

2019

①

$$\text{a) } X(0, -1, 0, 1) \quad Y(1, 0, -1, 0)$$

$$= \sqrt{(-0)^2 + 1^2 + 1^2 + 1^2}$$

$$= \underline{\underline{2}}$$

$$\text{b) F-measure} = \frac{2 \times P \times R}{P + R}$$

$$= \frac{2 \times 0.5 \times 0.6}{0.5 + 0.6}$$

$$= \underline{\underline{0.1}}$$

$$= \frac{0.6}{1.1} = \underline{\underline{\frac{6}{11}}}$$

(e) Yes,  $\{abc\}$  will also be infrequent because according to apriori principle, if an itemset is infrequent, all its supersets are infrequent

(d) i) Eliminate the missing values : we may lose some important attribute req. for analysis?

ii) Ignore the missing values - may affect our analysis

iii) Estimate the missing values - using KNN we can estimate the missing values

Q) Precision : The closeness of the measurement with each other is called precision.

Bias : The variation of the measurement with the actual value being measured. It is calculated by taking the diff between ~~& taking~~ the mean of all the values and the actual value being measured.

f) Variable transformation is the data preprocessing technique which is used to standardize the variables within a given range eg 0-1. If we have a mix of smaller values & larger values, we can apply variable transformation and bring them all within a given range. We can deal with outliers as well.

Methods :-

i) Normalization or standardization :

Using the Z-score method .

$$Z = \frac{x - \text{mean}}{\sigma}$$

2) Min-Max normalization

$$x' = \frac{x - \text{min}}{\text{max} - \text{min}} (\text{new\_max} - \text{new\_min}) + \text{new\_min}$$

(g)

$$\text{Support } (X \rightarrow Y) = \frac{4}{5}$$

$$\text{Conf } (X \rightarrow Y) = \frac{3}{4}$$

Yes, we can derive the support of  $Y \rightarrow X$  but not the confidence, as we need the number of occurrences of  $Y$ .

(h) Advantages:-

- 1) Whole dataset can be used for testing as well as training.



Disadvantages:

- 1) In leave one out, one record is used as test record & rest all as training records.
- 2) So, it is very computationally expensive to calculate.
- 3) It has very high variance.



Agglomerative

(i)

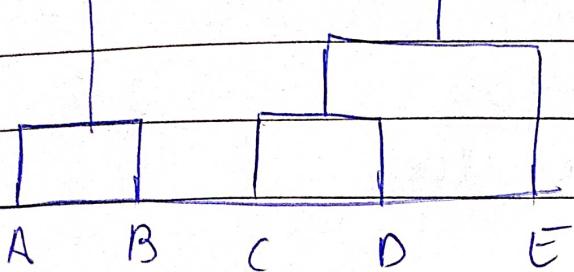
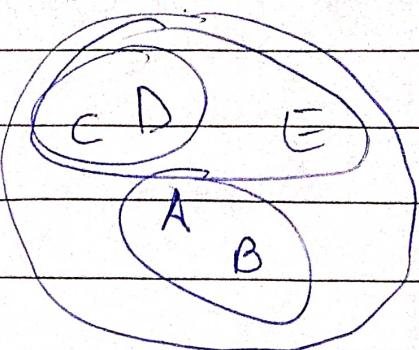
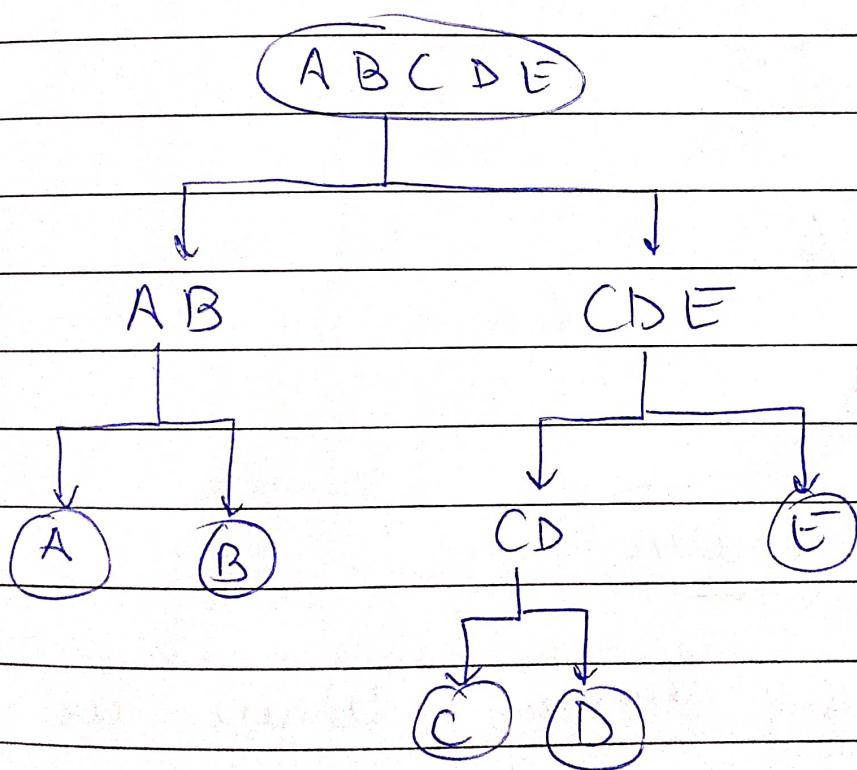
(i) Bottom up approach

Divisive

(ii) Top down approach

(ii) Singleton cluster at first & then successively combine until all clusters are merged into single cluster.

(ii) Splitting the cluster until we get the singleton clusters.

DendogramDivisive

(i))

\*\* (Asymmetric)

$$(k) \text{ Accuracy} = \frac{400+300}{1000} = 0.7$$

$$\text{Sensitivity} = \frac{400}{400+100} = 0.8$$

$$\text{Specificity} = \frac{300}{300+200} = 0.6$$

$$\text{TPR} = 0.8$$

$$\text{FPR} = 0.4$$

## SECTION-B

②

b) \*\* (metric)

c) Unsupervised learning :- Where we don't feed in any data and the algorithm on its own classifies the records into diff groups on the basis of their characteristics. Clustering is unsupervised learning eg - customer segmentation.

③ (a) Overall Gini index

(i)

$4+$        $6-$

$$5 \quad 1 - \left[ \frac{(4)^2}{105} + \frac{(83)^2}{105} \right]$$

$$= 1 - \left[ \frac{4}{25} + \frac{9}{25} \right]$$

$$\boxed{\frac{12}{25}}$$

A<sub>10</sub>

$7T$        $3F$   
 $4+$        $3-$        $3-$

$$T \quad 1 - \left[ \frac{(4)}{7} \right]^2 + \left[ \frac{(3)}{7} \right]^2$$

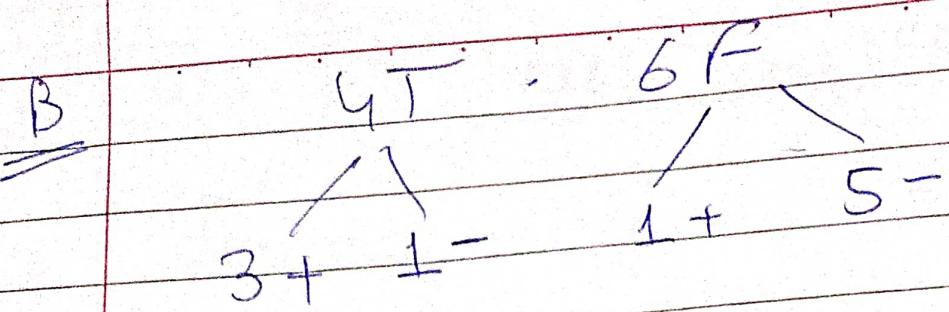
$$= 1 - \left[ \frac{16}{49} + \frac{9}{49} \right] = \boxed{\frac{24}{49}}$$

P  $\rightarrow$  ~~0~~ 0

20

$$\text{Weighted} = \frac{24}{49} \times \frac{12}{35} = \boxed{\frac{12}{35}} = 0.342$$

25



$$5 \quad 1 - \left[ \left( \frac{3}{4} \right)^2 + \left( \frac{1}{4} \right)^2 \right]$$

$$= 1 - \left[ \frac{9}{16} + \frac{1}{16} \right] = \frac{6}{16} = \boxed{\frac{3}{8}}$$

$$10 \quad 1 - \left[ \left( \frac{1}{6} \right)^2 + \left( \frac{5}{6} \right)^2 \right]$$

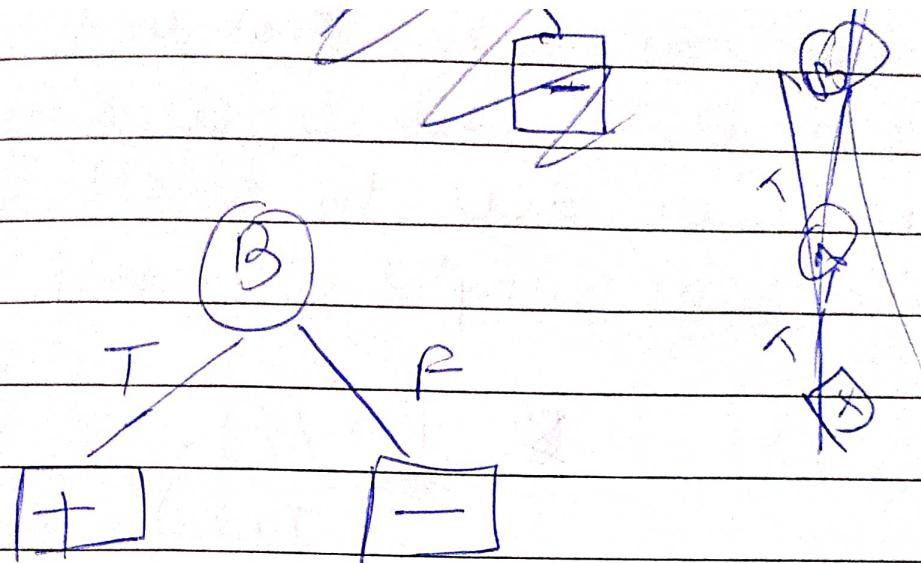
$$= 1 - \left[ \frac{1}{36} + \frac{25}{36} \right] = \boxed{\frac{5}{18}}$$

$$15 \quad \text{Weight of } \frac{3}{28} \times \frac{4}{10} + \frac{8}{318} \times \frac{6}{10} = \frac{3}{20} + \frac{1}{6}$$

$$= \frac{3}{20} + \frac{1}{6}$$

$$= \frac{18 + 20}{120} = \frac{38}{120} = \boxed{\frac{19}{60}} = 0.316$$

(ii) So, the decision tree induction would choose B.



(iv) 2 instances are misclassified by resulting decision tree

(b) KNN is classified as lazy learner because there is no need to build a model in KNN and it just learns all the training set but any record out of training set comes, it is inefficient to apply

(ii) (a) Exhaustive rules: Every record triggers at least one rule out of the ruleset given. If not exhaustive, we would not be able to classify my record into any class. Resolved using default class

(b) Progressive Sampling: Here, we start with a small sample size & progresses into a large sample size till an optimal size is obtained. Here, we don't need to fix the sample size at start and we can avoid the problems like high variance & low accuracy which can occur due to small sample size and large sample size resp.

$$(c) \quad P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)}$$

$$\Rightarrow P(Y|X) = \frac{\sum_i P(X_i|Y) \cdot P(Y)}{P(X)}$$

Assumption used by Naive Bayes classifier is of conditional independency.

$\Rightarrow X$  is said to be conditionally independent of  $Y$  given  $Z$  when,

$$P(X|Y, Z) = P(X|Z)$$

(5) (c) (ii) F-3 itemsets

F-1 itemset

$\{1, 2, 3\}$

$\{1\}$

$\{1, 2, 4\}$

$\{2\}$

$\{1, 2, 5\}$

$\{3\}$

$\{1, 3, 4\}$

$\{4\}$

$\{1, 3, 5\}$

$\{5\}$

$\{2, 3, 4\}$

$\{2, 3, 5\}$

$\{3, 4, 5\}$

Using  $F_{k-1} \times F_1$

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

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

(ii) Using  $F_{k-2} \times F_{k-2}$

Here, first  $k-2$  elements should be same.

=  $\{1, 2, 3, 4\}, \{1, 2, 3, 5\}, \{1, 3, 4, 5\},$

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

(b)

2 bits

Binary form

Awful	00
Poor	01
OK	10
Good	11

4 bits

Binary form

Awful	0000	1000
Poor	0100	
OK	0010	
Good	0001	

⑥

$$\{e\} = 8 \Rightarrow 80\%.$$

$$\{b, d\} = 2 \Rightarrow 20\%.$$

$$\{a, d\} = 4 \Rightarrow 40\%.$$

$$\{b, d, e\} = 2 \Rightarrow 20\%.$$

20

 $\{e\}$  and  $\{a, d\}$  are frequent(ii)  $\{b, d, e\}$ 

25

$\{b, d\} \rightarrow \{e\}$

Conf =  $100\%$

$\{b, e\} \rightarrow \{d\}$

Conf =  $2/4 = 50\%$

$\{d, e\} \rightarrow \{b\}$

$2/5 = 40\%$

$\{e\} \rightarrow \{b, d\}$

$2/8 = 25\%$

$\{d\} \rightarrow \{b, e\}$

$2/6 = 33.3\%$

$\{b\} \rightarrow \{d, e\}$

$2/6 = 33.3\%$

5

Strong rules  $\Rightarrow$

$\{b, d\} \rightarrow \{e\}$

10

b) Nominal - These attributes are used to distinguish the records  
eg - ex. color, customer ID,

15 Ordinal : These attributes are used when we want to rank or order the records  
eg - grades, positions in race

⑦ A 2018 done

20

25