

Classification

Chapt - 8

Date :

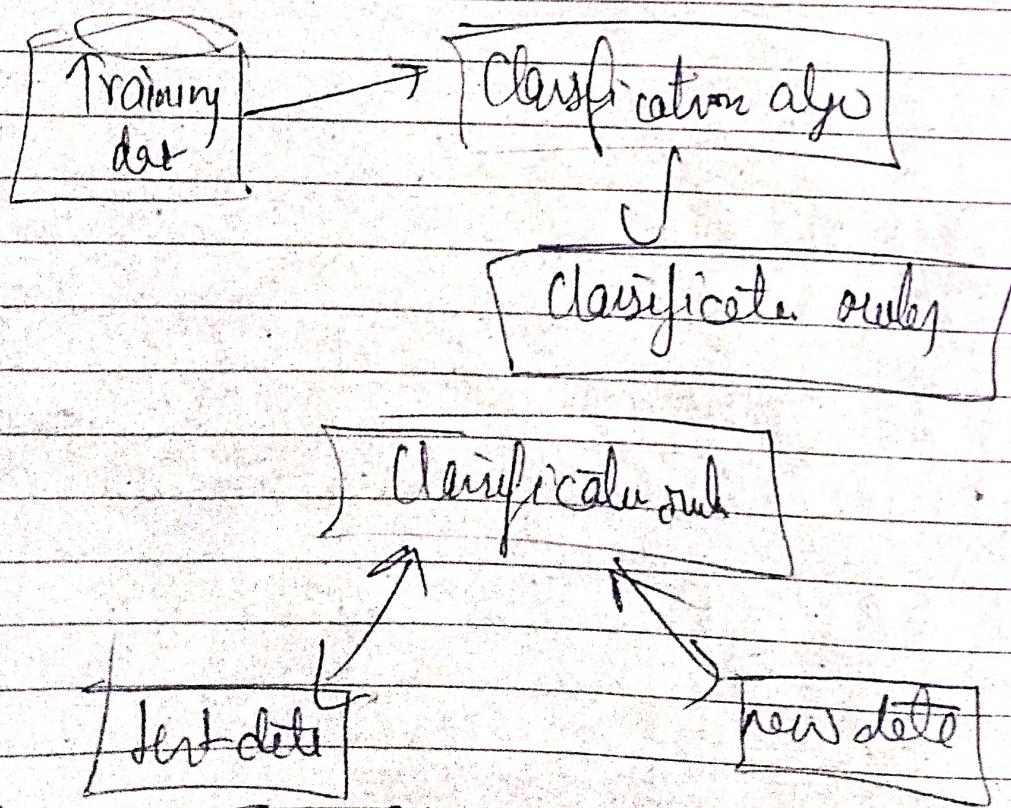
Page No.:

It is a form of data analysis that extracts model desired important data classes.

Data Classifier \rightarrow (i) Learning step
(ii) Classification step

a) Learning \rightarrow Training data are analyzed by a classification algo.

b) Classification \rightarrow Test data are used to estimate the accuracy of your classifier rules



Decision Tree Induction:

It is the learning of decision tree from class-labeled training tuples.

Attribute Selection Measures

It is a heuristic for selecting the splitting criterion that "best" separate a given data partition, D , of class-labeled training tuples into individual classes.

Information Gain

ID3 uses info gain as its attribute selection measure.

$$\text{Info}(D) = - \sum_{i=1}^m P_i \log_2(P_i)$$

$\text{Info}(D)$ is known as Entropy

$$\text{Entropy} = - \sum_{i=1}^m P_i \log_2(P_i)$$

$$= -\frac{P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

$$\text{Info-Gain}(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \text{Entropy}(D)$$

$|D_j|$ acts as the weight of j th partition

$$\text{Gain}(A) = \text{Info}(I) - \text{Info-Gain}(D)$$

$$= \text{Entropy} - \text{InfoGain}$$

tells how much
would be gained
by branching on A

Entropy \rightarrow Uncertainty associated with given info

Info Gain \rightarrow Information Gained

ID3 algo.

Day	Weather	Temp	Humidity	Play
D1	Sunny	Hot	High	N
D2	Sunny	Hot	High	N
D3	Cloudy	Hot	High	Y
D4	Rainy	Wld	Normal	Y
D5	Cloudy	Cool	Normal	Y

$$(1) \text{ Entropy (Enty)} \{ -3, -2 \} = -\frac{P}{P+N} \log_{10} \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_{10} \left(\frac{N}{P+N} \right)$$

$$= -\frac{3}{5} \log_{10}(3) - \frac{2}{5} \log_{10}(2)$$

$$(2) -\frac{P}{P+N} \log_{10} \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_{10} \left(\frac{N}{P+N} \right)$$

$$\log_{10}(2)$$

$$= 3.321 \left(-\frac{3}{5} \log_{10}(0.16) - \frac{2}{5} \log_{10}(0.04) \right)$$

$$= 9.721 / (0.133 + 0.159)$$

$$= 0.969$$

$$(i) \text{ Entropy (Weather) } \{ \text{Sunny} \{ 0, -2 \} \} = 3.321 \left[\frac{-0}{2} \log(0) - \frac{2}{2} \log(1) \right] \\ = 0$$

$$\text{Entropy (Cloudy) } \{ 1, 0 \} = \frac{1}{2} \log(1) - \frac{0}{2} \log(0) = 0$$

$$\text{Entropy (Rain)} \{ 2, 0 \} = 0$$

$$I_{\text{H}}(\text{Week}) = E(\text{entire}) - \frac{2 \times 0}{5} - \frac{1 \times 0}{5} - \frac{2 \times 0}{5}$$

$$I_{\text{H}}(\text{Week}) = 0.969$$

(ii) Temp

$$E(\text{H}(0)) \{ 1, -2 \} = 3.321 \left(\frac{-1}{3} \log\left(\frac{1}{3}\right) - \frac{2}{3} \log\left(\frac{2}{3}\right) \right) \\ = 3.321 (0.159 + 0.117) \\ = 0.9165$$

$$E(\text{H}(1)) \{ 1, 0 \} = 0$$

$$E(\text{H}(2)) \{ 1, 0 \} = 0$$

$$I_{\text{C}_1}(\text{Temp}) = 0.969 - \frac{3 \times 0.9165}{5} = 0.419$$

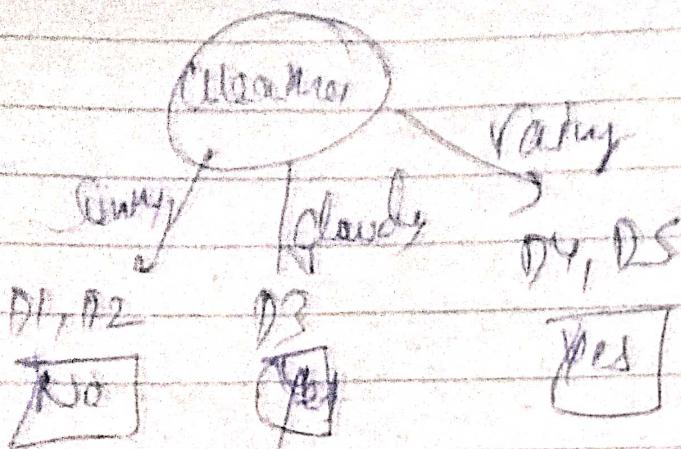
(iii) Humidity

$$E(\text{H}(0)) \{ 1, -2 \} = 0.9165$$

$$E(\text{H}(1)) \{ 2, 0 \} = 0$$

$$\underline{I_{\text{C}_2} = 0.419}$$

Maximum Entropy = 0.91



Gini Index

Purity of Data Set

$$Gini(D) = 1 - \sum_{i=1}^m (P_i)^2$$

$$Gini_A(D) = \frac{|D_1|}{D} Gini(D_1) + \frac{|D_2|}{D} Gini(D_2)$$

$$\Delta Gini_A(D) = Gini(D) - Gini_A(D)$$

$$Gini_A(D) = \frac{10}{14} Gini(D_1) + \frac{4}{14} Gini(D_2)$$

Ex	Day	Weather	Temp	Play
D1	Sunny	Hot	No	
D2	Sunny	not	No	
D3	Cloudy	Hot	Yes	
D4	Rainy	Hot	Yes	
D5	Rainy	Cool	Yes	

$$\text{Gini(Entire)} \{3, -2\} = 1 - 2p_i^2$$

$$= 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0.48$$

Gini (Weak)

$$\text{Gini(Sunny)} \{0, -2\} = 1 - \left(\frac{0}{2}\right)^2 - \left(\frac{2}{2}\right)^2 = 0$$

$$\text{Gini(Cloudy)} \{1, 0\} = 0$$

$$\text{Gini(Rainy)} \{2, -0\} = 0$$

$$\text{Gini(Weather)} = \frac{2 \times 0}{5} + \frac{1 \times 0}{5} + \frac{2 \times 0}{5} = 0$$

Gini(People)

$$\text{Gini(Hot)} \{1, -2\} = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 1 - 0.111 - 0.444 \\ = 0.445$$

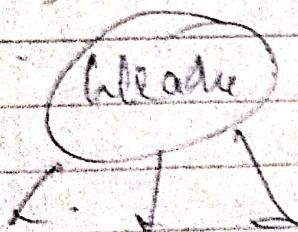
$$\text{Gini(Mild)} \{1, 0\} = 0$$

$$\text{Gini(Cool)} \{1, 0\} = 0$$

$$\text{Gini(People)} = \frac{3 \times 0.445}{5} + \frac{1 \times 0}{5} + \frac{0 \times 1}{5} = 0.267$$

Gini(Weather) < Gini(People)

less important



Day	Weather	Temp	age	income	Shirt	Grocery	buy
D1	Sunny	Hot	Y	H	No	F	No
D2	Sunny	Hot	M	H	No	E	No
D3	Cloudy		M	H	No	F	Y
D4	Rainy		S	M	No	F	Y
D5	Rainy		S	L	Yes	F	Y
D6	Rainy		S	L	Yes	E	N
D7	Cloudy		M	L	Yes	E	Y
D8	Sunny		Y	H	No	F	N
D9	Sunny		Y	L	Yes	F	Y
D10	Rainy		S	M	Yes	F	Y
D11	Sunny		Y	M	Yes	E	Y
D12	Cloudy		M	H	No	E	Y
D13	Cloudy		M	H	Yes	F	Y
D14	Rainy		S	M	No	E	N

$$E(\text{Entire}) \{9, -5\} = \frac{-9}{14} \log\left(\frac{9}{14}\right) - \frac{5}{14} \log\left(\frac{5}{14}\right) = 0.34$$

Intercell Age

$$E(\text{Young}) \{2, -3\} = 0.92$$

$$E(\text{Middle}) \{4, -3\} = 0$$

$$\Sigma (\text{Senior}) \{3, -2\} = 0.92$$

$$IC_1 = 0.34 - \frac{5}{14} \times 0.92 - \frac{4}{14} \times 0 - \frac{5}{14} \times 0.97 = 0.246$$

Cart Algo

Age

Young	5	Yes	2	
		No	3	→ 97% Acc.
Middle	4	Yes	4	
		No	0	
Senior	5	Yes	3	
		No	2	

Attribute

Rules

(error (err/total))

Age

Youth → NO	2/5	Total
Middle → Yes	0/4	0/4
Senior → Yes	2/5	

Bank Income

High	4	Yes	2
		No	2

Rules

Error

Total

High → Yes/No

3/4

High	6	Yes	4
		No	2

Med → Yes

2/6

5/14

High	9	Yes	3
		No	6

Low → Yes

1/4

Blurred

Yes	7	Yes	5
		No	1
No	7	Yes	5
		No	4

Rules

Error

Total

Yes → Yes

1/2

9/14

No → No

3/2

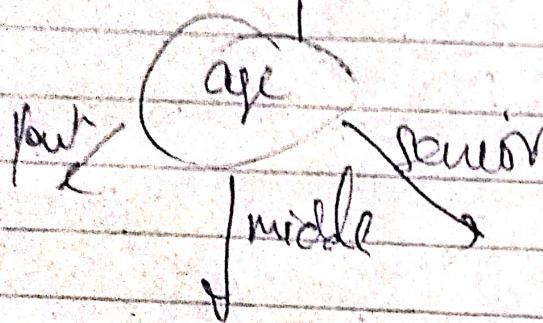
Crabdt

Page No.:

			Rules	Error
Pain	8	4.0	6	1
		No	2	0/8
Excellent	6	Yes	3	0/6
		No	3	0/6

$$\text{Error} = \frac{\text{No. of misclassifications}}{\text{Total number of samples}}$$

	Rules	Error	Total	
Age	Youth \rightarrow no Middle \rightarrow yes Senior \rightarrow yes	2/8 0/4 2/8	4/14	Same
Studies	Yes \rightarrow yes no \rightarrow no	2/2 2/2	4/4	direct <u>o</u>



Tree Pruning

To Reduce Overfitting

- Two approaches \rightarrow (i) Pre pruning
(ii) Post pruning

Pre Pruning \rightarrow halting its construction early

Upon halting, the node becomes a leaf. This leaf may fall the most frequent class among the subset tuple or the probability distribution of those tuple

When constructing a tree, measures such as statistical significance, information gain, Gini index, and so on, can be used to check the goodness of split.

If partitioning the leaf at a node would result in a split that falls below a prespecified threshold, then further partitioning is halted.

However choosing an appropriate threshold is very difficult.

High threshold \rightarrow oversimplified trees

Low threshold \rightarrow very little simplification

Various measure \rightarrow max-depth, max feature etc

"Re-Pruning is used on Big Data set"

Post Pruning \rightarrow remove subtree from a fully grown tree

A subtree at a given node is pruned by removing its branches and replacing it with a leaf. The leaf is labeled with the most frequent class among all subtree being replaced.

"It is used on Small Data Set"

CART \rightarrow Cost Complexity pruning algo \rightarrow no of leaves, error rate.

C4.5 \rightarrow pessimistic pruning \rightarrow overfitting

Bayesian Classifiers

Bayesian classifiers are statistical classifiers. They can predict class membership probabilities such as the probability that a given tuple belongs to a particular class.

→ Bayes' Theorem

→ Naive Bayesian classifier

Bayes' Theorem

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

$$P(A|B) \times P(B) = P(B|A) P(A)$$

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Likelihood

H → hypothesis
D → Data

$$P(H|D) = \frac{P(D|H) \times P(H)}{P(D)}$$

Posterior

marginal

Prior

$P(H) \rightarrow$ Prior \rightarrow Probabilities of hypotheses before giving evidence

Naive Bayes Classifier

Determine if weather = sunny, Temp = cool, Humidity
= high, Windy = strong?

$$P(\text{Yes}) = \frac{9}{14} = 0.64 \quad P(\text{No}) = \frac{5}{14} = 0.36$$

Conditional Probability of each feature

Weather	Yes	No	Humidity	Yes	No
Sunny	Sunny & Yes	Sunny & NO	high	3/9	4/5
	Total Yes	Total NO	Normal	6/9	1/5
	= 2/9	3/5			
Cloudy	4/9	0/5	Temp	Yes	No
Rainy	3/9	2/5	Hot	4/9	2/5
			Mild	4/9	2/5
			Cool	3/9	1/5

Windy	Yes	No
Strong	3/9	3/5
Weak	6/9	2/5

$$V_{NB} = \arg \max_{V_j \in \{\text{Yes}, \text{No}\}} P(V_j) \prod_i P(a_i | V_j)$$

$$V_{NB}(\text{Yes}) = P(\text{Yes}) P(\text{Sunny} | \text{Yes}) P(\text{Cool} | \text{Yes}) P(\text{High} | \text{Yes}) P(\text{Strong} | \text{Yes})$$

$$= \frac{9}{14} \times \frac{2}{9} \times \frac{3}{8} \times \frac{3}{9} \times \frac{3}{9} = 0.053$$

$$V_{NB}(\text{no}) = P(\text{no}) P(\text{Sunny}/\text{no}) P(\text{Cool}/\text{no}) P(\text{High}/\text{no}) \\ P(\text{Snowy}/\text{no})$$

$$= \frac{5}{14} \times \frac{3}{5} \times \frac{1}{5} \times \frac{4}{5} \times \frac{3}{5} = 0.026$$

To Overcome Zero Probability Problem

Laplace Smoothing

$$P_k(x|y) = \frac{\text{Count}(x,y) + k}{\text{Count}(y) + k |X|}$$

$k \rightarrow$ represents smoothening parameter > 0

$|X| \rightarrow$ no of dimension / features \rightarrow Weaker
 { Sun, Rain, Cloud } \rightarrow

$$P_{k=1}(\text{Outlook} = \text{Cloudy}/\text{no}) = \frac{\text{Count}(\text{Outlook} = \text{Cloudy}/\text{no}) + k}{\text{Count}(\text{no}) + 1 \times 3} \\ = \frac{1}{8}$$

Ex. Fruit = {Yellow, Sweet, Long}

Fruit	Yellow	Sweet	Long	Total
Orange	350	450	0	650
Banana	400	300	350	400
Orange	50	100	80	150
Total	800	850	400	1200

$$P(\text{Yellow} | \text{Orange}) = P(\text{Orange} / \text{Yellow}) \times P(\text{Yellow}) \\ P(\text{orange})$$

$$= \frac{350}{800} \times \frac{800}{1200} = 0.53$$

$$\frac{350}{1200}$$

$$P(\text{Sweet} | \text{Orange}) = P(\text{Orange} / \text{Sweet}) \times P(\text{Sweet}) = \frac{450}{850} \times \frac{850}{1200} \\ P(\text{orange}) \quad \frac{450}{1200}$$

$$= 0.69$$

$$P(\text{Long} | \text{orange}) = P(\text{Orange} / \text{Long}) \times P(\text{Long}) \\ P(\text{orange}) = \frac{0}{1200} = 0$$

$$P(\text{Fruit} | \text{orange}) = P(Y|O) P(S|O) \times P(L|O) = 0$$

Similarly

$$P(\text{Fruit L} | \text{banana}) = P(Y|B) P(S|B) P(L|B) = 1 \times 0.25 \times 0.9 \\ = 0.225$$

$$P(\text{Fruit} | \text{other}) = P(Y|\text{other}) P(S|\text{other}) P(L|\text{other}) = 0.122$$

∴ greatest \rightarrow fruit = banana

Rule Based Classification

Using IF - Then Rules

↓
Rule antecedent
or
Precondition → Rule Consequent

Rule Conflict

Conflict Resolution Strategy to figure out which rule gets to fire and assign its class predictor to X.

(i) Size Ordering → Assign the highest priority to the triggering rule that has the "toughest" requirements, where toughness is measured by the rule antecedent size.

(ii) No Rule Ordering →

↳ a) Class based → Classes are sorted in order of decreasing "importance" such as by decreasing order of prevalence.

b) Rule Based → Organized into class priority lists according to some measure & rule quality such as accuracy, coverage, etc.

$$\text{Coverage } (R) = \frac{n_{\text{covers}}}{|D|}$$

$$\text{Accuracy } (R) = \frac{n_{\text{correct}}}{n_{\text{covers}}}$$

Rule Induction Using a Sequential Covering algo

Rule_Set = \emptyset

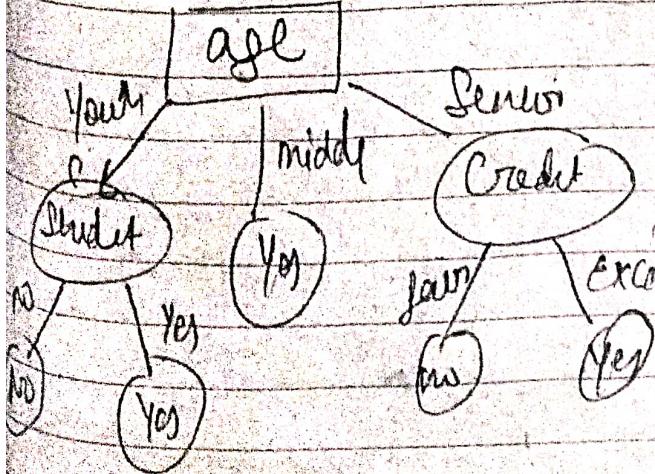
for each class c do

 Rule = Learn_One_Rule (D , att_valu, c)
 remove Rule covered by Rule from D
 Rule_Set = Rule_Set + Rule;

Until terminally condition

End for
return Rule_Set;

Rule Extraction from Decision Tree



R1: If age = Young and student = no
 then buy Computer = no

R2: If age = Young and student = yes
 then buy Computer = yes

R3: If age = middle and buy
 Computer = yes

R4: If age = senior and Credit = fair
 then buy Computer = no

R5: If age = senior and Credit = excellent then buy Computer = yes

Characteristics of Rules

- ① Mutually Exclusive \rightarrow we cannot have rule conflicts here because now the two rules will triggered for the same tuple.
- ② Exhaustive \rightarrow there is one rule for both, and any tuple can satisfy ~~any~~ for each possible attribute value combination.

Rule Quality Measures

FOIL \rightarrow first order induction learning

$$\text{① FOIL score} = \text{pos}' \log_2 \left(\frac{\text{pos}'}{\text{pos}' + \text{neg}'} \right) - \text{neg}' \log_2 \left(\frac{\text{neg}'}{\text{pos}' + \text{neg}'} \right)$$

$$\text{② Likelihood-Ratio} = 2 \sum_{j=1}^m f_j \log \left(\frac{f_j}{e_j} \right)$$

Metrics for Evaluating Classifier Performance

		Predicted		Total
		Yes	No	
Actual	Yes	TP	FN	P
	No	FP	TN	N
Total	P	N		P+N

MeasureFormula

accuracy, true positive rate

$$\frac{TP + TN}{P+N}$$

error rate, misclassified rate

$$\frac{FP + FN}{P+N}$$

sensitivity, true positive rate,

$$\frac{TP}{P}$$

Recall

$$P = (FP + FN)$$

Specificity, true negative rate,

$$\frac{TN}{N} = \frac{TN}{(FP + TN)}$$

Precision

$$\frac{TP}{TP+FP}$$

 f, f_1, f_2 score

2 x Precision x Recall

harmonic mean of precision

Recall

Recall

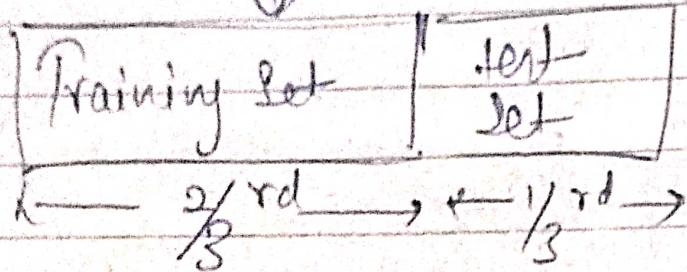
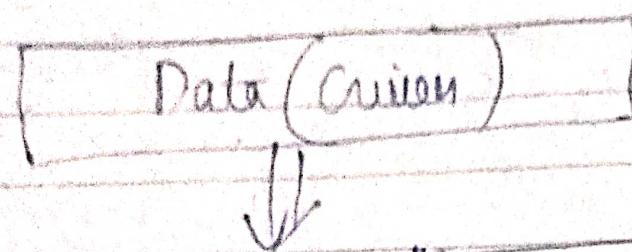
 β , $\beta \rightarrow$ non negative

$$(1+\beta^2) \times \text{Precision} \times \text{Recall}$$

$$\beta^2 \times \text{Precision} \times \text{Recall}$$

Holdout Method

The given data are randomly partitioned into two independent sets, a training set and test set.



Training set is used to derive the model

Test set is used to estimate the model's accuracy

The estimate is pessimistic because only a portion of the initial data is used to derive the model.

Random Subsampling

It is a variation of holdout method in which holdout method is repeated 'k' times.

The overall accuracy \rightarrow average of each accuracy

CROSS VALIDATION

K-FOLD Validation

Data

\Rightarrow

$D_1 | D_2 | D_3 | D_4 | D_5 \dots | D_k$

$D_1 | D_2 | D_3 \dots | D_k$

test
data set

training
data set

$D_1 | D_2 | D_3 | \dots | D_k$

↓
test
dataset

→ training, Val set

and so on.

Accuracy = $\frac{\text{Overall no. of correct classifications}}{\text{Total no. of samples in initial}}$

Leave - One - Out

It is special case of K-fold Cross Validation where K is set to the number of initial samples.

Only one sample is "left out" at a time for the test set.

Stratified K-fold

The folds are stratified so that the class distribution of the samples in each fold is approx the same as that in initial dat.

Half richtig 10 fold Cross Validation \rightarrow accuracy bias
low bias and variance

Bootstrap

The Bootstrap method samples from given training tuples uniformly with replacement. That is each time a tuple is selected, it is equally likely to be selected again and re-added to the training set.

ROC Curve Receiver Operating Characteristic

TPR vs FPR

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

↓
TP

↓
FP

N

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

True Positive

$$P \rightarrow P$$

False Positive

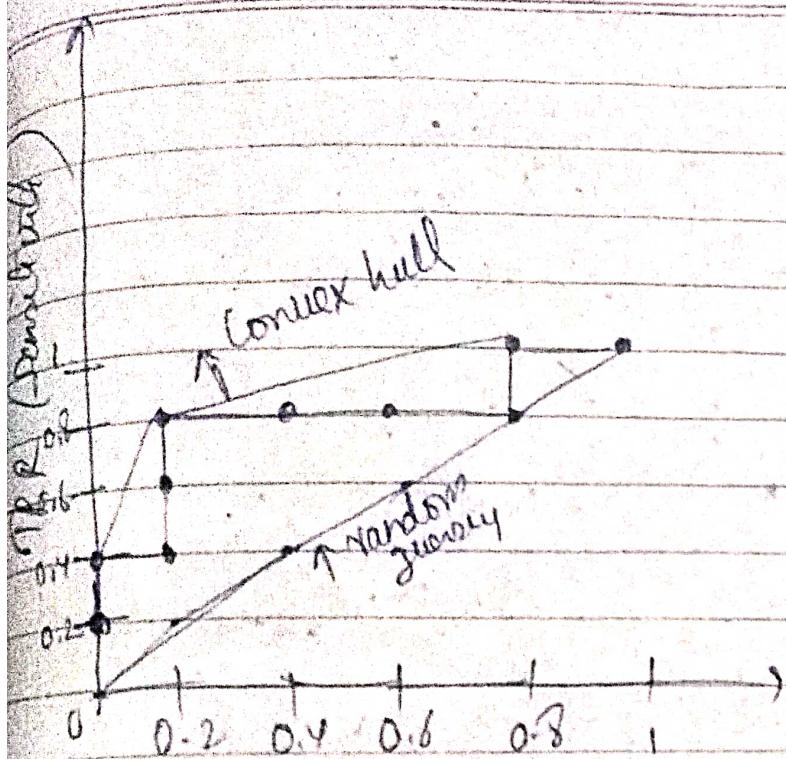
$$N \rightarrow P$$

→ Total Positive

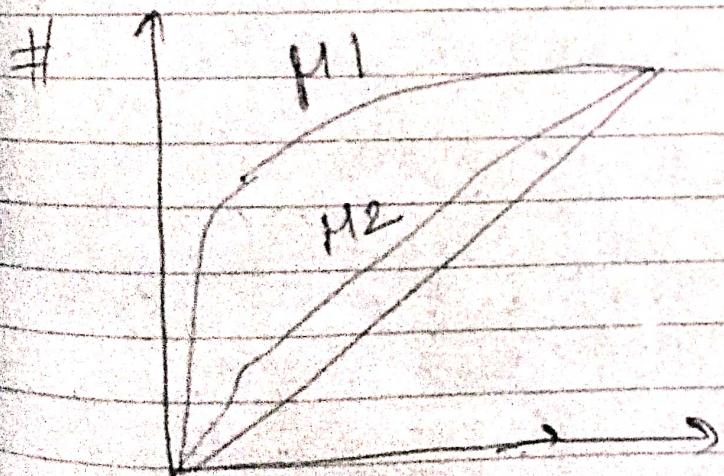
→ Total Negative

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

Ex Table #	Class	Prob	TP	FP	$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$	$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$
1	$P \rightarrow P$	0.90	1	0	$\frac{1}{1+0} = 1$	$0 = 0$
2	$P \rightarrow P$	0.80	1+2=3	0	$\frac{1}{1+2} = 0.33$	$0 = 0$
3	$P \rightarrow N$	0.80	2	0+1=1	$\frac{2}{2+1} = 0.66$	$0 = 0$
4	$P \rightarrow P$	0.60	2+1=3	1	$\frac{2}{2+1} = 0.66$	$0 = 0$
5	$P \rightarrow P$	0.55	3+1=4	1	$\frac{3}{3+1} = 0.75$	$0 = 0$
6	$P \rightarrow N$	0.54	4	1+1=2	$\frac{4}{4+2} = 0.66$	$0 = 0$
7	$P \rightarrow N$	0.53	4	2+1=3	$\frac{4}{4+3} = 0.53$	$0 = 0$
8	$P \rightarrow N$	0.51	4	3+1=4	$\frac{4}{4+4} = 0.5$	$0 = 0$
9	$P \rightarrow P$	0.50	2+1=3	4	$\frac{2}{2+3} = 0.4$	$0 = 0$
10	$P \rightarrow N$	0.40	5	4+1=5	$\frac{5}{5+5} = 0.5$	$0 = 0$



FPR ($1 - \text{specificity}$)

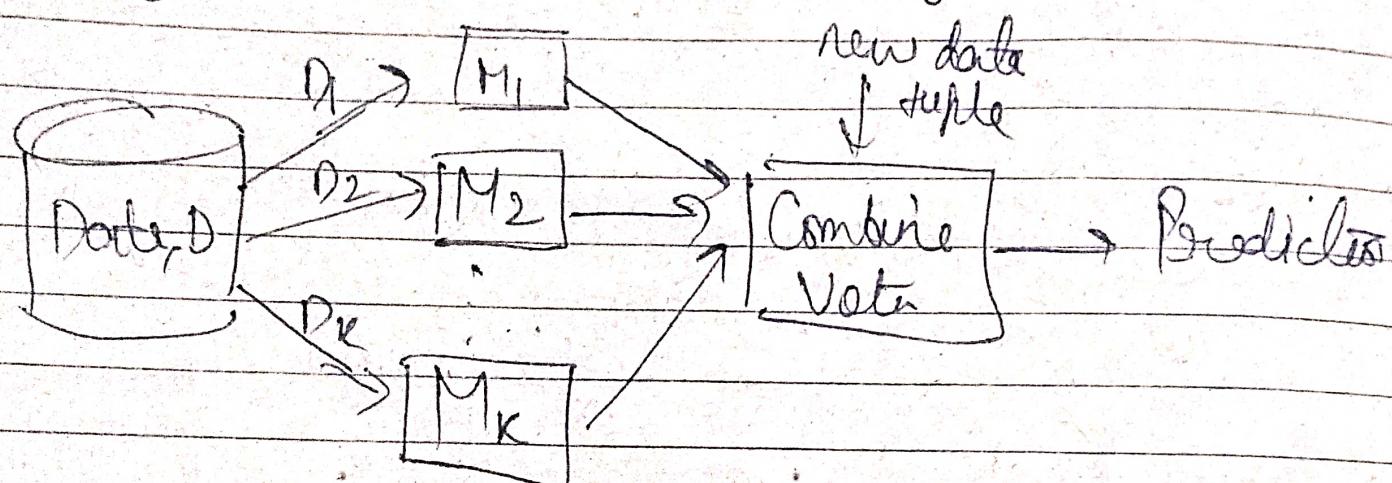


The ROC Curve is to the diagonal line, the less accurate the model is,

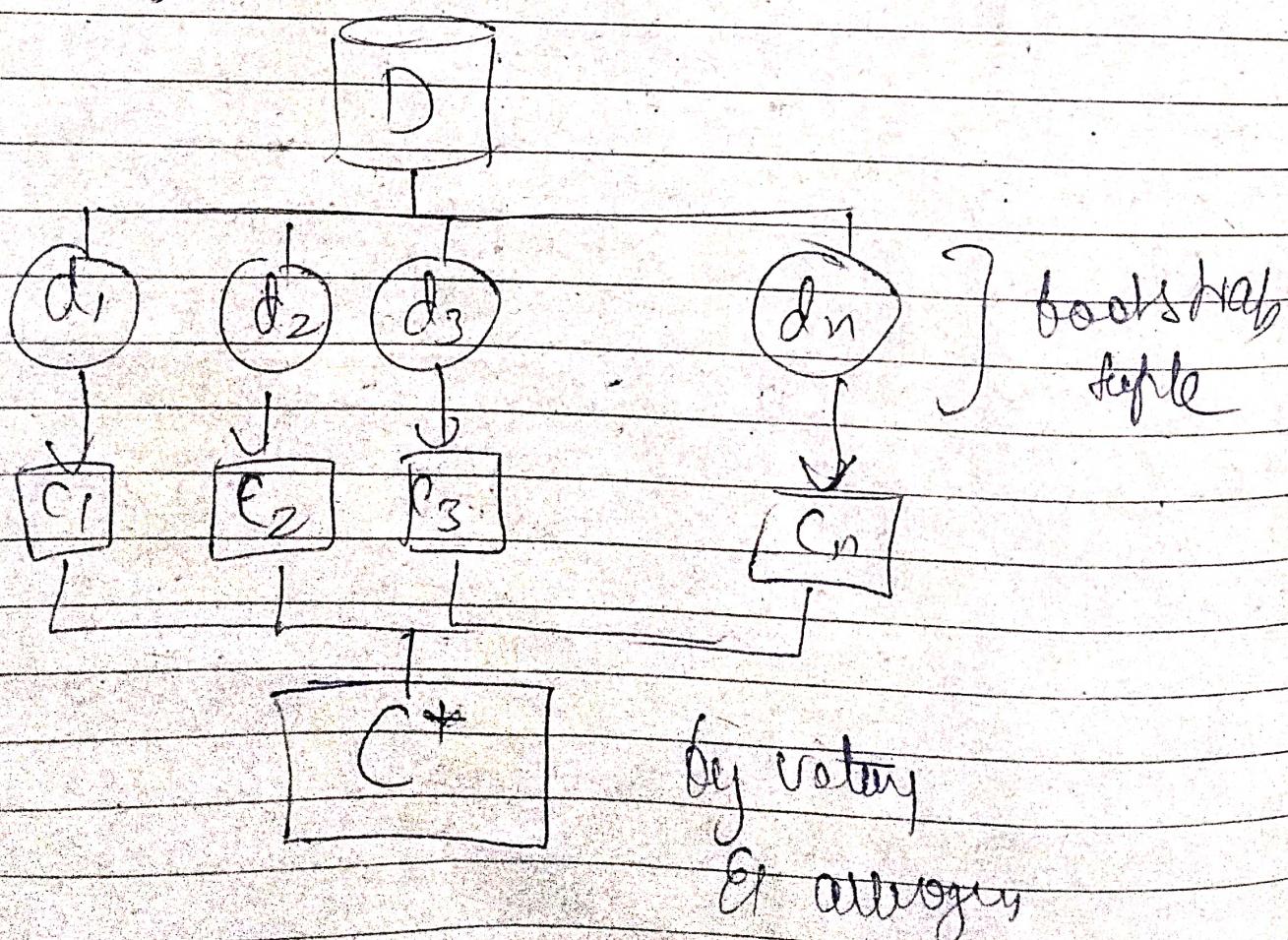
Thus M1 is more accurate

Ensemble Learning ~~Machine Learning~~

It is a supervised learning technique used in machine learning to improve overall performance by combining the prediction from multiple models.



Bagging \rightarrow Bootstrap Aggregation



Given set D , of n tuples, bagging works as follows.
For iteration i ($i=1, 2, \dots, K$). a training set, D_i of n tuples is sampled with replacement from original set of tuple D .

The bagged classifier often has significantly greater accuracy than a single classifier derived from D , the original training data.

It will not be considerably worse and is more robust to the effects of noisy data and overfitting.

The increased accuracy occurs because the composite model reduces the variance of the individual classifiers.

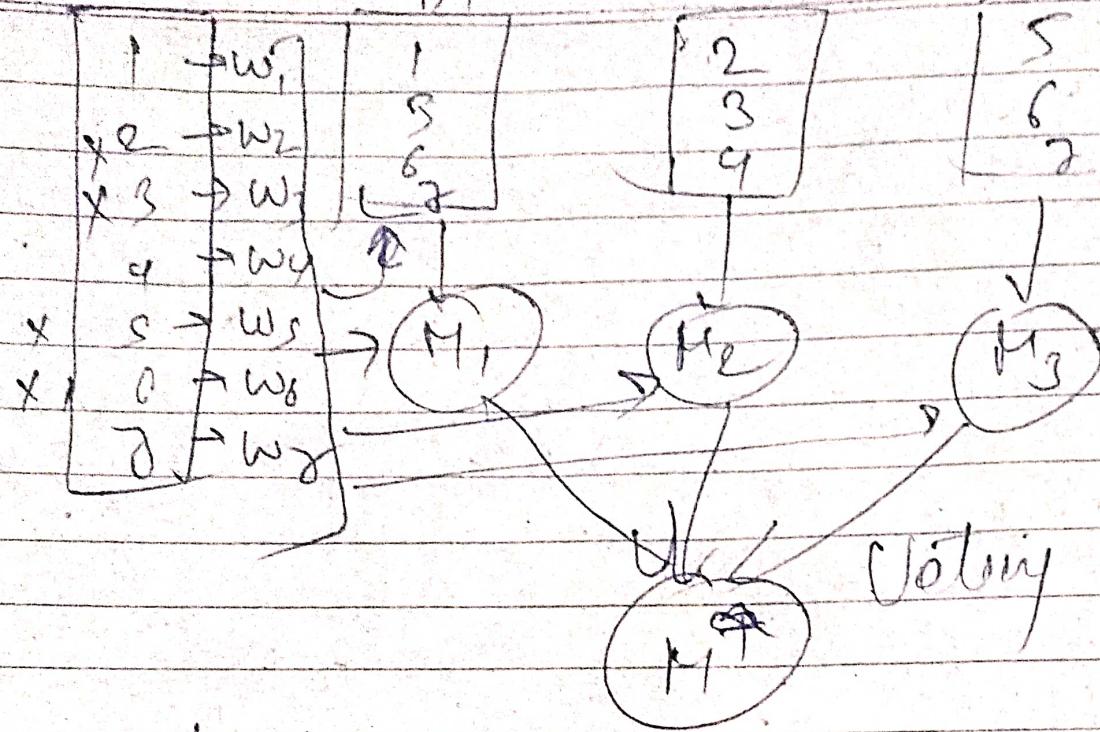
Boosting

Weights are also assigned to each training tuple. A series of K classifiers is iteratively learned. After a classifier, H_i is learned, the weight are updated to allow the subsequent classifier H_{i+1} to "pay more attention" to the training tuple that are misclassified by H_i . The final boosted classifier H^* , combine the votes of each individual classifiers, where the weight of each classifier's vote is a function of its accuracy.

Dataset

D₁

D₂



misclassified table should give more priority to go to next subsequent dataset. D₂

~~Not best~~

True Label	Predict Probable	TP	FP	TPR & FPR
N → 1	0.3	0	1	0
P → 0	0.6	0	1	0
P → 1	0.3	0	2	0
N → 0	0.4	0	2	0
P → 1	0.9	1	2	1/3

For 0.9 threshold
 ≈

True	Predicted	threshold(0,9)	(0,2)	(0,3)
1	0.8	0	1	1
0	0.6	0	0	1
1	0.2	0	1	1
0	0.4	0	0	1
1	0.9	1	1	1

(i) Threshold 0,9

True	Predicted	TP	FP	TPR	FPR
1	0	0	0	0	0
0	0	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0
1	1	1	0	1/3	0

$$TPR = \frac{TP}{TP+FN} = \frac{1}{1+2} = \frac{1}{3}$$

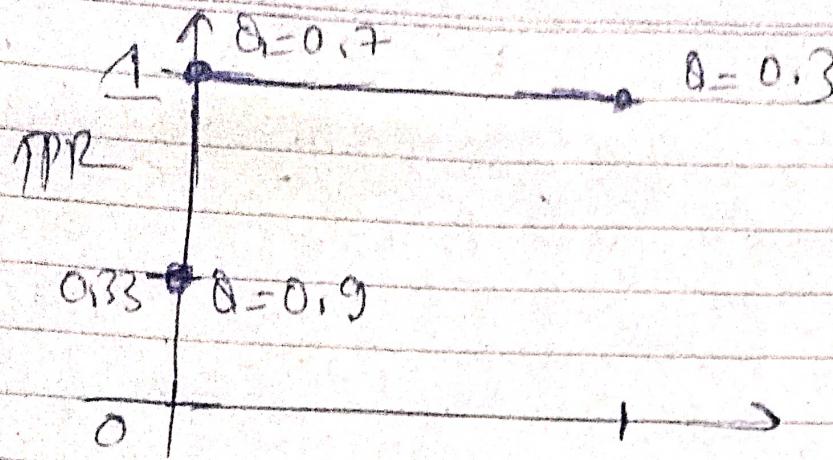
$$FPR = \frac{FP}{FP+TN} = \frac{0}{2} = 0$$

$$(ii) \text{ Threshold } (0,2) \quad TPR = \frac{TP}{TP+FN} = \frac{3}{3} = 1$$

$$FPR = \frac{FP}{FP+TN} = 0$$

$$(iii) \text{ Threshold } (0,3) \quad TPR = \frac{TP}{P} = \frac{3}{3} = 1$$

$$FPR = \frac{FP}{N} = \frac{2}{3} = \frac{2}{3}$$



Adaboost

Adaboost assign each training tuple an equal weight of $\frac{1}{d}$. Generating K classifier for the ensemble repeat K round through the start of the algo. In round i , the tuples from D are sampled to form a training set D_i , g_i is used,

$$\text{error}(h_i) = \sum_{j=1}^d w_j$$

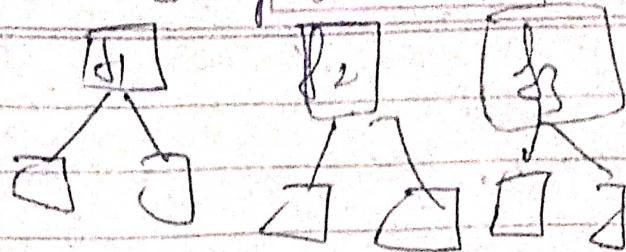
$$\text{Performance} = \frac{1}{2} \log \left(\frac{1 - \text{error}(h_i)}{\text{error}(h_i)} \right)$$

① Misclassified \rightarrow weight \uparrow

$$\text{New weight} = \text{old weight} \times e^{\text{Performance}}$$

② Classified \rightarrow weight \downarrow

$$\text{New weight} = \text{old weight} \times e^{-\text{Performance}}$$

$f_1, f_2, f_3 \text{ O/P}$ 

Random Forest

Random forest is a commonly used machine learning algorithm.

A Random forest is an ensemble learning method where multiple decision trees are constructed and then they are merged to get a more accurate prediction.

Random forest became popular because of its ease of use and flexibility in handling both classification and regression problems.

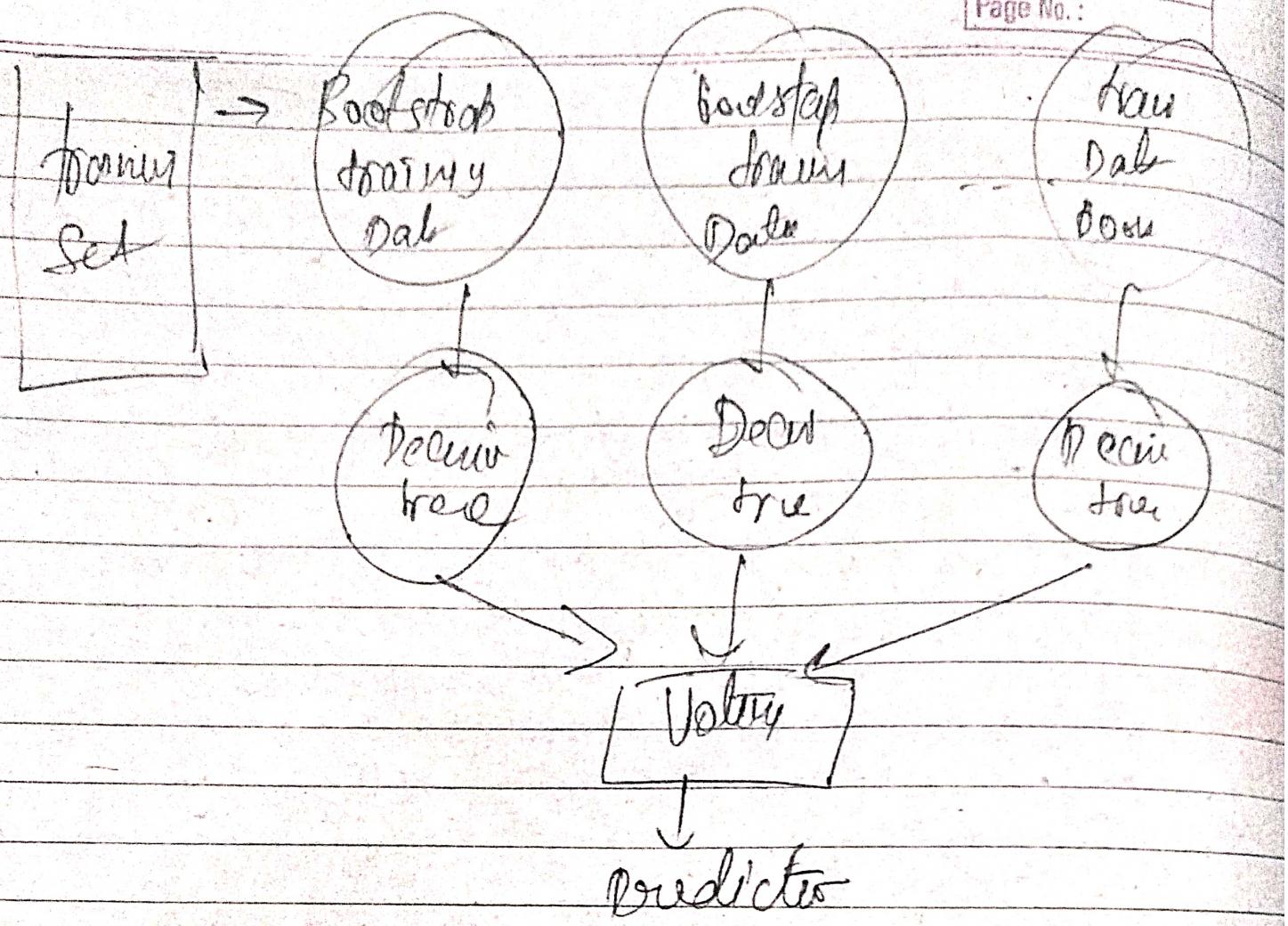
Steps

Build Random forest

a) If the number of examples in the training set is N , take a sample of n examples at random - but with replacement from the original set.

b) If there are ' M ' input variables, m variables are selected at random out of the M and the best split on these m is used to split the node. The value of m is held constant during the generation of the various trees in forest.

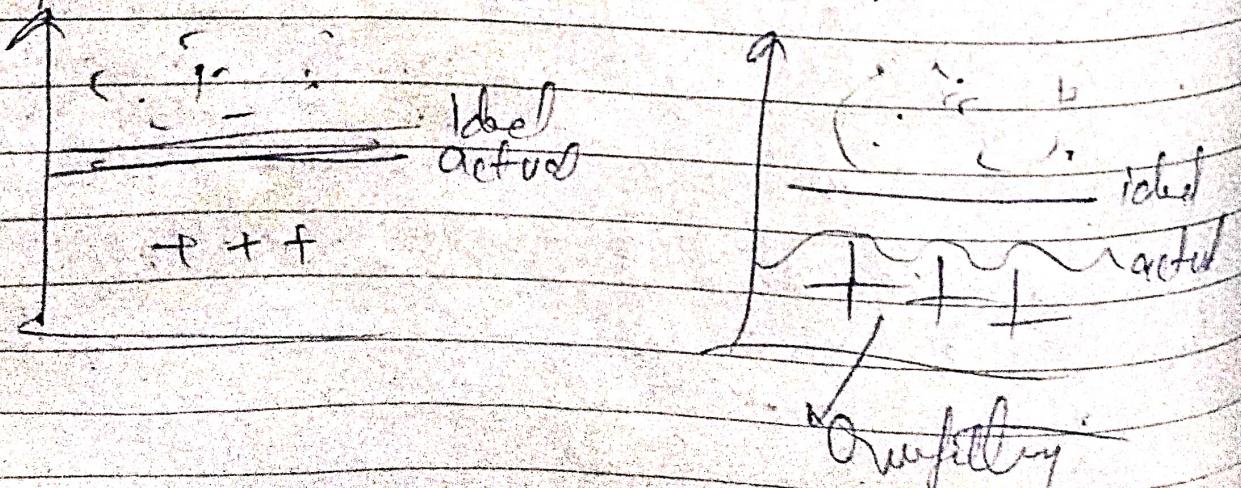
For new test data points, find the prediction of each decision tree, and assign the new data points to the category that wins the majority vote.



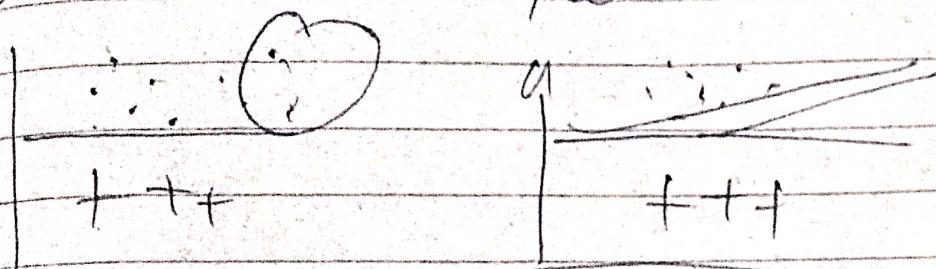
Random Forest are Comparable in accuracy to adaboost.
 Yet are more robust to errors and outliers

Imbalance data

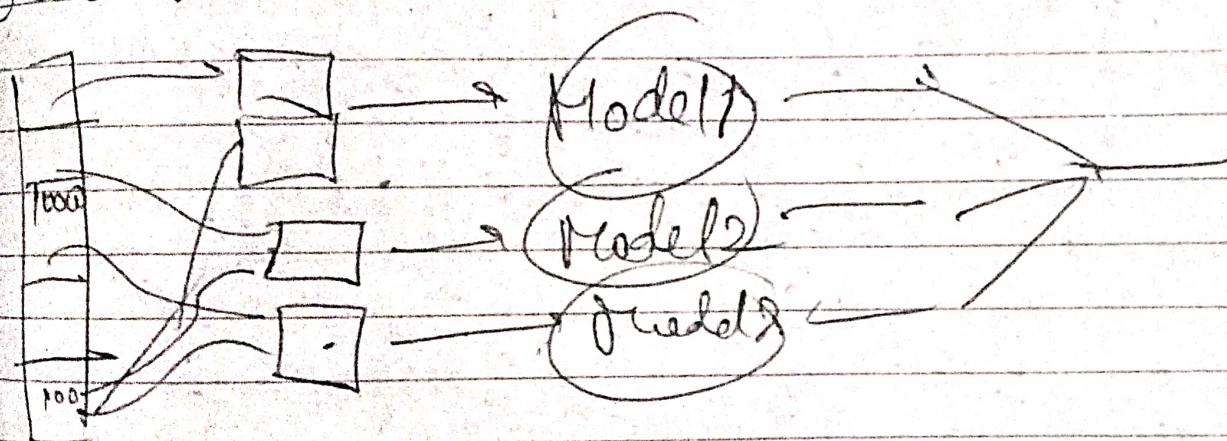
- ① Over sampling → resampling the positive tuple so that the resulting training set contains an equal number of positive and negative tuple



① Under Sampling → decreasing the number of negative sample. If randomly eliminate sample from the majority class until there are equal in number.



② Ensemble



③ threshold mining

Chapt 10 Data Mining

Clustering → Unsupervised Learning because the class label information is not present.

The following are typical requirements of clustering →

- (1) Scalability
- (2) Ability to deal with different types of attributes
- (3) Discovery of clusters with arbitrary shape
- (4) Requirement for domain knowledge to determine input parameters
- (5) Deal with noisy data
- (6) Incremental clustering and insensitivity to input order
- (7) Capability of clustering high dimensionality data
- (8) Constraint Based Clustering
- (9) Interpretability and usability

The following are orthogonal aspects with which clustering methods can be compared.

- (1) Partitioning Criteria
- (2) Separation of cluster
- (3) Similarity measure
- (4) Clustering space

Basic Clustering Methods

(1) Partitioning Method → Partitioning method
Construct K partitions of the data, where each partition represents cluster $k \in K$.

→ K-means Clust
K-medoids Clust

MethodsPartitioning methodsGeneral Characteristics

- Find mutually exclusive clusters of spherical shape
- Distance based
- May use mean or medoid as cluster center
- Effective for small to medium size data sets

Hierarchical Methods

- Clustering is a hierarchical downward process
- Cannot correct erroneous merges or splits
- May incorporate other techniques like microclustering or consider object "linkage"

Density based Methods

- Can find arbitrarily shaped clusters
- Cluster are dense regions w.r.t. object in shape that are separated by low-density regions
- Cluster density :- Each point has a minimum number of points within its "neighborhood".
- May filter out outliers

Grid based

- Use a multi-resolution grid data structure
- Fast processing time

K means Clustering

Point	Centroid (2,3)	Centroid (6,5)	Class
(2,3)	1	3.6	1
(3,3)	1	1.4	1
(5,4)	3.16	1.4	2
(6,5)	4.4	0	2
(9,8)	$\sqrt{65} = 8.0$	$\sqrt{11} = 3.3$	2
(10,8)	$\sqrt{41} = 6.4$	5	2

New cluster

$$1 \left(\frac{5}{2}, \frac{6}{2} \right)$$

$$2 \rightarrow \left(\frac{30}{9}, \frac{24}{9} \right)$$

$$\therefore (2.5, 6)$$

Point	Centroid 1 (2,3,1)	Centroid 2 (3,5,1)	old	new.
(2,3)	0.5	.	1	1
(3,3)	0.5	.	1	1
(5,4)	$\sqrt{2.25}$	$\sqrt{10.25}$	2	1
(6,5)	$\sqrt{16.25}$	$\sqrt{3.25}$	2	2
(9,8)	.	.	2	2
(10,8)	.	.	2	2

New Cent

$$1 (3.33, 3.73) \quad 2 (8.33, 6.66)$$

K-Medoids

Partitioning around Medoids or the K-medoids algorithm is a partitioning clustering algorithm which is slightly modified from the K-means algo.

In K-means algo, they choose means as the centroids but in the K-medoids, data point are chosen to be medoids

x	y	$C_1 = (3, 4)$	$C_2 = (7, 4)$	cluster
2	6	3	7	1
3	4	0	9	1
3	8	4	8	1
4	7	4	6	1
6	2	5	3	2
6	4	3	1	2
7	3	5	1	2
7	4	4	0	2
8	5	6	2	2
7	6	6	2	2

find distance Manhattan Dist = $|x_1 - x_2| + |y_1 - y_2|$

Q: $(2, 6), (3, 4), (3, 8), (4, 7)$

$(2, 6) (3, 4) (3, 8) (4, 7)$

Q: $(6, 2) (5, 9) (2, 1) (7, 4) (8, 5) (2, 1)$

$(6, 2) (6, 4) (2, 3) (2, 4) (2, 1) (8, 5)$

$$\text{Cost}(c_i, x) = \sum_{j \neq i} |c_j - x_j|$$

$$\begin{aligned} \text{Total Cost} &= \text{Cost}((3, 4), (2, 6)) + \text{Cost}((7, 4), (3, 8)) + \\ &\quad \text{Cost}((3, 4), (4, 2)) + \text{Cost}((7, 4), (6, 2)) + \\ &\quad \text{Cost}((7, 4), (6, 4)) + \text{Cost}((7, 4), (2, 2)) + \\ &\quad \text{Cost}((2, 4), (8, 5)) + \text{Cost}((2, 4), (2, 1)) \\ &= 1 + 2 + 0 + 4 + 1 + 3 + 1 + 2 + 1 + 0 + \\ &\quad 0 + 1 + 1 + 1 + 0 + 2 = 20 \end{aligned}$$

Step 3 Again I choose random new point others
than 2 medians

O: (7, 2)

X	Y	C1 = (3, 4)	O = (2, 3)	cluster
2	6	3	9	C1
3	4	0	5	C1
3	8	9	9	C1
4	2	4	2	C1
6	2	5	2	O
6	4	3	2	O
2	3	5	0	O
8	4	4	1	O
8	5	6	3	O
2	6	6	3	O

C1: {(2, 6), (3, 4), (3, 8), (4, 2)}

O: {(6, 2), (6, 4), (2, 3), (2, 4), (8, 5), (2, 1)}

$$\begin{aligned}
 \text{Cost} &= \text{Cost}((3,4), (2,6)) + \text{Cost}((3,4), (3,7)) + \text{Cost} \\
 &\quad ((3,4)(4,10) + (\text{out}((2,3)(6,2))) + \text{Cost}((2,3)(6,4)) \\
 &\quad + \text{Cost}((2,3)(2,4)) + \text{Cost}((2,3)(8,5)) + \text{Cost}((7,2)(2,1)) \\
 &= 1 + 2 + 0 + 4 + 1 + 3 + 1 + 1 + 1 + 1 + 1 + 2 \\
 &\quad \underline{\underline{0+2}} = 22
 \end{aligned}$$

Cost of Swapping

$$\begin{aligned}
 \text{Current Cost} &> \text{Previous Cost} \\
 22 &> 20
 \end{aligned}$$

Do no swapping is needed C2 with 0

$$\begin{aligned}
 \text{So cluster find are } C1: & \quad \{(2,6)(\underline{3,4})(7,2), (4,10) \\
 C2: & \quad \{(6,2)(6,4), (2,3)(2,4) \\
 & \quad (8,5), (2,1)\}
 \end{aligned}$$

"The K-medoids is less more robust than K
mean, as the presence of noise and outlier values
a medoid is less influenced by outlier or some
extreme values than a mean"

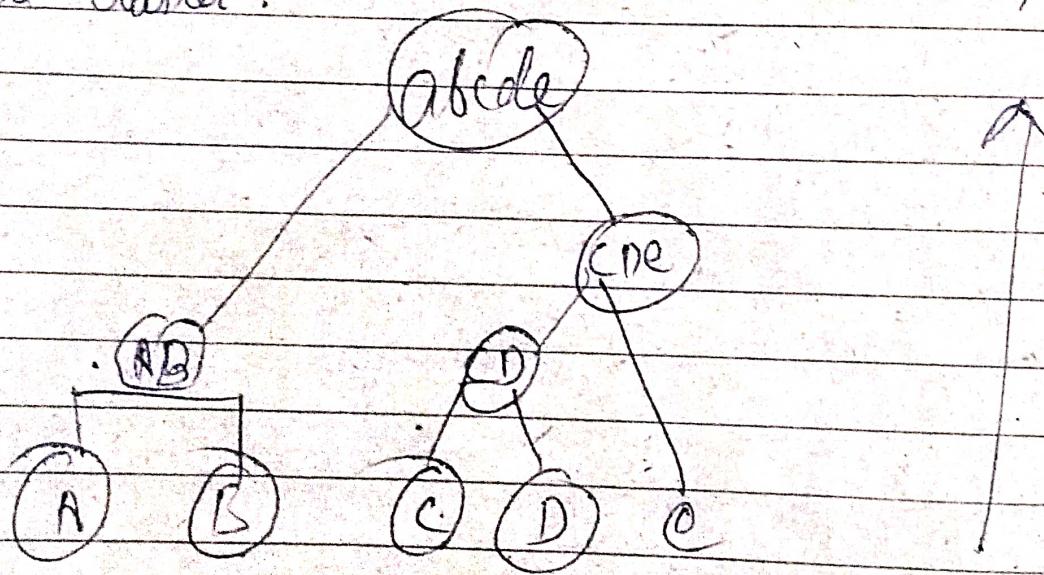
but K-medoids take $O(\underline{k(n-k)^2})$

Hierarchical Methods

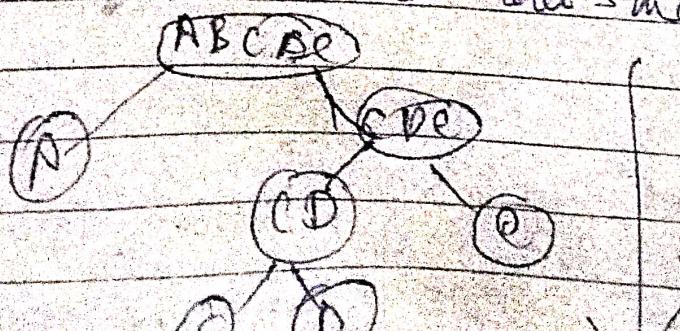
Working by grouping data objects into a hierarchy or "tree" of clusters.

Agglomerative Versus divisive

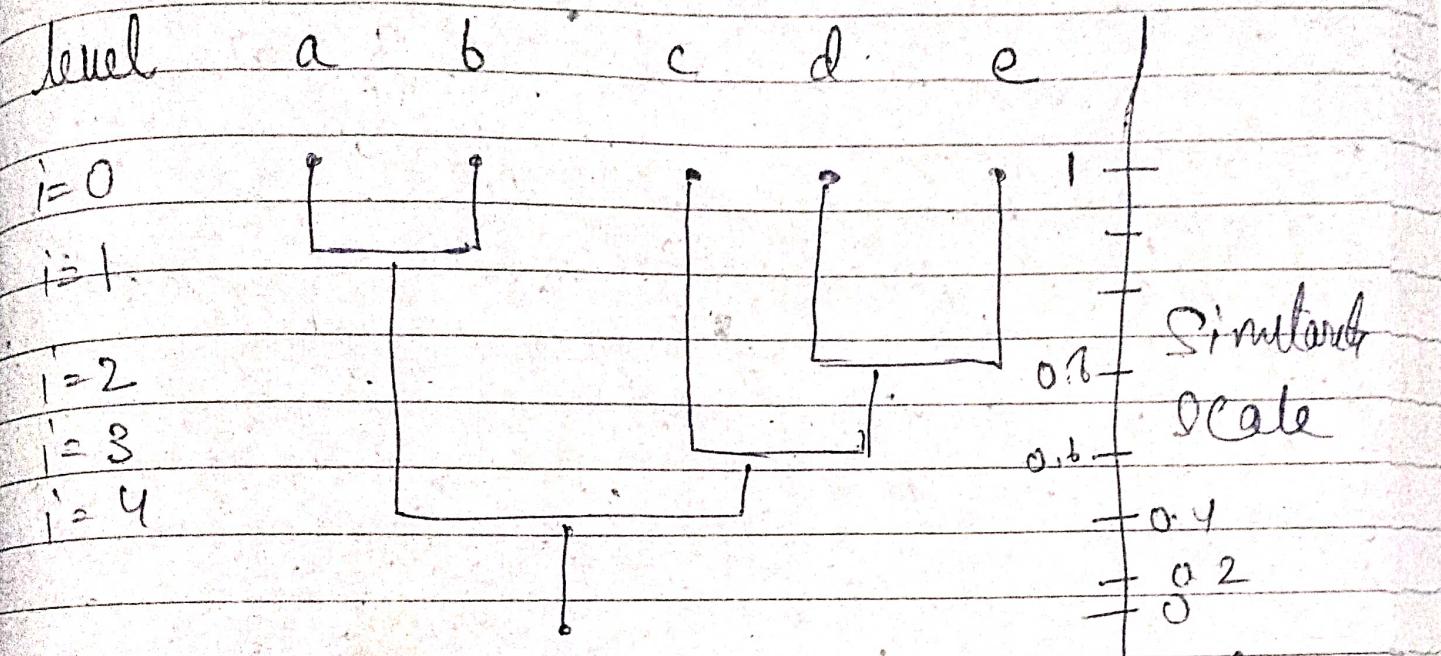
An agglomerative hierarchical method used bottom up strategy. It typically starts by letting each object form its own cluster and iteratively merges clusters into larger and larger clusters, until all the objects are in a single cluster.



An divisive method used a top down approach. It starts by placing all objects in one cluster, which is the hierarchy root. It then divides the root cluster into several smaller subclusters, and recursively partition those clusters into smaller ones.



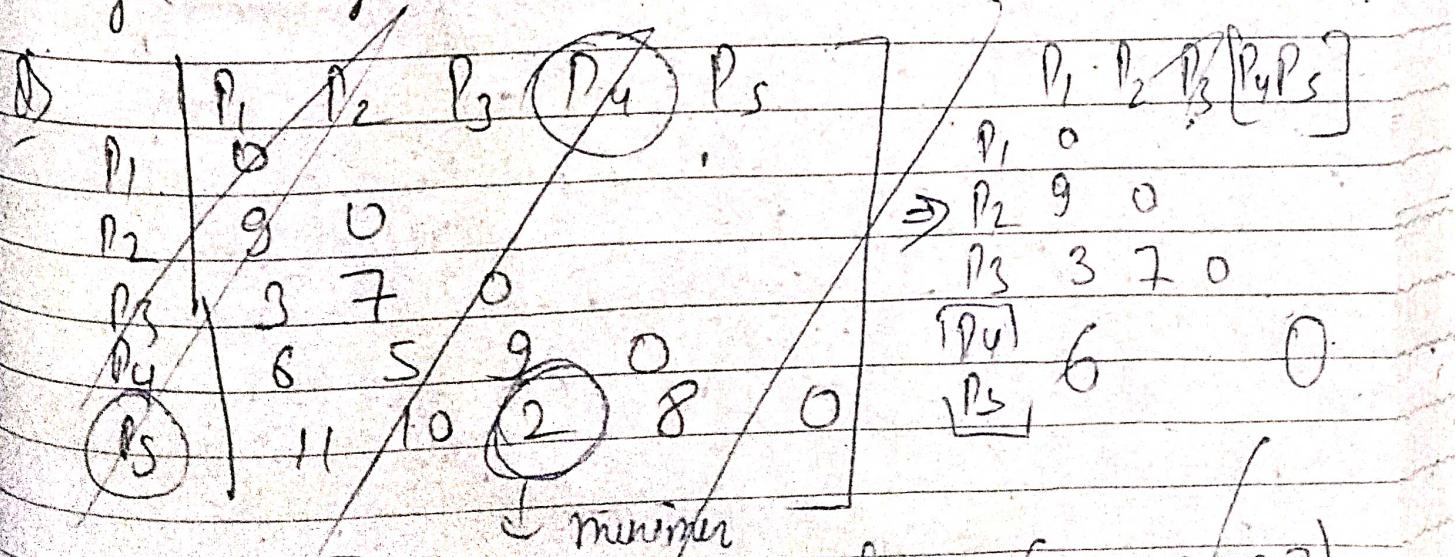
Dendrogram representation.



A Challenges within divisive method \rightarrow how to partition a large cluster into several smaller ones

Agglomerative

Single linkage \rightarrow minimum



P_4, P_5

minimum

distance $(P_1, [P_3, P_4])$

$$= \min(d_{14}(P_1, P_3), d_{14}(P_1, P_4))$$

Clark

$$= \min(9, 6) = 6$$

	P_1	P_2	$\boxed{P_3}$	P_4	P_5	
P_1	0					
P_2	9	0				
P_3	3	7	0			
P_4	6	5	9	0		
P_5	11	10	(2)	8	0	

	P_1	P_2	$\boxed{P_3 P_5}$	P_4	
P_1	0				
P_2	0				
P_3	0	3	7	0	
P_4	6	5	8	0	

minimum

$$\text{distave} (P_1, [P_3, P_5]) = \min (\text{dist}(P_1, P_3), \text{dist}(P_1, P_5)) \\ = \min (3, 11) = 3$$

Single linkage

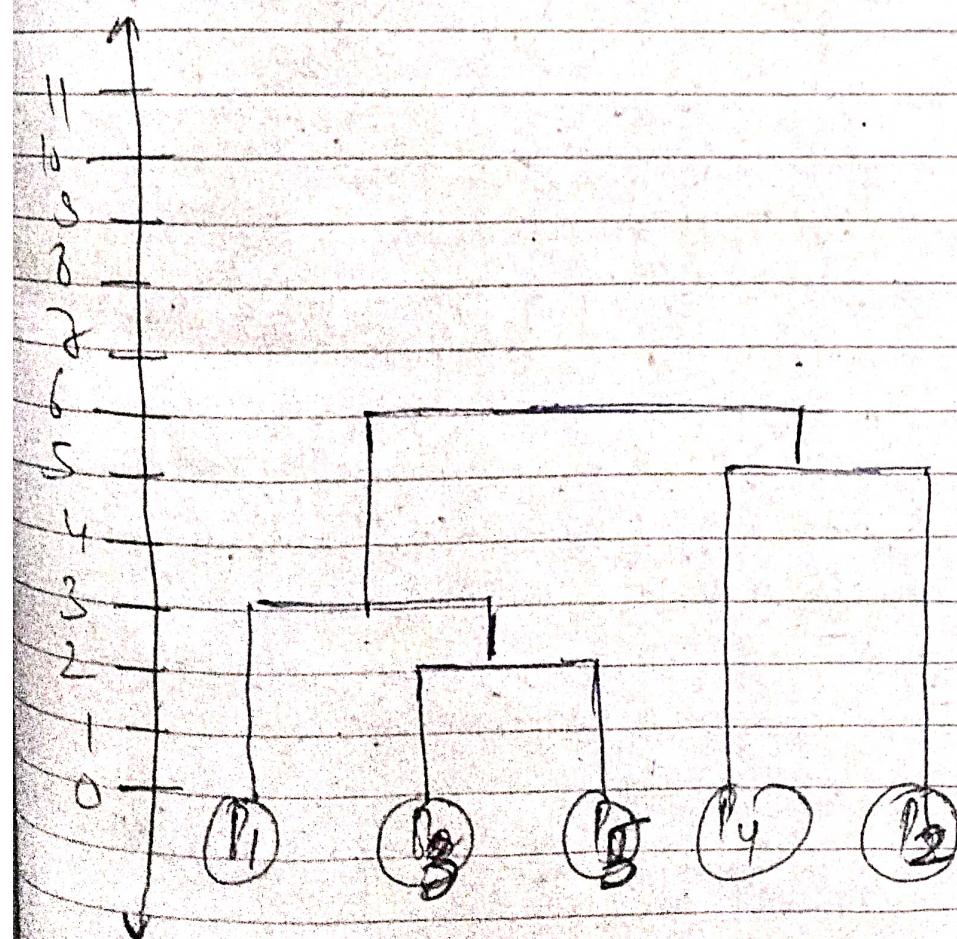
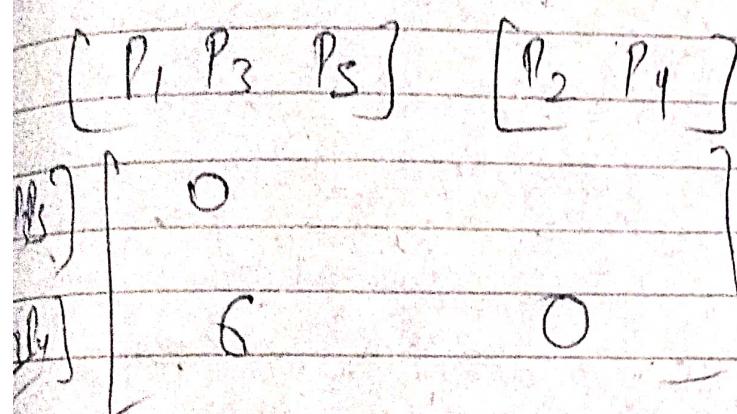
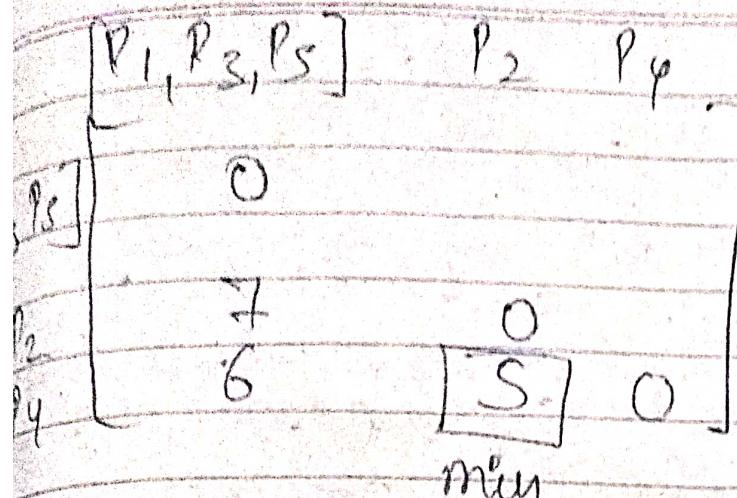
$$\text{dist}(P_2, [P_3, P_5]) = \min (\text{dist}(P_2, P_3), \text{dist}(P_2, P_5)) \\ = \min (7, 10) = 7$$

$$\text{dist}(P_4, [P_3, P_5]) = \min (\text{dist}(P_4, P_3), \text{dist}(P_4, P_5)) \\ = 9, 8 = 8$$

	P_1	P_2	$[P_3 P_5]$	P_4	
P_1	0				
P_2	9	0			
P_3	0	3	7	0	
P_4	6	5	8	0	

min min

$[P_1, P_3 P_5]$ now cluster



Agglomeration where we minimize deviations \rightarrow minimum
Spanning Tree
also

Soft Clustering

On the other hand in soft clustering each data point belongs to a cluster with a certain probability also known as membership value

FCM (fuzzy c-means clustering) algorithm is an example of soft clustering