

DATA MINING & ANALYSIS:

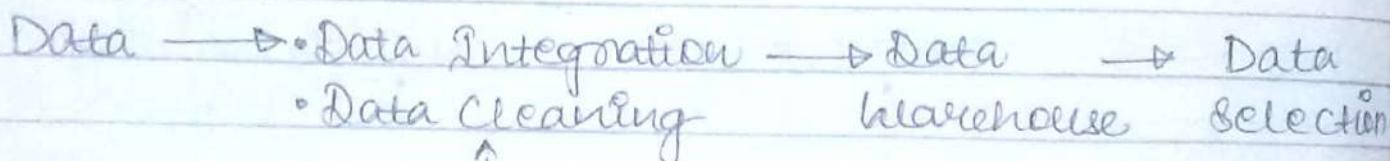
* Course Content :

- Unit I -
 - Data Preprocessing (08 hrs)
 - Frequent Pattern Mining (08 hrs)
- Unit II -
 - Classification techniques (08 hrs)
 - Cluster Analysis (08 hrs)
- Unit III -
 - Advanced Mining Techniques. (08 hrs)

CHAPTER-1

DATA PREPROCESSING

* KDD Process-



* Data Warehouse:

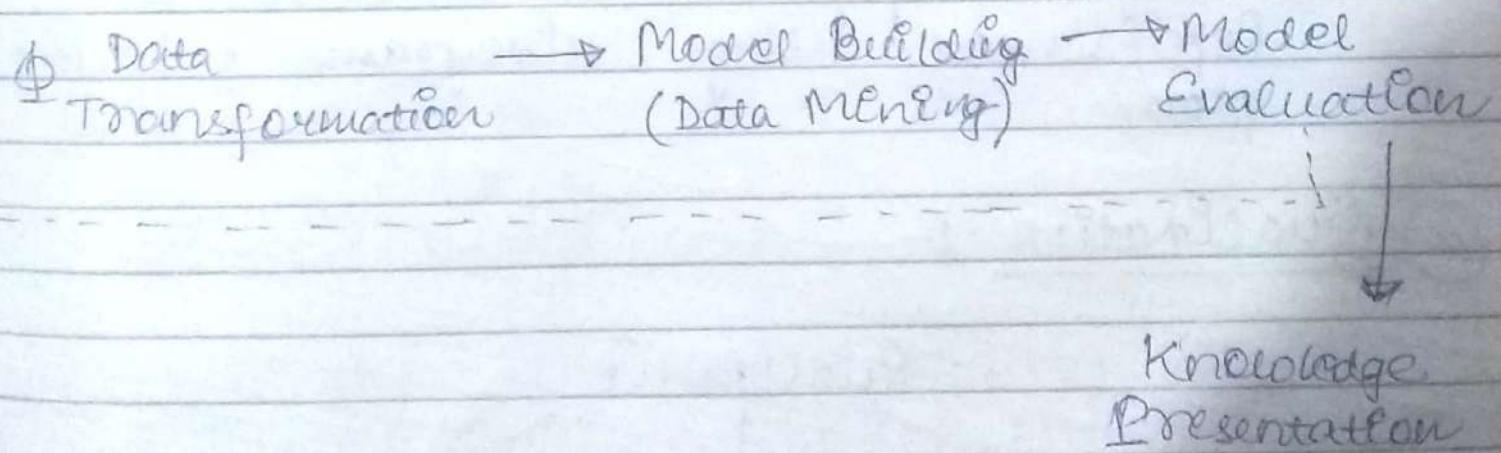
It is a repository of information collected from multiple sources, stored under a unified schema.

- Steps in KDD process - Data Mining (+1)
Reduction (+1)

Points

* APPLICATIONS OF DATA MINING:

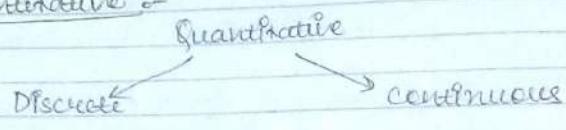
- a). Retail Industry :- Analyze Buying Patterns of customers.
- b). Banking Sector :- Fraud detection, Customer Retention
- c). Healthcare Industry :- Effective treatments, Best practices, Fraud detection
- d). Social & Information Networks :- New programs, new products / processes or services.



* DATA TYPES:

Types of Data :- a). Quantitative (Numerical)
b). Qualitative (Categorical)

(1) Quantitative :-



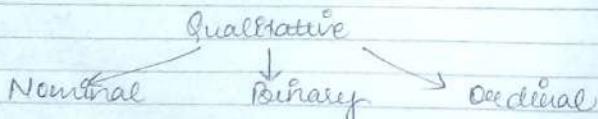
* Discrete

- takes integer values.
- E.g.: no. of people
- They represent the count.
- Represented using Bar Charts & Histograms.

* Continuous

- takes real nos. (can take ∞ no. of values)
- E.g.: ht. of person, distance covered.
- Represented using Histograms & Scatter plots.

(2) Qualitative :-



* Nominal :

- variables are labeled using symbols or words or letters
- E.g.: name of a person, city, color of eyes, etc.
- Represented using Pie charts usually.

* Ordinal :

- A well-defined order among categories.
- E.g.: size \rightarrow small, medium, large
- income
- Represented using bar graphs / bar charts

* Binary :

- takes only 2 values.
- E.g.: pass/fail, male/female, +ve/-ve
- Represented using pie charts.

Ques Select categorical data types -

- Color of a shirt.
- Finish order in a race.

Ex - Gender (Categorical, binary)

Age (Categorical, ordinal, when you snicker)

Teeth/Mouth (Numerical, continuous)

Item purchased (Categorical, numerical)

Quantity (Numerical, Discrete).

* DATA PREPROCESSING: INTRODUCTION

- Why Data Preprocessing is required?
- Raw (dirty) data is not adequate for analyzing.

Reasons:
• Quality
• Dimensionality { Too many attributes }

* Preprocessing of data assess the quality of the data. There are certain measures to assess quality of data.

Measures:

- Completeness, Consistency
- Timeliness, Accuracy
- Reliability
- Interpretability

1). Completeness:

- Not recorded, Missing Values

2). Consistency:

- Discrepancy Discrepancy in code or names.

3). Accuracy:

- degree to which recorded/noted value of data is correct or wrong.

- 4). Timeliness: indicates timely availability (updated, most recent data).
- 5). Interpretability: It indicates the ease with user can understand the data.
E.g. values, meaning, info. about that etc.
- 6). Deliverability: indicates the trustworthiness of the data. How trustable the data are correct?

* Preprocessing Techniques :-

Major Tasks in Data preprocessing:-

- Data Cleaning
- Data Integration
- Data Reduction
- Data Transformation

1). Data Cleaning: fills in the missing values & smoothes the noisy data.

2). Data Integration: It integrates the data sets by dissolving the issues during integration.

3). Data Reduction: It reduces the dimensionality of the data.

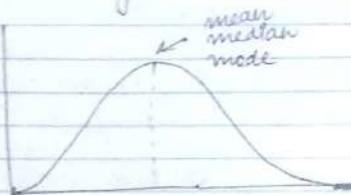
4). Data Transformation: It transforms the set of values of attributes to new set of values of attributes whenever required.

* DATA CLEANSING :

- Incomplete Data (Missing Values)
- Noisy Data

* Techniques to handle missing values:

- 1). Ignore the tuple : deleting those tuples whose values are missing.
- 2). Use a global constant: FNU in the missing values with either -x or ~~the~~ or unknown.
- 3). Use central tendency : the missing values are filled with mean, when the attributes are normally distributed this way.



Symmetric data

If attribute values are either + very skewed or - very skewed, then we replace the missing values with (median)

If the missing attribute is a categorical data, then that value is filled with the most frequent occurring value in the column.
i.e. Mode

{ mean : normally distributed
median : skewed
mode : categorical }

4). Fill in the missing values manually.

5). Use most probable value :

• Regression • Decision Tree

Quiz: Techniques to handle Missing Values :

• Regression .

• Central Tendency of the attribute.

* NOISY DATA

Noise is a random error or variance in a measured variable.

E.g. salary → -105 K.

Techniques to handle Noisy Data :

- 1). Binning
- 2). Regression
- 3). Outlier Analysis

1). Binning:

- a) Bin Means Method
- b) Bin Median Method
- c) Bin Boundary Method

E.g.: Data for price (in dollars): 7, 3, 19, 14, 23, 31, 24, 38, 33.

a).
Step① - Sort the data in ascending order.
3, 7, 14, 19, 23, 24, 31, 33, 38.

Step② - Distribute them equally into bins.

Bin 1: [3, 7, 14] Bin 2: [19, 23, 24] Bin 3: [31, 33, 38]

mean of 1st bin : 8

③ Replace all the values of the bin with its mean.

∴ Bin 1: [8, 8, 8] Bin 2: [22, 22, 22] Bin 3: [34, 34, 34]

This is Bin Means Method. →

b). Bin Median Method:

In this method, all the data in the bin with the median.

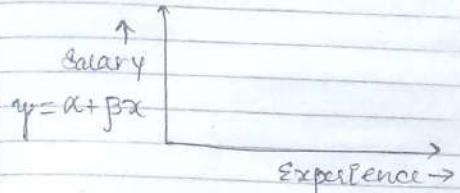
∴ Bin 1: [7, 7, 7]
Bin 2: [23, 23, 23]
Bin 3: [33, 33, 33]

c). Bin Boundary Method :

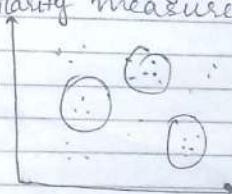
The Intermediate value is replaced by its closest boundary.

∴ Bin 1: [3, 3, 14] Bin 2: [19, 24, 24] Bin 3: [31, 31, 38]

2). Regression: Data can be smoothed by fitting the data to a regression function



3). Outlier Analysis: Uses clustering techniques, where in we group the data based on the similarity measures. The objects outside the clusters are called noise & are eliminated.



2). Non-Parametric Methods:

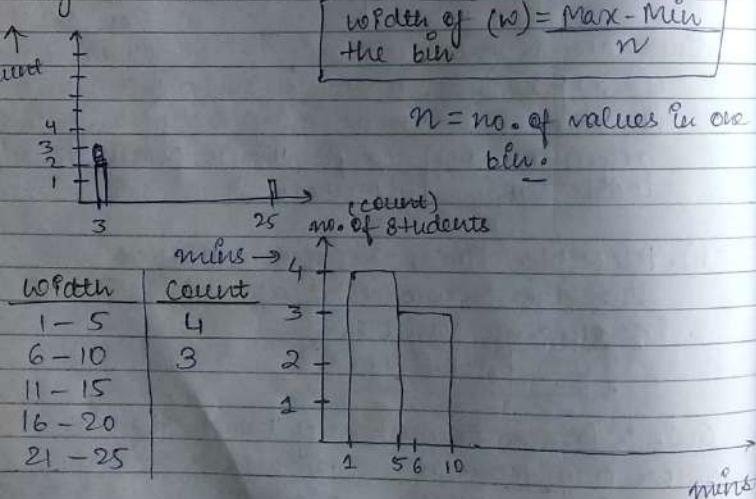
a). Histograms: It reduces the data by storing the average (sum) for each bin.

Bin is a range of values or a single value & the occurrence will be shown.

E.g.: The given data is the time (in mins) given by students for watching the video.

3, 4, 7, 11, 18, 3, 19, 22, 23, 25, 2, 10, 21, 14, 8, 11, 17, 23, 23, 14.

Singleton Bucket.



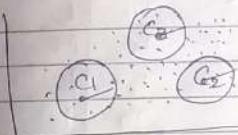
b). Clustering: unsupervised method (data mining technique).

• Store cluster representation (e.g. centroid & diameter) only.

The no. of centroids taken, depends on the no. of clusters. If we want 2 clusters, we will be having 2 centroids, etc.

Let $x_1 = (3, 4)$, $x_2 = (5, 6)$, $x_3 = (10, 12)$ In

find centroids c_1, c_2, c_3 if we want 3 clusters



Each circle will be having diameter.

So, instead of storing the entire data we can store centroids & diameters only.

c).
 c_1, d_1
 c_2, d_2
 c_3, d_3

Quiz: Clustering in Data Reduction :-

- Partitions data into clusters
- Store only cluster representation (Centroid & Diameter).

c). Sampling:

- Simple Random Sampling (SRS)
 - Samples are taken randomly i.e. every tuple has equal probability of getting selected.
- SRS without replacement (SRSWOR)
- SRS with replacement (SRSWR)
- Stratified Sampling

* SRSWOR:

Let we want two samples of size 3.

<u>S₁</u>	<u>S₂</u>
T ₂	T ₁
T ₄	T ₇
T ₆	T ₉

T ₁	L
T ₂	M
T ₃	L
T ₄	L
T ₅	M
T ₆	L
T ₇	H
T ₈	H
T ₉	H
T ₁₀	H

Here, tuples are not repeated.

* SRSWR:

<u>S₁</u>	<u>S₂</u>
T ₂	T ₂
T ₄	T ₂
T ₆	T ₉

Tuples can be repeated.

T₂, T₄, T₆ are randomly chosen for sample S₁. However, they are not taken out from the tuples set for they can be again selected for sample S₂.

* Stratified Sampling:

Let us suppose 'Income' attribute has non-linear values i.e. Low(L), Medium(M) & High(H).

1st strata : T₁, L

T₃, L

T₄, L

T₆, L

2nd strata : T₂, M

T₅, M

T₈, H

T₉, H

T₁₀, H

3rd strata : T₇, H

T₈, H

NOW, let us suppose that we need to select 50% from each strata.

From

1st strata → T₃
T₄

2nd strata → T₅

3rd strata → T₈
T₉

Stratified
Sample (S.S)

NOW, we see that let it be SRSWOR, SRSWR or S.S. the population size is reduced based on the sample size.

It reduces the data by selecting the data based on the sample size.

Ques: A selected object is not removed from population what type of sampling is it?
SRS with replacement (SRSWR)

* DATA TRANSFORMATION:-

- It is a pre-processing technique in which attribute values are mapped to a new set of values i.e., replaced by a new value.
- It is required because it leads the data mining process into a more efficient one and also the patterns obtained will be more easy to understand.

Data Transformation Techniques:-

- Smoothing
- Normalization
- Concept hierarchy
- Aggregation
- Discretization
- Generation

a). Smoothing:

Bin Means Method:

$$\begin{array}{r} 3+7+2 \\ \hline 3 \\ = 4 \end{array} \quad [3, 7, 2] \quad [4, 4, 4]$$

b). Aggregation:

Summarise the attribute value into one.

J	100	Q1 = 600	Half-yearly
F	200	Q2 = 1200	
M	300	Q3 = 1800	
A	400	Q4 = 2400	
;			Yearly
D			Half yearly

c). Normalization: gives / assigns equal weight to all the attributes

Normalization Methods:

- Min-Max Normalization
- Z-Score normalization
- Normalization by decimal scaling

* Min-Max Normalization:-

maps min value to '0' & maximum value to '1' & other values lying b/w will be mapped to any value b/w 0 to 1.

$$v' = \frac{(v - \text{min}_A)}{(\text{max}_A - \text{min}_A)}$$

Eg Income (Before Normalization):-

$$\begin{array}{ccc} 12,222 & 15,000 & 85,000 \\ \downarrow & & \downarrow \\ \therefore v' = \frac{(12,222 - 12,222)}{(85,000 - 12,222)} (1-0) & v' = \frac{(85,000 - 12,222)}{(85,000 - 12,222)} (1-0) \\ = 0 & & = 1 \end{array}$$

$$v' = \frac{(15,000 - 12,222)}{(85,000 - 12,222)} = 0.038$$

Ques: Which normalization is useful when min & max values of att are unknown?
 → Z-score Normalization Technique.

* Z-score Normalization: This method is used when the min & the max value of the attribute is not known.

$$v' = \frac{v - \mu_A}{\sigma_A}$$

• μ_A = mean of att. A
 • σ_A = std. dev. of A

E.g.

$$v' = \frac{73600 - 54000}{16000} = 1.225$$

- If value of attribute = mean of att. then v' (value of att) will become 0.
- If $v >$ mean of att. then $v' = +ve$.
- If $v <$ mean of att. then $v' = -ve$.

* Decimal Scaling Method:

- moves the decimal points of the att. values
- depends on the maximum absolute value of A.

$$v' = v / 10^T$$

T = no. of digits in v.

$$E.g. v = -986$$

$$v' = \frac{-986}{10^3} = -0.986$$

$$v = 917$$

$$v' = \frac{+917}{10^3} = 0.917$$

d). Discretization: The data values are divided into intervals & labelling them.

E.g.: age : 20. 80

20..35, 36..55, 56..80

young

middle

senior

This can be done by equal frequency or equal width method.

Thus, interval labels replace actual data values.

e). Concept Hierarchy Generation:

Replaces low level concepts by higher level concepts.

country

↑

state

↑

city

↑

street

Here, street can be replaced by the name of the country only

CHAPTER 2

FREQUENT PATTERN MINING

* Vdeo 1: Basic Concepts - Frequent Items, Support & Confidence

* Frequent Pattern Mining (Association rule mining): It is a process of finding frequent patterns i.e. the relevant hidden pattern from a database e.g. transactional database, relational database or any other repositories working together.

* Frequent Patterns: They are the item sets, subsequences or substructures that occur together.

* Frequent Item sets means, the item sets that are purchased together by the customer.

- Sub sequence → DNA Analysis
- Substructures → Protein Structure Analysis
- They represent important properties of the data.

* Frequent Pattern Mining

- Finding frequent patterns
- Deriving Association rules

Support & Confidence:

P.T.O

SUPPORT & CONFIDENCE

Trans	Items Purchased	occurrence
1	Milk, Bread, Coffee	Milk : 3
2	Milk, coffee	Bread : 2
3	Milk, Butter	Coffee : 2
4	Bread, Egg	Butter : 1

. we say item is frequent, when the occurrence of that item \geq threshold. This threshold is known as Support.

e.g.: Min-Sup = 2 (bt)

∴ frequent items are milk, bread, coffee.

Item set:

milk, bread : 1

milk, coffee : 2 ✓ (frequent item set).

∴ Association rule can be set up for frequent item set i.e. for milk, coffee

milk \Rightarrow coffee

Now, support (milk, coffee)

support (milk U coffee) = 2 = 50% = 2/4

Show many milk & coffee are occurring together in the transaction?.

$$\text{Confidence} = \frac{s(\text{milk U coffee})}{s(\text{milk})} = \frac{2}{3} = 66.6\%$$

• Support = 50% shows that the 50% of items milk & coffee are purchased together.

• Confidence = 66.6% indicates that 66% customer who purchased milk, also purchased coffee.

P.T.O

\therefore association Rule : milk \Rightarrow coffee ($s=2, c=66.6\%$)

Quiz: what do you mean by support (A)?

$$\text{Support}(A) = \frac{\text{No. of transactions containing } A}{\text{Total no. of transactions}}$$

* Video 2: Apriori Algorithm :-

* APRIORI ALGORITHM (Frequent Item Set Mining method)

* Apriori Property :-

1) If an itemset is frequent, then all of its subsets must also be frequent.

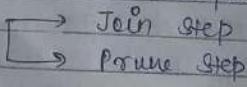
E.g. If AB is frequent then {A}, {B}, {AB} all are frequent.

2) If an itemset is infrequent, then all of its supersets must be infrequent too.

E.g. If {A} or {B} is infrequent, then their superset i.e. {AB} will also be infrequent.

Apriori Algorithm :-

It uses a 'levelwise search approach' in which 'k' itemsets are used to find out 'k+1' frequent item sets.



levelwise Search :- d steps

Join step: C_k is generated by joining L_{k-1} with itself.

Prune step: Any $(k-1)$ -item set that is not frequent cannot be a subset of a frequent k-item set.

E.g. Tid		Items Purchased	Min-sup=2
1		Milk, Bread, coffee	
2		Butter, Bread, Egg	
3		Milk, Butter, Bread, Egg	
4		Butter, Egg	

Step 1 : Generate C_1 & 1 frequent itemset L_1

C_1		L ₁	
Itemset	Support Count	Itemset	Sup-count
M	2	M	2
B	3	B	3
C	1	$\rightarrow (\leq 2) \therefore$ Ignored	B4
Bu	3		E3
E	3		

by going L_1 with L

Step 2 : Generate C_2 & 2 frequent itemset L_2

C_2		L ₂	
Itemset	Freq-count	Itemset	Sup-count
Milk, Bread	2	Milk, Bread	2
Milk, Butter	1	\rightarrow pruned	Bread, Butter
Milk, Egg	1		Bread, Egg
Bread, Butter	2		Butter, Egg
Bread, Egg	2		
Butter, Egg	3		

Step 3 : Generate C_3 & 3 freq itemset L_3 by going with L_2 .

<u>C_3</u>	<u>Item-set</u>	<u>Sup-count</u>	<u>(M, Bu)</u>
M, Br, Bu	In frequent	$\because (M, Br) \cdot (\underline{Br}, Bu)$	is Infrequent
M, Br, Egg	In frequent.		
Br, Bu, Egg			

↳ First Apply and Apriori Property.

<u>C_3</u>	<u>Item-set</u>	<u>sup-count</u>	<u>L_3</u>	<u>Item-set</u>	<u>Sup-count</u>
Br, Bu, Egg	2 ✓	→	Br, Bu, Egg	2	

Step 4 : Generate C_4 by going L_3 with L_3 .

$$C_4 = \emptyset \quad (\because \text{In } L_3 \text{ no 4 datasets})$$

\therefore Frequent Item Set = { Br, Bu, Eg } with Sup count = 2

Quiz : Let the size 2-frequent item sets obtained from Apriori Algorithm be $\{A, B\}$, $\{A, C\}$, $\{A, D\}$, $\{B, C\}$, $\{B, E\}$ and $\{C, E\}$. Select all the potential size- 3 item sets from the following :

$\{A, B, D\}$

$\{B, C, E\}$

$\{A, B, C\}$

$\{A, B, E\}$

* Video 1: Closed & Maximal Frequent Itemset

* Downward Closure Property: Any subset of a frequent itemset must also be frequent.

$$\{A \ B \ C\} \rightarrow \{\{A\} \ \{B\} \ \{C\} \ \{AB\} \ \{BC\} \ \{CA\}$$

$$\begin{matrix} A, B, C \\ N=3 \end{matrix} \rightarrow ?$$

$$\{ABC\}$$

• For N distinct items, frequent itemset will be $2^N - 1$.

* Closed & Maximal Frequent Itemsets:

MAXIMAL FREQUENT ITEMSET :

T.Id	Items Purchased
1	Milk, Bread, coffee
2	Butter, Bread, Egg
3	Milk, Butter, Bread, Egg
4	Butter, Egg.

Using Apriori Algorithm, we get : Min. Support = 2

Step 1: 1 frequent Itemset	Step 2: 2 freq Itemset	Step 3: 3 freq Itemset			
Itemset	Sup-count	Itemset	Sup-count	Itemset	Sup-count
{M}	2	{M, Br}	2		
{Br}	3	{Br, Bu}	2	{Br, Bu, Eg}	2
{Bu}	3	{Br, Eg}	2	{M, Br, Eg}	1
{Eg}	3	{Bu, Eg}	3	{M, Bu, Eg}	1

- An itemset is called as maximal frequent itemset if, its supersets are infrequent.

Sog: {Br, Bu, Eg} \leftarrow maximal freq. & max. size

Drawback

Maximal freq item sets do not have any information about the support count of the subsets.

Because of this it may give some redundant information.

PRESUENT

- * CLOSED ITEM SET: A frequent itemset is also known as closed frequent item set if its superset doesn't have same support count as of it, or greater sup. count.

$\text{fM}_4 : 2$ $\text{fM}, \text{Br}_4 : 2$

\therefore not a closed freq. function.

$\{B_2, B_3\}$ Maximal, \therefore Superset closed.
 $\{M, B_2\}$ \therefore Not closed.
 $\{B_1, B_2\}$:
 $\{B_1, B_3\}$:
 $\{B_2, B_3\}$:
 $\{\}$:

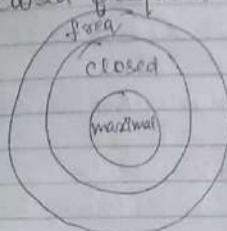
$f \in \mathcal{C}^1$: 3 $f_{\text{ext}, \text{eq}} \in \mathcal{C}^1$: 2 \therefore Not closed

§ Bui, Eq 3:3 because of § Bui, eq 4-3

NETC

A maximal frequent Hemset can be a closed one, but closed cannot be a maximal one.

- Thus, there exists a relationship b/w frequent closed frequent & maximal freq. itemset.



\therefore Maximal, Closed & freq
(Subset)

* Video 8: Generating Association Rules from frequent k-Hemets

Steps

- 1). For each frequent itemset I , generate all non-empty subsets of I . - for $\{B_3, B_4, E_3\}$ -
 $\Rightarrow \{\{B_3, B_4\}, \{B_3, E_3\}, \{B_4, E_3\}, \{B_3, B_4, E_3\}\}$

- 2). For every non-empty subset S of L , output the rule " $S \rightarrow (1-S)$ " if: $\frac{SC(L)}{SC(S)} > \text{min-confidence}$

Rule 1: $\{B_3, B_4\} \Rightarrow E_2$ // The customer who purchase B_3, B_4 also purchase E_2 .
 $\text{Supp}_{\text{cont}}(B_3, B_4, E_2) = \frac{2}{8} \times 100$
 $\text{Supp}_{\text{cont}}(B_3, B_4) = 100\%$.

Let us specify min-confidence = 70%
∴ This rule is a Strong Rule.

Rules: $\{B_1, B_2\} \Rightarrow \{B_1, B_2\}$
 $\text{Sup.-cont}(B_2, B_1, \text{tg}) = \frac{2}{3} \times 100 = 66.67\%$
 $\text{Sup.-cont}(B_1, \text{tg})$
 1. Not strong/Not interesting.

* Video 3: FP Growth Algorithm:-
 ↳ Frequent Item Set Mining Method

* Drawbacks of Apriori :-

- Generates large candidate elements if items in the database is large.
- Needs multiple scans of database to know the sup-cont.
- So, the problem with Apriori is space & time. These drawbacks can be overcome by FP GROWTH ALGORITHM.

FP GROWTH ALGORITHM:-
 ↳ Adv: Generate only frequent items from the infrequent items.
 (Divide & Conquer Strategy)

- It compresses the database into a frequent pattern tree.
- Divides the compressed database into a set of conditional databases.

- Each conditional database is associated with one frequent item.

E.g.: min-sup = 3

Algorithm:-

- First, scans the database to find sup-cont of each item.

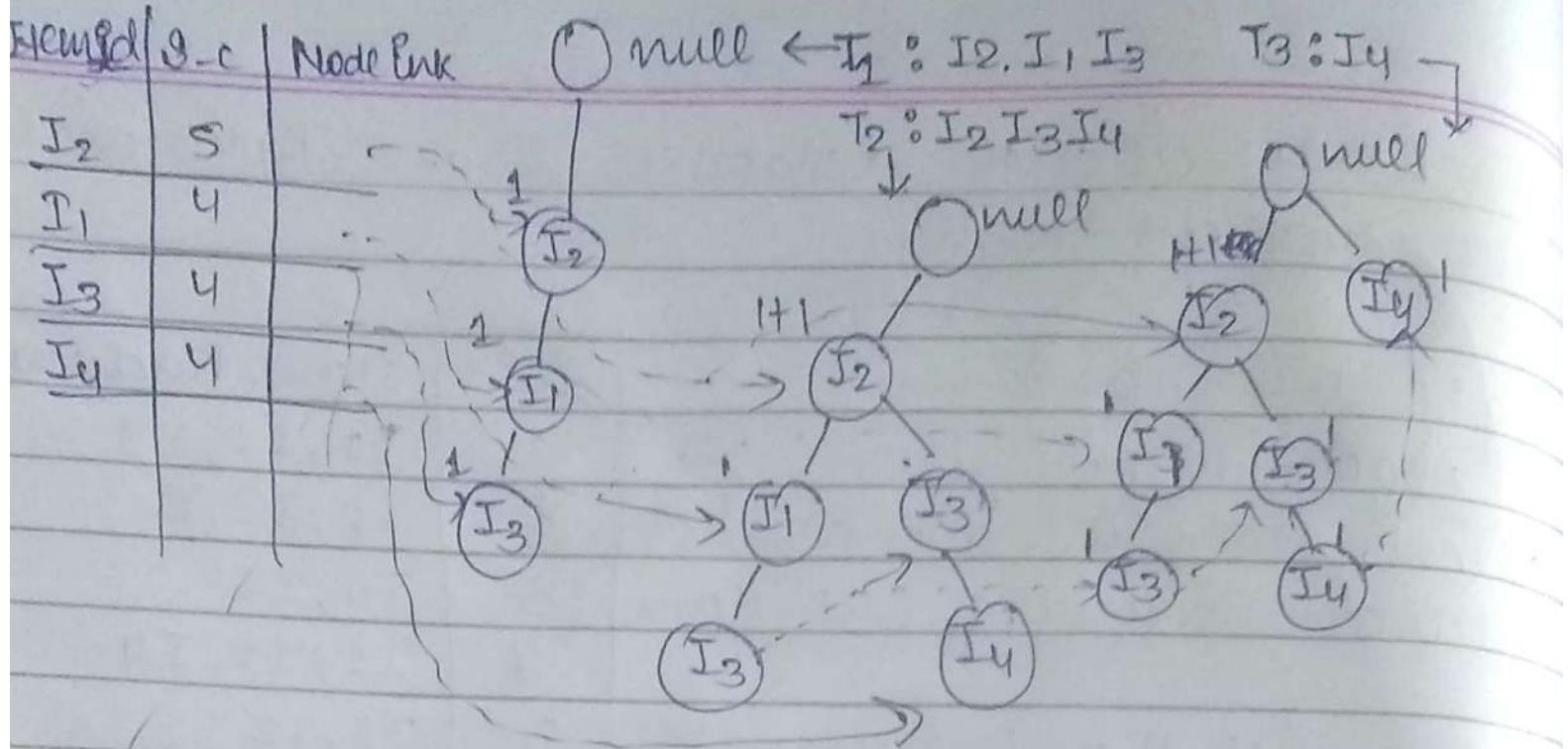
- Then, list them in descending order by eliminating the items whose sup-cont is less than min-supcount.

e.g. I2:5, I1:4, I3:4, I4:4

- Now, algorithm will reorganise the tables:-

Tid	Items Purchased	Items Purchased (sorted)
1	I1, I2, I3	I2, I1, I3
2	I2, I3, I4	I2, I3, I4
3	I4, I5	I4
4	I1, I2, I4	I2, I1, I4
5	I1, I2, I3, I5	I2, I1, I3
6	I1, I2, I3, I4	I2, I1, I3, I4

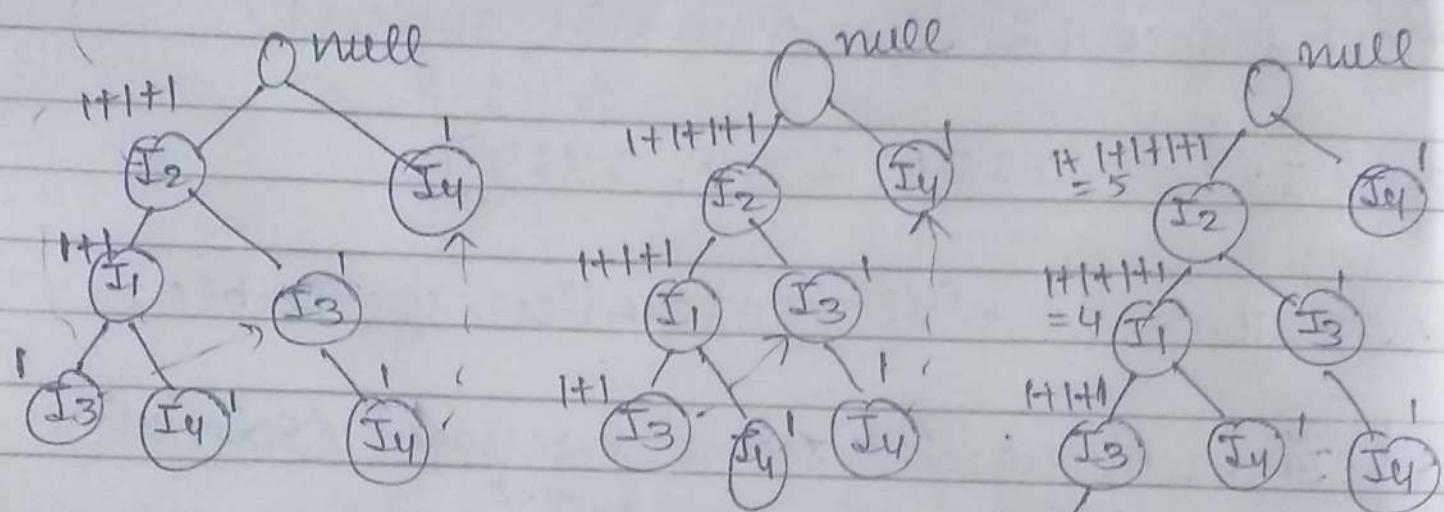
Header Table:			
Header	Sup.-cont	Node-link	
I2	5		
I1	4		
I3	4		
I4	4		



T₄: I₂, I₁, I₄

T₅: I₂, I₁, I₃

T₆: I₂, I₁, I₃, I₄



* Now, finding frequent itemsets from FP Tree :- min-sup = 3

Item	Conditional Pattern Base	Conditional FP Tree	Freq Pattern
I ₄	I ₂ I ₁ I ₃ : 1, I ₂ I ₁ : 1, I ₂ I ₃ : 1, I₂: 1 (Root node)	I ₂ : 3, I ₁ : 2, I ₃ : 2 (< 3)	I ₂ I ₄ : 3
I ₃	I ₂ I ₁ : 3, I ₂ : 1	I ₂ : 4, I ₁ : 3	I ₂ I ₃ ; I ₁ I ₃ : 3, I ₁ I ₂ I ₃ : 3
I ₁	{I ₂ : 4}	{I ₂ : 4}	{I ₂ I ₃ : 4}

Data Mining

Unit - 1

Data Pre-Processing - Chapter - 1

KDD Process



Data

May be
structured
or
Not

attribute values are
transformed to new
set of values.

(Data transformation)

[True - 1 false - 0]

* If evaluation results
are ^{not} upto mark then
these steps must be
repeated.

Needed bcs data might be
available in diff sources

→ Data Integration

&
Data cleaning

Needed bcs "null"
values or "duplicates"
values may be
present in data.

Data stored in well
defined
Data warehouse.

Data selection

&
Data transformation

Relavent attributes

must be selected out of
n no. of attributes.

(Data selection)

Model building

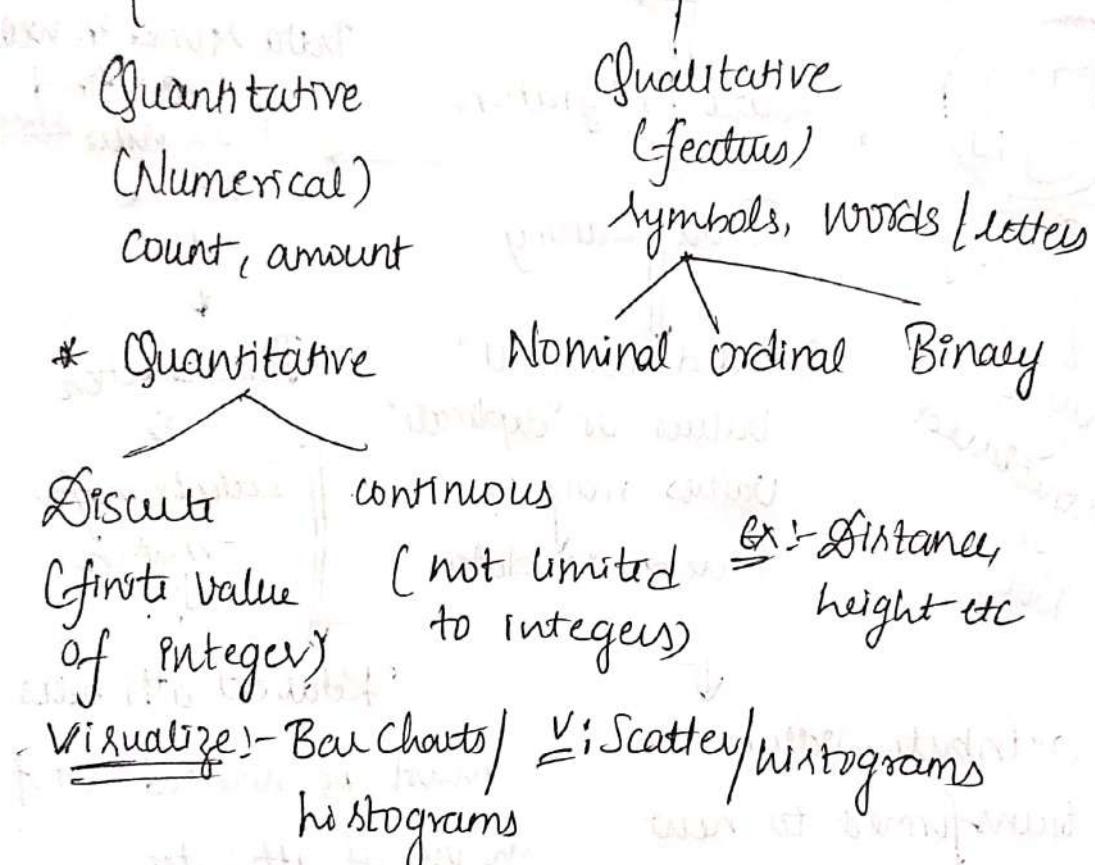
(Data mining)

Diff class, clustering
& pattern designing

Model evaluation
using diff metrics

Applications:- Retail sector, Banking sector, healthcare Industry & social and information networks etc.

Data types



* Nominal!- Labeled using symbols / words / letters

Ex:- Name of place, colour of skin etc

V:- Piecharts (percentage of male & female)

* Ordinal!- Order among categories

Ex:- small, medium

large etc.

Good

(or) Avg, worst

Grades: S, A, B, C, D

V:- Bar graph

Binary: Takes only 2 values

↳ Pass, fail etc.

↳ Pie charts.

Data Preprocessing

Why: Raw (dirty) data is not adequate for analysis.

Reasons: Quality (Not good quality).

Dimensionality (Too many attributes)

What: Assess the quality of data.

Measures

Not recorded, Missing values

- 1) Completeness, consistency. Disturbances happened
- 2) Timeliness, Accuracy (if or not) during data integration
- 3) Timeliness, Believability (two tables merging etc)
- 4) Interpretability.

* Accuracy:- Is recorded value correct / not.

* Timeliness:- Timely availability of data, mean must available / updated data.

* Interpretability:- How easily data is understand

* Believability:- The degree at which user can accept data.

Major tasks in Data Processing \rightarrow Dirty to quality data.

- * Data Cleaning (fills missing values & duplicates removed)
- * " Integration (Resolves issues during integration)
- * " Reduction (dimension reduction)
- * " Transformation (attribute value changes) to new set of values

* Data cleaning

- * Deals with incomplete data \rightarrow values are missing

EId	Name	Gender
1	-	-
2	A	-
3	B	F
4	C	M
5	D	M

Methods :- 1. Ignore the tuple.

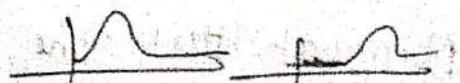
Ex:- tuple 1 is deleted/ignored.

2. Use a global constant

Ex:- Filling missing values with $-\infty$ or Unknown

3. Use central tendency :-

Mean :- Normally distributed data \approx

Median :- Skewed 

Mode:- If data is categorical
like replace Gender of tuple 2 with M
mode.

4. Fill in the missing values manually.

5. Use most probable value:

* regression

* Decision tree

Noisy Data

* It is a random error or covariance in a measured variable.

tId	Exper	salary
1	16	1000
2	10	-20
3	200	4000
4	15	8000

means

medium

boundary.

Methods:-

1. Binning
2. Regression
3. Outlier Analysis.

Binning

* First sort data

* Partition into equal bins

(equal freq) & (equal depth)

* Then apply means/m/bs

Ex:- 7 3, 9 19 14 23 31 24 38 33

SOF:- 3, 7, 14, 19, 23, 24, 31, 33, 38

Bin 1:- [3 7 14] Bin 2:- [19 23 24] Bin 3 = $\frac{9}{3} = 3$

\Rightarrow Mean:- [8, 8, 8] [22 22 22] [34, 34, 34] in each bin

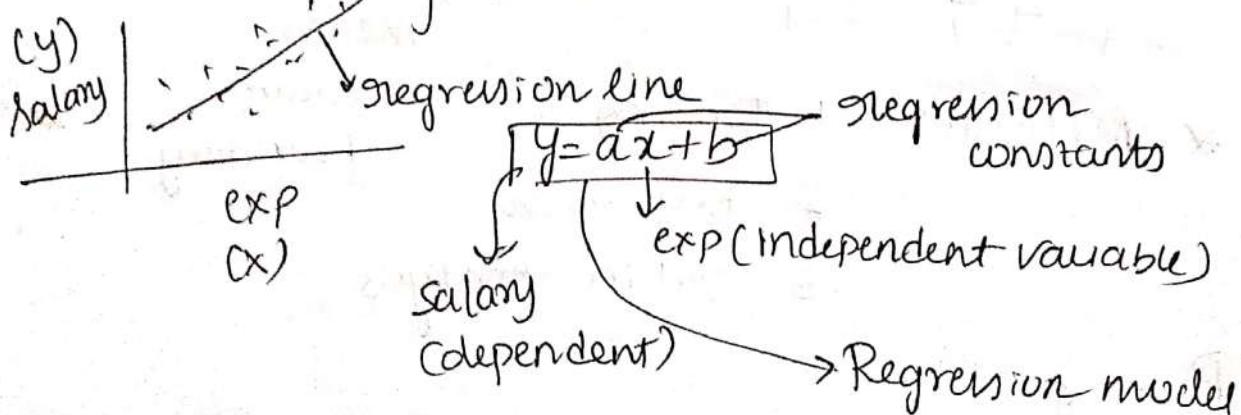
\Rightarrow Median:- [7 7 7] [23 23 23] [33 33 33]

\Rightarrow Boundary [3 3 14] [19 24 24] [31 31 38]

↓
Closest boundary

Value for
intermediate
values.

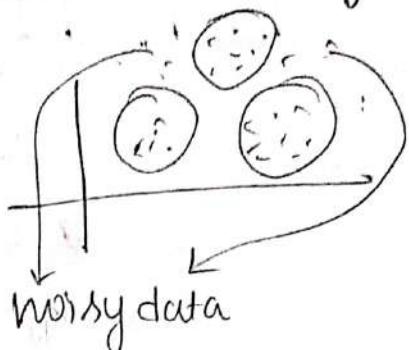
* Regression: Data can be smoothed by fitting the data to a regression function



Turing this we
replace / predict noisy
value.]

Outlier Analysis

- * We group data points based on similarity measures.
- * Here noisy data is identified & eliminated



Data integration

why:- When data required is available in different sources.

what:- Preprocessing technique, which integrates the data from diff sources into one to provide single view of data. & it also resolves the problems occurred during integration.

- problems / issues
- ▷ Schema integration.
[entity identification, Prob]
 - 2) Data value conflict Detection & resolution (DVC D & R)
 - 3) Redundancy.

▷ EI Prob

→ Objects that represents to the same real world entity are represented in diff way in diff tables.

for ex:-

CIId	Name/N	
1	1	1
2	2	2

C-Id	Name	N,
1	1	1
2	2	2

Name but ~~C-Id & CIId~~

* C-Id & CIId are represented as diff here which are same actually.

Sol:- Use metadata.



which has Name, datatype, meaning range of values allowed, & null rules to handle (or null & blank values) for attribute.

* So here metadata tells C-Id & C-Id are same, hence EI prob is resolved.

(ii) DVCD & R

↳ data is same but some attributes are represented in different units like Rs & \$ for money.

Sol:-

Use metadata

Redundancy

→ Derivable attribute

→ same attribute with diff names.

* Ex:- Age, DOB
 ↓ ↓
 in 1st in 2nd
 table table

* This occurs bcs there is relationship b/w attributes

Solution:-

* So we have to identify R.s b/w them

(1) correlation/covariance Analysis

↳ for numerical

(2) Chi square test

↳ for categorical.

Correlation

$$\rho_{AB} = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{(n-1) \sqrt{A} \sqrt{B}}$$

value of att in i^{th} tuple means
 \bar{A} \bar{B}

-0.95
 (strong -ve)

if ρ {
 < 0 -ve
 > 0 +ve
 = 0 No relation (independent)

$-1 < \rho < +1$

* It gives "direction & strength"

Covariance

$$\text{cov}(AB) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$$

-0.95
 (-ve)

* It gives only "direction"

if ρ {
 < 0 -ve
 > 0 +ve
 = 0 0nd

+0.2
 (+ve)

* Chi-square test (categorical data)

Null hypothesis - A & B are independent (default)

$$\chi^2 = \sum \frac{(\text{observed} - \text{expected})^2}{\text{expected}}$$

Example

	Male	Female	Sum (row)
Coffee	200 (90)	250 (360)	450
Tea	100 (210)	950 (840)	1050
Sum (col)	300	1200	1500

$$e_{11} = \frac{300 \times 450}{1500} = 90 \text{ ex}$$

$$\chi^2 = \sum_{i=1}^n \frac{(\text{ob}-\text{ex})^2}{\text{ex}}$$

$$= \frac{(200-90)^2}{90} + \frac{(250-360)^2}{360} + \frac{(100-210)^2}{210} +$$

$$\chi^2_{\text{cal}} = \frac{(950-840)^2}{840} = 240.07 \quad \begin{array}{l} \text{Reject Null} \\ \text{Hypothesis} \end{array}$$

$$\text{DOF} = (2-1)(2-1) = 1 \quad \chi^2_{\text{cal}} > \chi^2_{\text{obs}}$$

$$\text{Sig} = 0.001 \rightarrow \chi^2_{\text{obs}} = 10.828 \quad \therefore \text{Related}$$

Data Reduction (DR)

- Huge data analysis diff & inefficient hence DR is needed (bcs of redundancy, noise / null values)
- It is used to reduce volume of data, so pattern extraction will be quicker and efficient.
- Results will be same as the pattern extraction done without DR.

Methods

- 1) Dimensionality Reduction
 - removes redundant or irrelevant attr
- 2) Numerosity
 - we used attribute subset selection method (Dimer R)
(CASS)
 - Numerosity :- Alternate Representation of data
 - Parametric method
 - Non " "
 - 3) ASS - heuristic attribute selection method
 - 1. Stepwise forward selection
 - 2. " " backwards elimination
 - 3. Best combined attribute selection & elimination. (BCASE)

Stepwise F8

Attributes :- {A1, A2, A3, A4, A5, A6, A7, A8}

selection threshold :- '5'

Procedure

- ▷ { }
2) { A4 }
- 3) { A4 A7 }
- 4) { A4 A7 A2 }
- 5) { A4 A7 A2 A6 }
- 6) { A4 A7 A2 A6 A3 }

\downarrow $\brace{ \quad }$

$n=5$

S imp att

Step-wise B.C

same eliminates from least
produce

$\triangleright \{A_1, A_2, A_3, A_4, A_5, A_6, A_7, A_8\}$ $\xrightarrow{\text{to feature extraction}}$

27 Jun 11 Agy

B CASE

* In this both forward selection & backward elimination occurs simultaneously.

Procedure

1) { A₄ added, A₈ eliminated }

2) { A₄, A₇ added, A₁ " }

3) --- similar way.

* In this way, there 3 methods reduces no. of att based on threshold.

Xlumosity - R (X|R)

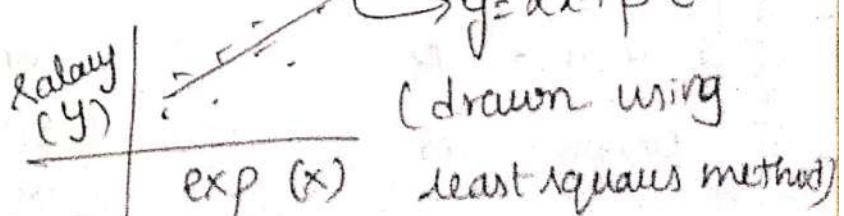
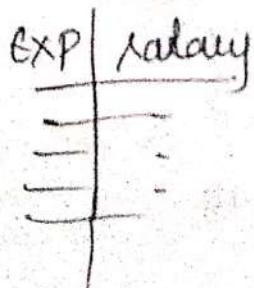
▷ Parametric { linear-R
 { mult. linear-R

2) Non " { histograms
 clustering
 Sampling.

(Parametric value
regression
coefficients)

Parametric

1) linear-R



* Only RC (α, β) are retained

(ii) Multiple L.R

* In this $\alpha_1, \alpha_2, \alpha_3, \dots, \beta$ are retained.

$$\text{where } Y = \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \dots + \beta$$

This reduces volume by storing parametric values of data model.

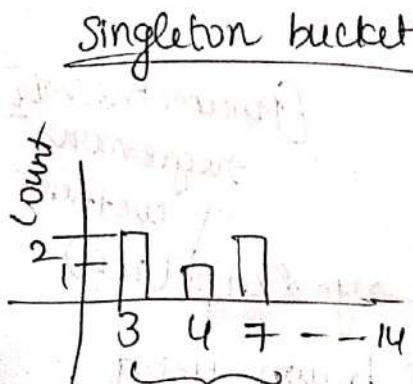
Non Parametric methods

Histogram :- By storing $\text{avg}(\text{sum})$ for each bin

Bin - single values or range of values

* Here Bin & occurrences of values will be shown

3, 4, 7 - 3 - 7 - 14

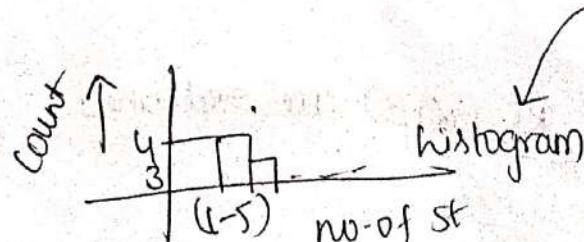


Equal width method

$$w = \frac{\text{max} - \text{min}}{n} \quad n \rightarrow \text{no. of bins}$$

$$w = 5$$

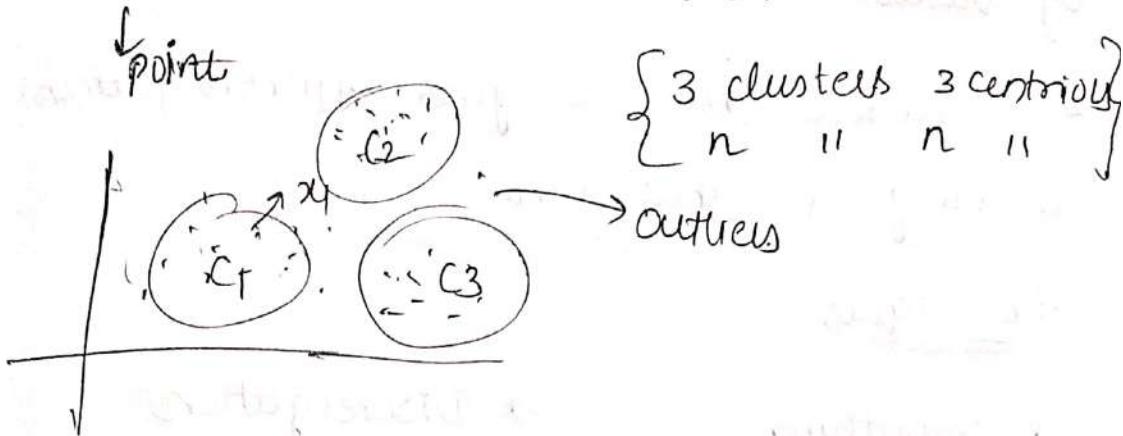
Width	Count
1-5	4
6-10	3
- -	1
20-25	10



Clustering :- Store cluster parameters

Eg centroid & diameter only

$$x_1 = (3, 4), x_2(5, 6) \dots x_n()$$



* Here we are storing (centroid & diameter) values instead of all data points.

Sampling → Reduces data by sample size.

* Simple random sampling (SRS)

* SRS without replacement (SRSWOR)

* " with " (SRSWR)

* Stratified sampling (SS)

SRSWOR

2 samples

with size 3

SRSWR
same

S_1	S_2
T_2	T_2
T_4	T_7
T_6	T_9

S_1	S_2
T_2	T_2
T_4	T_7
T_6	T_9
T_8	

T_1	L	
T_2	M	
T_3	L	
T_4	L	
T_5	M	
T_6	L	
T_7	H	
T_8	H	
T_9	H	
T_{10}	H	

Income

SD. state
 T_2

SS

$S_1 = T_1$

T_3
 T_4
 T_6

$S_2 = T_2$
 T_5
 M

$SD. = T_2$

Data Transformations

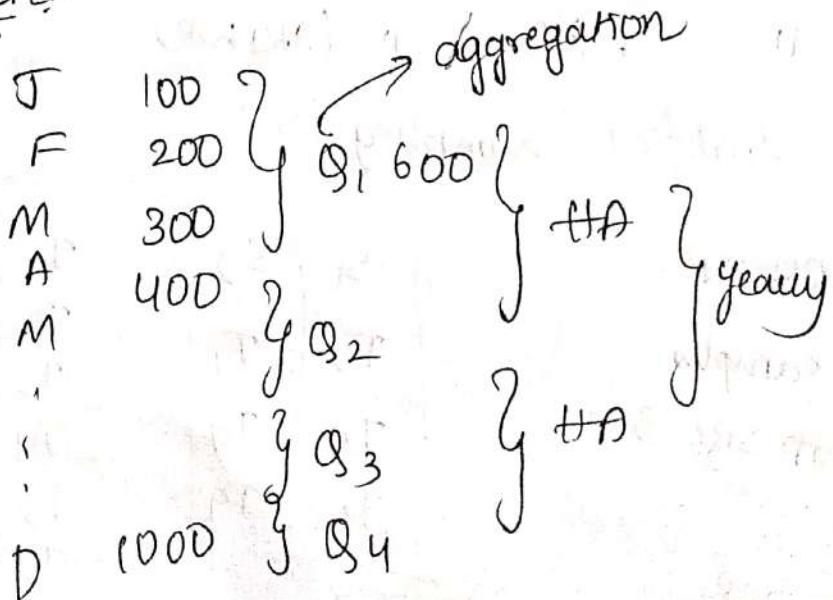
- Attribute values are mapped to new set of values.
- Needed bcs it gives efficient patterns & easy to understand techniques

- * Smoothing
- * Discretization
- * Aggregation
- * Concept hierarchy
- * Normalization
- * Generation

Smoothing :- Using Binning (removes spikes)

$$[3 \ 7 \ 2] \xrightarrow{\text{Mean-B}} [4 \ 4 \ 4]$$

Aggregation:



Normalization

* Gives equal weights to all attr

Methods

* 1. Min-Max-N

2. Z-score-N

3. Normalization by decimal scaling.

Min-Max-N

→ Maps values to "1" or "0" or any values b/w 0 & 1

→ min=0 max=1 (rest in b/w 0 & 1)

$$V' = \frac{(v - \text{min}_A)}{(\text{max}_A - \text{min}_A)}$$

Income

$$12222 \text{ (min)} = 0$$

$$15000 \text{ (b/w value)} = \frac{(15000 - 1222)(1-0)}{(15000 - 1222)} = 0.038$$

$$85000 \text{ (max)} = 1 = \frac{(85000 - 1222)(1-0)}{(85000 - 1222)}$$

2) Z-score (used when min & max values are not known)

$$V' = \frac{v - \text{mean}_A}{\sigma_A}$$

$$V' = \frac{73600 - 54000 \text{ (mean)}}{16000 \text{ (S.D.)}} = 1.225$$

if $V=U$ then $V^I=0$

if $V>U$ then $V^I=V$

if $V < U$ then $V^I=-V$.

Normalization by decimal scaling

→ Moving the decimal point of values of attr A

→ Depends on the max absolute value of A

$$V^I = V / 10^J$$

Ex $-986 \rightarrow +918$

$$= V^I \Rightarrow \frac{-986}{10^3} = -0.986$$

\downarrow
No. of digits

$$\downarrow \quad \downarrow \quad \downarrow$$
$$+ \frac{918}{10^3} = 0.918$$

Data Discretization

→ Data values are divided into intervals

age : $20, 30 \dots 40 \dots 80$
 $\underbrace{\quad\quad}_{\text{young}} \quad \underbrace{\quad\quad}_{\text{middle}} \quad \underbrace{\quad\quad}_{\text{senior}}$

{ equal freq
equal width

→ It labels data &
replaces data values
of actual data

Concept Hierarchy Generation

→ lower level concepts are mapped to higher level.

→ Replaces lower-l.c to H.l.c.

Ex:- Country

|
State

|
City

|
Street (replaced by country)

* Causation:- Causality mean two at A & B

have a cause & effect relationship b/w them

→ A & B appear at same time, one after other or wont appear together etc.

* Correlation:- strength & direction of relation
ship b/w A & B

LP

- 3) 200, 400, 800, 1000, 2000, 2200,

Min-Max Normalization

$$\text{min} = 200 \quad \text{max} = 2200$$

$$\text{new min} = 0 \quad \text{new max} = 1$$

$$V_i' = \frac{(V_i - \text{min})(\text{new max} - \text{min})}{(\text{max} - \text{min})}$$

$$V_1' = \frac{(200 - 200)(1 - 0)}{(2200 - 200)} = 0$$

$$V_5' = \frac{(2000 - 200)(1)}{2000}$$

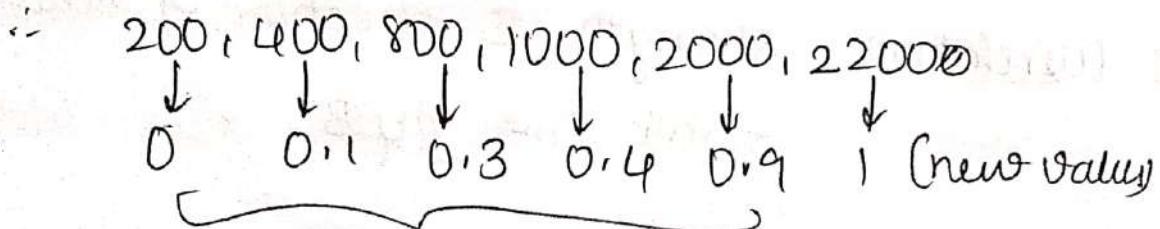
$$V_2' = \frac{(400 - 200)(1 - 0)}{2000} = 0.1$$

$$V_6' = \frac{(2200 - 200)(1)}{2000}$$

$$V_3' = \frac{(800 - 200)(1)}{2000} = 0.3$$

$$V_4' = \frac{(1000 - 200)(1)}{2000} = 0.4$$

$$= 1$$



Values obtained after normalization

⇒ 200, 400, 800, 1000, 2000, 2200

$n=2$ (no of bins)

(i) Equal width partitioning

$$\Rightarrow W = \frac{\max - \min}{n} = \frac{2200 - 200}{2} = \frac{1000}{2}$$

Dividing data into bins!-

For bin 1:

$$\Rightarrow 1000 + 200 = 1200 \quad [\text{include data from min } (x_1 + \min) \text{ to } \frac{1200}{x_2} \text{ in bin 1}]$$

∴ Bin 1: [200, 400, 800, 1000]

For bin 2:

$$\Rightarrow 1200 + 1000 = 2200 \quad [\text{from } 1200 \text{ to } 2200] \\ [\frac{x_2 + x_1}{2}]$$

∴ Bin 2: [2000, 2200]

∴ New bins are [200, 400, 800, 1000] [2000, 2200]

(ii) Equal Frequency Binning

→ All bins will have equal freq (count of n)

$n=2 \therefore$ Total count / n $\Rightarrow 6/2 = 3$ in each bin

∴ Bins are [200, 400, 800] [1000, 2000, 2200]

Bin-1 Bin-2

2) Data discretization techniques involves many, some of them are:

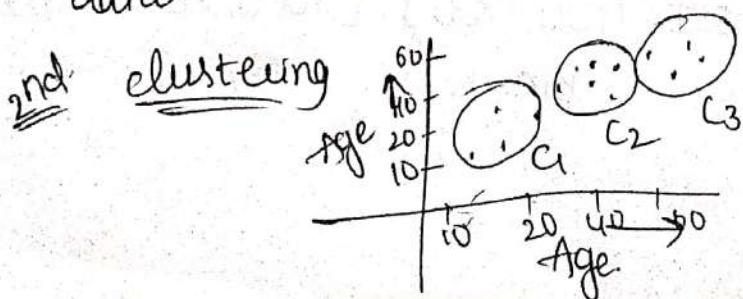
- (i) histogram analysis
 - (ii) Binning
 - (iii) Correlation analysis
 - (iv) Clustering "
 - (v) Decision tree "
 - (vi) Equal width partitioning
 - (vii) Equal depth "
- } explain them (only 2)
} and explain data transformation

Ex:- $\rightarrow 10, 11, 13, 14, 18, 20, 25, 40, 48, 60, 65, 80$
Ages input data [Before discretization]

1st We will discretize the data using labels

(i) Young (ii) Middle (iii) Old age.
↑ age. age ↓ age
 $[10 \ 11 \ 13 \ 14 \ 18 \ 20] \quad [25, 40, 48] \quad [60, 65, 80]$

Here we did just manually,
Using (vi) & (vii) etc also we can do for given
data.



(3) 10/20

Krishna

Frequent Pattern Mining - Chapter-2

(FPM)

FPM :- It is process of finding frequent patterns i.e. relevant hidden patterns from database.

Frequent Patterns :- Itemsets / sub sequences or substructures that appear together in DB.

They represent very imp patterns / properties of data.

Methods to find

- Finding frequent patterns
- Drawing Association rules.

Support & Confidence

Tid	Items Purchased
1.	Milk, Bread, Coffee
2.	Milk, Coffee
3	Milk, Butter
4	Bread, Egg.

threshold

items.
frequent, ↗
Milk = 3 ✓
Bread = 2 ✓
Coffee = 2 ✓
Butter = 1 X
egg = 1 X
if support-min = 2

Support :- We say that item is frequent only when the occurrence of item in Transaction DB is more or equal to threshold of support.

Two items frequent item net

which occur

Milk, Bread - 1 X frequently in data

Milk, coffee - 2 ✓ (comparing with

Bread, coffee - 1 X threshold)

∴ Milk, coffee frequent itemset of 2

Association rule

Milk \rightarrow coffee

Support of this rule = (Milk U coffee)

$$\Rightarrow S(M \cup C) = 2$$

Relative support $\Rightarrow \frac{S(M \cup C)}{\text{Total No. of transaction}} = \frac{2}{4} = 50\%$

Confidence = $\frac{S(M \cup C)}{S(M)} = \frac{2}{3} = 66.6\%$

* Milk \rightarrow coffee (S=2 (50%), C=66.6%).

Support- It indicates that 50% of the transactions Milk & coffee are purchased together.

Confidence- It indicates that 66% of users who purchased milk has purchased coffee.

* These 2 measures are used to find interesting rules of rule.

Apriori Algorithm

- * Method to find frequent itemsets even we have FP growth Algorithm.

Apriori Property

- * If an itemset is frequent, then all of its subsets must be frequent.
 $\{ABC\} \rightarrow$ frequent then $\{A\}$ $\{B\}$ $\{C\}$ are frequent
- * If an IS is infrequent, then all of its supersets are infrequent too.
 $\{AB\} \rightarrow$ infrequent then $\{ABC\}$ $\{AC\}$ are infrequent.

Algorithm

- * General wine search.

Join step :- C_k is generated by joining L_{k-1} with itself.

Prune step :- Any $(k-1)$ item set that is not frequent cannot be a subset of a frequent k -itemset.

Ex:-

T-id	Items Purchased
1	Milk, Bread, coffee
2	Butter, Bread, Egg
3	Milk, Bread, Butter, Egg
4	Butter, Egg

Min support = 2

Step 1 Generate C_1

C_1

Itemset	Support count
M	2
Br	3
C	1
Bu	3
E	3

L_1

Item	S-C
M	2
Br	3
Bu	3
E	3

We prune

Step



Step 2 C_2 (Joint C_1 with L_1)

C_2

Itemset	Sc
M Br	2
M Bu	1
M E	1
Br Bu	2
Br Eg	2
Bu Eg	3

Prune C_2



L_2

Itemset	Sc
M Br	2
Br Bu	2
Br E	2
Bu E	3

Step-3 (join L_2 with L_2)

C_3 (candidate table)

L_3 (3 freq item set)

Item	SC	
M Br Bu	1 X	(M Bu) not there
M Br Eg	1 X	(M E) not there
Br Bu Eg	2	

→ Use apriori algorithm here
if any of these subsets are
infrequent then that particular
set is infrequent.

Cloud & Maximal frequent Item set

* When we generate FIS with the
threshold set by user to support, there
will be problem if min-support value
is low, bcs we will get many no-of
subsets.

Downward closure Property

* Any subset of a frequent itemset must
also be frequent.

$$\{ABC\} - \{A\} \{B\} \{C\}$$

For N distinct items,

FDS will be 2^{N-1}

$$\{AB\} \{BC\} \{CA\}$$

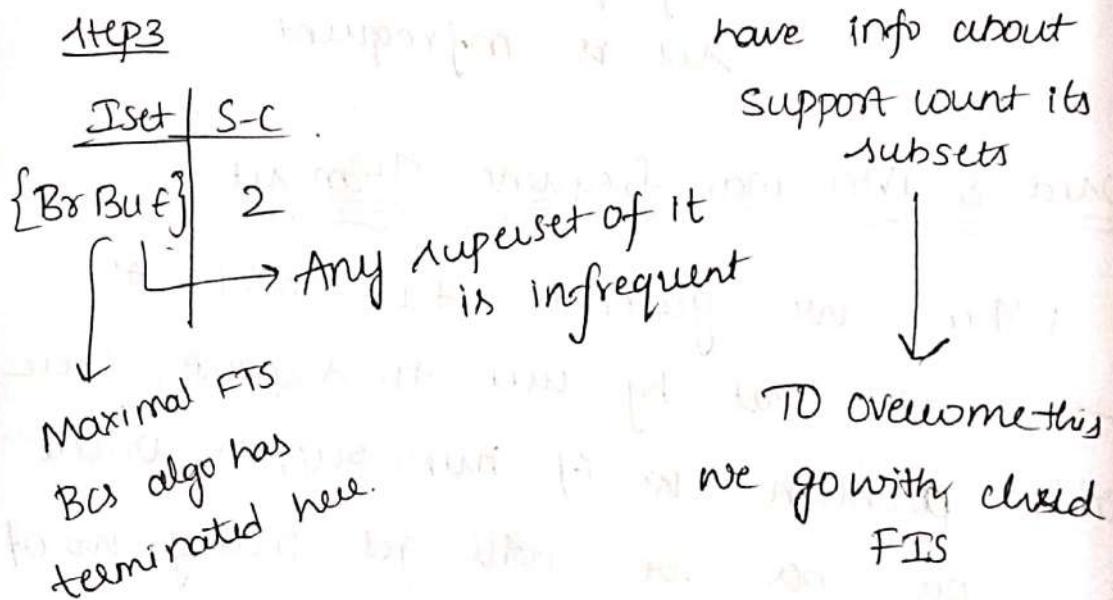
$$2^3 - 1 = 7$$

* if N is large. is very large like 100
 $2^{100}-1$ will be very huge number, hence
to overcome this problem we have sol
as closed & Maximal frequent set.

↓
Downward closure property.

Maximal FIS:- An FIS is called as maximal
FIS if its supersets are infrequent.

Consider the last ex



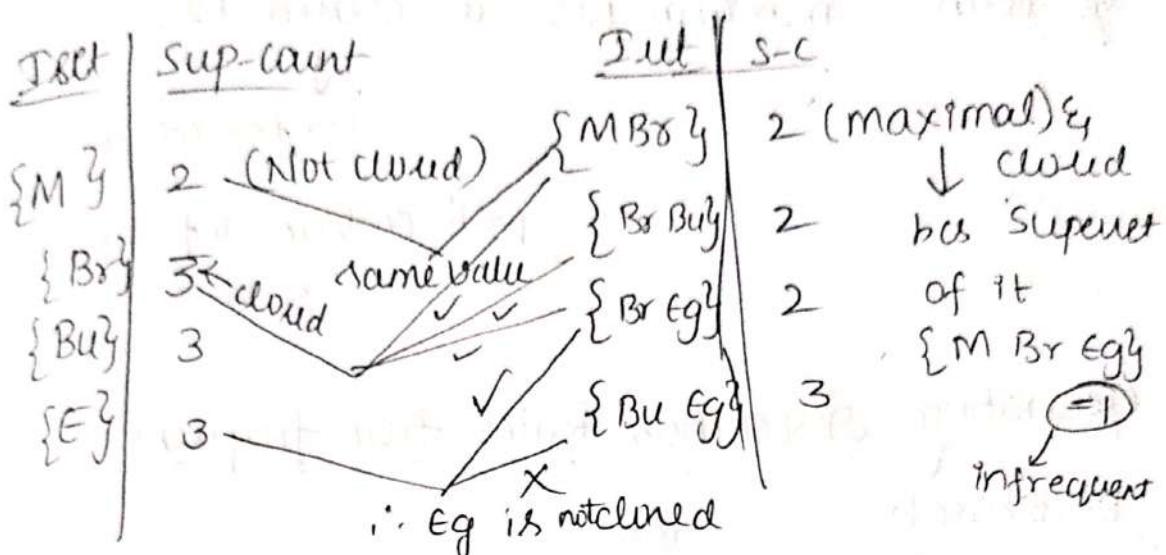
Closed FIS:- An FIS is called as closed Fis if its superset don't have same support count of it.

consider the last ex

Step 1

1

2



* If any one of superset is having same or greater value of support count, then it can't be closed.

$\{M\}, \{E\}$ are not closed, $\{Br\}$ is closed \rightarrow bcs its supersets s-c is not equal & less than its s-c.

* Maximal FIS can be closed but closed is not maximal.

Relation b/w them

* Closed \in FIS
* Maximal \in closed & FIS

Maximal \subset closed \subset FIS



Q. When we have many subsets with given min-support using apriori algorithm, we can go with maximal FIS or closed FIS

↓
preferred as it will contain info about its subsets.

Generating Association Rules from frequent k-itemsets

(AR)

Ex

T-id	Items Purchased
1.	Milk, Br, cof
2.	Bu, Br, Eg
3.	M, Bu, Br, Eg
4	Bu, Eg

Using apriori Algo, we found FIS as $\{Bu, Br, Eg\}$

Rules of generating AR

- * For each FIS I , generate all nonempty subsets of I

$$\{Br, Bu, Eg\} \quad \{Br, Bu\} \quad \{Br, Eg\}$$

$$\{Br, Bu, Eg\} \quad \{Br, Bu\} \quad \{Br, Eg\}$$

* For every non empty subset of S of I , O/P

the rule " $S \rightarrow \{I-S\}$ " if $\frac{Sc(I)}{Sc(S)} \geq \text{mincof}$
↓ ↓
Subset remaining subsets. then
we can
generate rule.

Strong Rules

$$|\text{Min-cof} = 70\%|$$

R1: $\{Br, Bu\} \Rightarrow Eg$ condition to check: $\frac{Sc(Br, Bu, Eg)}{Sc(Br, Bu)}$

$$\Rightarrow \frac{2}{2} \approx 100 = 100\%$$

∴ Strong Rule. $\checkmark \geq 70\%$

R2: $\{Bu, Eg\} \Rightarrow \{Bu\} \Rightarrow \frac{2}{Sc(Bu, Eg)} = \frac{2 \times 100}{3} = 66.67\% \leq 70\%$
∴ Not strong rule. \times

Other Method of finding FIS

FP growth {Frequent Pattern growth} Algorithm

Drawbacks of Apriori

- * Generates large candidate sets if I is in the database is large.
- * Needs multiple scans of database.

To overcome this drawback, we go with FP growth algorithm.

FP GM

- * Divide & conquer strategy
- * It compresses the database into a frequent pattern tree.
- * It divides the compressed database into a set of conditional databases.
- * Each conditional database is associated with one frequent item.

Example Min-Sup=3

T-id	Items purchased
1	I ₁ , I ₂ , I ₃
2	I ₂ , I ₃ , I ₄
3	I ₄ , I ₅
4	I ₁ , I ₂ , I ₄
5	I ₁ , I ₂ , I ₃ , I ₅
6	I ₁ , I ₂ , I ₃ , I ₄

Step 1

I₂: 5, I₁: 4, I₃: 4, I₄: 4 (in desc order)

* We have reorganise them acc to desc order of their sc.

Step-2

Tid	I.P
1	I ₂ I ₁ I ₃
2	I ₂ I ₃ I ₄
3	I ₄
4	I ₂ I ₁ I ₄
5	I ₂ I ₁ I ₃
6	I ₂ I ₁ I ₃ I ₄

Step-3

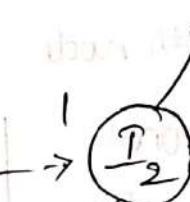
FPTree

header table

Item Id	sc	Node link
I ₂	5	-
I ₁	4	-
I ₃	4	-
I ₄	4	-

①

{ 3 }



T₁: I₂, I₁, I₃

null



null

T₃: I₄

null

T₄: I₂, I₁, I₄

null

T₅: I₂, I₁, I₃

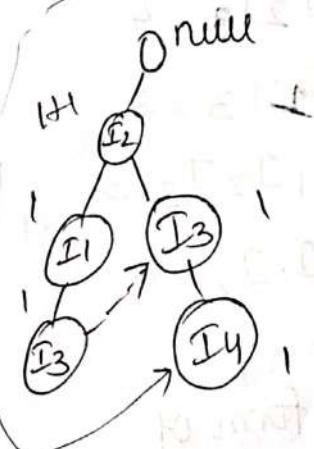
null

T₆: I₂, I₁, I₃, I₄

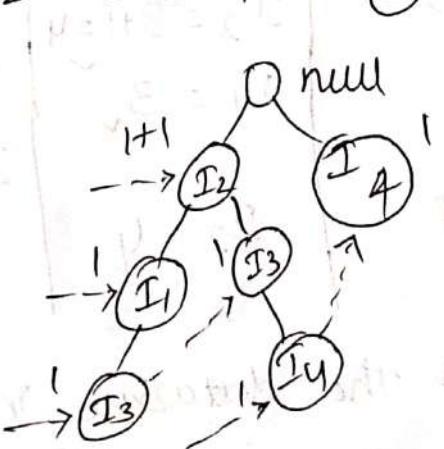
null

②

T₂: I₂, I₃, I₄



③ T₃: I₄



null

T₄: I₂, I₁, I₄

null

T₅: I₂, I₁, I₃

null

T₆: I₂, I₁, I₃, I₄

null

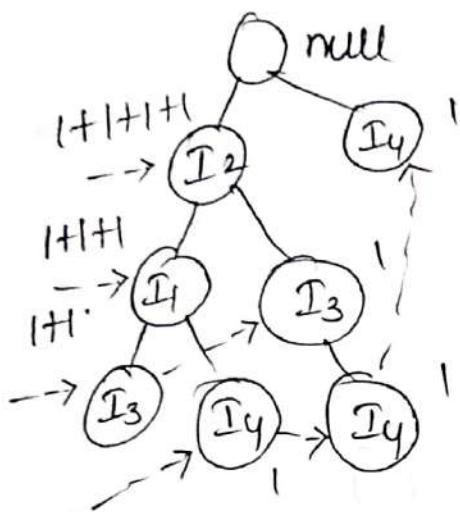
⑤ T₅: I₂, I₁, I₃

null

⑥ T₆: I₂, I₁, I₃, I₄

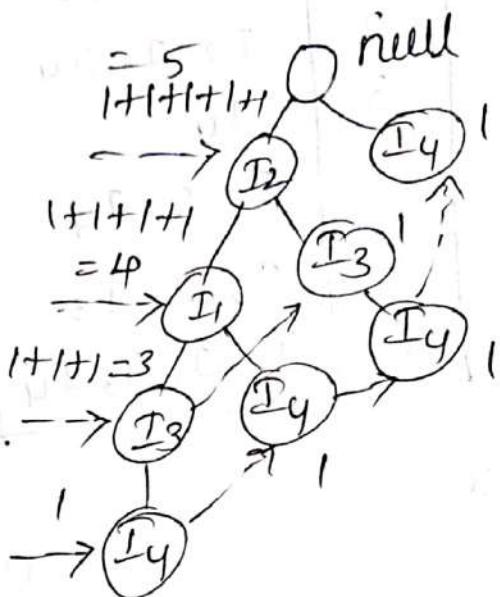
null

⑤ $I_5: I_2 I_1 I_3$



start from leaf ~~path~~ node

⑥ $I_6: I_2, I_1, I_3, I_4$



$\min-S = 3$

Item	Conditional Pattern Base	Conditional FP Tree	Freq Patterns
I_4	$I_2 I_1 I_3 : 1,$ $I_2 I_1 : 1$ $I_2 I_3 : 1$	$I_2 = 3 \checkmark$ $I_1 = 2 \times$ $I_3 = 2 \times$	$I_2 I_4 = 3$
I_3	$I_2 I_1 : 3$ $I_2 : 1$	$I_2 = 3 + 1 = 4 \checkmark$ $I_1 = 3 \checkmark$	$I_2 I_3 = 4$ $I_1 I_3 = 3$
I_1	$I_2 : 4$	$I_2 = 4$	$I_1 I_2 I_3 = 3$ $I_1 I_2 = 4$

* It compresses the database in form of FP Tree

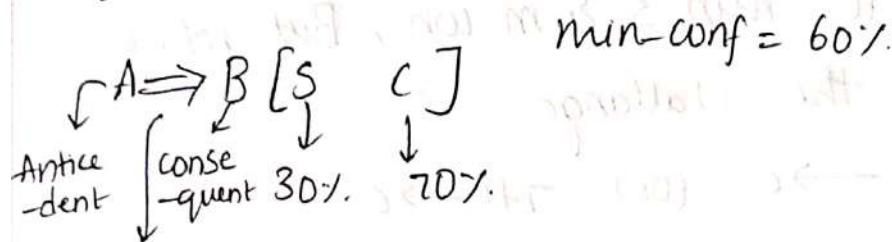
* It divides FP Tree into conditioned DB from which we extract frequent patterns

* It generates only frequent Items.

Limitations of Support & Confidence framework

Misleading Association Rules

$\{AB\}$ min-sup = 20% value



Strong rule ~~in reality~~

* In some cases this rule might not be strong rule though it satisfies min-sup & min-conf.
So, let's solve that problem.

Example

	Horlicks	Not Horlicks	Sum (Total)
coffee	350	800	1150
Not coffee	90	10	100
Sum (Total)	440	810	1250

R1:

$H \Rightarrow C$ [$S=28\%$, $C=79\%$] ✓ strong rule.

$$S = \frac{S(H \wedge C)}{\text{Total} \cdot T}$$

$$= \frac{350 \times 100}{1250} = 28\%$$

$$C = \frac{S(H \wedge C)}{S(H)}$$

$$= \frac{350 \times 100}{440} = 79\%$$

R2: $H \Rightarrow \neg C$ [$S=7.2\%$, $C=20.4\%$] ✗ not strong rule.

$$S = \frac{S(H \wedge \neg C)}{T \cdot T} = 7.2\%, \quad C = \frac{S(H \wedge \neg C)}{S(H)} = \frac{90}{440} \times 100 = 20.4\%$$

$R_3 \vdash H \Rightarrow C [S = 64\%, C = 98.7\%]$

✓ strong rule

* Here R_1 & R_3 both are strong rules according to min-S & m-won, But which is cut is the challenge

$$H \rightarrow C \text{ (or) } \neg H \rightarrow C$$

Here the problem occurred due to conf

* To overcome this, we go with '2' rules

LIFT & Chi-square test

$$\text{Lift : Lift}(H,C) = \frac{C(H \rightarrow C)}{S(C)}$$

$$\begin{aligned} R_1 & \quad S(C) \\ &= \frac{S(H \cup C)}{S(H)} \\ &= \frac{S(H \cup C)}{S(H)S(C)} \end{aligned}$$

$$\begin{aligned} \text{Lift}(H \cup C) &= \frac{93}{S(H \cup C)} = 0.86 \\ R_3 & \quad S(C) \\ &= \frac{90}{S(H)S(C)} \\ &= \frac{90}{1250} \end{aligned}$$

$$\begin{aligned} \text{Lift} & \left\{ \begin{array}{l} \geq 1 - \text{Independent} \\ \leq 1 - \text{correlation is there} \\ < 1 - \text{not correlated} \end{array} \right. \\ & \frac{440}{1250} \times \frac{100}{1250} = 2.55 \end{aligned}$$

$\therefore R_1$ is not strong ($0.86 < 1$)

R_3 is strong ($2.55 > 1$)

Chi-square

$$\chi^2_{\text{cal}} = \sum \frac{(\text{Obs} - \text{Exp})^2}{\text{Obs}}$$

$\chi^2_{\text{cal}} = 0$ (Independent)

$\chi^2_{\text{cal}} > 0$ (+ve or -vely correlated)

* First calculate expected values for all obs
then calculate dof then χ^2_{cal}

if $\chi^2_{\text{cal}} > \chi^2_{\text{obs}}$ (from table)

→ Reject H₀ (null hypothesis)

which states they are independent

if Obs = 350

and Exp = 440

then χ^2 is 0 / mean 350-440 is -ve
∴ -vely correlated.

if Obs = 90

Exp = 20 $90-20 = +ve$ ∴ +vely correlated.

* Both Lift & Chi-square are correlated measures.

* They overcome the problems of Support & confidence.

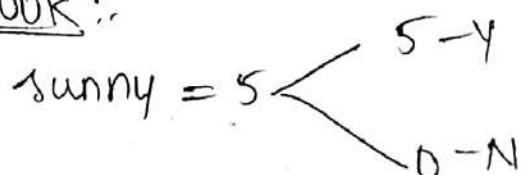
$$15) \quad Y_{10} = 7(P) \\ N_0 = 3(n)$$

$$\log_2 x = \frac{\ln x}{\ln 2}$$

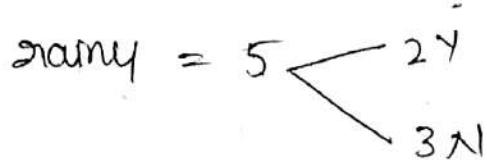
$$I_D = \frac{-7}{10} \log_2 \left(\frac{7}{10} \right) - \frac{3}{10} \log_2 \left(\frac{3}{10} \right) \\ = -0.7 \times (-0.5145) - 0.3 \times (-1.7369) \\ \approx 0.8813 \text{ J/J}$$

$$\log_2 \left(\frac{7}{10} \right) \\ = \frac{\ln(0.7)}{\ln(2)}$$

Outlook:



$$\sum \frac{P_i + n_i}{P+n} = \frac{P_i + n_i}{P+n} I(p_i, n_i) +$$



$$\frac{n_i + p_i}{P+n} I(p_i, n_i)$$

$$I_{outlook} = \frac{(5+0)}{10} I(5,0) + \frac{(2+3)}{10} I(2,3) \\ \left(\frac{7+3}{10} \right) \quad \downarrow \quad \left(\frac{7+3}{10} \right)$$

$$= \frac{5}{10} \left[- \underbrace{\frac{5}{5} \log_2 \frac{5}{5}}_0 - 0 \right] + \frac{5}{10} \left[- \frac{2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right]$$

$$\approx 0.2 [-0.4 \times -1.3219 - 0.6 \times -0.7369]$$

$$\approx 0.4855$$

$$Gain_{outlook} = I_D - I_{outlook} = 0.8813 - 0.4855 = 0.3958$$

$$G_{outlook} = 0.3958$$

Com

company

$$\text{big} = 3 \begin{cases} 1Y \\ 0N \end{cases}$$

$$\text{Med} = 4 \begin{cases} 3Y \\ 1N \end{cases}$$

$$\text{NO} = 3 \begin{cases} 1Y \\ 2N \end{cases}$$

$$I_{\text{company}} = \frac{3}{10} I(3, 10) + \frac{4}{10} I(3, 1) + \frac{3}{10} I(1, 2)$$

$$= \frac{3}{10} + 0.4 \left[-\frac{3}{4} \log_2 \left(\frac{3}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right]$$

$$+ \frac{3}{10} \left[-\frac{1}{3} \log \left(\frac{1}{3} \right) - \frac{2}{3} \log \left(\frac{2}{3} \right) \right]$$

$$= 0.4 \left[-0.75 \times -0.41503 - 0.25 \times -2 \right] +$$

$$0.3 \left[-0.33 \times -1.5849 - 0.66 \times -0.5849 \right]$$

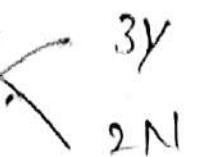
$$= 0.3245 + 0.2727$$

$$= 0.5972 = 0.6$$

$$\boxed{\text{Gain}_{\text{company}} = 0.2813}$$

sailboat

Small = 5 

Big = 5 

$$\begin{aligned} I_{\text{sailboat}} &= \frac{5}{10} I(4,1) + \frac{5}{10} I(3,2) \\ &= \frac{5}{10} \left[\frac{-4}{5} \log_2 \left(\frac{4}{5} \right) - \frac{1}{5} \log_2 \left(\frac{1}{5} \right) \right] + \\ &\quad \frac{5}{10} \left[\frac{-3}{5} \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} \log_2 \left(\frac{2}{5} \right) \right] \\ &= 0.5 [-0.8 \times 0.3219 - 0.2 \times -2.3219] + \\ &\quad 0.5 [0.4421] \\ &= 0.8465 \end{aligned}$$

$$I_{\text{Gain}} = 0.0348$$

So, the outlook has higher gain value. \therefore it becomes node

Outlook

Sunny

Rainy

Company	Sailboat	Class
big	small	yes
med	small	"
med	big	"
no	small	"
big	big	"

done

Company	Sailboat	Class
no	small	no
med	small	yes
big.	big	yes
no	big	no
med	big	no

Consider this table
now

$$Yes = 2(P) \quad P+n=5$$

$$No = 3(n)$$

$$ID = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log\left(\frac{3}{5}\right) = 0.9709$$

company

$$\begin{aligned} I_{comp} &= \frac{2}{5} I(0,2) + \frac{1}{5} I(1,0) \\ &\quad + \frac{2}{5} I(1,1) \\ &= 0.4 \end{aligned}$$

$$\begin{aligned} no &= 2 \quad 0 \quad 1 \\ big &= 1 \quad 1 \quad 0 \\ med &= 2 \quad 1 \quad 0 \end{aligned}$$

$$I_{gain} = 0.5709$$

sailboat

Small = 2 $\begin{cases} 1Y \\ 1N \end{cases}$

big = 3 $\begin{cases} 1Y \\ 2N \end{cases}$

$$\begin{aligned}
 I_{\text{sailboat}} &= \frac{2}{5} I(X_1|1) + \frac{3}{5} I(X_1|2) \\
 &= \frac{1}{5}(1) + \frac{3}{5} (0.2753) \\
 &= 0.4 + 0.1636 \\
 &= 0.9510.
 \end{aligned}$$

$$G_{\text{sailboat}} = 0.0199$$

