw))$

$\min d \{ d(yu), d(yv), d(yw) \}$   
(3, 3, 5)

uvwx	uvwx	y	
u	0		
y		0	

$d(y, uvwx)$

$\min \{ d(yu), d(yv), d(yw), d(yx) \}$

$\min (3, 3, 5, 3)$

