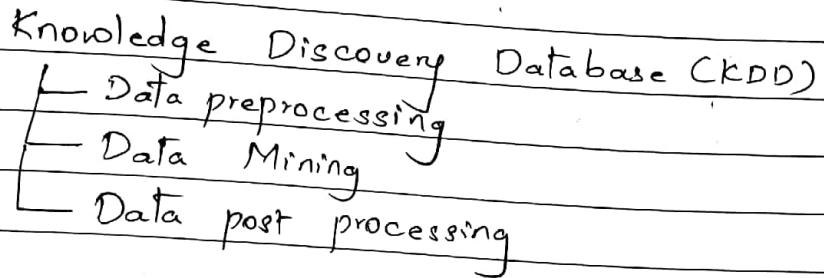


DATA MINING



data cleaning
data transformation
data integration
data selection

- 1) pattern evaluation
- 2) " filtering
- 3) evaluation
- 4) cluster visualization

Traditional Data Analysis

- Scalability
- Out of core algorithm
- Dimensionality
- heterogeneous

KDD - converting raw data into useful info.

Applications of DM

- Business
- Medicine, Science & Engineering

Challenges

- Scalability
- High Dimensionality
- Heterogeneous & Complex Data
- Data ownership & distribution
- Non-traditional analysis

Origins of Data Mining

Data Mining

- Database & DataWarehouse
 - Statistics - Sampling
 - Information Retrieval - find of imp. topics
 - = Parallel Computing
 - ML learning
 - supervised learning → classification
 - unsupervised " → clustering (grouping data)
 - semi-supervised " → requesting ML ip on labels.
 - active
- Explanatory Variable & Class labels.

Data Mining Task

<u>Predictive modelling</u>				<u>Descriptive Modelling</u>
create				↳ association mining → associated values ↳ laptop (because?)
explanatory variables to derive				↳ Clustering (grouping closely related things)
target variable				↳ anomaly detection (significantly different data) (credit card, IP's)
Animal name	Animal type	Habitat	Class	
Elephant	Terrestrial	Forest	Herbivore	

↳ Classification → used for discrete data
 ↳ Regression → continuous data
 (floating pt values)
 (temp, pressure, wind)

Pattern interestingness

- valid on new data → potentially useful
- easily understood
- novel, if it validates a hypothesis

→ useful.

→ interesting patterns represents knowledge

Objective measures of pattern

Objective Interestingness - based on structure of discovered patterns

- ↳ Rule support
- ↳ $P(X \cup Y)$

It is a percentage of transaction that contains $X \cup Y$
It is defined as $P(X \cup Y)$

↳ Confidence : tells association b/w either X or Y

→ It accesses the degree of certainty of a detected association.

It is defined by $P(Y|X)$

Conditional probability

↳ Accuracy

Defines the percentage of transaction correctly classified by the rule. If $X \rightarrow Y$
Note : Here both X and Y should be satisfied

↳ Coverage

Percentage of transaction covered by the rule
Note : Only X should be satisfied.

Subjective interestingness : based on user beliefs in data

↳ unexpected : contradicts his belief/pattern

↳ expected : validates his belief

↳ actionable : gives info. on which you can take action taking the measurement

Data Characterisation / Data Discrimination

Target class:

Summarization of general features which describes the target class

Ex: Customer profiles (Credit rating, occupation)

Data Discrimination: Finding the discriminating features,

- class description
- concept "

Data

- raw facts or observation.

Data Object / (Record) / Vector

Data set: Collection of data object which are categorized by certain characteristics

- Attribute / Feature / Characteristics
- Measurement Scale
- Measurement Process

Types of attributes

↳ based on number properties

↳ Distinctness = \$ ≠

↳ Order <, >

↳ addition +, -

↳ multiplication *, /

Categorical / Qualitative Attributes

distinctness: \hookrightarrow nominal : distinguishes one data from another

order: \hookrightarrow ordinal: " one data object from another
 \rightarrow Product rating.

Numerical / Quantitative Attributes

additive \hookrightarrow Interval attribute: finding difference is useful e.g. temp

add & multip \hookrightarrow Ratio " : e.g. age

\rightarrow Based on no. of values

\hookrightarrow discrete attribute: it has finite no. of values or
 countably in finite no. of values
 e.g. rollno,

\hookrightarrow Continuous attribute: contains floating pt. values
 e.g. stock price, weather

\hookrightarrow Binary attributes :- 1 or 0

\hookrightarrow asymmetric attribute :- non-zero values are imp.
 e.g. medical profiles like diabetes measure

Characteristics of Data Set

\hookrightarrow Dimensionality (attri): Curse of dimensionality

\rightarrow if no. of dimensions \rightarrow 1000, it is difficult to process data.

\hookrightarrow Sparsity \rightarrow feature of dataset where only non-zero values are stored.

\hookrightarrow Resolution \rightarrow value of data object varies with different levels of resolution.

e.g. satellite images.

e.g. data varies according to time.

Types of Data Set

(I) i) Record Data: A ^{fixed} data set with fixed no. of ~~fixed~~ ^{most} fields.

1) Transaction / Market Basket Data

- Each record is Transaction
- Each attribute is item

T ₁	milk, eggs
T ₂	Cheese, milk, flour.
T ₃	Cake, detergent
T ₄	jam, bread, milk

Trans Id	milk	eggs	cheese	flour	Detergent	Cake	Jam	Bread
T ₁	1	1	0	0	0	0	0	0
T ₂	1	0	1	1	0	0	0	0

Record Data

2) Data Matrix: used for numerical data & each pt is plotted as multidimensional ~~set~~ ^{point}

row → data object

column → attributes

Sepal length	Sepal width	Petal length	Petal width
7.5	3.5	2.5	3.5

Sparse Data Matrix

↳ asymmetric attribute

eg: document term matrix : represented by 1 or 0 + term frequency

Doc No	Term 1	Term 2	Term 3
1	3	1	0

→ Here only non-zero values are important.

II

Graph Data

(i) node → data object → we show link b/w diff. data object
 e.g. web pages over another where each object is a node

(ii) Displaying the relationship b/w sub-objects

e.g. representⁿ of chemical compound.

→ Graph → for one data object

↳ each data object → contains sub object
 nodes → sub objects



Ordered Data

attributes → linked to time / space

1] Sequential Data :- should have 1 attri which is related to time
 e.g. login access.

Login Time	Logout Time	Username
10:00 AM	11:00 AM	user1

2] Sequence data :- order of sequence is important.
 e.g. genetic info (plants) (contains seq. of individual entities)

3] Temporal data / time series data :- capturing data according to the time, we are not storing the time

↳ Temporal auto correlation :- rainfall in June is diff. from rainfall in Aug

4] Spatial Data :- data which is represented as latitude & longitude
(area, position)

↳ Spatial - auto correlation :- values which are close in space is similar

Data Cleaning \rightarrow Data Quality

aspects affect data quality \rightarrow Measurement error : it results from measurement process
 ↳ Data Collection : human errors (missing data...)

Noise :- Random variation in the attribute.

↳ attribute noise :- missing, incorrectly entered

↳ class noise :- wrong class label

Artifact :- deterministic distortion of data

Techniques to handle noise :-

- (1) Binning (Data Smoothing)
- (2) Regression
- (3) Outliers.

(1) Binning

* Unsupervised binning :- target class info is not available

→ equal width binning

→ equal depth binning : no. of element frequency same
 or equal frequency binning

* Supervised binning

→ entropy based binning

Equal width binning

→ divides ^{data} into k ^{intervals} bins which will have equal width

Equal Depth binning

- sort the data

eg: 0, 4, 12, 16, 18, 24, 26, 28 ~~30~~ $k=3$.

bin1 0, 4, 12

bin2 16, 18, 18

bin3 24, 26, 28, 28

* Data Smoothing Technique

(1) bin means. → Find mean & replace the values as mean

bin1 5, 5, 5

bin2 17, 17, 17

bin3 26 26 26

(2) bin boundaries - Closest value

bin1 0 0 12

bin2 16 16 18

bin3 24 24 28

(3) bin median - Find median

bin1 4 4 4

bin2 16 16 16

bin3 26 26 26

* Equal Width binning :-

$$W = \frac{\max - \min}{k}$$

Intervals

[min, min + W)

[min + W, min + 2W)

:

[min + (k-1)W, min + kW]

- eg: $0, 4, 12, 16, 16, 18, 24, 26, 28$
 $k=3$

binning is used

→ noise handling

→ discretization

$$\omega = \frac{28-0}{3} = \lceil 9.33 \rceil \approx 10$$

$$[0, \min + \omega]$$

$$[0, 10] \quad [10, 20]$$

$$[10, 20) \quad [12, 16, 16, 18]$$

$$[20, 30] \quad [24, 26, 28]$$

bin means

$$\text{bin 1 } [2 \ 2]$$

$$\text{bin 2 } [16 \ 16 \ 16 \ 16]$$

$$\text{bin 3 } [26 \ 26 \ 26] \cancel{\text{etc}}$$

bin median

$$\text{bin 1 } [2 \ 2]$$

$$\text{bin 2 } [16 \ 16 \ 16 \ 16]$$

$$\text{bin 3 } [26 \ 26 \ 26]$$

bin boundaries

$$\text{bin 1 } [0 \ 4]$$

$$\text{bin 2 } [12 \ 16 \ 18 \ 18]$$

$$\text{bin 3 } [24 \ 24 \ 28]$$

→ Supervised Binning :- When target class info is available

→ Entropy based binning :- measures impurity measure

$$\text{entropy}(c) = -\sum P_{ij} \log_2 P_{ij}$$

$j = \text{no. of class}$
 $i^{\text{th}} \text{ interval}$

$$P_{ij} = \frac{m_{ij}}{m_i}$$

$m_{ij} \rightarrow$ no. of record having class j
in i^{th} interval

$m_i \rightarrow$ no. of record in the i^{th} interval

$$\text{net entropy} = \sum w_i e_i$$

$$w_i = \frac{m_i}{m}$$

$m \rightarrow$ total no. of records.

Gain = Total entropy of data - net entropy (attribute)

Hours Studied	A or Test
4	N
5	Y
8	N
12	Y
15	Y

$$\text{entropy (D)} = - \left[\left(\frac{3}{5} \right) \log_2 \left(\frac{3}{5} \right) + \left(\frac{2}{5} \right) \log_2 \left(\frac{2}{5} \right) \right]$$

$$= 0.97099 \approx 0.971$$

Take adjacent values (4, 5) \rightarrow find average $= (4+5)/2$

$$= 4.5$$

		Y	N
	A or Test	Lower than 4.5	
≤ 4.5	0	1	
> 4.5	3		1

$$\text{entropy } (\leq 4.5) = - \left[\left(\frac{1}{4} \right) \log_2 \left(\frac{1}{4} \right) + \left(\frac{3}{4} \right) \log_2 \left(\frac{3}{4} \right) \right]$$

$$= 0$$

$$\text{entropy } (> 4.5) = - \left[\left(\frac{1}{4} \right) \log_2 \left(\frac{1}{4} \right) + \left(\frac{3}{4} \right) \log_2 \left(\frac{3}{4} \right) \right]$$

$$= 0.811$$

$$\text{net entropy} = \frac{m_i}{m} = \left(\frac{1}{5} \right) \times 0 + \left(\frac{4}{5} \right) \times 0.811$$

$$= 0.6488$$

$$\text{Gain} = \text{entropy}(D) - \text{net entropy}(\text{for } 4.5)$$

$$= 0.971 - 0.6488$$

$$= 0.328$$

		Entropy
Split value (4.5)	Gain = 0.328	0.6488
" value (6.5)	Gain = 0.232	0.9538
" value (10)	Gain = 0.4172	0.5538
" value (13.5)	Gain = 0.191	0.8

Entropy bin 1 ≤ 10 4.5, 8
 bin 2 > 10 12, 15

→ For split value (6.5)

	Y	N	
≤ 6.5	1	1	
> 6.5	2	1	

$$\text{entropy } (< 6.5) = - \left[\left(\frac{1}{2} \right) \log_2 \left(\frac{1}{2} \right) + \left(\frac{1}{2} \right) \log_2 \left(\frac{1}{2} \right) \right]$$

$$= 1$$

$$\text{entropy } (> 6.5) = - \left[\left(\frac{2}{3} \right) \log_2 \left(\frac{2}{3} \right) + \left(\frac{1}{3} \right) \log_2 \left(\frac{1}{3} \right) \right]$$

$$= 0.8899$$

$$\text{net entropy} = \frac{m_i}{m} - \frac{2}{5} \times 1 + \frac{2}{5} \times 0.8899 = 0.93394$$

$$\begin{aligned}\text{Gain} &= \text{entropy}(D) - \text{net entropy}(\text{for } 6.5) \\ &= 0.971 - 0.93394 \\ &= 0.03706\end{aligned}$$

→ For split value (10)

	Y	N	
≤ 10	1	2	
> 10	2	0	

$$\begin{aligned}\text{entropy}(\leq 10) &= - \left[\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3} \right] \\ &= 0.9182\end{aligned}$$

$$\begin{aligned}\text{entropy}(> 10) &= - \left[\frac{2}{2} \log_2 \frac{2}{2} + \frac{0}{2} \log_2 \frac{0}{2} \right] \\ &= 0\end{aligned}$$

$$\begin{aligned}\text{net entropy} &= \frac{m_i}{m} = \frac{3}{5} \times 0.9182 + \frac{2}{5} \times 0 = 0.55092 \\ &\approx 0.551\end{aligned}$$

$$\begin{aligned}\rightarrow \text{Gain} &= \text{entropy}(D) - \text{net entropy}(\text{for } 10) \\ &= 0.971 - 0.9182 \\ &\approx 0.0529\end{aligned}$$

→ For split value (13.5)

	Y	N	
≤ 13.5	2	2	
> 13.5	1	0	

$$\text{entropy}(\epsilon=13.5) = - \left[\frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{1} \log_2 \frac{2}{4} \right]$$

$$= 1$$

$$\text{entropy}(\epsilon > 13.5) = - \left[\frac{1}{1} \log_2 \frac{1}{1} + \frac{0}{1} \log_2 \frac{0}{1} \right]$$

$$= 0$$

$$\text{net entropy} = \frac{m_i}{m} = \frac{4}{5} \times \frac{1}{2} + \frac{1}{5} \times 0$$

$$= 0.8$$

$$\text{Gain} = \text{entropy}(ID) - \text{net entropy}(\text{for } 13.5)$$

$$= 0.971 - 0.8$$

$$= 0.171$$

Noise Reduction

(1) binning

(2) Regression : use to predict the value of $y = mx + c$

y = ^{should} contain noise

x = should not contain noise

Linear Regression

Interpolation

If you are trying to predict the y within known range

Interpolation: points within the curve

Extrapolation: points outside the curve

Std devial

std devial ^{new}

mean

x	y	xy	x^2
1	2	2	1
2	5	10	4
3	3	9	9
4	8	32	16
5	7	35	25
$\sum x = 15$	$\sum y = 25$	$\sum xy = 88$	$\sum x^2 = 55$

$$y = mx + c$$

for 'c' equation, $\sum y = n * c + m * \sum x$

for 'm' equation, $\sum xy = c * \sum x + m \sum x^2$

$$25 = 5c + 15m \quad (1)$$

$$88 = 15c + 55m \quad (2)$$

$$c = 1.1$$

$$m = 1.3$$

$$y = 1.3x + 1.1$$

$$y = \beta_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \dots + \alpha_n x_n$$

(3) Outlier : pt which doesn't fit
region containing density

- precision : closeness of measurement taken to one another
- bias : systematic variation of measurement from known value
- accuracy : closeness of given measurement to the true value.

Measurements = 1.015,

0.990

Known value = 1 kg

1.013

0.986

1.001

Mean = $\frac{9.986}{5} = 1.001$

$$\begin{aligned}\text{Bias} &= \text{mean} - \text{known value} \\ &= 1.001 - 1 \\ &= 0.001\end{aligned}$$

Precision = std deviation

$$= \sqrt{\frac{1}{N} \sum_{i=1}^n (\bar{x}_i - \text{mean})^2}$$

$$\begin{aligned}= \sqrt{\frac{1}{5} \left[(1.015 - 1.001)^2 + (0.990 - 1.001)^2 + (1.013 - 1.001)^2 \right.} \\ \left. + (0.986 - 1.001)^2 + (1.001 - 1.001)^2 \right]} \\ = 0.0117\end{aligned}$$

→ bias: diff. b/w measurement & known value

Outliers

Anomalous attribute

Anomalous data object e.g. if characteristics of one object is different from other data objects

Techniques to handle missing value

→ Missing values:-

→ info was not collected

→ Attribute is not applicable for all records.



Techniques to handle missing values

1) Eliminate data object / attribute

2) Estimate missing values.

→ Find avg

Continuous & estimate

	Jan	Feb	March
31	32	33	

71	81	
81	93	1.3

Categorical In nominal →
most commonly
occurring attrb.

3) Ignore missing value (which performing data mining)

Inconsistent Data (Invalid data)

Eg. -ve height

Duplicate data (Multiple entries for same object)

- Same entries : keep only one.

- different entries : keep most correct entry.

AcctNo	CustName	Address	balance
--------	----------	---------	---------

1003	John	Boston	1000
1003	John	Houston	3000

analyse which is
correct

- Similar objects are not duplicates.

Deduplicating

Removal of duplicate entries from a data set

Issues related to application (Desirable properties)

→ Timeliness (Dataset should be upto date)

eg. Century old rainfall data will produce wrong pattern.

→ Relevance (Dataset should be useful for application)

→ Knowledge about data (Description of the data)

Variable transformation

$$\log_{10}(x) \quad e^x$$

Normalisation Censure data fits within a range
Standardization

* Z-score normalisation (Range of attributes not known)

$$V_i' = V_i - \mu_x$$

$$\sigma_x$$

μ_x → mean of attribute x

σ_x → std deviation of attribute x

$$\mu_x = \frac{V_1 + V_2 + \dots + V_n}{n}$$

V_1, V_2, \dots, V_n - value of attribute

n - no. of attribute 'n'

$$\sigma_x = \sqrt{\frac{1}{n} (V_1 - \mu_x)^2 + (V_2 - \mu_x)^2 + \dots + (V_n - \mu_x)^2}$$

* min-max normalisation (Range of attribute is known)

$$V_i' = \frac{V_i - \text{min}_A}{\text{max}_A - \text{min}_A} (\text{newmax}_A - \text{newmin}_A) + \text{newmin}_A$$

* Decimal Scaling

$$V_i' = \frac{V_i}{10^j}$$

$j \rightarrow$ smallest integer such that $\max(|V_i'|) < 1$

$j = \log_{10}(\max(v)) \leftarrow \dots$ to nearest integer

Q. 24, 29, 31, 32, 32 32, 37, 42, 42, 44 44, 64, 68

$$\text{newmin} = -2$$

$$\text{newmax} = 8$$

$$\mu_x = 40.07$$

$$\begin{aligned}
 (24 - 40.07)^2 &= 258.244 \\
 (29 - 40.07)^2 &= 122.54 \\
 (31 - 40.07)^2 &= 82.26 \\
 (32 - 40.07)^2 \times 3 &= 195.37 \\
 (37 - 40.07)^2 &= 9.42 \\
 (42 - 40.07)^2 &= 4.44 \\
 (44 - 40.07)^2 &= 30.88 \\
 (64 - 40.07)^2 &= 572.64 \\
 (68 - 40.07)^2 &= 780.08 \\
 &\quad \underline{2058}
 \end{aligned}$$

$$\sqrt{x} = \sqrt{\frac{1}{13} \times (2058)} = 12.58$$

$$\sqrt{x} = \sqrt{158.374} = 12.58$$

new value,

$$(24) = \frac{24 - 40.07}{12.58} = -1.2774$$

$$\text{new value (29)} = \frac{29 - 40.07}{12.58} = -0.8799$$

$$\begin{aligned}
 \text{new value (31)} &= \dots = -0.7209 \\
 (32) &= \dots = -0.6414 \\
 (37) &= \dots = -0.2440 \\
 (42) &= \dots = 0.1534 \\
 (44) &= \dots = 0.3124 \\
 (64) &= \dots = 1.9022 \\
 (68) &= \dots = 2.2201
 \end{aligned}$$

Using decimal scaling

$$\max(v) = 68$$

$$j = \log_{10}(68) = 1.83 \approx 2$$

$$\text{new value (24)} = \frac{24}{10^3} - \frac{24}{10^2} = 0.24$$

$$\begin{aligned}
 (81) &= 0.31 \\
 (82) &= 0.32 \\
 (87) &= 0.37 \\
 (42) &= 0.42 \\
 (44) &= 0.44 \\
 (64) &= 0.64 \\
 (68) &= 0.68
 \end{aligned}$$

Using min-max

~~$$\text{new value}(24) = \frac{24 - 24}{(68 - 24)} (8 - (-2)) + (-2) = -2$$~~

~~$$\text{new value}(29) = \frac{29 - 24}{(68 - 24)} (8 - (-2)) + (-2) = -0.8636$$~~

~~$$\text{new value}(31) = -0.4090$$~~

min-max normalization

$$\text{old min} = 24 \quad \text{new min} = -2$$

$$\text{old max} = 68 \quad \text{new max} = 8$$

$$\text{new value} = \frac{\text{old value} - \text{old min}}{\text{old max} - \text{old min}} (\text{new max} - \text{new min}) + \text{new min}$$

$$\text{new value}(24) = \frac{24 - 24}{68 - 24} (8 - (-2)) + (-2) = -2$$

$$\text{new value}(29) = \frac{29 - 24}{68 - 24} (8 - (-2)) + (-2) = -0.86$$

$$\text{new value}(31) = \frac{31 - 24}{68 - 24} (8 - (-2)) + (-2) = -0.4090$$

$$\text{new value}(32) = \frac{32 - 24}{68 - 24} (8 - (-2)) + (-2) = -0.18$$

$$\text{new value}(37) = \frac{37 - 24}{68 - 24} (8 - (-2)) + (-2) = 0.95$$

$$\text{new value (42)} = \frac{42-24}{44} (10)-2 = 2.09$$

$$\text{new value (44)} = \frac{44-24}{44} (10)-2 = 2.5$$

[2, 8]

$$\text{new value (64)} = \frac{64-24}{44} (10)-2 = 7.09$$

$$\text{new value (68)} = \frac{68-24}{44} (10)-2 = 8$$

3

* Data Preprocessing

1] Aggregation : Combining one or more data together.

→ Reduce the no. of attribute values.

Adv

→ less
memory

	Storename	Item bought	Purchasedate	Item Amt	Categorical: Store location Item bought
	Chicago	Watch	7/07/19	800	
	Chicago	Wallet	17/7/19	900	
	Chicago	Penset	20/7/19	200	
	Boston	Vase	25/7/19	1200	
	Boston	Cutlery set	27/7/19	2500	

Store location	Item	Purchasedate	Item Amt
Chicago	Giftset	7 - 20/7/19	1900
Boston	Household Items	25 - 27/7/19	3700

→ Advantage

- Less time, Less memory
- offers a high level view of data.

→ Disadvantages

- Some data analyses → require low level view of data.

2] Sampling

(picking up subject from the given set)

Nyquist theor

or Nyquist

Sample representativeness

Property of original set & sample set should be more or less same.

↳ Simple random sampling without replacement (SR SWOR)

↳ " " " with " (SR SWOR)

↳ Cluster "

↳ Stratified "

↳ Progressive (Adaptive) " :- Each iteration, addition done → checks

accuracy
use when do not know actual sample size

not removing → →
→ SR SWOR :- chances of picking any set is same

removing & → do not replace with original data set.

→ Cluster Sampling :- Use clusters to pickup a sample

e.g. Given 11 Clusters

↳ choose 3 clusters SCM

↳ Stratified : correct grouping the data

ex:

Youth	Middle Age	Senior Citizens	
4	8	2	

Sample size = 7

$$- \text{ Youth} = \frac{4}{14} \times 7 = 2$$

$$\text{no. of samples to be chosen from each group} = \left(\frac{\text{no. of objects in the group}}{N} \right) \times \text{sample size}$$

$$\text{Middle age} = \frac{8}{14} \times 7 = 4$$

(no. of samples)

$$\text{Senior Citizen} = \frac{2}{14} \times 7 = 1$$

Sample :

Youth
Youth
Middle Age
"
"
"
Senior Citizen

Ex:	Branch	Comp	Civil	Production	Automobile	IT	ETC
No	52	38	36	46	44	64	

Sample Size = 60

N = 280

$$\begin{aligned}
 - \text{ Comp} &= \frac{52}{280} \times 60 = 11.1 \approx 11 & - \text{ Automobile} &= \frac{46}{280} \times 60 = 9.85 \approx 10 \\
 - \text{ Civil} &= \frac{38}{280} \times 60 = 8.14 \approx 8 & - \text{ IT} &= \frac{44}{280} \times 60 = 9.42 \approx 9 \\
 - \text{ Production} &= \frac{36}{280} \times 60 = 7.71 \approx 8 & - \text{ ETC} &= \frac{64}{280} \times 60 = 13.71 \approx 14
 \end{aligned}$$

3) Dimensionality Reduction

Linear Algebra techniques : transforming old attr. to new attr.

↳ Principal Component Analysis

10 attr. to 3 attr.

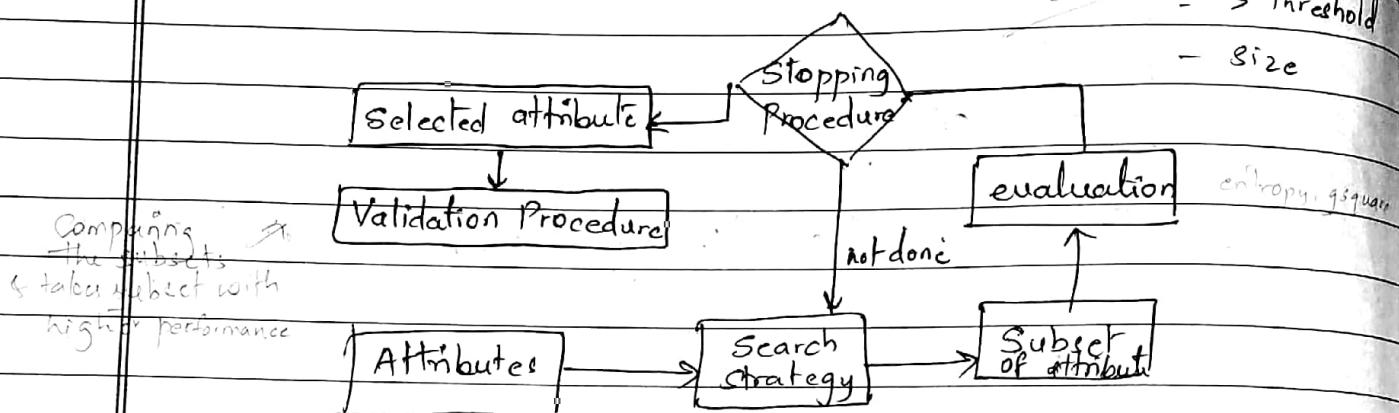
↳ Singular Value Decomposition : $A A^T$

4) Attribute subset Selection

↳ Irrelevant features

- ↳ redundant features : age & birthdate
eliminate one
- embedded approach : decides which attr. should be kept & which should be discarded.
- filter " :- imp., not imp
- Wrapper " :- (black-box) decides which subset should be used & which should be discarded

Architecture for Feature Subset Selection



			Principal Component Analysis
X	Y	Z	
7	8	13	↙
8	9	14	
11	13	18	
13	10	8	
Xmean	Ymean	Zmean	

$$\text{Cov} = \frac{\sum (x - \bar{x})(y - \bar{y})}{n}$$

$$\text{Cov}(x,y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{n}$$

$$\text{Covariance} = \begin{bmatrix} \text{cov}(x,x) & \text{cov}(x,y) & \text{cov}(x,z) \\ \text{cov}(y,x) & \text{cov}(y,y) & \text{cov}(y,z) \\ \text{cov}(z,x) & \text{cov}(z,y) & \text{cov}(z,z) \end{bmatrix} \rightarrow A$$

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

↓

quadratic equation

$$\lambda_1, \lambda_2, \lambda_3 = 0$$

$$\text{eg. } \lambda_2 > \lambda_1 > \lambda_3$$

$n \rightarrow k$
A dimension to k dimension
Where $k < n$

- Stepwise forward selection :- (empty set initially) best attr is selected & added to the subset (method: gain...)
- Stepwise backward selection: (Set has all elements)
Remove the worst attr from the set.
- Combination of both: Putting best & removing worst
- decision tree:
 - Internal nodes
 - leaf nodes.

* Feature Weighting

- 3) *
- * Feature Creation: understanding the data
 - ↳ feature extraction: analyse the data set creating new attributes
 - ↳ mapping data to a new space transferring data (changing the range)
 - ↳ feature construction: e.g. given birthdate we can find age new attr

Identify sub-structure & create this substructure as attribute

- 6) *
- Discretisation: Converting continuous values into discrete
 - (doesn't use class info)
 - ↳ Unsupervised discretization:
 - Equal width
 - Equal depth / frequency

range
quant

eg X	10 ≤ x < 20	low
	"	"
13	"	"
19	"	"
21	20 ≤ x < 30	medium
22	"	"

- ↳ Supervised discretization: use class info

eg: of Unsupervised

bin1	11, 13, 19, 21	low
bin2	22, 26, 33, 37	high

Discrete $10 \leq x \leq 21$

$22 \leq x \leq 37$

* Supervised discretization: Entropy based binning)

- find split values
- Grain

*) Binarisation

- If it is continuous \rightarrow it has to be converted to discrete and then discrete to binary.

Categorical: Excellent $\rightarrow 5 \rightarrow 101$
values

Very good $\rightarrow 4 \rightarrow 100$

Good $\rightarrow 3 \rightarrow 011$

Fair $\rightarrow 2 \rightarrow 010$

bad $\rightarrow 1 \rightarrow 001$

Product Name	Excellent	Verygood	good	Fair	bad
prod 1	1	0	0	0	0
prod 2	0	0	0	0	1

asymmetric: 0 or 1

Drawback:

- Wastage of space
- Dimensionality \uparrow goes

10/08/19

Dissimilarities between Data Objects

Minkowski distance

$$d(x, y) = \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{1/r} \quad \text{--- (1)}$$

If $r = 1$,

$$d(x, y) = \sum_{k=1}^n |x_k - y_k| \rightarrow \begin{array}{l} \text{CityBlock / Manhattan distance} \\ / L_1 \text{ Metric} \end{array}$$

different data objects

If $r = 2$

$$d(x, y) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2} \rightarrow \begin{array}{l} \text{Euclidean distance} \\ / L_2 \text{ norm} \end{array}$$

If $r = \infty$,

$$d(x, y) = \lim_{r \rightarrow \infty} \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{1/r} \rightarrow \begin{array}{l} \text{Supremum} \\ \text{distance /} \\ L_\infty \text{ or } L_\infty \text{ metric} \end{array}$$

* Euclidean distance properties

positivity $d(x, y) \geq 0$
 for (x, y)

Symmetry - $d(x, y) = d(y, x)$

Triangle inequality - $d(x, z) \leq d(x, y) + d(y, z)$

Similarities between data object (for binary data)

$$\text{Simple Matching Coefficient (SMC)} = \frac{f_{00} + f_{11}}{f_{00} + f_{11} + f_{01} + f_{10}}$$

Note: Jaccard Coefficient. only ~~suitable~~ work for ^{suitable} ~~binary~~ assymetric ^ attributes.

$$\text{Jaccard Coefficient} = \frac{f_{11}}{f_{11} + f_{01} + f_{10}}$$

Cosine Similarity (Primary used for document data)

$$\cos(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

$$\begin{aligned} x \cdot y &= \sum_{k=1}^n x_k \cdot y_k \\ &= x_1 y_1 + x_2 y_2 + \dots + x_n y_n \end{aligned}$$

$$\|x\| = \sqrt{\sum_{k=1}^n x_k^2}$$

$$\|y\| = \sqrt{\sum_{k=1}^n y_k^2}$$

Q. $x = (0, -1, 0, 2)$ $y = (2, 0, -1, 0)$

Find the distance using Euclidean distance formula.

⇒

$$d(x, y) = \sqrt{(0-2)^2 + (-1-0)^2 + (0+1)^2 + (2-0)^2}$$

$$d(x, y) = \sqrt{4+1+1+4}$$

$$d(x, y) = \sqrt{10}$$

$$\text{City block} = |0-2| + |-1-0| + |0+1| + |2-0| = 2+1+1+2 = 6$$

Cosine Similarity

$$\cos(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

$$\|x\| = \sqrt{0^2 + (-1)^2 + 0^2 + 2^2} = \sqrt{1+4} = \sqrt{5}$$

$$\|y\| = \sqrt{2^2 + 0^2 + (-1)^2 + 0^2} = \sqrt{4+1} = \sqrt{5}$$

$$x \cdot y = (0 \times 2) + (-1) \times 0 + 0 \times (-1) + 2 \times 0$$

$$x \cdot y = 0$$

$$\cos(x, y) = \frac{0}{\sqrt{5} \cdot \sqrt{5}} \\ = 0$$

Q. $x = (1, 1, 0, 1, 1, 0, 0, 1)$

$y = (0, 0, 0, 1, 1, 1, 0, 1)$

$f_{00} = 2$

$f_{11} = 3$

$f_{01} = 1$

$f_{10} = 2$

$$SMC = \frac{f_{00} + f_{11}}{f_{00} + f_{11} + f_{01} + f_{10}}$$

$$SMC = \frac{2+3}{2+3+1+2} - \frac{5}{8} = 0.6$$

$$Jaccard = \frac{3}{3+1+2} = \frac{3}{6} = \frac{1}{2} = 0.5$$

Q.4) Given the foll. data set perform min-max normalization
(newmin = 0 & newmax = 1)

① Z-score

② Decimal Scaling

$$\rightarrow 200, 300, 400, 600, 1000$$

a)

$$\text{New Value} = \frac{\text{old value} - \text{oldmin}}{\text{oldmax} - \text{oldmin}} (\text{newmax} - \text{newmin}) + \text{newmin}$$

$$\text{oldmin} = 200, \quad \text{oldmax} = 1000$$

$$\text{New value}(200) = \frac{200 - 200}{1000 - 200} (1 - 0) + 0 = 0$$

$$\begin{aligned}\text{New value}(300) &= \frac{300 - 200}{1000 - 200} (1 - 0) + 0 \\ &= 0.125\end{aligned}$$

$$\begin{aligned}\text{New value}(400) &= \frac{400 - 200}{1000 - 200} (1 - 0) + 0 \\ &= 0.25\end{aligned}$$

$$nv(500) = 0.5$$

$$nv(1000) = 1$$

b) Z-score normalization

$$v' = \frac{v - \mu}{\sigma}$$

$$\text{Mean}(\mu) = 500$$

$$\text{Standard deviation}(\sigma) = 282.84$$

$$SD = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$= \sqrt{\frac{1}{5} [(200 - 500)^2 + (300 - 500)^2 + (400 - 500)^2 + (600 - 500)^2 + (1000 - 500)^2]}$$

$$= 282.84$$

$$v' = \frac{200 - 500}{282.84} = -1.0607$$

decimal scaling

$$v' = \frac{v}{10^j}$$

$$j = \log_{10}(\max(v)) = \log_{10}(1000) = 3$$

$$v' = \frac{v}{10^3} = \frac{200}{10^3} = 0.2$$

$$v'(300) = \frac{300}{10^3} = 0.3$$

Q. $5, 10, 11, 13, 35, 50, 55, 72, 92, 204, 215$

$k = 3$.

Perform equal width & equal depth & apply data smoothing techniques for the foll.

\Rightarrow E10B

$$w = \frac{\max - \min}{k} = \frac{215 - 5}{3} = 70$$

$[5, 75)$ $5, 10, 11, 13, 35, 50, 55, 72$

$[75, 145)$ $92,$

$[145, 215]$ $204, 215$

bin mean

bin1 $31.37, 31.37, 31.37, 31.37, 31.37, 31.37, 31.37, 31.37$

bin2 92

bin3 $209.5 \quad 209.5$

bin median

bin1 24, 24, 24, 24, 24, 24, 24, 24

bin2 92

bin3 209.5 209.5

bin boundaries

bin1 5, 5, 5, 5, 5, 72, 72, 72

bin2 92

bin3 204 215

EDB

Bin1 5, 10, 11, 13

Bin2 35, 50, 55, 72

Bin3 92, 204, 215

bin mean

bin1 9.95, 9.95, 9.95, 9.95

bin2 53, 53, 53, 53

bin3 170.3, 170.3, 170.3

bin median

bin1 10.5, 10.5, 10.5, 10.5

bin2 52.5, 52.5, 52.5, 52.5

bin3 170.3, 170.3, 170.3

bin boundary

bin1 5, 13, 13, 13

bin2 35, 35, 72, 72

bin3 92, 215, 204

UNIT - I**PCA**

hours (h)	marks (m)
2	1
9	4
5	0
7	6
9	2

$$h_{\text{mean}} = 5.2$$

$$m_{\text{mean}} = 2.6$$

$h - h_{\text{mean}}$	$m - m_{\text{mean}}$	$(h - h_{\text{mean}})(m - m_{\text{mean}})$
-3.2	-1.6	5.12
-2.2	1.4	-3.08
-0.2	-2.6	0.52
1.8	3.4	6.12
3.8	-0.6	-2.28
		sum = 6.4

↑
data adjust

$$\text{covariance } (h, m) = \frac{6.4}{(5-1)} = 1.6$$

$$\text{Var}(h) = 8.2 = \frac{1}{n-1} \sum (h - h_{\text{mean}})^2$$

$$\text{Var}(m) = 5.8$$

Construct covariance matrix

$$\begin{bmatrix} \text{cov}(h,h) & \text{cov}(h,m) \\ \text{cov}(m,h) & \text{cov}(m,m) \end{bmatrix}$$

$$A = \begin{bmatrix} 8.2 & 1.6 \\ 1.6 & 5.8 \end{bmatrix}$$

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

$$\begin{vmatrix} 8.2 & 1.6 \\ 1.6 & 5.8 \end{vmatrix} \begin{vmatrix} \lambda & 0 \\ 0 & \lambda \end{vmatrix} = 0$$

$$\begin{vmatrix} 8.2 - \lambda & 1.6 \\ 1.6 & 5.8 - \lambda \end{vmatrix}$$

$$\lambda^2 - 14\lambda + 45 = 0 \rightarrow \text{Use quadratic formula}$$

$$\lambda = 5, \lambda = 9$$

$$\text{eigen values} = [5, 9]$$

$$= [9, 5]$$

Solve for eigen vectors

$$\begin{bmatrix} 0.8944 \\ 0.4472 \end{bmatrix} \quad \begin{bmatrix} 0.4472 \\ 0.8944 \end{bmatrix}$$

$$\begin{bmatrix} 0.8944 \\ 0.4472 \end{bmatrix} \quad \frac{5}{2} \times \frac{2}{1} = \frac{5}{1} \text{ matrix}$$

X data adjust

\downarrow
1 attribute is removed.

* Person's Correlation

$$S(x,y) = \frac{\text{covariance}(x,y)}{\text{standard-deviation}(x) * \text{standard-deviation}(y)}$$

n - entire population

$n-1$ - assume sample data set (doesn't represent entire population)

$$\text{Covariance}(x,y) = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})$$

$$\text{Standard-deviation}(x) = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^2}$$

$$\text{Standard-deviation}(y) = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (y_k - \bar{y})^2}$$

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k$$

$$\bar{y} = \frac{1}{n} \sum_{k=1}^n y_k$$

$$1. \quad x = (0, -1, 0, 2)$$

$$y = (2, 0, -1, 0)$$

$$\begin{aligned} \text{Covariance}(x,y) &= \frac{1}{4-1} \left[(0-0.25)(2-0.25) + (-1-0.25)(0-0.25) \right. \\ &\quad \left. + (0-0.25)(-1-0.25) + (2-0.25)(0-0.25) \right] \\ &= -0.0833 \end{aligned}$$

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k$$

$$= \frac{1}{4} \times [0 + (-1) + 0 + 2] = \frac{1}{4} = 0.25$$

$$\bar{y} = \frac{1}{n} \sum_{k=1}^n y_k$$

$$= \frac{1}{4} \times [2 + 0 + (-1) + 0] = \frac{1}{4} = 0.25$$

$$S(x,y) = \text{Covariance}(x,y)$$

standard-deviation(x) * standard-deviation(y)

$$\begin{aligned} \text{standard-deviation}(x) &= \sqrt{\frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^2} \\ &= \sqrt{\frac{1}{4-1} [(0-0.25)^2 + (-1-0.25)^2 + (0-0.25)^2 + (2-0.25)^2]} \\ &= 1.258 \end{aligned}$$

$$\text{Standard-deviation}(y) = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (y_k - \bar{y})^2}$$

$$= \sqrt{\frac{1}{4-1} [(2-0.25)^2 + (0-0.25)^2 + (-1-0.25)^2 + (0-0.25)^2]} \\ = 1.258$$

$$S(x,y) = \frac{-0.0833}{1.258 \times 1.258}$$

$$= -0.0526 //$$

$$x = (1, 1, 0, 1, 1, 0, 0, 1)$$

$$y = (0, 0, 0, 1, 1, 0, 1)$$

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k$$

$$= \frac{1}{8} \times [1+1+0+1+1+0+0+1]$$

$$\bar{x} = 0.625$$

$$\bar{y} = \frac{1}{n} \sum_{k=1}^n y_k$$

$$= \frac{1}{8} \times [0+0+0+1+1+1+0+1]$$

$$\bar{y} = 0.5$$

$$\text{Covariance } (x, y) = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})$$

$$= \frac{1}{8-1} \times [(1-0.625) \times (1-0.5) + (1-0.625) \times (0-0.5) +$$

$$+ (1-0.625) \times (1-0.5) + (1-0.625) \times (0-0.5) + (0-0.625) \times (1-0.5) +$$

$$+ (0-0.625) \times (0-0.5)]$$

$$= 0$$

$$= \frac{1}{8-1} \times [(1-0.625) \times (0-0.5) + (1-0.625) \times (0-0.5) +$$

$$+ (0-0.625) \times (0-0.5) + (1-0.625) \times (1-0.5) +$$

$$(1-0.625) \times (1-0.5) + (0-0.625) \times (1-0.5) +$$

$$(0-0.625) \times (0-0.5) + (1-0.625) \times (1-0.5)]$$

$$= \frac{1}{7} \times [-0.1875 + -0.1875 + 0.3125 + 0.1875 +$$

$$+ 0.1875 + -0.3125 + 0.3125 + 0.1875]$$

$$= 0.07142$$

$$\begin{aligned} \text{Standard deviation } (x) &= \sqrt{\frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^2} \\ &= \sqrt{\frac{1}{7} \times [(1-0.625)^2 + (1-0.625)^2 + (0-0.625)^2 + \\ &\quad + (1-0.625)^2 + (1-0.625)^2 + (0-0.625)^2 + \\ &\quad (0.625)^2 + (1-0.625)^2]} \\ &\approx 0.5219 \end{aligned}$$

$$\begin{aligned} \text{Standard deviation } (y) &= \sqrt{\frac{1}{n-1} \sum_{k=1}^n (y_k - \bar{y})^2} \\ &= \sqrt{\frac{1}{7} \times [(0-0.5)^2 + (0-0.5)^2 + (0-0.5)^2 + \\ &\quad (1-0.5)^2 + (1-0.5)^2 + (1-0.5)^2 + (0-0.5)^2 + \\ &\quad + (1-0.5)^2]} \\ &= 0.5345 \end{aligned}$$

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

Standard deviation (x) * Standard deviation (y)

$$= 0.07142$$

$$0.5219 \times 0.5345$$

$$= 0.2560$$

26/08/19 * Extended Jaccards Coefficient / (Tanimoto Coefficient)

$$EJC(x,y) = \frac{x \cdot y}{\|x\|^2 + \|y\|^2 - x \cdot y}$$

* Bregman Divergences : Used to measure loss or distortion

$$D(x,y) = \phi(x) + \phi(y) - (\nabla \phi(y), (x-y))$$

↳ Vector difference between (x,y)

* Issues with similarity

- 1) → attribute with different scales (ranges) without correlation
↳ normalise data
- attribute with different scale with correlation
↳ Mahalanobis distance → captures correlation
Mahalanobis

Mahalanobis distance

$$\text{Mahalanobis}(x,y) = (x-y) \sum^{-1} (x-y)^T$$

↳ Inverse of Covariance matrix $\Sigma [x,y]$

- 2) → heterogeneous attributes : attrib. having diff. types of values
 - find similarity b/w each attrib.
 - check if it is asymmetric attrib
(if it is assy. & both places have $(0,0)$, should not be consider)

Similarities of heterogeneous Objects

- 1] for k^{th} attribute, complete similarity $s_k(x,y)$ in range $[0,1]$
- 2] define an indicator variable s_k

$s_{ik} = 0 \rightarrow$ if k^{th} attrib → asymmetric and both obj's have '0' value

or

one obj → have missing value

$= 1$ otherwise

w_k weight

3) Compute overall similarities

$$\text{Similarity}(x, y) = \frac{\sum_{k=1}^n w_k s_k(x, y)}{\sum_{k=1}^n w_k}$$

→ Using weights:

$$s(x, y) = \frac{\sum_{k=1}^n w_k s_k(x, y)}{\sum_{k=1}^n w_k}$$

or

Using minkowski

$$d(x, y) = \left(\sum_{k=1}^n w_k |x_k - y_k|^r \right)^{1/r}$$

* Summary statistics : measures to find diff features in dataset
data exploration

Frequency & Mode

	Temp
Day1	low
Day2	high
Day3	medium
Day4	low
5	medium

Frequency = no. of objects with attrib value 'k'

freq(low) = $\frac{2}{5} = 0.4$

freq(medium) = $\frac{2}{5} = 0.4$

freq(high) = $\frac{1}{5} = 0.2$

mode = 0.4

Mode is the value having the highest frequency

Percentile

Given an ordinal or continuous attribute x , and a no. p between 0 and 100, the p th percentile x_p is a value of x such that $p\%$ of the observed values of x are less than x .

	0.1 x total no. of values	average	no. of values	exp
1	10%	$10/100 = 0.1$	$= 0.1 \times 10 = 1$	$x_p = 1.5$
2	20%	$20/100 = 0.2$	$= 0.2 \times 20 = 2$	$x_p = 2.5$
3	30%	0.3	$= 0.3 \times 30 = 3$	3.5
4	40%	0.4	4	4.5
5	50%	0.5	5	5.5
6	60%	0.6	6	6.5
7	70%	0.7	7	7.5
8	80%	0.8	8	8.5
9	90%	0.9	9	9.5
10	100%	1.0	10	10.5

Srno	x	Srno	x
1	43	16	85
2	54	17	87
3	56	18	88
4	61	19	89
5	62	20	93
6	66	21	95
7	68	22	96
8	69	23	98
9	69	24	99
10	70	25	99
11	71		
12	72		
13	77		
14	78		
15	79		

$$10\% = 0.1 \quad 0.1 \times 25$$

$$= 2.5$$

$$\approx 3$$

$$\text{average} = \frac{56+61}{2} \quad x_p = 58.5$$

3 values < $x_p(58.5)$

$3/25 = 12\%$ values.

$$80\% = 0.8 \quad 0.8 \times 25$$

$$= 20$$

$$\text{average} = \frac{93+95}{2} \quad x_p = 94$$

20 values < $x_p(94)$
 $20/25 = 80\%$

$$60\% = 0.6 \quad 0.6 \times 25$$

$$= 15$$

$$\text{average} = \frac{79+85}{2} \quad x_p = 82$$

$$15 \text{ values} < x_p(82) \\ 15/25 = 60\%$$

Measures of locations

↳ mean

↳ median

$$\text{mean} = \frac{1}{k} \sum_{i=1}^k x_i$$

$$\begin{aligned} \text{median} &= x_{(r+1)} && - \text{odd} \\ &= \frac{x_{(r+1)} + x_{(r)}}{2} && - \text{even} \end{aligned}$$

trimmed mean

$$[1, 2, 3, 4, 5, 90]$$

$$p = 40\%$$

$$(p)1\% = 20\%$$

$$20\% \text{ of } 6 \text{ values} = 1.2 \quad x_1$$

$$\underline{\underline{0.20 \times 6}} = 1.2$$

trimmed data set $2, 3, 4, 5$

trimmed mean = 3.5

Measures of spread range and variance

$$\text{range} = \max(x) - \min(x)$$

$$\text{Variance } [Sx^2] = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})^2$$

mean absolute deviation
 $= \text{median}(|x_1 - \bar{x}|, |x_2 - \bar{x}|, \dots, |x_n - \bar{x}|)$

absolute average deviation $= \frac{1}{m} \sum_{i=1}^m |x_i - \bar{x}|$

interquartile range $= x_{75\%} - x_{25\%}$

$$\begin{array}{cccc} 75\% = 0.75 & 0.75 \times 25 & 25\% = 0.25 & 0.25 \times 25 \\ & 18.75 & & = 6.25 \\ & \approx 19 & & \approx 6 \\ \text{average} = \frac{89+93}{2} & x_p = 91 & \text{average} = \frac{66+68}{2} & x_p = 67 \end{array}$$

$$\begin{array}{ll} 19 \text{ values} < x_p(91) & 8 \text{ values} < x_p(67) \\ \therefore 19/25 = 76\% & 8/25 = 32\% \\ \therefore 91 - 67 = 24\% \end{array}$$

Multivariate Summary Statistics

$\bar{x} = (\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_n)$

$$S_{ij} = \text{covariance}(x_i, x_j) \\ = \frac{1}{m-1} \sum_{k=1}^m (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j)$$

$\therefore x_{ki} \rightarrow$ value of x_i attrib for k^{th} object
 $x_{kj} \rightarrow$ " " x_j " " .. "

$$x_{ij} = \text{correlation}(x_i, x_j) \\ = \text{covariance}(x_i, x_j) / S_i S_j$$

classmate
Date _____
Page _____

OLAP \rightarrow extracts the data

classmate
Date _____
Page _____

$S_i, S_j \rightarrow$ variance of x_i, x_j \rightarrow standard deviation

Skewness

OLAP & Data Warehousing

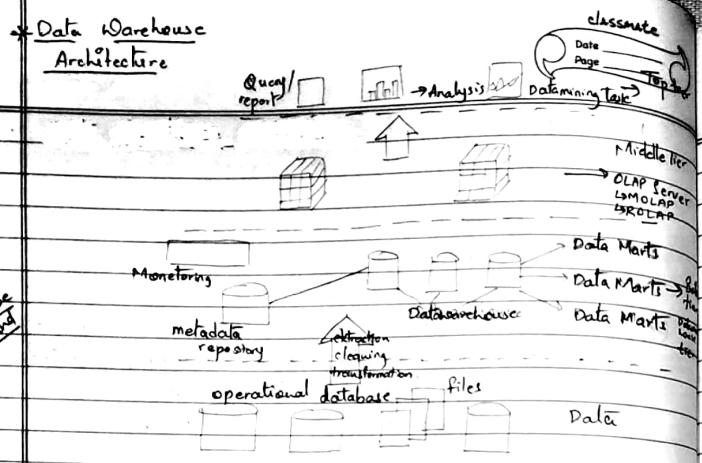
(*) Warehouse \rightarrow designed to store large amount of data
 \rightarrow A data warehouse is subject oriented, integrated, time variant and non-volatile collection of data in support of the management decision making process.

Integrated \rightarrow Integrator / wrappers \rightarrow datawarehouse has its own update driver: data is first integrated, transformed, given to datawarehouse

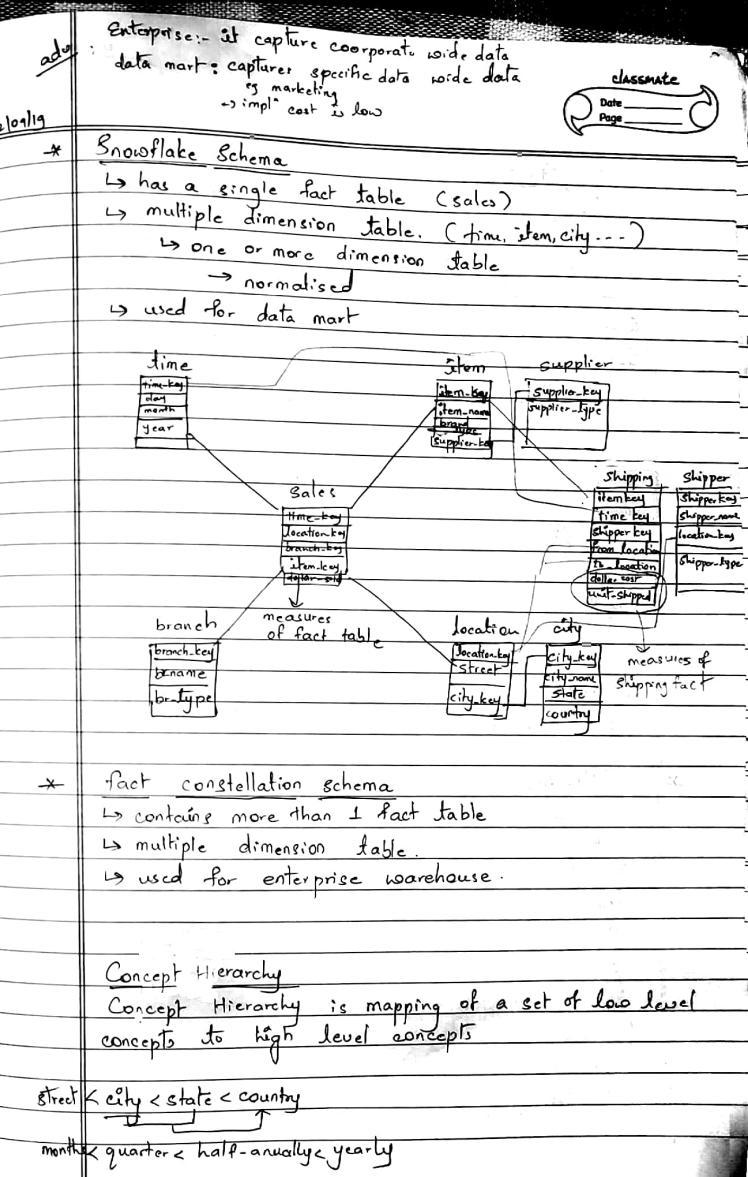
\rightarrow Data is transformed & then stored into datawarehouse schema (semantic data store)

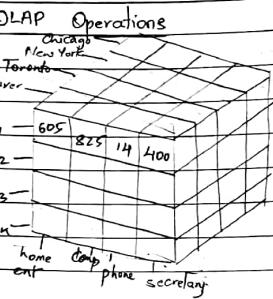
- customer profile
- repositioning of product
- sources of profit

	OLAP	vs	OLTP
Users &	\rightarrow market oriented		\rightarrow customer oriented
System Orientation			
Data Contents	\rightarrow Data from historic perspective		\rightarrow Data from current database transaction
Database Design	\rightarrow specialized datawarehouse schema		\rightarrow User ER based design
View	\rightarrow multiple views		\rightarrow current detail/transactional
Access pattern	\rightarrow complex read-only query		\rightarrow consist of short atomic transaction

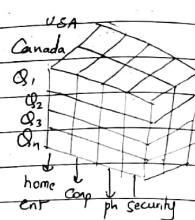


ongoing Absent

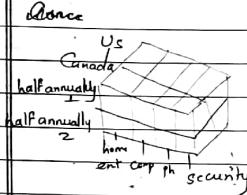




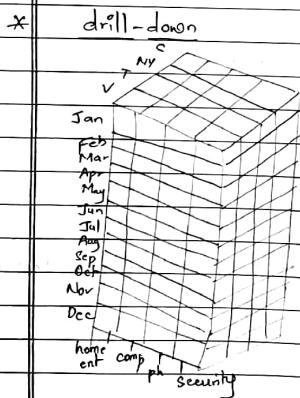
Rollup



The rollup operation performs aggregation on a data cube by either climbing up a concept hierarchy or by dimension reduction.



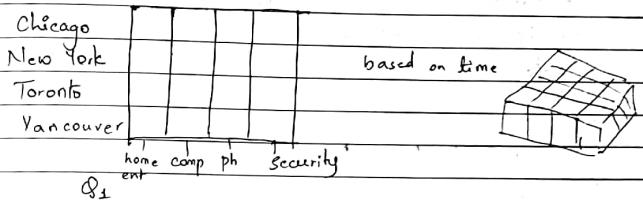
If rollup is performed using dimension reduction, one or more dimensions are removed from the given cube.



classmate
Date
Page

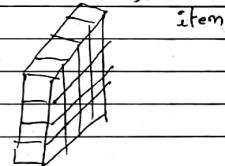
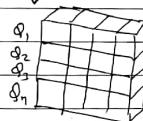
- A drill down operation navigates from less detailed data to more detailed data.
- It can be realized by stepping down a concept hierarchy for a given dimension or by introducing more dimension.

Slice Operation (Selection of data based on one dimension)



- > Slice performs the selection on one dimension of a given cube resulting in a sub-cube.

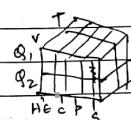
slice based on one location

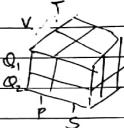


Dice Operation

Dice Operation defines a sub-cube by performing a selection on two or more dimension.

Selection based on time & entity location.





* Pivot Operation

- Rotating the axes of the diff. dimensions
- It's a visualisation operation that rotates the data axis with views to provide an alternate data representation.

he			
c			
ph			
s	C	NY	T

* drill across

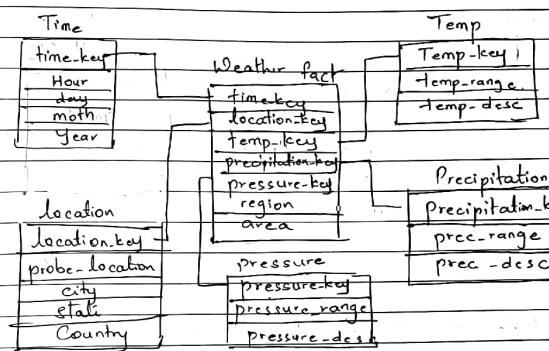
it refers queries across more than one fact table.

* drill through

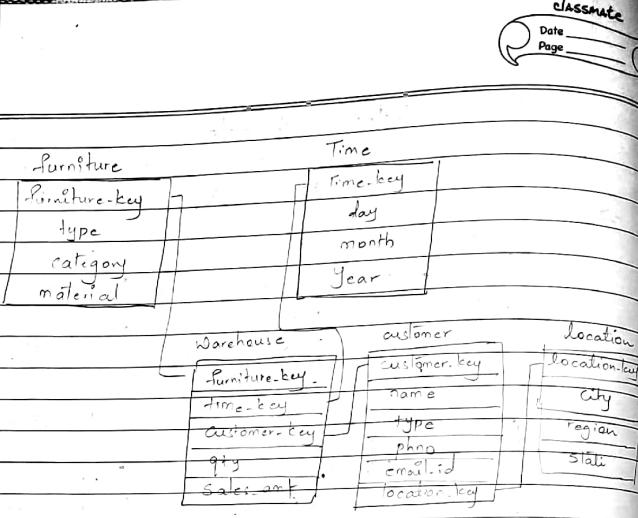
is an operation that drills through the bottom level of the data cube down to its back-end relational tables.

Q. Design a datawarehouse for a regional weather bureau. The weather bureau has about 1000 probes which are scattered throughout various land & ocean location in the region to collect basic weather data such as air-pressure, temp, precipitation at each hour. All data are sent to the central station.

which is collected such data for over 10 years. Your design should facilitate efficient query & all line analytical processing



Q. Design a datawarehouse for the wholesale furniture company. The datawarehouse has to allow to analyse the companies situation atleast w.r.t. to furniture, customers & time. Moreover the company needs to analyse the furniture w.r.t. to its time (Chair, table, Wardrobe etc.) forms a category (Kitchen, living room, office, etc) and material (wood, marble, etc). If also needs to analyse the customer w.r.t. to their special location by considering cities, regions & states. The company is interested in learning the quantity of furniture units sold & total sales amount generated.



* Data Generalisation using Attribute Oriented Induction
 Data Generalisation summarises the data by replacing low level concepts/values with high level concepts or by reducing the dimensions.
 → If we cannot generalise then remove the dimensions.

Generalisation can be done in 2 ways.

↳ Characterisation :- for one collection

↳ Discrimination :- for more collections

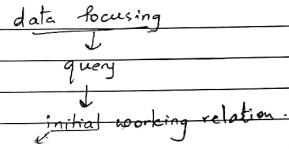
→ Characterization provides the concise and succinct summarization of a given data collection.

→ Discrimination provides description comparing two or more data collection

Q. Why OLAP → not used in data generalisation.
 → can't handle complex types
 depends on end user domain knowledge.

* Attribute Oriented Induction
 query → results are obtained & then generalised

↳ data focusing
 ↳ Attribute Relevance analysis
 ↳ correlation analysis
 ↳ entropy, gain, gini index



Attribute Removal

If we cannot generalise then remove the attributes
 eg. rollno, Aadhar number, etc.

Attribute Generalisation

Attribute threshold :- decide how much generalisation should be done
 Attribute threshold :- generalise until the threshold value
 relational threshold or generalising the tuple

→ Output of this is → prime generalised relation.

e.g.
 select name, gender, major, birthplace, birth_date, residence, phone, gpa
 from student
 where station in { "MSL", "MA", "MBA" }

- name → attribute removal
- gender → retain
- major → Science, Arts, Business Mgmt
- birth-place birthplace → city → state → country.
- birth-date → age → age ranges
- residence → city → state → country
- phone → attribute removal
- gpa → grades

i) date focusing: info.
ii) query → result is called as
iii) initial working relation

Gpa score: 3.97, classmate
Date: _____
Page: _____

Data Generalisation using Attribute Oriented Induction

Initial Working Relation							
name	gender	major	birth-place	birth-date	residence	phone	gpa
Jim	M	CS	Vancouver, Canada	12-8-76	3511, st	6537-008	3.97
Larsa	F	Phys	Seattle, USA	8-12-76	RJ Austin	659-5233	3.83

Attribute threshold
Generalisation → Attribute Removal
Relation threshold
Attribute Generalisation

Output of Generalisation → prime generalisation.

→ major: generalised into broad disciplines
birth-place: state < country
birth-date: age < age-range / age-groups
residence: city < state < country
gpa: grades.

Generalisation	gender	major	birth-country	age-group	residence-country	grade	count
	M	Sc	US	20-25	US	First	16
	F	Eng	Canada	20-25	Canada	First	20

Graduate Student Table

Candler	Major	Count
M	Sc	20
F	Sc	30
M	Eg	10
F	Eng	40
		100

UG Student table

Gender	Major	Count
M	Sc	15
F	Sc	40
M	Engg	10
F	Engg	10
		75

Q. Explain attribute oriented induction. Consider the following tables, from the given data, compute the gain of gender, compute the gain of major & analyse which attribute is more relevant.

⇒ Target Class = Graduate, UG

$$\text{Entropy}(D) = - \left[\left(\frac{100}{175} \right) \times \log_2 \left(\frac{100}{175} \right) + \left(\frac{75}{175} \right) \times \log_2 \left(\frac{75}{175} \right) \right]$$

$$= 0.985$$

or

$$\text{Gender}$$

$$\text{Gender} = M$$

$$\text{entropy} = - \left[\left(\frac{30}{55} \right) \times \log_2 \left(\frac{30}{55} \right) + \left(\frac{25}{55} \right) \times \log_2 \left(\frac{25}{55} \right) \right]$$

$$= 0.9940$$

$$\text{Gender} = F$$

$$\text{entropy} = - \left[\left(\frac{40}{120} \right) \times \log_2 \left(\frac{40}{120} \right) + \left(\frac{80}{120} \right) \times \log_2 \left(\frac{80}{120} \right) \right]$$

$$= 0.979$$

$$\text{net entropy (gender)} = \frac{55}{175} \times (x_1) + \frac{120}{175} \times (x_2)$$

$$= x_3 =$$

$$= 0.983$$

classmate
Date _____
Page _____

$$\text{Gain(Gender)} = \text{entropy}(D) - x_3$$

$$= 0.985 - 0.983$$

Major

$$\text{Major} = \text{Sc}$$

$$\text{entropy} = - \left[\left(\frac{50}{105} \right) \times \log_2 \left(\frac{50}{105} \right) + \left(\frac{55}{105} \right) \times \log_2 \left(\frac{55}{105} \right) \right]$$

$$= 0.998$$

Major = Engg

$$\text{entropy} = - \left[\left(\frac{50}{70} \right) \times \log_2 \left(\frac{50}{70} \right) + \left(\frac{20}{70} \right) \times \log_2 \left(\frac{20}{70} \right) \right]$$

$$= 0.863$$

$$\text{net entropy (major)} = \frac{105}{175} \times 0.998 + \frac{70}{175} \times 0.863$$

$$= 0.944$$

$$\text{Gain(Major)} = 0.985 - 0.944$$

$$= 0.041$$

goal

Attribute Oriented Induction
→ Descriimation / Class Comparison

- 1] Data Collection
→ query

→ partition data → target class : high risk diabetes
set of contrasting classes : low & medium risk

* generalising data for more than one data collection \rightarrow classification

CLASSMATE

Date _____

Page _____

- 2] Dimension relevance analysis : finding most imp. attributes using impurity measure & correlation analysis.
- 3] Synchronous Generalisation : generalising more than one data collection simultaneously.
 - \hookrightarrow prime target class relation.
 - \hookrightarrow prime contrasting class relation.
- 4] Data presentation.
can be presented in
 - \hookrightarrow charts, tables, etc.

e.g. find the characteristics of students pursuing UG & PG courses.

20/09/19

MODULE 3

Classification : used for discrete values.

model / learning function

explanatory variable \rightarrow target variable.
attribute set X .

- A fn that map explanatory variable to target variable is called learning function.
- explanatory variables :- describes the target class.

Training Set : class labels are given.

Test Set : class labels are not given.

* Decision Tree

types of nodes
 \rightarrow root node

\rightarrow internal nodes : 1 incoming & many outgoing edges.
 \rightarrow leaf node : 1 incoming & no outgoing edges.

* Hunt's Algorithm

D_t = training set

$y = \{y_1, y_2, \dots, y_n\}$

\rightarrow set of class labels

1] $D_t \rightarrow$ leaf node

2] $D_t \rightarrow$ more than one class label
use attribute test condition

\hookrightarrow create child node for each outcome of the chosen attribute (D_{tj})

\hookrightarrow partition data among child nodes.

\hookrightarrow recursively apply the algorithm for each child node.

Test Conditions :

- 1] If child node is empty:
 - \hookrightarrow depends on the parent node.
 - \hookrightarrow creates a leaf node \rightarrow with majority class of parent node.

2] Identical attribute values.

Habitat	Mammal	Class	\rightarrow we cannot split
Land	yes	herbivore	
Land	yes	carnivore	
Land	yes	herbivore,	

\hookrightarrow create a leaf node

\hookrightarrow majority class of given child node.

Design issues \rightarrow decision tree

\rightarrow attribute test condition

\rightarrow an objective measure : used to check the performance of test condition.

\rightarrow stopping condition : when all records have same class label

a) when identical attribute, we have atleast one class label in majority.

* Performance of Classification model

\rightarrow confusion matrix.

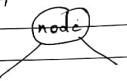
	No. of records classified as '0'	No. of records classified as '1'
Class 0	f_{00}	f_{01}
Class 1	f_{10}	f_{11}

$$\text{Accuracy} = \frac{f_{00} + f_{11}}{f_{00} + f_{11} + f_{01} + f_{10}}$$

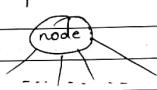
$$\text{Error rule} = \frac{f_{01} + f_{10}}{f_{00} + f_{11} + f_{01} + f_{10}}$$

* Types of splits

- Binary Split



- Multway Split

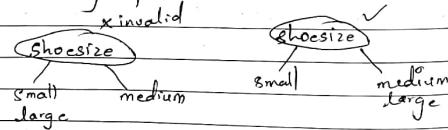


* Nominal Attributes

↳ binary / multway split

Ordinal Attribute: order should be preserved.

↳ binary / Multway Split



* Continuous Attribute

↳ binary / Multway Split

$A > v$ or $A \leq v$

or

classmate
Date _____
Page _____

classmate
Date _____
Page _____

range query
 $v_i \leq A < v_{i+1}$

Measures of determine best split

impurity measures

$$\rightarrow \text{entropy} = - \sum p(c|t) \log_2 (p(c|t))$$

Probability of a given class at note t

$$\rightarrow \text{gini index} = 1 - \left[\sum (p(c|t))^2 \right]$$

$$\rightarrow \text{classification error} = 1 - \max_t [p(c|t)]$$

Nominal Attribute

car type (N1)	(N2)	CN1
Sports, Luxury		
C0	9	1
C1	7	3

$$\text{Gini index}, N_1 = 1 - \left(\frac{9}{16} \right)^2 - \left(\frac{7}{16} \right)^2$$

$$N_2 = 0$$

$$N_1 = 0.49218$$

$$\text{Gini index}, N_2 = 1 - \left(\frac{1}{4} \right)^2 - \left(\frac{3}{4} \right)^2$$

$$N_2 = 0.375$$

$$\text{Weighted Gini index} = \frac{16 \times 0.49218 + 4 \times 0.375}{20}$$

$$= 0.468$$

	Sports N ₃	Family Luxury (N ₄)
C0	8	2
C1	0	10

$$\text{Gini Index } N_3 = 1 - \left(\frac{8}{10}\right)^2 - \left(\frac{2}{10}\right)^2$$

$$N_3 = 0$$

$$\text{Gini Index } N_4 = 1 - \left(\frac{2}{12}\right)^2 - \left(\frac{10}{12}\right)^2$$

$$= 0.2977$$

$$\text{Weighted Gini Index} = \frac{8}{20} \times 0 + \frac{12}{20} \times 0.2977$$

$$= 0.1666$$

Q	Node n.1	Count
	class 0	1
	class 1	5

$$\text{entropy } (CN_1) = - \left[\left(\frac{1}{6}\right) \log_2 \left(\frac{1}{6}\right) + \left(\frac{5}{6}\right) \log_2 \left(\frac{5}{6}\right) \right]$$

$$= 0.650$$

$$\text{Gini index } (CN_1) = 1 - \left(\frac{1}{6}\right)^2 - \left(\frac{5}{6}\right)^2$$

$$= 0.298$$

$$\text{Classification error} = 1 - \max \left(\left(\frac{1}{6}\right), \left(\frac{5}{6}\right) \right)$$

$$= 1 - \max (0.1667, 0.833)$$

$$= 1 - 0.833$$

$$= 0.167$$

Annual Income	default borrowers
60 <= 63	no
70	no
75	no
85	yes
90	yes
95	yes
100 > 102	yes
120	no
125	no
220	no

$$\text{avg} = \frac{60+70}{2} = 65$$

- least entropy value is the split point

- Continue the same procedure
+ Sort if not sorted

	yes	no
<= 65	0	1
> 65	3	6

$$\text{entropy } (<= 65) = - \left[\left(\frac{0}{3}\right) \log_2 \left(\frac{0}{3}\right) + \left(\frac{1}{3}\right) \log_2 \left(\frac{1}{3}\right) \right]$$

$$= 0$$

$$\text{entropy } (> 65) = - \left[\left(\frac{3}{9}\right) \log_2 \left(\frac{3}{9}\right) + \left(\frac{6}{9}\right) \log_2 \left(\frac{6}{9}\right) \right]$$

$$= 0.9182$$

$$\text{Weighted entropy} = \left(\frac{1}{10}\right) \times 0 + \left(\frac{9}{10}\right) \times 0.9182$$

$$= 0.8263$$

Gain : Difference of impurity measure

$$\Delta I(\text{parent}) = \sum_{j=1}^k \frac{N(v_j)}{N} I(v_j)$$

$N(v_j)$: No. of records associated with childnode v_j

N : No. of records at parent node

No. of attribute values.

Using ID3 (Iterative Dichotomiser)

Gender	height	Class
F	1.6	Short
F	1.6	Short
F	1.7	Short
M	1.7	Short
F	1.75	Medium
F	1.8	Medium
F	1.8	Medium
M	1.85	Medium
F	1.88	Medium
F	1.9	Medium
M	1.95	Medium
M	2.0	Medium
M	2.1	Tall
M	2.2	Tall

Q. Divide the height attribute into ranges as follows:

(0, 1.6]

(1.6, 1.7]

(1.7, 1.8]

(1.8, 1.9]

(1.9, 2.0]

(2.0, 5.0]

construct the decision tree for the following data

(0, 1.6] OR height ≤ 1.6

Gender	height	Class
F	[0, 1.6]	Short
F	(0, 1.6]	Short
F	(1.6, 1.7]	Short
M	(1.6, 1.7]	Short
F	(1.7, 1.8]	Medium
F	(1.7, 1.8]	Medium
F	(1.7, 1.8]	Medium
M	(1.8, 1.9]	Medium
F	(1.8, 1.9]	Medium
F	(1.8, 1.9]	Medium
M	(1.9, 2.0]	Medium
M	(1.9, 2.0]	Medium
M	(2.1, 5.0]	Tall
M	(2.1, 5.0]	Tall

$$\text{entropy}(D) = - \left[\left(\frac{4}{15} \right) \log_2 \left(\frac{4}{15} \right) + \left(\frac{9}{15} \right) \log_2 \left(\frac{9}{15} \right) + \left(\frac{2}{15} \right) \log_2 \left(\frac{2}{15} \right) \right] = 1.3381$$

$$\text{Cgender} \quad \begin{matrix} \text{short} & \text{medium} \\ \text{M} & \end{matrix}$$

$$\text{entropy}(\text{Cgender} = M) = - \left[\left(\frac{1}{6} \right) \log_2 \left(\frac{1}{6} \right) + \left(\frac{3}{6} \right) \log_2 \left(\frac{3}{6} \right) + \left(\frac{2}{6} \right) \log_2 \left(\frac{2}{6} \right) \right] = 1.4591$$

$$\text{Cgender} \quad \begin{matrix} \text{short} & \text{medium} \\ F & \end{matrix}$$

$$\text{entropy}(\text{Cgender} = F) = - \left[\left(\frac{3}{9} \right) \log_2 \left(\frac{3}{9} \right) + \left(\frac{6}{9} \right) \log_2 \left(\frac{6}{9} \right) + \left(\frac{0}{9} \right) \log_2 \left(\frac{0}{9} \right) \right] = 0.9182$$

net entropy / weighted entropy = $\frac{6}{15} \times 1.459 + \frac{9}{15} \times 0.9182$
 (gender) = $= 1.1345$

Gain(gender) = entropy(D) - weighted entropy(gender)
 $= 1.1338 - 1.1345$
 $= 0.00354$

Height
 entropy(Height) = $(0, 1.6] = -\left[\left(\frac{2}{2}\right) \log_2 \left(\frac{2}{2}\right) + \left(\frac{0}{2}\right) \log_2 \left(\frac{0}{2}\right) + \left(\frac{0}{2}\right) \log_2 \left(\frac{0}{2}\right) \right]$
 $= 0$

entropy(Height = $(1.6, 1.7]$) = $-\left[\left(\frac{2}{2}\right) \log_2 \left(\frac{2}{2}\right) + \left(\frac{0}{2}\right) \log_2 \left(\frac{0}{2}\right) + \left(\frac{0}{2}\right) \log_2 \left(\frac{0}{2}\right) \right]$
 $= 0$

entropy(Height = $(1.7, 1.8]$) = $-\left[\left(\frac{2}{3}\right) \log_2 \left(\frac{2}{3}\right) + \left(\frac{3}{3}\right) \log_2 \left(\frac{3}{3}\right) + \left(\frac{0}{3}\right) \log_2 \left(\frac{0}{3}\right) \right]$
 $= 0$

entropy(Height = $(1.8, 1.9]$) = $-\left[\left(\frac{4}{4}\right) \log_2 \left(\frac{4}{4}\right) + \left(\frac{1}{4}\right) \log_2 \left(\frac{1}{4}\right) + \left(\frac{0}{4}\right) \log_2 \left(\frac{0}{4}\right) \right]$
 $= 0$

entropy(Height = $(1.9, 2.0]$) = $-\left[\left(\frac{0}{2}\right) \log_2 \left(\frac{0}{2}\right) + \left(\frac{2}{2}\right) \log_2 \left(\frac{2}{2}\right) + \left(\frac{0}{2}\right) \log_2 \left(\frac{0}{2}\right) \right]$
 $= 0$

entropy(Height = $(2.0, 5.0]$) = $-\left[\left(\frac{0}{10}\right) \log_2 \left(\frac{0}{10}\right) + \left(\frac{0}{10}\right) \log_2 \left(\frac{0}{10}\right) + \left(\frac{2}{10}\right) \log_2 \left(\frac{2}{10}\right) \right]$
 $= 0$

Weighted entropy(Height) = $\left(\frac{2}{15}\right) \times 0 + \left(\frac{2}{15}\right) \times 0 + \left(\frac{3}{15}\right) \times 0 + \left(\frac{4}{15}\right) \times 0 + \left(\frac{2}{15}\right) \times 0 + \left(\frac{0}{15}\right) \times 0$
 $= 0$

Gain(Height) = entropy(D) - weighted entropy(Height)
 $= 1.1338 - 0$
 $= 1.1338$

Since gain(Height) is highest \rightarrow split based on height

	Gender	Car Ownership	TravelCost	IncomeLevel	TransportationMode
M	0	Cheap	low	bus	
M	1	Cheap	medium	bus	
F	1	Cheap	medium	train	
F	0	Cheap	low	bus	
M	1	Cheap	medium	bus	
M	0	Standard	medium	train	
F	1	Standard	medium	train	
F	1	expensive	high	car	
M	2	expensive	medium	car	
F	0	expensive	high	car	

Entropy(D) = $-\left[\left(\frac{4}{10}\right) \log_2 \left(\frac{4}{10}\right) + \left(\frac{3}{10}\right) \log_2 \left(\frac{3}{10}\right) + \left(\frac{3}{10}\right) \log_2 \left(\frac{3}{10}\right) \right]$
 $= 1.5709 = 1.571$

Gender

$$\text{entropy}(\text{gender} = \text{M}) = - \left[\left(\frac{1}{5} \right) \log_2 \left(\frac{1}{5} \right) + \left(\frac{2}{5} \right) \log_2 \left(\frac{2}{5} \right) + \left(\frac{2}{5} \right) \log_2 \left(\frac{2}{5} \right) \right]$$

$$= 1.5229$$

Gender (gender = M)

$$\text{entropy}(\text{gender} = \text{M}) = - \left[\left(\frac{3}{5} \right) \log_2 \left(\frac{3}{5} \right) + \left(\frac{1}{5} \right) \log_2 \left(\frac{1}{5} \right) + \left(\frac{1}{5} \right) \log_2 \left(\frac{1}{5} \right) \right]$$

$$= 1.3709$$

Ownership

$$\text{entropy}(\text{ownership} = 0) = - \left[\left(\frac{2}{3} \right) \log_2 \left(\frac{2}{3} \right) + \left(\frac{1}{3} \right) \log_2 \left(\frac{1}{3} \right) + \left(\frac{1}{3} \right) \log_2 \left(\frac{1}{3} \right) \right]$$

$$= 0.9182$$

net entropy / weighted entropy = $\frac{5}{10} \times 1.3709 + \frac{5}{10} \times 1.5229$
(gender)

$$= 1.447$$

Gain(gender) = entropy(D) - net entropy(gender)

$$= 1.571 - 1.447$$

$$= 0.1246$$

Car Ownership

$$\text{entropy}(\text{cos} = 0) = - \left[\left(\frac{2}{3} \right) \log_2 \left(\frac{2}{3} \right) + \frac{1}{3} \log_2 \left(\frac{1}{3} \right) + \frac{0}{3} \log_2 \left(\frac{0}{3} \right) \right]$$

$$= 0.9182$$

$$\text{entropy}(\text{cos} = 1) = - \left[\left(\frac{2}{5} \log_2 \left(\frac{2}{5} \right) + \frac{2}{5} \log_2 \left(\frac{2}{5} \right) + \left(\frac{1}{5} \right) \log_2 \left(\frac{1}{5} \right) \right] \right]$$

$$= 1.5219$$

$$\text{entropy}(\text{cos} = 2) = - \left[\frac{0}{2} \log_2 \frac{0}{2} + \frac{0}{2} \log_2 \frac{0}{2} + \frac{2}{2} \log_2 \frac{2}{2} \right]$$

$$= 0$$

Travel Cost

$$\text{entropy}(\text{TC} = \text{cheap}) = - \frac{5}{10}$$

net entropy = $\frac{3}{10} \times 0.918 + \frac{5}{10} \times 1.5229 + 0$

$$= 1.0364$$

Gain(COS) = entropy(D) - net entropy(cos)

$$= 1.571 - 1.0364$$

$$= 0.5346$$

Travel Cost

$$\text{entropy}(\text{TC} = \text{cheap}) = - \left[\frac{4}{5} \log_2 \left(\frac{4}{5} \right) + \frac{1}{5} \log_2 \left(\frac{1}{5} \right) + \frac{0}{5} \log_2 \left(\frac{0}{5} \right) \right]$$

$$= 0.7219$$

$$\text{entropy}(\text{TC} = \text{standard}) = - \left[\left(\frac{0}{2} \right) \log_2 \left(\frac{0}{2} \right) + \frac{2}{2} \log_2 \frac{2}{2} + \frac{0}{2} \log_2 \frac{0}{2} \right]$$

$$= 0$$

$$\text{entropy}(\text{TC} = \text{expensive}) = - \left[\frac{0}{2} \log_2 \frac{0}{2} + \frac{0}{2} \log_2 \frac{0}{2} + \frac{2}{2} \log_2 \frac{2}{2} \right]$$

$$= 0$$

net entropy(TravelCost) = $\frac{5}{10} \times 0.7219 + \frac{2}{10} \times 0 + \frac{2}{10} \times 0$

$$= 0.36095 = 0.361$$

Gain(TravelCost) = entropy(D) - net entropy(TC)

$$= 1.571 - 0.36095$$

$$= 1.21005 = 1.210$$

Income Level

$$\text{entropy (IL = low)} = - \left[\frac{2}{2} \log_2 \frac{2}{2} + \frac{0}{2} \log_2 \frac{0}{2} + \frac{0}{2} \log_2 \frac{0}{2} \right]$$

$$= 0$$

$$\text{entropy (IL = medium)} = - \left[\frac{2}{6} \log_2 \frac{2}{6} + \frac{3}{6} \log_2 \frac{3}{6} + \frac{1}{6} \log_2 \frac{1}{6} \right]$$

$$= + 1.459$$

$$\text{entropy (IL = high)} = - \left[\frac{0}{2} \log_2 \frac{0}{2} + \frac{0}{2} \log_2 \frac{0}{2} + \frac{2}{2} \log_2 \frac{2}{2} \right]$$

$$= 0$$

$$\text{net entropy (IncomeLevel)} = \frac{2}{10} \times 0 + \frac{6}{10} \times 1.459 + \frac{2}{10} \times 0$$

$$= 0.8754$$

$$\text{Gain (IncomeLevel)} = \text{entropy}(D_0) - \text{net entropy (IncomeLevel)}$$

$$= 1.571 - 0.8754$$

$$= 0.6956$$

		travel cost		
		Gender	Carownership	Incomelevel
Car	expensive	M	0	low
Train	standard	M	1	medium
Bus	cheap	F	1	medium
		F	0	low
		M	1	medium

2nd Iteration

$$\text{entropy (D}_1\text{)} = - \left[\frac{1}{5} \log_2 \frac{1}{5} + \frac{4}{5} \log_2 \frac{4}{5} \right] = 0.722$$

entropy (gender)

$$\text{entropy (gender = F)} = - \left[\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right]$$

$$= 1$$

$$\text{entropy (gender = M)} = - \left[\frac{3}{3} \log_2 \frac{3}{3} + \frac{0}{2} \log_2 \frac{0}{2} \right]$$

$$= 0$$

net $\text{entropy(gender)} = \frac{2}{5} \times 1 + \frac{3}{5} \times 0 = 0.4$

gain(gender) = $\text{entropy}(D_0) - \text{net entropy(gender)}$
 $= 0.722 - 0.4 = 0.322$

Carownership

$$\text{entropy (CO = 0)} = - \left[\frac{2}{2} \log_2 \frac{2}{2} + \frac{0}{2} \log_2 \frac{0}{2} \right]$$

$$= 0$$

$$\text{entropy (CO = 1)} = - \left[\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3} \right]$$

$$= 0.9182$$

$$\text{net entropy (CarOwnership)} = \frac{2}{5} \times 0 + \frac{3}{5} \times 0.9182$$

$$= 0.5509$$

gain (CarOwnership) = $\text{entropy}(D_1) - \text{net entropy(CarOwnership)}$
 $= 0.722 - 0.5509$
 $= 0.1711$

classmate
Date _____
Page _____

$$\text{Income level entropy (IL=low)} = - \left[\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3} \right]$$

$$= 0$$

$$\text{entropy (IL=medium)} = - \left[\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3} \right]$$

$$= 0.9182$$

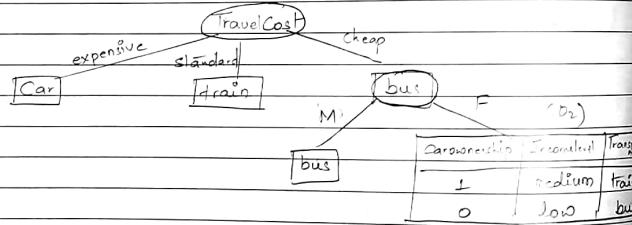
$$\text{net entropy (IncomeLevel)} = \frac{2}{5} \times 0 + \frac{3}{5} \times 0.9182$$

$$= 0.5509$$

✓ gain (Incomelevel) = entropy(D) - net entropy (Incomelevel)

$$= 0.722 - 0.5509$$

$$= 0.1711$$



3rd iteration

$$\text{entropy (D}_2\text{)} = - \left[\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right]$$

$$= 1$$

Carownersip

$$\text{entropy (CO=0)} = - \left[\frac{1}{1} \log_2 \frac{1}{1} + \frac{0}{1} \log_2 \frac{0}{1} \right] = 0$$

classmate
Date _____
Page _____

$\left\{ \begin{array}{l} \text{Cosine Similarity (Document Data)} \\ \text{Any association in FP Tree, Apriori, vertical Data Format} \\ \text{Decision Tree (C4.5, Entropy)} \end{array} \right.$

$$\text{entropy (CO=1)} = - \left[\frac{1}{1} \log_2 \frac{1}{1} + \frac{0}{1} \log_2 \frac{0}{1} \right]$$

$$= 0$$

$$\text{net entropy (Carownership)} = 0$$

$$\text{gain(Carownership)} = \frac{\text{entropy (CO=0)}}{\text{net entropy}} = \frac{1-0}{1-0} = 1$$

Income Level

$$\text{entropy (IL=medium)} = - \left[\frac{1}{1} \log_2 \frac{1}{1} + \frac{0}{1} \log_2 \frac{0}{1} \right]$$

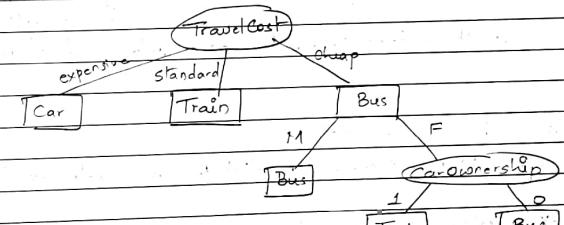
$$= 0$$

$$\text{entropy (IL=low)} = - \left[\frac{1}{1} \log_2 \frac{1}{1} + \frac{0}{1} \log_2 \frac{0}{1} \right]$$

$$= 0$$

$$\text{net entropy (IncomeLevel)} = 0$$

$$\text{gain (Incomelevel)} = 1-0 = 1$$



$$\text{Gain ratio} = \frac{\text{Gain}}{\text{Split Info}}$$

$$\text{Split Info} = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \frac{|D_j|}{|D|}$$

$v \rightarrow$ no. of outcomes

$D_j \rightarrow$ records associated with each split

Alternative for Gain Ratio is Binary Split.

* Characteristics of a decision tree

- non-parametric approach
- NP complete problem (Non-deterministic Polynomial)
- Computationally inexpensive
- easy to interpret → small decision tree.
- it is a top-down recursive partitioning approach
- data fragmentation :- handle by using gain ratio, binary splits, pruning.
- it divides attribute space into disjoint regions (with defined classes)

* Decision Border

↳ rectilinear → (parallel to co-ordinate axis)

* Oblique decision tree

Attribute test condition with more than one attribute.

* Constructive Induction

- Only one attribute

→ uses composite attribute.

* The impurity measure does not affect performance of Decision Tree.

* → quite robust to noise and redundant data.

Outlook	Temperature	humidity	Wind	PlayTennis
Sunny	hot	high	weak	no
Sunny	hot	high	strong	no
Overcast	hot	high	weak	yes
Rain	mild	high	weak	yes
Rain	cool	normal	weak	yes
Rain	cool	normal	strong	no
Overcast	cool	normal	strong	yes
Sunny	mild	high	weak	no
Sunny	cool	normal	weak	yes
Rain	mild	normal	weak	yes
Sunny	mild	normal	strong	yes
Overcast	mild	high	strong	yes
Overcast	hot	normal	weak	yes
Rain	mild	high	strong	no

$$\text{entropy}(D) = - \left[\frac{9}{14} \log_2 \left(\frac{9}{14} \right) + \left(\frac{5}{14} \right) \log_2 \left(\frac{5}{14} \right) \right]$$

$$= 0.9402$$

Outlook

$$\text{entropy}(\text{Outlook} = \text{sunny}) = - \left[\left(\frac{2}{5} \right) \log_2 \left(\frac{2}{5} \right) + \left(\frac{3}{5} \right) \log_2 \left(\frac{3}{5} \right) \right]$$

$$= 0.9705$$

$$\text{entropy}(\text{Outlook} = \text{overcast}) = - \left[\frac{4}{14} \log_2 \frac{4}{14} + \frac{0}{14} \log_2 \frac{0}{14} \right]$$

$$= 0$$

$$\text{entropy}(\text{Outlook} = \text{rain}) = - \left[\frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5} \right]$$

$$= 0.9709$$

classmate
Date _____
Page _____

$$\text{net entropy} = \frac{5}{14} \times 0.970 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.970$$

$$= 0.6928$$

$$\text{gain}(\text{Outlook}) = \text{entropy}(D) - \text{net entropy}(\text{Outlook})$$

$$= 0.940 - 0.6928$$

$$= 0.248.$$

$$\text{Split Info} = - \left[\left(\frac{5}{14} \right) \log_2 \left(\frac{5}{14} \right) + \left(\frac{4}{14} \right) \log_2 \left(\frac{4}{14} \right) + \frac{5}{14} \log_2 \frac{5}{14} \right]$$

$$= 1.577$$

$$\text{Gain ratio} = \frac{\text{Gain}(\text{Outlook})}{\text{Split Info}(\text{Outlook})} = \frac{0.248}{1.577} = 0.157$$

Temperature

$$\text{entropy}(T = \text{hot}) = - \left[\frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{4} \log_2 \frac{2}{4} \right]$$

$$= 1$$

$$\text{entropy}(T = \text{mild}) = - \left[\frac{4}{6} \log_2 \frac{4}{6} + \frac{2}{6} \log_2 \frac{2}{6} \right]$$

$$= 0.918$$

$$\text{entropy}(T = \text{cool}) = - \left[\frac{3}{4} \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \frac{1}{4} \right]$$

$$= 0.811$$

$$\text{net entropy(Temp)} = \frac{4}{14} \times 1 + \frac{6}{14} \times 0.918 + \frac{4}{14} \times 0.811$$

$$= 0.9108 = 0.91108$$

classmate
Date _____
Page _____

$$\text{gain(Temp)} = 0.940 - 0.9108 = 0.0292 = 0.030$$

$$\text{splitInfo(Temp)} = -\left[\frac{4}{14} \log_2 \frac{4}{14} + \frac{6}{14} \log_2 \frac{6}{14} + \frac{4}{14} \log_2 \frac{4}{14} \right]$$

$$= 1.55665$$

$$\text{Gain ratio} = \frac{0.030}{1.55665} = 0.019.$$

→ Humidity

$$\text{entropy(H=high)} = -\left[\frac{3}{7} \log_2 \frac{3}{7} + \frac{4}{7} \log_2 \frac{4}{7} \right]$$

$$= 0.98522$$

$$\text{entropy(H=normal)} = -\left[\frac{6}{7} \log_2 \frac{6}{7} + \frac{1}{7} \log_2 \frac{1}{7} \right]$$

$$= 0.59167$$

$$\text{net entropy(Humidity)} = \frac{7}{14} \times 0.98522 + \frac{7}{14} \times 0.59167$$

$$= 0.788$$

✓ gain(Humidity) = 0.940 - 0.788 = 0.15155

$$\text{splitInfo(Humidity)} = -\left[\frac{7}{14} \log_2 \frac{7}{14} + \frac{7}{14} \log_2 \frac{7}{14} \right]$$

$$= 1$$

$$\text{Gain ratio} = \frac{0.15155}{1} = 0.15155$$

Wind

$$\text{entropy(Wind=Weak)} = -\left[\frac{6}{8} \log_2 \frac{6}{8} + \frac{2}{8} \log_2 \frac{2}{8} \right]$$

$$= 0.8112$$

$$\text{entropy(Wind=Strong)} = -\left[\frac{3}{6} \log_2 \frac{3}{6} + \frac{3}{6} \log_2 \frac{3}{6} \right]$$

$$= 1$$

$$\text{net entropy(Wind)} = \frac{8}{14} \times 0.8112 + \frac{6}{14} \times 1$$

$$= 0.89211$$

← gain(Wind) = 0.940 - 0.89211 = 0.04789

$$\text{splitInfo(Wind)} = -\left[\frac{8}{14} \log_2 \frac{8}{14} + \frac{6}{14} \log_2 \frac{6}{14} \right]$$

$$= 0.98522$$

$$\text{Gain ratio} = \frac{0.04789}{0.98522} = 0.0486$$

03/08/19

Classification

classmate

Date _____

Page _____

- training error: ~~how~~ no. of records which are incorrectly classified during training phase
 - Exact value for training value.
 - generalization error: errors which occurs during Testing phase
 - only be approximated
- | | |
|---|-------------------------------------|
| → <u>Model overfitting</u> (More records) | - Model Underfitting (less records) |
| → low training error | → high training error |
| → relatively high generalization error. | → high generalization error |

→ Model overfitting takes place b'coz of noise, lack of representative samples, multiple comparison procedure (multiple test conditions)

→ Model overfitting occur when you have detailed decision tree.

Techniques to find approximate generalisation error

(1) Resubstitution Error

→ training error is considered as approx. value for generalⁿ error

(2) pessimistic estimate

↪ sum of training error + penalty term

$$e_g(T) = \frac{\sum_{i=1}^k [e(t_i) + \eta(t_i)]}{\sum_{i=1}^k n(t_i)}$$

$$= \frac{e(T) + \Omega(T)}{N_t}$$

training error = 4

no. of nodes = 7

no. of training records = 24 $\Omega = 0.5$

$$\therefore = \frac{e(T) + \Omega(T)}{N_t} = \frac{4 + 7 \times 0.5}{24}$$

* Minimum Description Length $\xrightarrow{\text{OR}}$ Occam's Razor / Principle of Parsimony

Cost = Cost(model) + cost(data|model)

$\xrightarrow{\text{OR}}$ Choose model which is smaller / simpler one

(b) Statistical Correction: The generalization is upperbound of the statistical correction applied to the training error.

$$\text{Upper (N,i.e., } \alpha) = e + \frac{Z_{\alpha/2}^2}{2N} + Z_{\alpha/2} \sqrt{\frac{e(1-e)}{N}} + \frac{Z_{\alpha/2}^2}{4N^2}$$

$\xrightarrow{\text{confidence level}}$

Using Validation set

break training into 2 parts \rightarrow a few records are kept for testing

optional

* Evaluating Performance of Classifier
 \rightarrow holdout method

Data set
 ——————
 training set
 ——————
 Test set

Random Subsampling:
 K times.

Cross-Validation

- α fold validation: divide into training & test set, perform validation + for next iteration swap
- K fold validation: - data divide into K partition: one part for testing
- leave one approach. ($K=N$)
 N : size of the data set
- $N-1 \rightarrow$ used for training
 N^{th} record \rightarrow testing

Note: Continuous data can be represented as range query

$$\therefore c \leq A \leq d$$

$$A < V \text{ or}$$

$$A \geq V$$



* Rule Based Classification

$A \rightarrow y$

↑ consequent

antecedent /
Precondition

$$A = C_1 \vee C_2 \vee C_3 \dots \vee C_n$$

→ Rule gets triggered when antecedent satisfies the condition

$$A \rightarrow y$$

- e.g.: $r_1 : (\text{Gives birth} = \text{no}) \wedge (\text{Aerial Creature} = \text{yes}) \rightarrow \text{Birds}$
- $r_2 : (\text{Gives birth} = \text{no}) \wedge (\text{Aquatic Creature} = \text{yes}) \rightarrow \text{Fishes}$
- $r_3 : (\text{Gives birth} = \text{yes}) \wedge (\text{Body Temp} = \text{warm blooded}) \rightarrow \text{Mammals}$
- $r_4 : (\text{Gives birth} = \text{no}) \wedge (\text{Aerial Creature} = \text{no}) \rightarrow \text{Reptiles}$
- $r_5 : (\text{Aquatic Creature} = \text{semi}) \rightarrow \text{Amphibians}$

$|D|$: No. of records in a data set

$|A \cap Y|$: No. of records satisfying the antecedent &

and consequent

$|A|$: No. of records satisfying the antecedent

$$\therefore \text{Accuracy} = \frac{|A \cap Y|}{|A|}$$

$$\text{Coverage} = \frac{|A|}{|D|}$$

r3:

$$|A| = 5$$

$$|A \cap Y| = 5$$

$$|D| = 15$$

$$\therefore \text{Accuracy} = 5/5 = 100\%$$

$$\text{Coverage} = 5/15 = 33.33\%$$

r1:

$$|A| = 2$$

$$|A \cap Y| = 2$$

$$|D| = 15$$

$$\text{Accuracy} = 2/2 = 100\%$$

$$\text{Coverage} = 2/15 = 13.33\%$$

r2:

$$|A| = 2$$

$$|A \cap Y| = 2$$

$$|D| = 15$$

$$\text{Accuracy} = 2/2 = 100\%$$

$$\text{Coverage} = 2/15 = 13.33\%$$

r4:

$$|A| = 8$$

$$|A \cap Y| = 3$$

$$|D| = 15$$

$$\text{Accuracy} = 3/8 = 37\%$$

$$\text{Coverage} = 3/15 = 20\%$$

r5:

$$|A| = 4$$

$$|A \cap Y| = 2$$

$$|D| = 15$$

$$\text{Avg} = 2/4 = 50\%$$

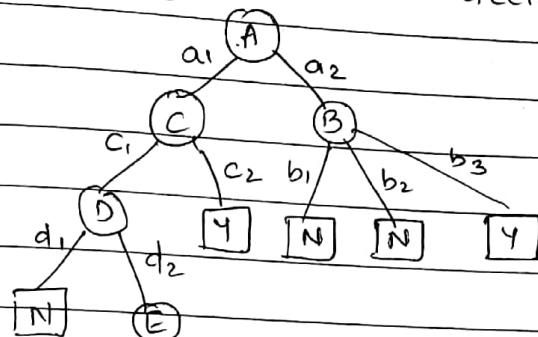
$$\text{Coverage} = 2/15 = 13.33\%$$

Continuation of decision tree

1) Pre-prunning :- finds gain & if gain < threshold, it halts the growing of tree. (Doesn't allow to construct the entire tree)

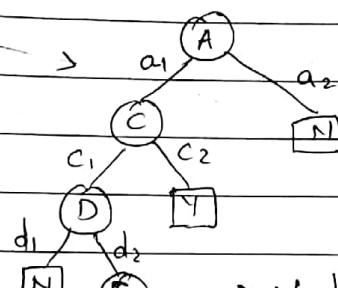
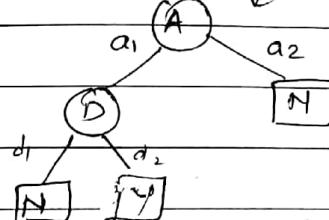
2) Post-prunning : Allow to construct the entire tree

prunning :-
 is done to handle model overfitting
 → reducing branches in decision tree



- Subtree Replacement

- Subtree Raising



→ Node is replaced by a leaf node

→ Node is replaced by a subtree

* Mutually Exclusive Rules

The rules in a rule set are said to be mutually exclusive, if no two rules are triggered by the same set of records.

* Exhaustive Coverage

The rules in a rule set are said to have exhaustive coverage, if there is a rule for each combination of attribute value.

training error = 0

* Default Rule $r_d(x) \rightarrow y_d$

→ Majority class

→ When none of the rules are satisfied

* 2 Types of Rules

1] Ordered Rules: Rules which are ordered with priority

→ rules → ordered ^{according} → priority

↳ decision list

Adv of Ordering Rules:

2] Unordered Rules:

Types of Rule Ordering

→ Rule by Rule ordering

→ class by class ordering: rules are grouped having same class → then group the class

Rule based classifier

1] Direct method for rule extraction :- analyse → apply algo and extract rule.

2] Indirect method for rule extraction :- extracting from other classifier

1] Direct method

→ Sequential Covering Algorithm

→ D → dataset

→ Attributals → set of attribute & values.

1] Rule set = { }

2] for each class C do

3] repeat

4] Rule = Learn-One-Rule (D, AttrVal_s, C)

5] remove tuples from Discovering rule.

6] Ruleset = Ruleset + Rule

7] until terminating condition

8] end for

9] return Rule set

Find max. coug

For Rule s

+ve example

* Rule growing strategies

→ general to specific:

Contain empty set & go on adding

→ specific to general:

↳ choose one +ve example

↳ eliminates condition & check the accuracy.

e.g.: If it is a mammal, we need only body-temp and gives birth
remove has-legs X, aerial creature --.

* Rule Quality Measure

R ratio / Likelihood ratio

$$R = 2 \sum_{i=1}^k f_i \log_2 (f_i/e_i)$$

f_i is observed freq. of class i covered by the rule
and e_i is the expected freq. of class i covered
by the rule

$$r_i = (C_1 \wedge C_2 \rightarrow m_i)$$

$$r'_i = (C_1 \wedge C_2 \wedge C_3 \rightarrow m_i)$$

$$RC(r_i) \quad RC(r'_i)$$

10/10/19

Laplace Measure

$$= \frac{f_+ + 1}{n+k}$$

f_+ : no. of the tuples which are covered by the rule

n : total no. of tuple covered by the rule.

m-estimate

$$= \frac{f_+ + K p_+}{n+k}$$

 K : no. of classes p_+ : Prior probability for the classes

$$\text{FOIL Gain } (R_0, R_1) = p_1 \times \left(\log_2 \frac{p_1}{p_1 + n_1} \right) - \log_2 \left(\frac{p_0}{p_0 + n_0} \right)$$

 R_0 : Original rule R_1 : Modified rule p_1 : +ve examples covered by R_1 p_0 : the tuples covered by R_0 n_0 : no. of -ve tuples covered by R_0

Q. training set 60 +ve examples

100 -ve examples.

 $R_1 \rightarrow$ covers 50 +ve examples & 5 -ve examples. $R_2 \rightarrow$ covers 2 +ve examples & 0 -ve examples.

$$e_i = \frac{55 \times 60}{160} = 20.625$$

(for +ve examples)

$$e_i \text{ (for -ve examples)} = \frac{55 \times 100}{160} = 34.375$$

$$R = 2 \left[50 \times \log_2 \left(\frac{50}{20.625} \right) + 5 \times \log_2 \left(\frac{5}{34.375} \right) \right]$$

$$R = 99.9$$

r_2

$$e_i \text{ (for +ve examples)} = 2 \times 60 / 160 = 0.75$$

$$e_i \text{ (for -ve examples)} = 2 \times 100 / 160 = 1.25$$

$$R = 2 \left[-2 \times \log_2 \left(\frac{2}{0.75} \right) + 0 \times \log_2 \left(\frac{0}{1.25} \right) \right]$$

$$R = 5.66$$

Laplace Measure

$$(\text{for } r_1) = \frac{50+1}{50+2} = 0.894 = 89.4\%$$

$$(\text{for } r_2) = \frac{2+1}{2+2} = 0.75 = 75\%$$

$$P_t = 0.2$$

 r_1

$$\text{FOIL Gain } (R_0, r_1) = 50 \left[\log_2 \left(\frac{50}{50+5} \right) - \log_2 \left(\frac{60}{60+100} \right) \right]$$

modified rule

Original rule

$$\text{FOIL Gain } (r_1, r_2) = 2 \times \left[\log_2 \left(\frac{2}{2+0} \right) - \log_2 \left(\frac{50}{50+5} \right) \right]$$

Q. training set 100 +ve examples
 400 -ve examples

$r_1 \rightarrow$ covers 4 +ve examples & 1 -ve example.

$r_2 \rightarrow$ covers 30 +ve examples & 10 -ve examples.

$$k=2 \quad P_t = 0.2$$

 y_1

$$\text{m-estimate} = \frac{4 + 2 \times 0.2}{5 + 2} = 62.857\%$$

 y_2

$$\text{m-estimate} = \frac{30 + 2 \times 0.2}{40 + 2} = 72.380\%$$

- * RIPPER Algorithm (general to specific)
 - for datasets - with imbalanced class distributions,
 - suitable for noisy data sets
 - uses validation/prune set.

2 class problem

- ↳ select class with a majority $\xrightarrow{\text{Assumed}}$ default class,
- ↳ find rules for class (minor)

multiclass problem

$$y = \{y_1, y_2, \dots, y_n\}$$

$y_i \rightarrow$ less frequent class

$y_c \rightarrow$ default class, /max. freq.

$y_i =$ least majority class,

$y_c =$ max. majority class.

1st iteration

y_i

+ve examples are examples belonging to class y_i

-ve examples are examples not belonging to class y_i

$R_0: C \rightarrow y, \exists \text{ Empty Rule}$

FOIL gain (R_0, r_i)

$r_i \rightarrow \text{Rule Pruning}$

FOIL Prune value = $\frac{p-n}{p+n}$

e.g. $r_i(C_1 \wedge C_2 \wedge C_3) \Rightarrow y,$

FOIL Prune value (r_i)

$r_i'(C_1 \wedge C_2) \Rightarrow y,$

FOIL Prune value (r_i')

add rule to rule set

remove tuples covered by r_i

Stopping Condition

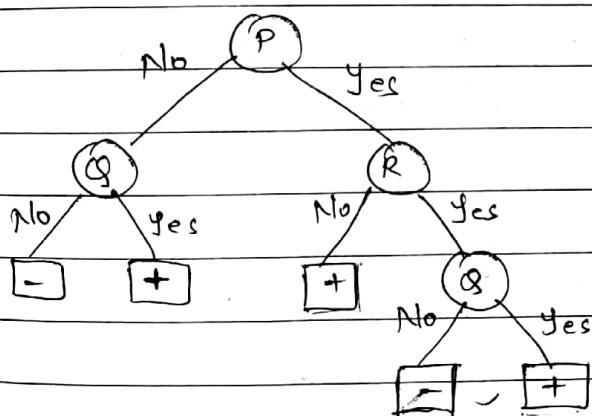
→ d bits ($d = 64$ bits)

→ If error rate $> 50\%$ on Prune set

* consequent = class

Indirect method for rule extraction

eg:



$P = \text{No} \wedge Q = \text{No} \Rightarrow -$

$P = \text{No} \wedge Q = \text{Yes} \Rightarrow +$

$P = \text{Yes} \wedge R = \text{No} \Rightarrow +$

$P = \text{Yes} \wedge R = \text{Yes} \wedge Q = \text{Yes} \Rightarrow +$

$P = \text{Yes} \wedge R = \text{Yes} \wedge Q = \text{No} \Rightarrow -$

$P = \text{No} \wedge Q = \text{No} \Rightarrow -$

$P = \text{Yes} \wedge R = \text{Yes} \wedge Q = \text{No} \Rightarrow -$

$P = \text{No} \wedge Q = \text{Yes} \Rightarrow +$

$P = \text{Yes} \wedge R = \text{No} \Rightarrow +$

$P = \text{Yes} \wedge R = \text{Yes} \wedge Q = \text{Yes} \Rightarrow +$

Reduce
Andysis

$$\begin{aligned}
 \therefore Q = \text{Yes} &\rightarrow + \\
 P = \text{Yes} \wedge R = \text{No} &\rightarrow + \\
 Q = \text{No} &\rightarrow -
 \end{aligned}$$

$\left\{ \text{Not} \right.$

* Characteristics

- The expensiveness of a rule base classifier is equivalent to the decision tree since the decision tree can be represented as a set of mutually exclusive & exhausted rules.
- They are used to produce descriptive models that are easier to interpret but gives the comparison performance to the decision tree classifier.
- The class based ordering approach is suitable for data sets with imbalanced class distribution.

* Nearest neighbour based classification

Rule classifier : class for tuples are assigned to test records

- searches for exact match then assign class.
- drawback! doesn't get exact match

* K-nearest neighbours

- called as Lazy-learning Technique : do not perform training

↳ majority voting

↳ inverse distance weighted approach

X-Axis Durability

$x_1 = \text{Acid Durability}$

7

7

3

1

 $x_2 = \text{Acid Strength}$

7

4

4

4

 $y = \text{Acid Quality}$

Bad

Bad

Good

Good

test record = $(x_1 = 3, x_2 = 7)$

x_1	x_2	y	dist
7	7	Bad	$d = 7-3 + 7-7 = 4$
7	4	Bad	$d = 7-3 + 4-7 = 7$
3	4	Good	$d = 3-3 + 4-7 = 3$
1	4	Good	$d = 1-3 + 4-7 = 5$

 $k=1$ $x_1=3, x_2=7 = \text{Good}$ $k=2$ $x_1=3, x_2=7 = \text{Bad}$ $k=3$ $x_1=3, x_2=7 = \text{Good}$

* Inverse distance weighted voting $y^* = \underset{(x_i, y_i \in D_2)}{\operatorname{argmax}} \sum w_i x_i I(v=y_i)$

for class bad

$$= \frac{1}{(4)^2} \times 1 + \frac{1}{(3)^2} \times 0 = 0.0625$$

for class good

$$= \frac{1}{(4)^2} \times 0 + \frac{1}{(3)^2} \times 1 = 0.1111$$

$$\therefore y^* = \underset{\text{greater}}{\operatorname{argmax}} (0.0625, 0.111) = \text{good}$$

12/10/19

 $k=1, k=3, k=4 \quad x' = 5$

	x	y	distance from x'
1)	0.1	-	$d = 4.9$
2)	0.7	+	4.3
3)	1	+	4
4)	1.6	-	3.4
5)	2.2	+	3
6)	2.5	+	2.5
7)	3.2	-	1.8
8)	3.5	-	1.5
9)	4.1	+	0.9
10)	4.9	+	0.1

 $k=1$ Record 10 is nearest (class = '+')

$$x' = 5 \Rightarrow y = +$$

$$\underline{k=3}$$

$$x' = 5 \Rightarrow y = +$$

Record 8 $\rightarrow -$

$$9 \rightarrow +$$

$$10 \rightarrow +$$
 using majority voting

 $k=4$

$$x' = 5 \Rightarrow y = +$$

Record 7 $\rightarrow -$

$$8 \rightarrow -$$

$$9 \rightarrow +$$

$$10 \rightarrow +$$

inverse distance weighted measure
for '+'

$$= \frac{1}{(1.8)^2} \times 0 + \frac{1}{(1.5)^2} \times 0 + \frac{1}{(0.9)^2} \times 1 + \frac{1}{(0.1)^2} \times 1$$

$$\downarrow \sqrt{}$$

$$= 101.23$$

for 'i'

$$= \frac{1}{(1.8)^2} \times 1 + \frac{1}{(1.5)^2} \times 1 + \frac{1}{(0.9)^2} \times 0 + \frac{1}{(0.1)^2} \times 0 \\ = 5.56$$

$$y^i = \text{argmax} \{ 101.23, 0.730 \}$$

$y^i = +$

Characteristics of k-nearest neighbour.

- K-nearest neighbour is part of a more general technique known as instance based learning which uses specific training instance to make predictions without having to maintain and abstraction (all model) derived from the data
- We do not require model building, hence classifying a test example is an expensive approach because we need to compute the proximity value individually b/w the test & the training example.
- Nearest neighbours make the predictions based on local information whereas decision tree are rule base classifier attempt to find the global model that fix the entire ip space.
- They can produce arbitrarily shape decision boundary.
- It's highly dependent on the proximity measures and can produce wrong predictions unless the correct proximity measures & data free processing steps are taken.

* CLUSTERING : grouping closely related data

- Unsupervised learning technique.

↳ In the context of understanding data clusters are potential classes & cluster analysis is the study of automatically finding the classes.

* Applications of clustering

→ information retrieval

→ climatic data

→ medicine

→ Utility & Summarisation

→ Compression: instead of using entire dataset it uses clusters
- reduce data

→ finding nearest neighbours:

* Types of Clusters

→ Well separated Clusters : are ideal clusters

- inter similarity should be small value

→ prototype based
 centroid : avg. of all points

medoid : one point taken from clusters

→ Graph based clusters : every object is node

→ density based clusters : if points are in high density \rightarrow cluster
: " " " " low \rightarrow outliers

→ shared / conceptual clusters : they cluster together if they have some common property.

* Types of Clustering Techniques:

→ Partitioning : k-means, k-medoid

→ hierarchical

↳ nested clusters



→ exclusive clustering :- single turn clusters

↳ one cluster contain only one record

→ overlapping :-

→ Fuzzy clustering : soft clustering ↳ 40% correct & 60% wrong

↳ membership function

0 - 1

→ Probabilistic Clustering:-

* Kmeans Clustering

1] $k \rightarrow$ no. of clusters

2] randomly choose k points from the data set.

3] Compute the similarities of remaining points from the

↳ k points

1) means are same

2) clustering ...

3) after does it convert

Then stop after certain no. of iterations

} Stopping condition.

Object	x	y
A	1	1
B	2	1
C	4	3
D	5	4

Initial centroids A (1,1) & B (2,1)

No. of clusters $k = 2$

Cluster 1 Cluster 2
distance from dist. from

Centroid A Centroid B

A	0	1
B	1	0
C	3.61	2.83
D	5	4.24

$$\text{dist}(A, B) = \sqrt{(1-2)^2 + (1-1)^2} \\ = 1$$

$$\text{distance}(A, C) = \sqrt{(2-1)^2 + (3-1)^2} \\ = 3.61$$

$$\text{dist}(B, C) = \sqrt{(4-2)^2 + (3-1)^2} \\ = 2.83$$

$$\text{dist}(A, D) = \sqrt{(5-1)^2 + (4-1)^2} \\ = 5$$

$$\text{dist}(B, D) = \sqrt{(5-2)^2 + (4-1)^2} \\ = 4.24$$

cluster 1 (A)

cluster 2 (B, C, D)

cluster mean for cluster 1 - A (1,1) \rightarrow Centroid 1

cluster mean for cluster 2 - $(\frac{2+4+5}{3}, \frac{1+3+4}{3})$

= 3.66, 2.66 \rightarrow Centroid 2

dist from dist. from

Centroid 1 Centroid 2

A 0 3.14

B 1 2.36

C 3.61 0.47

D 5 1.89

$$\text{dist}(A, \text{Centroid 2}) = \sqrt{(C_1 - 3.67)^2 + (1 - 2.67)^2} \\ = 2.07$$

$\text{dist}(B, \text{Centroid 2})$

$$= \sqrt{(2 - 3.67)^2 + (1 - 2.67)^2} \\ = 2.36$$

$\text{dist}(C, \text{Centroid 2})$

$$= \sqrt{(4 - 3.67)^2 + (3 - 2.67)^2} \quad \text{dist} \\ = 3.61$$

cluster 1 A, B

cluster 2 C, D

$$\text{cluster 1 mean} = \left(\frac{1+2}{2}, \frac{1+1}{2} \right) = (1.5, 1) \rightarrow \text{Centroid 1}$$

$$\text{cluster 2 mean} = \left(\frac{4+5}{2}, \frac{3+4}{2} \right) = (4.5, 3.5) \rightarrow \text{Centroid 2}$$

	Centroid 1	Centroid 2	$\text{dist}(A, \text{Centroid 1})$
A	0.5	4.50	$= \sqrt{(1 - 1.5)^2 + (1 - 1)^2}$
B	0.5	3.5	$= 0.5$
C	3.20	0.71	
D	4.61	0.71	

$$\text{dist}(A, 2) = \sqrt{(1 - 4.5)^2 + (1 - 3.5)^2} = 0.5$$

$$\text{dist}(C, 1) = \sqrt{(4 - 1.5)^2 + (3 - 1)^2} = 3.20$$

$$\text{dist}(D, 1) = \sqrt{(5 - 1.5)^2 + (4 - 1)^2} = 4.609$$

$$\text{dist}(D, 2) = \sqrt{(5 - 4.5)^2 + (4 - 3.5)^2} = 0.71$$

cluster 1 mean = (1.5, 1)
 cluster 2 mean = (4.5, 3.5)

Same clusters so stop

Q.

Object	x	y
1	20	10
2	30	20
3	30	30
4	35	35
5	40	40
6	50	45

Object 2 (30, 20) and
 Object 5 (40, 50)
 → initial centroids.

	Centroid 1	Centroid 2	
1	14.14	44.72	$\text{dist}(1, 1) = \sqrt{(20-30)^2 + (10-20)^2}$
2	0	31.62	= 14.14
3	10	22.36	$\text{dist}(1, \text{Cent}2) = \sqrt{(20-40)^2 + (10-50)^2}$
4	15.81	15.81	= 44.72
5	22.36	22.36	$\text{dist}(2, 1) = \sqrt{(30-30)^2 + (20-20)^2}$
6	32.01	11.18	= 0

$$\text{dist}(2, 2) = \sqrt{(30-40)^2 + (20-50)^2} = 31.62$$

$$\text{dist}(3, 1) = \sqrt{(30-30)^2 + (30-20)^2} = 10$$

$$\text{dist}(3, 2) = \sqrt{(30-40)^2 + (30-50)^2} = 22.36$$

$$\text{dist}(4, 1) = \sqrt{(35-30)^2 + (35-20)^2} = 15.81$$

$$\text{dist}(4, 2) = \sqrt{(35-40)^2 + (35-50)^2} = 15.81$$

$$\text{dist}(5, 1) = \sqrt{(40-30)^2 + (40-20)^2} = 22.36$$

$$\text{dist}(5, 2) = \sqrt{(40-40)^2 + (40-50)^2} = 10$$

$$\text{dist}(6, 1) = \sqrt{(50-30)^2 + (45-20)^2} = 32.01$$

$$\text{dist}(6, 2) = \sqrt{(50-40)^2 + (45-50)^2} = 11.18$$

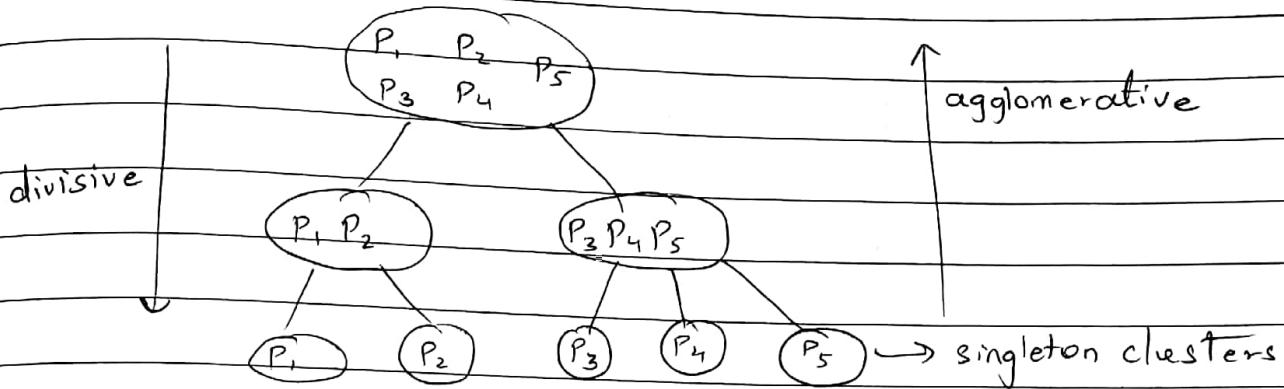
Cluster 1 (1, 2, 3, 4)

Cluster 2 (5, 6)

Hierarchical Clustering

- ↳ creation of nested clusters
- ↳ divisive
- ↳ agglomerative (AGNES) clustering

Co Dataset

Agglomerative ClusteringProximity

- ↳ single linkage (min)
- ↳ complete linkage (max)
- ↳ group average
- ↳ Ward's distance

Q.

Single Linkage (min)

Point	x	y
P ₁	0.40	0.53
P ₂	0.22	0.38
P ₃	0.35	0.32
P ₄	0.26	0.19
P ₅	0.08	0.41
P ₆	0.45	0.30

	P_1	P_2	P_3	P_4	P_5	P_6	
P_1	0	0.234	0.215	0.367	0.341	0.235	
P_2	0.234	0	0.143	0.194	0.143	0.243	
P_3	0.215	0.143	0	0.158	0.284	0.101	
P_4	0.367	0.194	0.158	0	0.284	0.219	
P_5	0.341	0.143	0.284	0.284	0	0.386	
P_6	0.235	0.243	0.101	0.219	0.386	0	

$$(P_1, P_2) = \sqrt{(0.40 - 0.22)^2 + (0.53 - 0.38)^2} = 0.234$$

$$(P_1, P_3) = 0.2158$$

$$(P_1, P_4) = \sqrt{(0.40 - 0.26)^2 + (0.53 - 0.19)^2} = 0.367$$

$$(P_1, P_5) = \sqrt{(0.40 - 0.08)^2 + (0.53 - 0.41)^2} = 0.341$$

$$(P_1, P_6) = \sqrt{(0.40 - 0.45)^2 + (0.53 - 0.30)^2} = 0.235$$

$$(P_2, P_3) = \sqrt{(0.22 - 0.35)^2 + (0.38 - 0.32)^2} = 0.143$$

$$(P_2, P_4) = \sqrt{(0.22 - 0.26)^2 + (0.38 - 0.19)^2} = 0.194$$

$$(P_2, P_5) = \sqrt{(0.22 - 0.08)^2 + (0.38 - 0.41)^2} = 0.143$$

✓ min dist betn P_3 and P_6

→ merge P_3, P_6

	P_1	P_2	$P_{3,6}$	P_4	P_5	
P_1	0	0.234	0.215	0.367	0.341	
P_2	0.234	0	0.143	0.194	0.143	
$P_{3,6}$	0.215	0.143	0	0.158	0.284	
P_4	0.367	0.194	0.158	0	0.284	
P_5	0.341	0.143	0.284	0.284	0	

$$\text{dist}(P_{3,6}, P_1) = \min(d_{3,1}, d_{6,1})$$

$$= \min(0.215, 0.235)$$

$$= 0.215$$

$$\text{dist}(P_{3,6}, P_2) = \min(d_{3,2}, d_{6,2})$$

$$= \min(0.143, 0.243)$$

$$= 0.143$$

$$\text{dist}(P_{3,6}, P_4) = \min(d_{3,4}, d_{6,4}) \\ = \min(0.158, 0.219) \\ = 0.158$$

$$\text{dist}(P_{3,6}, P_5) = \min(d_{3,5}, d_{6,5}) \\ = \min(0.284, 0.386) \\ = 0.284$$

min dist b/w P_2 & P_5 ,
merge P_2 & P_5

	P_1	$P_{2,5}$	$P_{3,6}$	P_4
P_1	0	0.234	0.215	0.367
$P_{2,5}$	0.234	0	0.143	0.194
$P_{3,6}$	0.215	0.143	0	0.158
P_4	0.367	0.194	0.158	0

$$\text{dist}(P_{2,5}, P_1) = \min(d_{2,1}, d_{5,1}) \\ = \min(0.234, 0.341) \\ = 0.234$$

$$\text{dist}(P_{2,5}, P_{3,6}) = \min(d_{2,3}, d_{2,6}, d_{5,3}, d_{5,6}) \checkmark_{\text{correct}} \\ = \min(0.143, 0.284) \\ = 0.143$$

$$\text{dist}(P_{2,5}, P_4) = \min(d_{2,4}, d_{5,4}) \\ = \min(0.194, 0.284) \\ = 0.194$$

min dis b/w $P_{2,5}$ & $P_{3,6}$
merge $P_{2,5}$ & $P_{3,6}$

	P_1	$P_{2,3,5,6}$	P_4
P_1	0	0.215	0.367
$P_{2,3,5,6}$	0.215	0	0.158
P_4	0.367	0.158	0

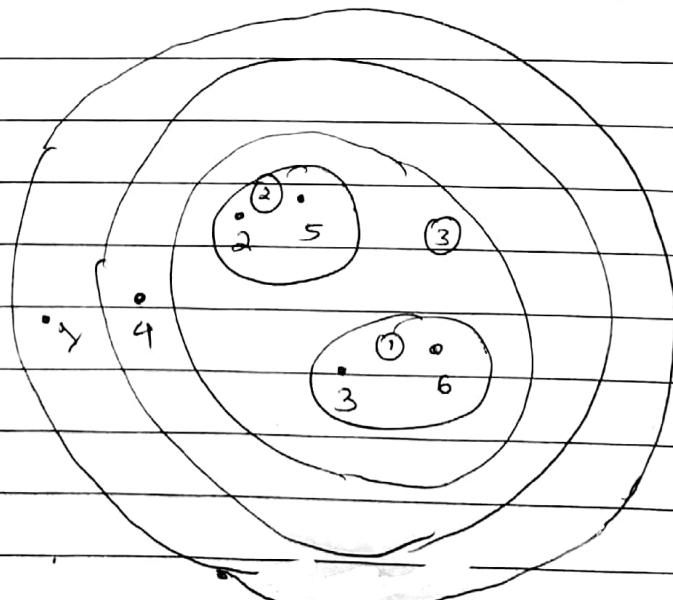
$$\begin{aligned} \text{dist}(P_{2,3,5,6}, P_1) &= \min(d_{2,1}, d_{3,1}, d_{5,1}, d_{6,1}) \\ &= \min(0.234, 0.215, 0.341, 0.235) \\ &= 0.215 \end{aligned}$$

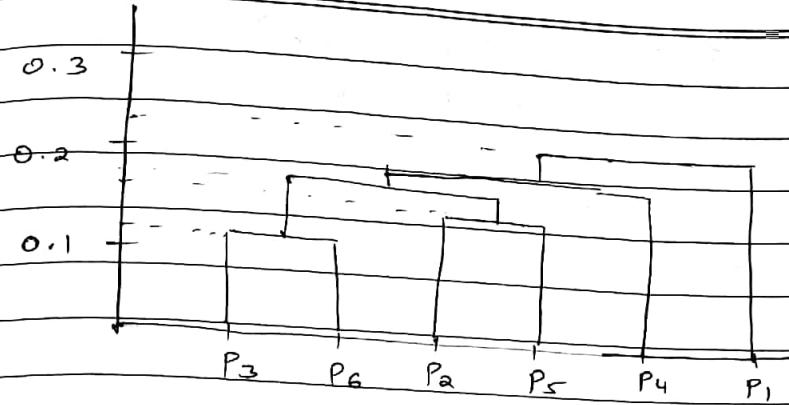
$$\begin{aligned} \text{dist}(P_{2,3,5,6}, P_4) &= \min(d_{2,4}, d_{3,4}, d_{5,4}, d_{6,4}) \\ &= \min(0.194, 0.158, 0.284, 0.219) \\ &= 0.158 \end{aligned}$$

min dist b/w $P_{2,3,5,6}$ & P_4

	P_1	$P_{2,3,4,5,6}$
P_1	0	0.215
$P_{2,3,4,5,6}$	0.215	0

$$\begin{aligned} \text{dist}(P_{2,3,4,5,6}, P_1) &= \min(d_{2,1}, d_{3,1}, d_{4,1}, d_{5,1}, d_{6,1}) \\ &= \min(0.234, 0.215, 0.367, 0.341, 0.235) \\ &= 0.215 \end{aligned}$$





Complete Linkage Problem

Point	x	y
P ₁	0.40	0.53
P ₂	0.22	0.38
P ₃	0.35	0.32
P ₄	0.26	0.19
P ₅	0.08	0.41
P ₆	0.45	0.30

soln	P ₁	P ₂	P ₃	P ₄	P ₅	P ₆
P ₁	0	0.24	0.22	0.37	0.34	0.23
P ₂	0.24	0	0.15	0.20	0.14	0.25
P ₃	0.22	0.15	0	0.15	0.28	0.11
P ₄	0.37	0.20	0.15	0	0.29	0.22
P ₅	0.34	0.14	0.28	0.29	0	0.39
P ₆	0.23	0.25	0.11	0.22	0.39	0

min dist P₃ & P₆

→ merge P₃ & P₆

	P ₁	P ₂	P _{3,6}	P ₄	P ₅
P ₁	0	0.24	0.23	0.37	0.34
P ₂	0.24	0	0.25	0.20	0.14
P _{3,6}	0.23	0.25	0	0.22	0.39
P ₄	0.37	0.20	0.22	0	0.29
P ₅	0.34	0.14	0.39	0.29	0

$$\begin{aligned}\text{dist}(P_1, P_{3,6}) &= \max(d_{1,3}, d_{1,6}) \\ &= \max(0.22, 0.23) \\ &= 0.23\end{aligned}$$

$$\begin{aligned}\text{dist}(P_2, P_{3,6}) &= \max(d_{2,3}, d_{2,6}) \\ &= \max(0.15, 0.25) \\ &= 0.25\end{aligned}$$

$$\begin{aligned}\text{dist}(P_4, P_{3,6}) &= \max(d_{4,3}, d_{4,6}) \\ &= 0.22\end{aligned}$$

$$\begin{aligned}\text{dist}(P_5, P_{3,6}) &= \max(d_{5,3}, d_{5,6}) \\ &= 0.39\end{aligned}$$

✓ min dist P_2 and P_5

Merge P_2 & P_5

	P_1	$P_{2,5}$	$P_{3,6}$	P_4
P_1	0	0.34	0.23	0.37
$P_{2,5}$	0.34	0	0.39	0.29
$P_{3,6}$	0.23	0.39	0	0.22
P_4	0.37	0.29	0.22	0

$$\begin{aligned}\text{dist}(P_1, P_{2,5}) &= \max(d_{1,2}, d_{1,5}) \\ &= 0.34\end{aligned}$$

$$\begin{aligned}\text{dist}(P_{2,5}, P_{3,6}) &= \max(d_{2,3,6}, d_{5,3,6}) \\ &= 0.39\end{aligned}$$

$$\begin{aligned}\text{dist}(P_4, P_{2,5}) &= \max(d_{4,2}, d_{4,5}) \\ &= 0.29\end{aligned}$$

✓ min dist P_4 & $P_{3,6}$

merge P_4 & $P_{3,6}$

	P_1	$P_{2,5}$	$P_{3,4,6}$	
P_1	0	0.34		
$P_{2,5}$	0.34	0	0.37	
$P_{3,4,6}$	0.37	0.39	0	

$$\text{dist}(P_1, P_{3,4,6}) = \max(d_{1,3}, d_{1,4}, d_{1,6}) \\ = 0.37$$

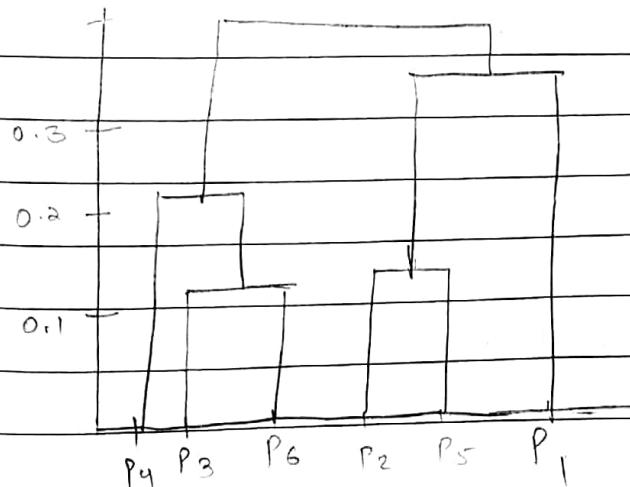
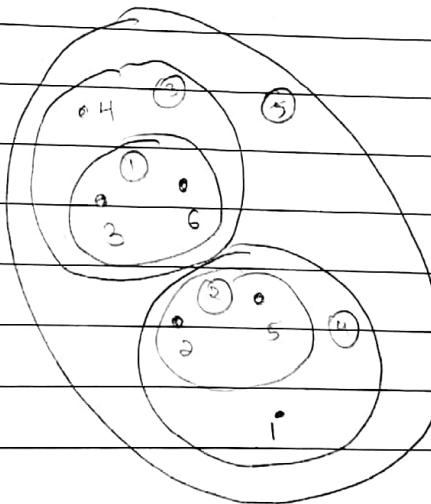
$$\text{dist}(P_{2,5}, P_{3,4,6}) = \max(d_{2,3}, d_{2,4}, d_{2,6}) \\ = 0.39$$

min dist P_1 and $P_{2,5}$

merge P_1 and $P_{2,5}$

	$P_{1,2,5}$	$P_{3,4,6}$
$P_{1,2,5}$	0	0.39
$P_{3,4,6}$	0.39	0

$$\text{dist}(P_{1,2,5}, P_{3,4,6}) = 0.39$$



* Group average method
 $\text{proximity } (C_i, C_j) = \frac{\sum_{x \in C_i} \sum_{y \in C_j} \text{proximity}(x, y)}{m_i \times m_j}$

m_i = size of cluster C_i

m_j = size of cluster C_j

Q.	Point	x	y
=	P ₁	0.40	0.53
	P ₂	0.22	0.38
	P ₃	0.35	0.32
	P ₄	0.26	0.19
	P ₅	0.08	0.41
	P ₆	0.45	0.36

	P ₁	P ₂	P ₃	P ₄	P ₅	P ₆
P ₁	0	0.24	0.22	0.37	0.34	0.23
P ₂	0.24	0	0.15	0.20	0.14	0.25
P ₃	0.22	0.15	0	0.15	0.28	0.11
P ₄	0.37	0.20	0.15	0	0.29	0.22
P ₅	0.34	0.14	0.28	0.29	0	0.39
P ₆	0.23	0.25	0.11	0.22	0.39	0

min dist $P_3 + P_6 \rightarrow$ merge $P_3 \& P_6$

	P ₁	P ₂	P _{3,6}	P ₄	P ₅
P ₁	0	0.24	0.225	0.37	0.34
P ₂	0.24	0	0.2	0.20	0.14
P _{3,6}	0.225	0.2	0	0.185	0.335
P ₄	0.37	0.20	0.185	0	0.29
P ₅	0.34	0.14	0.335	0.29	0

$$\text{dist}(P_1, P_{3,6}) = (0.22 + 0.23) / (1 \times 2) = 0.225$$

$$\text{dist}(P_2, P_{3,6}) = (0.15 + 0.25) / (1 \times 2) = 0.2$$

$$\text{dist}(P_4, P_{3,6}) = (0.15 + 0.22) / (1 \times 2) = 0.185$$

$$\text{dist}(P_5, P_{3,6}) = (0.28 + 0.39) / (1 \times 2) = 0.335$$

	min dist	bln	$P_2 \rightarrow P_5$	$\rightarrow \text{merge } P_2 \text{ & } P_5$
P_1	0	$P_{2,5}$	$P_{3,6}$	P_4
P_1	0	0.295	0.225	0.37
$P_{2,5}$	0.295	0		
$P_{3,6}$	0.225		0	0.185
P_4	0.37		0.185	0

$$\text{dist}(P_1, P_{2,5}) = (0.25 + 0.34) / (1 \times 2) = 0.295$$

10/10/14

Absentee Losses)

Hierarchical Clustering
 ↳ Ward's distance

SSE → sum of squared errors.

$$\text{Cluster SSE} = \sum_{i=1}^n (x_i - \bar{x})^2$$

\bar{x} → centroid of the cluster.

$$\text{Total SSE} = \sum_{j=1}^k \sum_{x_i \in C_j} (x_i - \bar{x})^2$$

k = no. of clusters

Increase in SSE - I_{AB}

$$I_{AB} = SSE_{AB} - (SSE_A + SSE_B)$$

$$SSE_A = \sum_{l=1}^{n_A} (a_i - \bar{a})^2$$

$$SSE_B = \sum_{l=1}^{n_B} (b_i - \bar{b})^2$$

$$SSE_{AB} = \sum_{l=1}^{n_{AB}} (y_i - \bar{y}_{AB})^2$$

a_i - data object belonging to cluster A

b_i - data object belonging to cluster B

y_i - data object belonging to cluster AB

\bar{a} - centroid of cluster A

\bar{b} - centroid of cluster B

\bar{y}_{AB} - centroid of cluster AB

AB = avg

	x	y	
p ₁	0.4	0.53	
p ₂	0.22	0.38	
p ₃	0.35	0.32	
p ₄	0.26	0.19	
p ₅	0.08	0.41	
p ₆	0.45	0.3	

Wards distance :

DBSCAN Problem

	x	y	
P ₁	2	2	min p ₁₃ = 3
P ₂	4	4	$\Sigma = 3$
P ₃	6	6	
P ₄	0	6	
P ₅	6	0	
P ₆	5	5	
P ₇	7	7	
P ₈	9	9	

	P ₁	P ₂	P ₃	P ₄	P ₅	P ₆	P ₇	P ₈
P ₁	0	2.83	5.65	4.47	4.47	4.24	7.07	9.90
P ₂	2.83	0	2.83	4.47	4.47	1.41	4.24	7.07
P ₃	5.65	2.83	0	6.00	6.00	1.41	1.41	4.24
P ₄	4.47	4.47	6.0	0	8.49	5.10	7.07	9.49
P ₅	4.47	4.47	6.0	8.49	0		7.07	9.49
P ₆	4.24	1.41	1.41	5.10	5.10	0	2.83	5.66
P ₇	7.07	4.24	1.41	7.07	7.07	2.83	0	2.83
P ₈	9.90	7.07	4.24	9.49	9.49	5.66	2.83	0

classmate
Date _____
Page _____

core point should have less no. of pts < min pts

	ϵ -neighbourhood	no. of points
P ₁ -	P ₁ , P ₂	2
P ₂ -	P ₁ , P ₂ , P ₃ , P ₆	4
P ₃ -	P ₂ , P ₃ , P ₆ , P ₇	4
P ₄ -	-	-
P ₅ -	-	-
P ₆ -	P ₂ , P ₃ , P ₆ , P ₇	4
P ₇ -	P ₃ , P ₇ , P ₆ , P ₈	4
P ₈ -	P ₈ , P ₇	2

P₄ & P₅ - noise points (neither core/border points)

P₁ → border point (< min pts) & lies in ϵ -neighbourhood of P₂

P₈ → border point (lies in ϵ -neighbourhood of P₇) & (< min pts)

P₁ is density reachable from P₂ (lies in ϵ -neighbourhood of P₂ (corepoint))

P₃ is density reachable from P₂ (lies in ϵ -neighbourhood of P₂ (corepoint))

P₆ is density reachable from P₂

P₂ is density reachable from P₃

P₇ is density reachable from P₃

hence P₂ & P₇ are density connected

P₃ → density reachable from P₇

P₈ → " " " P₇

P₃ & P₈ are density connected

So P₁, P₂, P₃, P₆, P₇, P₈ lie in the same cluster

→ P₄, P₅ → outliers/noisepts

Notes

Meas

- Q. Differentiate b/w Eigen learning & Lazy learning. (or)
Why is k-nearest neighbour called as a lazy classifier.
Compare its performance to a decision tree classifier
→ Eager learner:
Always build model.
Decision tree is better

- Q. Explain the different rule order scheme & list its advantages & disadvantages.
- Q. Explain the diff. voting strategies employed in k-nearest neighbours algorithm & illustrate each strategy.
→ majority & distance weighted approach.
- Q. Explain the learn one rule function with the suitable example.
- Q. Explain rule induction using sequential covering algorithm.
- Q. Define accuracy & coverage of the rule with respect to rule base classifier.
- Q. What is rule pruning, How it is carried out.
- Q. Explain the RIPPLE algorithm :- it is a direct algorithm.
- Q. Explain the diff. types of clusters & explain the diff. clustering techniques.
- Q. Explain k-means algorithm / static k-means.
- Q. State the DBSCAN algorithm.
- Definitions related to density based clustering:
density, reachability, density connected, core points,
Noise points

Clustering

Selecting centroids

~~cluster SSE~~ (sum of squared errors)

- Ideally, Euclidean distance is used to find distance of each point from its centroid.

Cluster SSE

$$\text{cluster SSE} = \sum_{i=1}^n (x_i - c_0)^2 \quad (c_0 \rightarrow \text{centroid of the cluster})$$

$x_i \in \text{cluster}$.

Total SSE

- Sum of all cluster SSEs.

$$\text{Total SSE} = \sum_{i=1}^k \sum_{j=1}^n (x_j - c_i')^2 \quad (c_i' \rightarrow \text{centroid of cluster } i)$$

$x_j \in c_i'$

for document data

$$\text{Total cohesion} = \sum_{i=1}^k \left(\sum_{j=1}^n \cosine(x_j, c_i') \right) \quad \begin{matrix} \text{cosine} \\ \text{similarity} \end{matrix} \quad \begin{matrix} \text{cluster} \\ \text{cohesion} \end{matrix}$$

$x_j \in c_i$

Issues with K-means clustering

- handling empty clusters

→ Outliers

use outlier detection.

Postprocessing with SSE

To minimise SSE

→ reduce no. of clusters

↳ disperse a cluster

↳ select a cluster with least SSE.

Redistribute points to remaining ~~points~~ clusters

↳ merge two clusters

↳ merge two clusters having the least SSE.

→ split the clusters

↳ split clusters having highest SSE.

Incrementally updating centroid



Strengths and Weaknesses

Lance William Formula

$$p(R, \Theta) = \alpha_A p(A, \Theta) + \alpha_B p(B, \Theta) + \beta p(A, B) + \gamma |p(A, \Theta) - p(B, \Theta)|$$

	α_A	α_B	β	γ
single linkage	$\frac{1}{2}$	$\frac{1}{2}$	0	$-\gamma_2$
complete linkage	$\frac{1}{2}$	$\frac{1}{2}$	0	$+\frac{1}{2}$
group average	$\frac{m_A}{m_A + m_B}$	$\frac{m_B}{m_A + m_B}$	0	0
Ward's	$\frac{m_A + m_B}{m_A + m_B + m_C}$	$\frac{m_B + m_C}{m_A + m_B + m_C}$	$-\frac{m_A}{m_A + m_B + m_C}$	$\frac{m_A + m_B + m_C}{m_A + m_B + m_C}$

- lack of global objective function.
- handling varying cluster size.
 - weighted approach
 - unweighted approach
- merge decisions are final. (can't reassign point to another cluster only merged)
- outliers

single / complete linkage → last iteration maybe eliminated as it is an outlier.

Drawback of hierarchical clustering is that, it is ~~strength~~ - computationally expensive because it has to find distance between one point to another.

DBSCAN

Strength

- Detect outliers

Weakness

- High dimensional data can't be handled.
- Cannot handle data of different density.
- Computationally expensive as we have to find distance between all points.

K-medoid clustering

dataset

$K=2$

		dist from medoid 1	dist from medoid 2
2	2	0 ✓	2
4	4	2	0
10	10	8	6
12	12	10	8
3	3	1 ✓	3
20	20	18	16
30	30	28	26
11	11	9	7
25	25	23	21

PT

$$\text{cluster 1} \quad 2 \quad 0 \quad \text{absolute error criterion}$$

for medoid 1 → $3 \quad 1 = \sum_{i=1}^K \sum_{j=1}^n |x_j - c_i'|$

$$\text{cluster 2} \quad 4 \quad 0$$

$$\text{for medoid 2} \rightarrow 10 \quad 6$$

$$12 \quad 8 = |2-2| + |2-3|$$

$$20 \quad 16 + |4-4| + |4-10| + |4-12| +$$

$$30 \quad 26 + |4-20| + |4-30| + |4-11|$$

$$11 \quad 7 + |4-25|$$

$$25 \quad 21 = 85$$

old medoids $2, 4$

new medoids $2, 20$

$m_1 \quad m_2$

medoid 4 is replaced

Absent - $\alpha = 1/10$
28/10

clustm1 clustm2

	0	18	cluster 1	2
4	2	-	16	4
10	8 ✓	10		10
12	10	8		3
3	1	17		11
20	18	0	cluster 2	12
30	20	10		20
11	9 ✓	9		30
25	23	5		25

pt dist

cluster 1. 2 0 absolute error criterion = 13
(new cost).

4 2

10 8

3 1

11 9

cluster 2 12 8

20 0

30 10

25 5

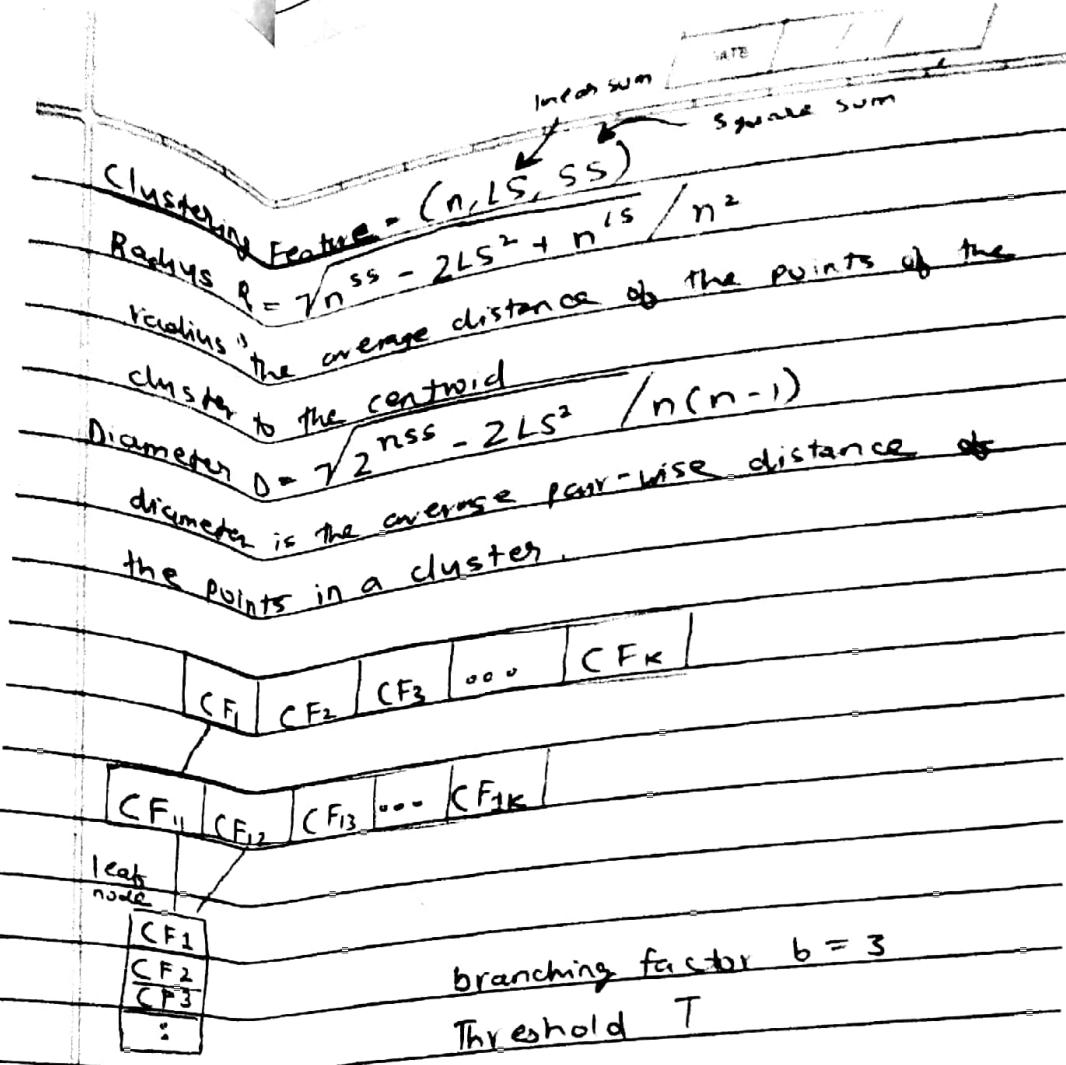
10/19

③ Hierarchical Clustering

BIRCH (Balanced Iterative Reducing Clustering using hierarchies).

Phase 1: An initial Clustering Feature tree is built (CF tree)

Phase 2: A conventional clustering algorithm such as k-means or k-medoid is applied to cluster the leaf nodes of the CF tree.



	A	B	C
cluster1	(2,5)	(3,2)	(4,3)
cluster2	(7,3)	(9,10)	(3,8)

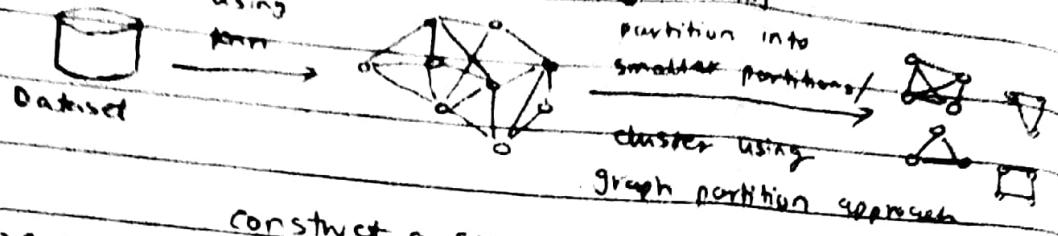
$$\begin{aligned}
 CF_1(\text{cluster } 1) &= (3, (2+3+4, 5+2+3), (2^2+3^2+4^2, \\
 &\quad 5^2+2^2+3^2)) \\
 &= (3, (9, 10), (29, 38))
 \end{aligned}$$

$$\begin{aligned}
 CF_2(\text{cluster } 2) &= (4, (7+9+3+2, 3+10+8+5), \\
 &\quad (7^2+9^2+3^2+2^2, 3^2+10^2+8^2+5^2)) \\
 &= (4, (21, 24), (143, 182))
 \end{aligned}$$

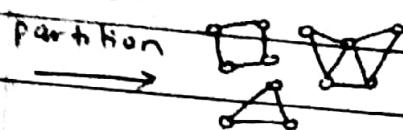
$$\begin{aligned}
 CF_1 + CF_2 &= (n_1+n_2, LS_1+LS_2, SS_1+SS_2) \\
 &= (3+4, (9, 10)+(21, 24), (29, 38)+(143, 182))
 \end{aligned}$$

$$\star (7, (30, 34), (194, 220))$$

Chameleon \rightarrow hierarchical algorithm



construct a sparse graph



Relative interconnectivity (RI)

$$RIC(C_i, C_j) = \frac{EC_{[C_i - C_j]}}{\frac{1}{2}(|EC_{C_i}| + |EC_{C_j}|)}$$

EC edge cut

Relative interconnectivity is defined as the absolute interconnectivity between C_i and C_j normalised w.r.t the internal connectivity between C_i and C_j

$EC_{[C_i - C_j]}$ \rightarrow absolute connectivity

\hookrightarrow edge cut for cluster containing $(C_i \text{ and } C_j)$

Absolute connectivity

Sum of the edges that connect vertices in C_i and C_j

EC_{C_i} \rightarrow internal connectivity

\hookrightarrow edge cut of min bisector

Internal connectivity \rightarrow sum of the edges that roughly partition cluster C_i and C_j into two roughly equal partitions.

Relative Closeness (RC)

Relative Closeness is defined as the absolute closeness between C_i and C_j normalised w.r.t the internal closeness between C_i and C_i .

$$RC = \frac{\bar{SEC}_{C_i \rightarrow C_j}}{\left(\frac{|C_i|}{|C_i| + |C_j|} \right) \bar{SEC}_i + \left(\frac{|C_j|}{|C_i| + |C_j|} \right) \bar{SEC}_j}$$

$\bar{SEC}_{C_i \rightarrow C_j}$ → average weight of the edges of the vertices that connect C_i to C_j

\bar{SEC}_i : → average weight of edges that belong to the min bisector of C_i and C_j .

Single linkage

	X	Y	Proximity Matrix					
P ₁	1	1	P ₁	P ₂	P ₃	P ₄	P ₅	P ₆
P ₂	1.5	1.5		0	0.71			
P ₃	5	5	P ₂	0	0.71	5.66	3.61	4.24
P ₄	3	4	P ₃		0	4.95	2.92	3.20
P ₅	4	4	P ₄	5.66	4.95	0	2.24	2.50
P ₆	3	3.5	P ₅	3.61	2.92	2.24	0	1.41
			P ₆	4.24	3.54	1.41	1.00	0.50
				3.20	2.50	2.50	0.50	1.12
								0

min dist b/w P₄ & P₆⇒ merge P₄ & P₆

	P ₁	P ₂	P ₃	P ₄ , P ₆	P ₅
P ₁	0	0.71	5.66	3.20	4.24
P ₂	0.71	0	4.95	2.50	3.54
P ₃	5.66	4.95	0	2.24	1.41
P ₄ , P ₆	3.20	2.50	2.24	0	1.00
P ₅	4.24	3.54	1.41	1.00	0

$$\text{dist}(P_1, P_4, P_6) = \min(d_{1,4}, d_{1,6})$$

$$= \min(3.61, 3.20) = 3.20$$

min dist b/w P₁ & P₂⇒ merge P₁ & P₂

	P ₁ , P ₂	P ₃	P ₄ , P ₆	P ₅
P ₁ , P ₂	0	4.95	2.50	3.54
P ₃	4.95	0	2.24	1.41
P ₄ , P ₆	2.50	2.24	0	1.00
P ₅	3.54	1.41	1.00	0

min dist b/w P₅ & P₄, P₆⇒ merge P₄, P₅, P₆

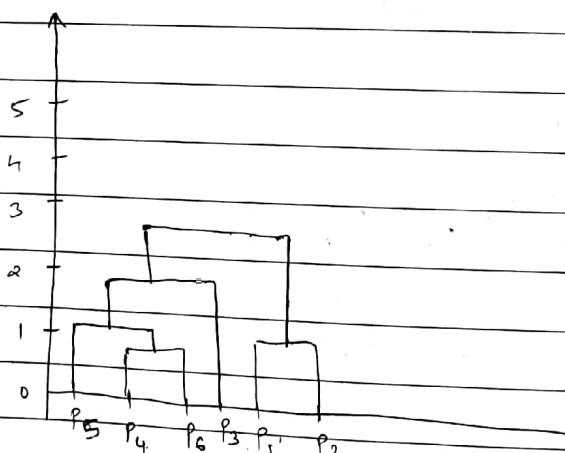
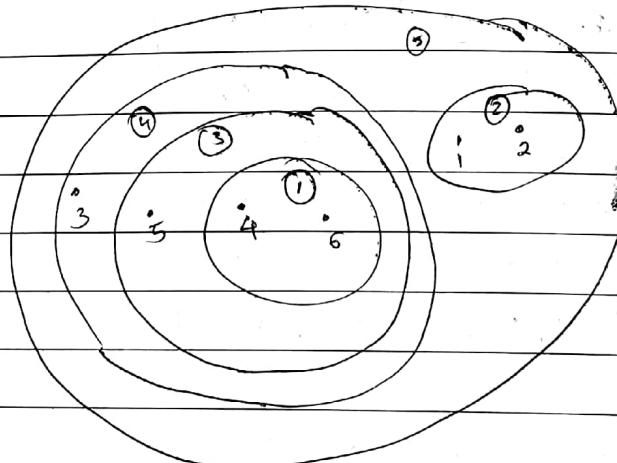
	P_1P_2	P_3	$P_4P_5P_6$	
P_1P_2	0	4.95	2.50	
P_3	4.95	0	1.41	
$P_4P_5P_6$	2.50	1.41	0	

~~min dist b/w P_3 & $P_4P_5P_6$~~

⇒ merge P_3 & $P_4P_5P_6$

$$\text{dist}(P_1P_2, P_3P_4P_5P_6) = 2.50$$

	P_1P_2	$P_3P_4P_5P_6$	
P_1P_2	0	2.50	
$P_3P_4P_5P_6$	2.50	0	



DBSCAN

$\epsilon = 3$

min. points = 3

Note: circle all values which have diff less than equal to ϵ

Meow!



Meow

Meooooow

FOOT
classmate

Date: ii/ii
Page: ii/iv

Consider same proximity matrix.

Dorae-moo

P₂, P₄: P₆

P₂ & P₃: P₄

ϵ -neighbourhood

no. of points

Type of pt

P ₁	P ₁ , P ₂	2	border point
P ₂	P ₁ , P ₂ , P ₄ , P ₆	4	Core pts
P ₃	P ₃ , P ₄ , P ₅ , P ₆	4	"
P ₄	P ₂ , P ₃ , P ₄ , P ₅ , P ₆	5	"
P ₅	P ₃ , P ₄ , P ₅ , P ₆	4	"
P ₆	P ₂ , P ₃ , P ₄ , P ₅ , P ₆	5	"

P₂ & P₄ are density reachable from P₄, hence P₂ & P₃ are density connected

P₂ & P₄ are density " " " P₆, " P₂ & P₄ " " "

P₂ & P₅ " " " " P₄ or P₆, " P₂ & P₅ " " "

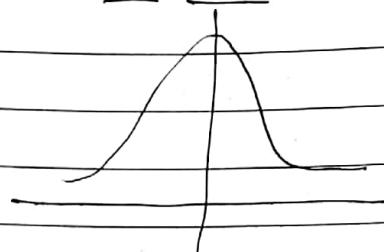
P₂ & P₆ " " " " P₄, " P₂ & P₆ " " "

o10a1q1

Statistical based outlier Detection

with

If the value of $|x| \geq c$ such that probability of $(|x|) \geq c = \alpha$, where c is a constant satisfies that condition, and α is a probability that a value for the given distribution may be mistakenly classified as outlier.

Multivariate Data

$$\text{Mahalanobis dist } (x, \bar{x}) = (x - \bar{x}) S^{-1} (x - \bar{x})^T$$

Multivariate data

↳ mixture model

* Likelihood / log likelihood approach

Likelihood is a measure of the extend to which the sample provides support for particular values of a parameter in parametric mode.

* Likelihood-based outlier detection

- 1] At time $t=0$, let $M_t \rightarrow$ contain all objects
 $A_t \rightarrow$ empty

$$LL_t(D) = LL(M_t) + LL(A_t) \rightarrow \log \text{likelihood of data set.}$$

- 2] for each point x that belongs to M_t do

- 3] Move x from M_t to A_t to produce new Datasets

$$A_{t+1} \beta M_{t+1}$$

- 4] Compute new loglikelihood of D $LL_{t+1}(D) = LL(A_{t+1}) + LL(M_{t+1})$
- 5] Compute difference $\Delta = LL_t(D) - LL_{t+1}(D)$
- 6) IF $\Delta > c$; c -threshold
- 7) x is an anomaly $\rightarrow M_{t+1} \& A_{t+1} \rightarrow$ become current data sets.
- 8) end if
- 9) end for.

Note:- It is suitable for data containing more than one distribution.

* Strengths & Weakness

\rightarrow not efficient for multivariate data / data set with mixed distribution.

* Proximity based Outlier Detection

$$\hookrightarrow O(m^2)$$

It checks whether an object is outlier based on the distance to k-nearest neighbour.

* Density based outlier detection

An outlier score of an object is the inverse of the density around the object.

Density of the given object is a reciprocal of the average distance to the k-nearest neighbours. & is given as

$$\text{density}(x, k) = \left(\frac{\sum_{y \in N(x, k)} \text{dist}(x, y)}{|N(x, k)|} \right)^{-1}$$

Mean! Mean Mean

$$\rightarrow \frac{\text{Average relative density } (\alpha, k)}{= \frac{\text{density } (\alpha, k)}{\sum_{y \in N(\alpha, k)} \text{density}(y, k) / |N(\alpha, k)|}}$$

UNIT 4:

Association Mining: How two strongly the attributes are related to each other.

eg: Market Basket Analysis

List of Customer transaction (D)

↳ Transaction T

I - set of items

$$T \subseteq I$$

Support : % of transaction containing a given itemset

Confidence is the probability that in given item set if item A is brought then item B is also brought. It verify the degree of certainty for given association A → B
Frequent itemset: frequency occurs in the itemset.

Eg. 10 transaction

40%

= 40% of 10 transaction

= 4

Support → min support

Support count is occurrence frequency of a given item set in a list of transaction.

$$\text{Support} = P(A \cup B)$$

$$\text{Confidence} = P(B|A) = \frac{\text{Sup-count}(A \cup B)}{\text{Sup-count}(A)}$$

$$A \rightarrow B$$

↳ min confidence.

① Generate frequent element - min support

② Generate strong association rules - min conflicts

Frequent subsequence \rightarrow common sequence
occurring find the common subsequence.

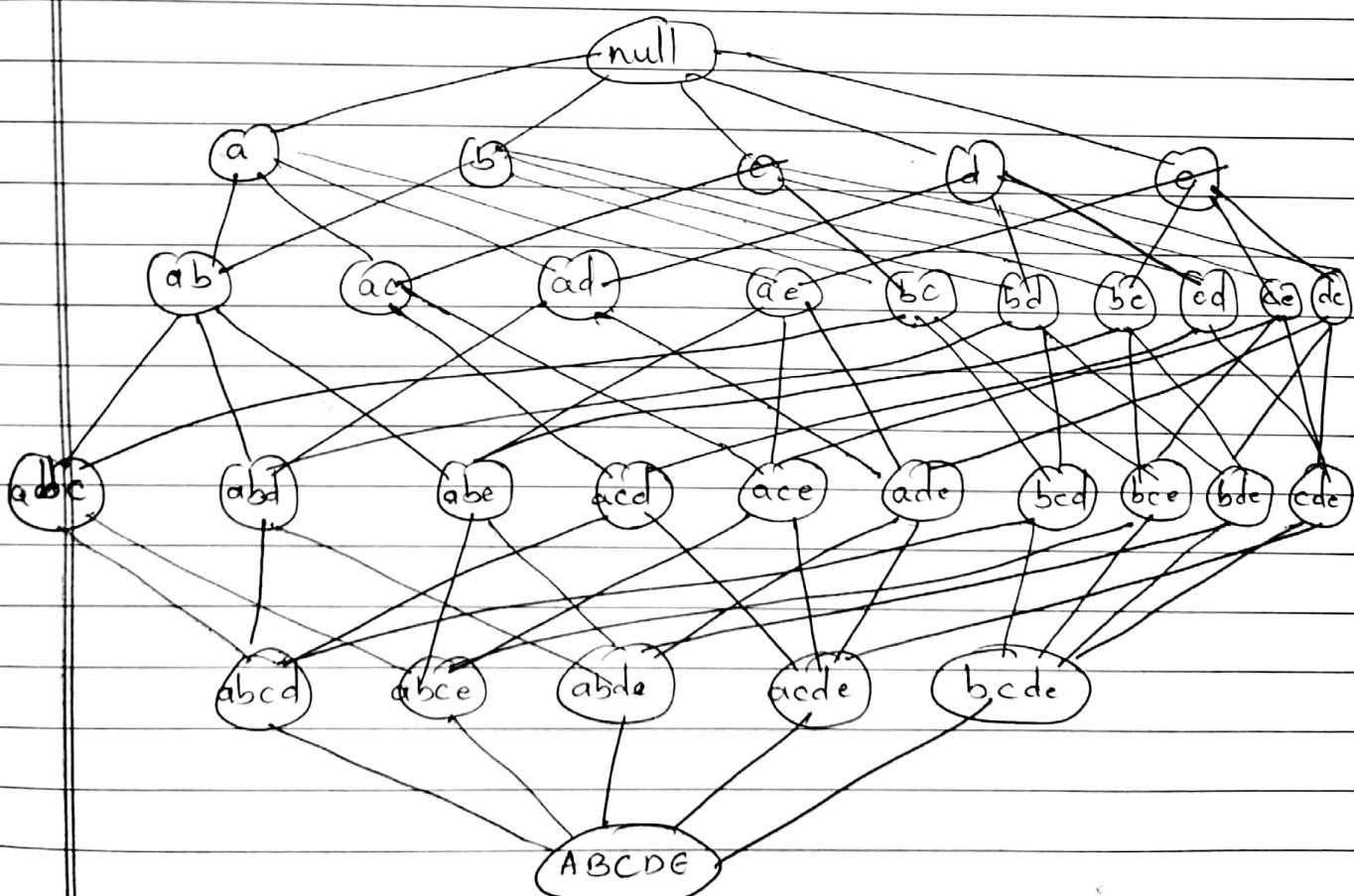
Frequent substructure \rightarrow tree, graph data

Maximal freq.

Is a item set where none of item set are immediate itemset are frequent.

Maximal freq

Every item above the border is frequent & every item below the border is not frequent.



26/08/19

Closed itemset : if none of immediate superset have the same support count as α

Closed frequent element

* Closed itemset:

If item set α is said to be closed itemset if none of the immediate supersets have the same support count as α

* Closed frequent Element:-

If item set α is a closed frequent itemset, if α is closed (none of its imm. supersets have the same support count as α) and α has a support count which is greater than equal to the minimum support count

Tid	Itemset
1	a,b,c
2	a,b,c,d
3	b,c,e
4	a,c, d,e
5	d,e

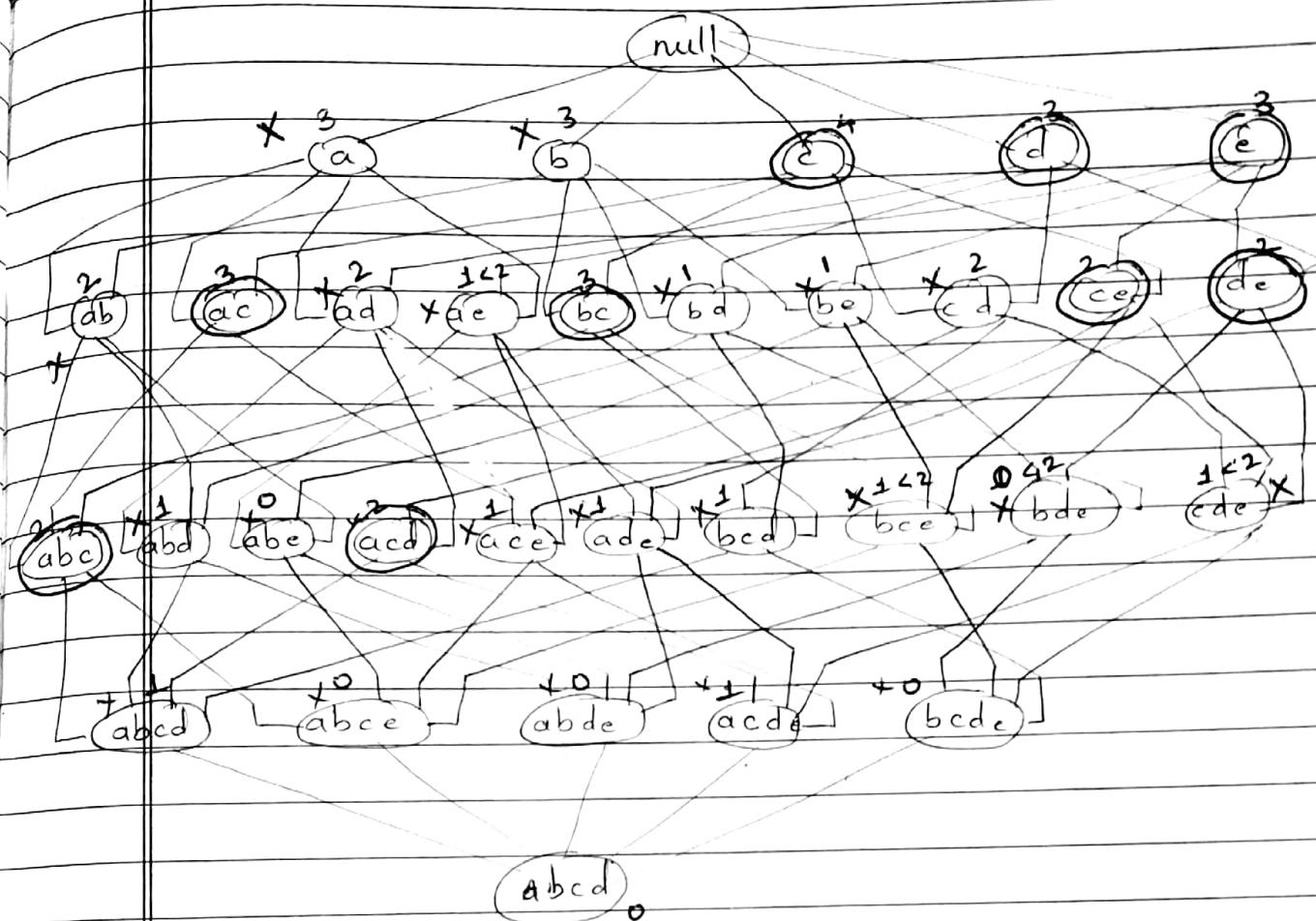
min sup = 40%

min sup-cnt = 40% of total transaction.

$$= 40\% \text{ of } 5$$

$$= \frac{40}{100} \times 5$$

$$= 2$$



$a = \{ab\}^3 \{ac\}^2 \{ad\}^2 \{ae\}^3 \times abd$ a closed freq. itemset

$c = \{ae\}^3 \{bc\}^3 \{cd\}^2 \{ce\}^2 \checkmark$ closed freq. itemset

Tid	Itemset	min sup-count = 3
		min confidence = 80%.
1	a, d, e	
2	a, b, c, e	
3	a, b, d, e	
4	a, c, d, e	
5	b, c, e	
6	b, d, e	
7	c, d	
8	a, b, c	
9	a, d, e	
10	a, b, e	

(A) Apriori^o
By Generation

1st iteration

Candidate set $C_1 \rightarrow$ frequent itemset L_1

Item	Sup-cnt	item	Sup-cnt
a	7	a	7
b	6	b	6
c	5	c	5
d	6	d	6
e	8	e	8

2nd iteration

$C_2 \rightarrow L_2$

Itemset	sup-cnt	itemset	sup-cnt
ab	4	ab	4
ac	3	ac	3
ad	4	ad	4
ae	6	ae	6
bc	3	bc	3
bd	2 X	be	5
be	5	ce	3
cd	2 X	de	5
ce	3		
de	5		

C_3

3rd iteration

L_2

Itemset	sup-cnt	Itemset	sup-cnt
abc	2 X	abc	
abd	1	abe	
abe	3	abe	3
acd	1	ade	
ace	2 X	ade	4
adc	4		
bce	2 X		

$$abc = \{ab\} \ \{bc\} \ \{ac\} \checkmark$$

$$C_4 \Rightarrow \phi \quad L_4 = \phi$$

frequent itemsets = $\{a, b, e\}$ $\{a, d, e\}$

2] Generate Association Rules.

$\{a, b, e\}$

$\{a\} \{b\} \{e\}$ $\{a, b\}$ $\{a, e\}$ $\{b, e\}$

$\{a\} \Rightarrow \{b, e\}$

$$\text{Confidence} = \frac{\text{sup.-cnt}(\{a\} \cup \{b, e\})}{\text{sup.-cnt}(\{a\})}$$

$$= \frac{\text{sup.-cnt}(a, b, e)}{\text{sup.-cnt}(a)}$$

$$= \frac{3}{7} = 0.42 = 42\% < \text{min confidence}$$

$\{b\} \Rightarrow \{a, e\}$

$$\text{Confidence} = \frac{\text{sup.-cnt}(a, b, e)}{\text{sup.-cnt}(b)} = \frac{3}{6} = 50\%$$

$< \text{min confidence}$

$\{e\} \Rightarrow \{a, b\}$

$$\text{Confidence} = \frac{\text{sup.-cnt}(a, b, e)}{\text{sup.-cnt}(e)} = \frac{3}{8} = 37.5\%$$

$< \text{min confidence}$

$\{a, b\} \Rightarrow \{e\}$

$$\text{Confidence} = \frac{\text{sup.-cnt}(a, b, e)}{\text{sup.-cnt}(a, b)} = \frac{3}{4} = 75\%$$

$< \text{min confidence}$

$\{a, e\} \Rightarrow \{b\}$

$$\text{Confidence} = \frac{\text{sup.-cnt}(a, b, e)}{\text{sup.-cnt}(a, e)} = \frac{3}{6} = 50\% < \text{min confidence}$$

$\{b\} \Rightarrow \{a\}$

Confidence = $\frac{3}{5} = 60\% < \text{min confidence.}$

$I = \{a, d, e\}$

$\{a\} \{d\} \{e\} \{ad\} \{ae\} \{de\}$

- $\{adj\} \Rightarrow c$

Confidence = 100%

- $\{de\} \Rightarrow a$

Confidence = 80%

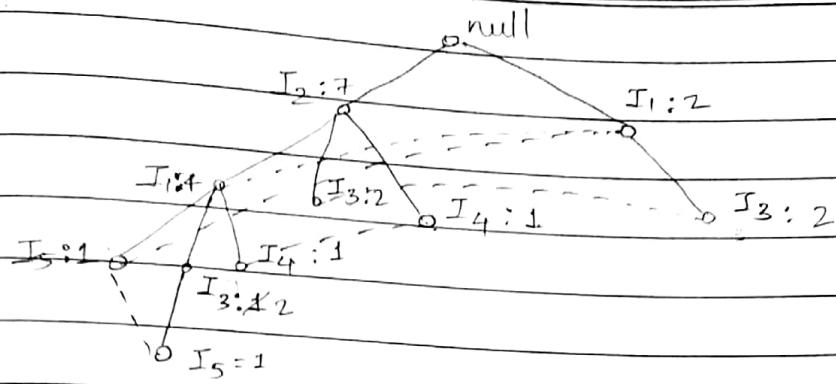
gloegig

FP tree

Tid	Items	Item	sup-ent	Item	sup-ent
1	$I_1 I_2 I_5$	I_1	6	I_2	7
2	$I_2 I_4$	I_2	7	I_1	6
3	$I_2 I_3$	I_3	6	I_3	6
4	$I_1 I_2 I_4$	I_4	2	I_1	2
5	$I_1 I_3$	I_5	2	I_5	2
6	$I_2 I_3$				
7	$I_1 I_3$				
8	$I_1 I_2 I_3 I_5$				
9	$I_1 I_2 I_3$				

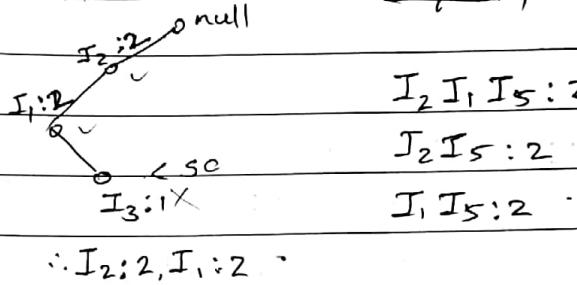
min sup = 2

Tid	Items	Tid	Items
1	$I_2 I_1 I_5$	9	$I_2 I_1 I_3$
2	$I_2 I_4$		
3	$I_2 I_3$		
4	$I_2 I_1 I_4$		
5	$I_1 I_3$		
6	$I_2 I_3$		
7	$I_1 I_3$		
8	$I_2 I_1 I_3 I_5$		



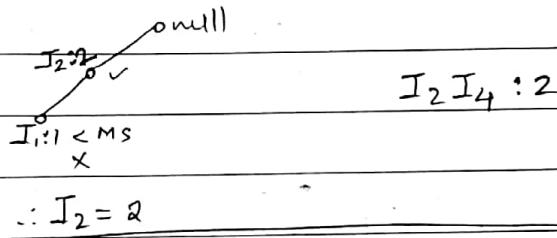
<u>Item</u>	<u>Conditional pattern base</u>	<u>Conditional FP-tree</u>	<u>frequent pattern</u>
<u>I5</u>	$I_2 I_1 : 1$		

$I_2 I_1 I_3 : 1$

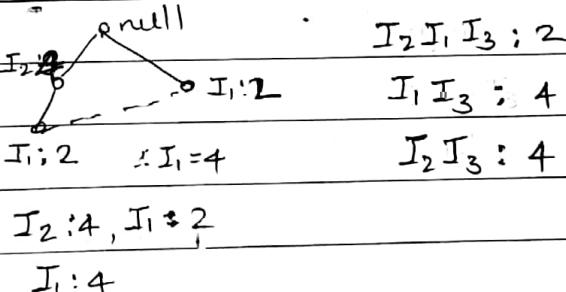


I4 $I_2 : 1$

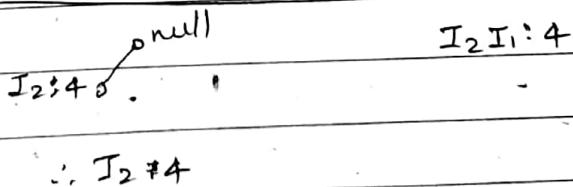
$I_2 I_1 : 1$



I3 $I_2 : 2$
 $I_2 I_1 : 2$
 $I_1 : 2$



I1 $I_2 : 4$



Q.	Tid	Items	①	Item	sup-cnt	②	Item	sup-cnt
=	t_1	$I_1 I_3 I_4$		I_1	3	\rightarrow	I_2	4
	t_2	$I_2 I_3 I_5$		I_2	4		I_3	4
	t_3	$I_1 I_2 I_3 I_5$		I_3	4		I_5	4
	t_4	$I_2 I_5$		I_4	1		I_1	3
	t_5	$I_1 I_2 I_3 I_5$		I_5	4		I_4	1

\therefore Exclude I_4 \because $I_4 < sc$

min-sup-cnt = 2

Tid	Items							
t_1	$I_3 I_1$							
t_2	$I_2 I_3 I_5$							
t_3	$I_2 I_3 I_5 I_1$							
t_4	$I_2 I_5$							
t_5	$I_2 I_3 I_5 I_1$							

- association \rightarrow find frequent sets

Classification

Hunts

types of splits (Binary, multiway ...)

Properties of decision tree

Impurity measures

gain & gain ratio

Mod 4

Association Mining

Using vertical data format

Tid	Itemset		
1	Strawberry, litchi, orange	l, o, s	
2	Strawberry, butter-fruit	b-f, s	
3	butter-fruit, vanilla	b-f, v	
4	strawberry, litchi, oranges	l, o, s	
5	Banana, Orange	b, o	
6	Banana	b	
7	banana, butter-fruit	b, b-f	
8	Strawberry, litchi, orange, apple	a, l, o, s	
9	apple, vanilla	a, v	
10	Strawberry, litchi	l, s	min sup-ct=3

①	Item	Transaction	②	Item	Transaction	
a	8, 9	X		b, b-f	7	X
b	5, 6, 7			b, l	-	X
b-f	2, 3, 7			b, o	5	X
l	1, 4, 8, 10			b, s	-	X
o	1, 4, 5, 8			b-f, l	-	X
s	1, 2, 4, 8, 10			b-f, o	-	X
v	3, 9	X		b-f, s	2	X
				l, o	1, 4, 8	
				l, s	1, 4, 8, 10	
				o, s	1, 4, 8	

③	Item	Transaction
	I, O, S	1, 418

frequent item set
f.o.s

1,0,5

Q. Transaction 1

Item

1	Bread, cheese, Eggs, Juice	B, C, E, J
2	Bread, cheese, Juice	B, C, J
3	Bread, Milk, Yogurt	B, M, Y
4	Bread, Juice, Milk	B, J, M
5	Cheese, Juice, Milk	C, J, M

The min-support is 50%.

$$\begin{aligned}
 \text{min-sup-ent} &= 50\% \text{ of total transaction} \\
 &= 50\% \text{ of } 5 \\
 &= \frac{50}{100} \times 5 \\
 &= 2.5 \approx 3
 \end{aligned}$$

(1)	Item	Transaction	(2)	Item	Transaction
	B	1, 2, 3, 4		B, C	1, 2, X
	C	1, 2, 5		B, J	1, 2, 4
	E	1, X		B, M	3, 4 X
	J	1, 2, 4, 5		B, J	
	M	3, 4, 5		C, J	1, 2, 5
	Y	3, 4, 5 X		C, M	5 X
				C, M	
				J M	4, 5 X

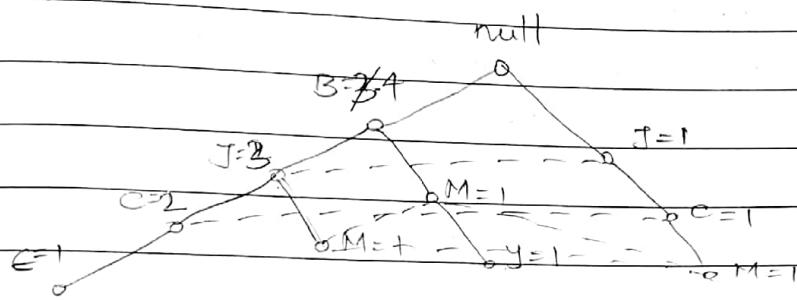
(3)	Item	Transaction
	B, C, J	1, 2 X

Frequent item set: .. B, J, C, J

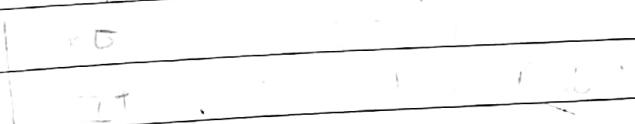
Q. FP Tree (Same question)

Item	Sup-cnt	Item	Sup-cnt
B	4	B	4
C	3	J	4
E	1	C	3
J	4	M	3
M	3	E	1
Y	1	Y	1

Tid	Item
1	B, J, C, E
2	B, J, C
3	B, M, Y
4	B, J, M
5	J, C, M



Item	Conditional pattern base	Conditional FP-tree	Frequent pattern
------	--------------------------	---------------------	------------------



FP Tree

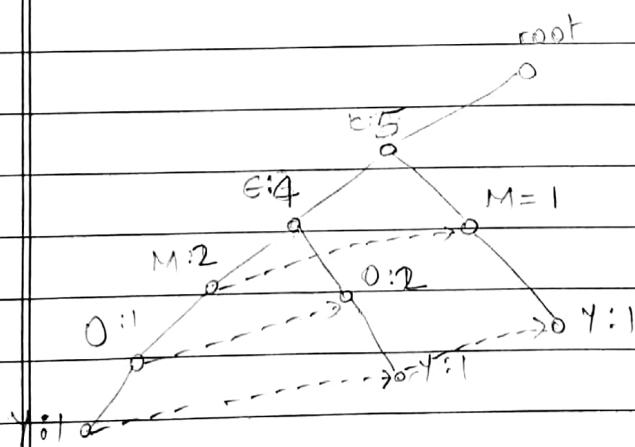
Tid	Items
1	M, O, N, K, E, Y
2	D, O, N, K, E, Y
3	M, A, K, E
4	M, D, C, K, Y
5	C, O, O, K, I, E

min sup = 60%

min sup-cnt = 60% of transaction
= 3

① Item	Sup-cnt	② Item	Sup-cnt
M	3 ✓	K	5
O	3 ✓	E	4
N	2 ✗	M	3
K	5 ✓	O	3
E	4 ✓	Y	3
Y	3 ✓		

③ Tid	Items
1	K, E, M, O, Y
2	K, E, O, Y
3	K, E, M
4	K, M, Y
5	K, E, O



Note: Y should not be the first item

Item	Conditional pattern base	Conditional FP tree	Frequent Pattern
Y	K, E, M, O : 1		
	K, E, O : 1		
	K; M : 1	<p style="text-align: center;">root</p> <p style="text-align: center;">K:3 E:2 M:1 O:1</p> <p style="text-align: center;">K:3, M:1</p>	
		<p style="text-align: center;">E:2</p> <p style="text-align: center;">M:2</p> <p style="text-align: center;">O:2</p>	