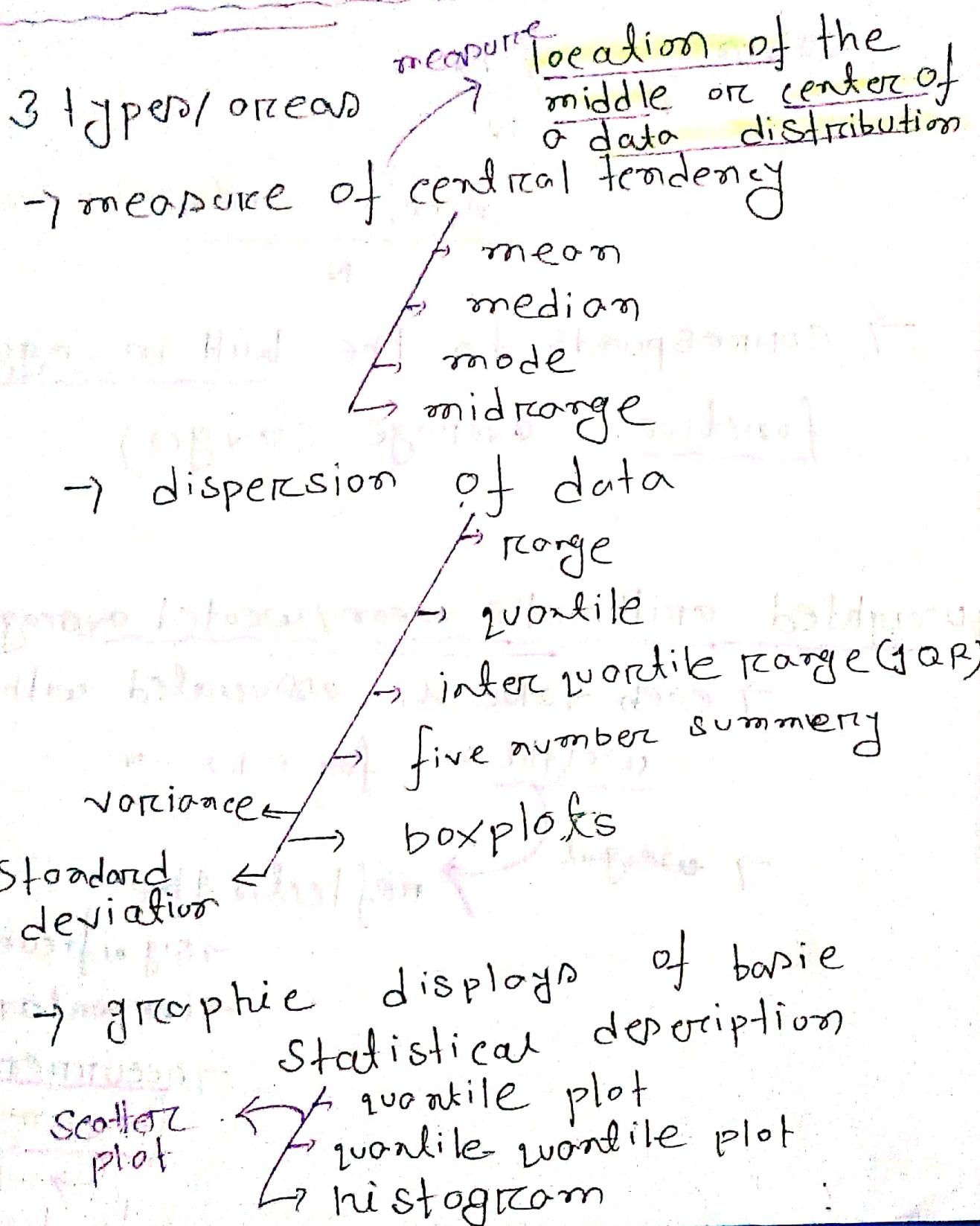


## Chapter 2

### Getting to your Data:

#### Basic Statistical Structure of Data:



## Mean / arithmetic mean

Let  $x_1, x_2, \dots, x_n$  be a set of observations.

the mean of this set of values,

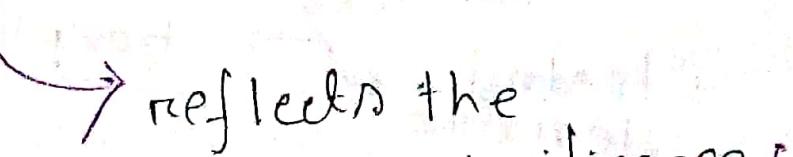
$$\bar{x} = \frac{\sum_{i=1}^n x_i}{N}$$

$$= \frac{x_1 + x_2 + \dots + x_N}{N}$$

→ corresponds to the built-in aggregate function - average (avg)

## Weighted arithmetic mean / weighted average:

→ each value  $x_i$  associated with a weight  $w_i$  for  $i=1, 2, \dots, N$ .

→ weight  reflects the

- significance
- importance
- occurrence
- frequency

attached to their respective

$$\bar{x} = \frac{\sum_{i=1}^N w_i x_i}{\sum_{i=1}^N w_i}$$

$$w_1x_1 + w_2x_2 + \dots + w_Nx_N$$

$$w_1 + w_2 + \dots + w_N$$

Problem of mean:

→ sensitivity to extreme / outlier values

outliers

Small number of extreme values  
can correct the mean

Sol<sup>a</sup>: trimmed mean

mean obtained after  
chopping off values at the  
high & low extreme.

→ remove top & bottom 2%

- Median: → expensive to compute
- better for skewed data → large number of observations
- middle value in a set of ordered data values.
- $N$  is odd
- $N$  is odd → middle value of the ordered set
- $N$  is even → two middlemost values & any values in between  
→ conventionally average of the two middlemost values

### For grouped data

- data are grouped in intervals  
in a.t. their  
→ data values  
→ frequency

$$\text{median} = L_1 + \left( \frac{\frac{N}{2} - (\sum f_{\text{prev}})_L}{f_{\text{req median}}} \right) \text{width.}$$

where,

$L_1$  = Lower boundary

$(\sum f_{\text{prev}})_L$  = sum of the frequencies of intervals that are lower than median"

$f_{\text{median}}$  = frequency of the median.  
width = width of the median interval.

## ■ Mode:

Mode for a set of data is

value that occurs most frequently in the set.

Unimodal: Data set with one mode

bimodal: " " two modes

trimodal: " " three "

multimodal: " " with two or more modes

for unimodal numeric data  
Relation among mean, mode, median: moderately skewed  
 $\text{mean} - \text{mode} \approx 3(\text{mean} - \text{median})$

## Midrange:

Average of the largest and smallest values in the set.

## Skewed data/asymmetric data

2 types

→ positively skewed

→ negatively skewed

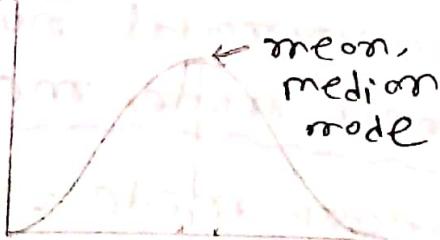
mode occurs at a value, smaller than the median ( $\underline{\text{mode}} < \text{median} < \text{mean}$ )

mode occurs at a "greater" than the median.

( $\text{mean} < \text{median} < \underline{\text{mode}}$ )

## Symmetric data:

$\text{mean} = \text{median} = \text{mode}$



## Measuring the Dispersion of Data:

- Range
- Quartiles
- Deciles
- Percentiles
- Interquartile Range

### Range:

Difference between the largest and smallest values.

$$\therefore \text{Range} := \max() - \min()$$

### Quartiles:

Points taken at regular intervals of data distribution.

dividing it into essentially equal size consecutive data sets.

★  $k$ -th quantile for a given distribution

is the value  $\underline{\mathcal{N}}$

such that  $\begin{cases} \rightarrow \text{at most } k/2 \text{ data values } \leq \mathcal{N} \\ \rightarrow \text{ " " } (2-k)/2 \text{ " " } > \mathcal{N} \end{cases}$

$\rightarrow k$  is an integer  $\boxed{0 < k < 2}$

$\rightarrow 2-1$  quantiles

■ Quartiles:

→ 4 quartiles / each of the four

distributed part of data is known

as quartile:

■ Percentiles:

→ 100 quantiles are commonly referred to as percentiles.

→ divide the data distribution into 100 equal-sized consecutive data.

## Box quartile range (QQR)

Distance between the first, and third quartiles is called QQR and defined as,

$$QQR = Q_3 - Q_1$$

→ cut off the lowest 25% of the data

→ 1st quartile ( $Q_1$ )

→ 2nd " ( $Q_2$ ) → it gives the center of the data distribution

→ 3rd " ( $Q_3$ ) → cut off the lowest 75% of data

## Rule of thumb

→ used for identifying suspected outliers

values falling at least

$$1.5 \times QQR$$

above the 3rd quartile  
below the first "

Outliers

## Five number Summary:

Consists of → median ( $Q_2$ )

→ the quartiles ( $Q_1$  &  $Q_3$ )

→ Smallest and largest individual observation

→ written in the order of

minimum,  $Q_1$ , Median,  $Q_3$ ,

maximum.

Boxplot → proper way of visualizing a distribution.

→ can be computed in  $O(n \log n)$  time

→ ends of the box are at quartiles

box length is  
the IQR

→ median is marked by

a line within the  
box

→ 2 lines outside the box extend to the

(whiskers)

→ Smallest / minimum

→ & largest / maximum

→ cut

extended to the extreme low and high observations only if

values are less than  $1.5 \times \text{IQR}$

→ Outliers are indicated by dot (•)

## Variance:

Variance of N observations.

$x_1, x_2, \dots, x_N$

for numeric attribute

$$\begin{aligned} \text{Variance } \sigma^2 &= \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \\ &= \left( \frac{1}{N} \sum_{i=1}^N x_i^2 \right) - \bar{x}^2 \end{aligned}$$

## Standard deviation:

Square root of the variance ( $\sigma^2$ )

→ denoted by  $\sigma$

## Properties

## Basic properties of Standard deviation ( $\sigma$ )

→ measures spread about the mean.

→ considered only when

mean is chosen  
as the measure of

- $\sigma = 0$ , there is no spread
- all observations have the same value.
- otherwise  $\sigma \rightarrow 0$

Chebyshev's inequality:



Shown that  $\text{at least } (1 - \frac{1}{k^2}) \times 100\%$  of the observations are no more than  $k$  standard deviations from the mean.

## Chapter-6

Mining Frequent Patterns,

Associations and Correlations:  
Basic Concepts and Methods

→ Frequent itemset mining leads to

discovery of association & correlations

among items in large transaction / relational datasets

→ discovery of interesting correlation relationships

can help in

→ catalog design

→ cross-marketing

→ customer shopping behavior analysis

## Example of frequent pattern mining:

market basket analysis

→ help retailers to develop

marketing strategies

by gaining insight  
into which items  
are frequently purchased

tog ether by  
customer

## Item set:

A set of items referred  
to as an itemset.

→ denoted by  $I = \{I_1, I_2, \dots, I_m\}$

## k-itemset:

An itemset that contain k-items

$$X = \{x_1, x_2, \dots, x_k\}$$

■ Occurrence frequency of an itemset / frequency

support count / count:

Number of transactions that contain the itemset.

→ also called absolute support

■ relative support / support  $(A \Rightarrow B)$

$$\text{Support } (A \Rightarrow B) \leq P(A \cup B)$$

probability that a transaction contains  $A \cup B$

■ confidence  $(A \Rightarrow B)$ :

$$= \frac{\text{Support}(A \cup B)}{\text{Support}(B)}$$

$$(\text{confidence } (A \Rightarrow B), e = P(B|A))$$

conditional probability that a transaction having  $B$  also contains  $A$

## frequent itemset:

A itemset  $J$  is frequent if  
 $J$ 's support is no less than  
minimum support threshold.

i.e  $\text{support}(J) \geq \text{minsup}(J)$

## Strong:

Rule that satisfied both  
a minimum support threshold

→ and a "confidence"

## Association Rule:

It is an implication of  
the form  $A \Rightarrow B$ , where  $A \neq \emptyset$ ,  
 $B \neq \emptyset$ ,  $A \cap B = \emptyset$

## Process of association rule mining:

two step process:

→ find all frequent itemsets

→ generate strong association rules

from the frequent itemset.

→ rule must satisfy

minimum support

.. confidence

## Closed itemset:

An itemset  $X$  is closed in a data set  $D$  if there exists no proper super itemset  $Y$

such that  $\text{sup-count}(Y) = \text{sup-count}(X)$

## Closed frequent itemset:

An frequent itemset  $X$  is a closed frequent itemset in set  $D$  if  
 $\rightarrow X$  is both closed and frequent in  $D$

## ■ maximal frequent itemset / (max-itemset)

→  $x$  is frequent

→ no super-itemset  $y$

such that  $x \subset y$

→  $y$  is frequent

## ■ proper super itemset:

$y$  is a proper super-itemset of  $x$ , if

$x$  is a proper sub-itemset of  $y$ . i.e.

$x \subset y$

every item of  $x$  is contained in  $y$ .

at least one item of  $y$  that is not in  $x$

## Algorithm to find frequent itemset:

1994

→ Apriori → R. Agrawal, R. Srikant

→ FP-growth → Han, Pei, Yin

→ vertical data format approach  
for forming frequent itemset

for Boolean association rules

## Apriori property / Downward closure property

★ ★ All nonempty subsets of

- frequent itemset must also be frequent.

## Apriori pruning principle:

- if there is any itemset which is infrequent,

its superset should not be generated

tested

## Steps of Apriori algorithm:

★ 2 steps procedure

→ The join step

→ The prune step

The join step: denoted by  $C_k$

→ find set of candidate

→ generated by joining

→  $C_k$  is generated by joining  $L_{k-1}$  with itself

The prune step:

→ Scan  $C_k$  & compare with min-sup

→ determine  $L_k$

→ frequent k-itemset

- ~~to~~
- A priori property is used

→ to reduce the  
size of  $C_k$

## Subset testing:

→ Any  $(k-1)$ -itemset that is not frequent, can't be subset of frequent  $k$ -itemset.

→ Any  $(k-1)$  subset of  $C_k$  is not frequent, then candidate can't be frequent.

→ so can be removed from  $C_k$

potential frequent

## ■ Frequent itemset $\rightarrow$ Association Rules:

- $\rightarrow$  for each frequent itemset  $I$ ,  
generate all nonempty subset of  $I$ .
- $\rightarrow$  for every nonempty subset  $S$  of  $I$ ,  
output the rule " $S \Rightarrow (I - S)$ "

if       $\frac{\text{support-count}(I)}{\text{support-count}(S)}$   $\geq \text{minsup}$

Support Count	Itemset
2	{A}
2	{B}
2	{C}
2	{AB}
2	{AC}
2	{BC}
1	{ABC}

## Example - G-3, 6-1

Given that,

TID	List of item-IDs
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

$$\text{min-Sup} = 2$$

$$\text{min-Cof} = 70\%$$

Step 1

⑥ Scan D for count of each candidate

Itemset	Sup-Count
{I1}	6
{I2}	7
{I3}	6
{I4}	2
{I5}	2

## Step 2

Compute candidate support with minimum support count

$L_1$

Itemset	Sup-count
{g1}	6
{g2}	7
{g3}	6
{g4}	2
{g5}	2
{g1, g2}	

## Step 3

Generate  $C_2$  from  $L_1$  and second for count of each.  $C_2 = L_1 \bowtie L_1$

$C_2$

Itemset	Sup-count
{g1, g2}	1.
{g1, g3}	1
{g1, g4}	1
{g1, g5}	2
{g2, g3}	9
{g2, g4}	1
{g2, g5}	2
{g3, g4}	0

Itemset	Sup-count
{g3, g5}	1
{g4, g5}	0

## Step-4

Compare candidate count with minimum count: 3

$L_2$

Gemset	Sup. count
{g1, g2}	1
{g1, g3}	1
{g1, g5}	2
{g2, g3}	1
{g2, g4}	2
{g2, g5}	2

## Step-5

Generate  $C_3$  from  $L_2$  and  
second D. for each Count

$C_3$

Gemset	Sup-count
{g1, g2, g3}	2
{g1, g2, g5}	2
{g1, g3, g4}	-
{g1, g3, g5}	-
{g1, g2, g4}	-
{g2, g3, g4}	-
{g2, g3, g5}	-

Itemset	SUP-count
{j2, j4, j5}	-

But According to apriori property

Last 4 sets are not be a member  
of frequent itemset.

### Step-6

Generate  $L_3$  from  $C_3$

$L_3$

Itemset	sup-count
{j1, j2, j3}	2
{j1, j2, j5}	2

~~Ans~~

### Step 7:

$$\text{Now } C_4 = L_3 \times L_3$$

$$= \{j1, j2, j3, j5\}$$

But it's subset {j1, j3, j5} is not frequent

So it is not frequent.

∴ Frequent itemset.

$$\{ \{g_1, g_2, g_3\}, \{g_1, g_2, g_5\} \}$$

Ans

Generating association rule:

For frequent itemset  $\{g_1, g_2, g_3\}$

Subsets are  $\{g_1\}, \{g_2\}, \{g_3\}$ ,  
 $\{g_1, g_2\}, \{g_1, g_3\}, \{g_2, g_3\}$

$$\therefore \{g_1, g_2\} \rightarrow \{g_3\}, \text{ con } \frac{\text{sup-count}\{g_1, g_2, g_3\}}{\text{sup-count}\{g_1, g_2\}}$$

$$= \frac{2}{4}$$

$$= 50\%$$

$$\{g_1, g_3\} \rightarrow \{g_2\}, \text{ con } \frac{\text{sup-count}\{g_1, g_2, g_3\}}{\text{sup-count}\{g_1, g_3\}}$$

$$\{g_2, g_3\} \rightarrow \{g_1\}, \text{ con } \frac{2}{4} = 50\%$$

$$\{g_2, g_3\} \rightarrow \{g_1\}, \frac{\text{sup-count}\{g_1, g_2, g_3\}}{\text{sup-count}\{g_2, g_3\}}$$

$$\frac{\text{sup-count}\{g_1, g_2, g_3\}}{\text{sup-count}\{g_2, g_3\}} = \frac{2}{4}$$

$$\{g_1\} \rightarrow \{g_2, g_3\}, \frac{\text{sup-count}\{g_1, g_2, g_3\}}{\text{sup-count}\{g_1\}}$$

$$= \frac{2}{6} = 33.3\%$$

$$\{g_2\} \rightarrow \{g_1, g_3\}, \frac{\text{sup-count}\{g_1, g_2, g_3\}}{\text{sup-count}\{g_2\}}$$

$$= \frac{2}{7} = 28.57\%$$

$$\{g_3\} \rightarrow \{g_1, g_2\}, \frac{\text{sup-count}\{g_1, g_2, g_3\}}{\text{sup-count}\{g_3\}}$$

$$= \frac{2}{6} = 33.33\%$$

As minimum confidence is 70%  
there is no association rule

for frequent itemset  $\{g1, g2, g5\}$

Subsets are,  $\{g1, \{g2\}, \{g5\}\}$

$\{g1, g2\} \{g1, g5\} \{g2, g5\}$

$$\{g1, g2\} \rightarrow \{g5\} \quad \text{confidence} = \frac{\text{sup\_count } \{g1, g2, g5\}}{\text{sup\_count } \{g1, g2\}}$$

$$= \frac{2}{4} = 50\%$$

$$\{g1, g5\} \rightarrow \{g2\} \quad \text{confidence} = \frac{2}{2} = 100\%$$

$$\{g2, g5\} \rightarrow \{g1\} \quad \text{confidence} = \frac{2}{2} = 100\%$$

$$\{g1\} \rightarrow \{g2, g5\} \quad \text{confidence} = \frac{2}{6} = 33.3\%$$

$$\{g2\} \rightarrow \{g1, g5\} \quad \text{confidence} = \frac{2}{7} = 28.57\%$$

$$\{g5\} \rightarrow \{g1, g2\} \quad \text{confidence} = \frac{2}{2} = 100\%$$

As minimum confidence threshold is 70%  
∴ association rules are,

$$\{g1, g5\} \rightarrow \{g2\}$$

$$\{g2, g5\} \rightarrow \{g1\}$$

$$\{g5\} \rightarrow \{g1, g2\}$$

Ans

Q 2017-1(c)

Given that, D,

$$|D|=9$$

Transaction ID	itemset
T <sub>1</sub>	{KA, D, B}
T <sub>2</sub>	{D, A, CE, B}
T <sub>3</sub>	{C, A, B, E}
T <sub>4</sub>	{B, A, D}

$$\text{min-sup} = 60\%$$

$$\text{min-conf} = 80\%$$

$$\therefore \text{min-sup} = 60\% = 9 \times \frac{60}{100} =$$

## Step-1

C<sub>1</sub>

Itemset	Sup-count	%Sup-count
{A}	4	100%
{B}	4	100%
{C}	2	50%
{D}	3	75%
{E}	2	50%
{K}	1	25%

## Step-2

L<sub>1</sub>, Comparing with minimum support count

itemset	sup-count
{A}	4
{B}	4
{D}	3

### Step-3

$L_1 \rightarrow e_2$

$$e_2 = L_1 \bowtie L_1$$

itemset	SUP-count	% SUP-count
{A,B}	9	100%
{A,D}	3	75%
{B,D}	3	75%

### Step-4

$L_2$ , comparing  $\emptyset$  with minimum supported

itemset	SUP-count
{A,B}	9
{A,D}	3
{B,D}	3

$\{A, B\} \in E^1 = \{\{A, B, A\}\}$

## Step-5

$$C_3 = L_2 \times L_2$$

itemset	SUP-count	% SUP-count
{A, B, D}	3	75%

## Step-6

$L_3$ , comparing with minimum support

itemset	SUP-count
{A, B, D}	3

$$\therefore C_3 = L_3 \times L_3$$

$$\subseteq \{A, B, D\} = L_3$$

$\therefore$  Frequent itemset is  $\{A, B, D\}$

generating association rule

frequent itemset is  $\{A, B, D\}$

$\therefore$  Subsets of are

$\{A\}, \{B\}, \{D\}, \{AB\}, \{AD\}$   
 $\{BD\}$

$\therefore \{A, B\} \rightarrow \{D\}$ , confidence  $= \frac{3}{84} = 75\%$

$\{A, D\} \rightarrow \{B\}$ ,  
" "  $= \frac{3}{3} = 100\%$

$\{B, D\} \rightarrow \{A\}$ ,  
" "  $= \frac{3}{3} = 100\%$

$\{A\} \rightarrow \{B, D\}$ ,  
" "  $= \frac{3}{4} = 75\%$

$\{B\} \rightarrow \{A, D\}$ ,  
" "  $= \frac{3}{3} = 100\%$

$\{D\} \rightarrow \{A, B\}$ ,  
" "  $= \frac{3}{3} = 100\%$

As minimum confidence threshold is

80% Association rules are,

$\{A, D\} \rightarrow \{B\}$  /  $\{D\} \rightarrow \{A, B\}$   
And

## Antimonotonicity property:

If a set cannot pass a test, all of its supersets will fail the same test as well.

## Finding interesting patterns:

→ support and confidence measure are insufficient

→ after filtering out uninteresting association rule

⇒ Spearman's correlation measure can be used.

## # Correlation rule

It is measured not only by its support & confidence but also by the correlation between items A and B.

2 types of correlation analysis

→ lift

→  $\chi^2$

lift:

Lift between occurrence of A and B is,

$$\text{lift}(A, B) = \frac{\text{PCA}(AB)}{\text{PCA}(A)\text{PCA}(B)}$$
$$= \frac{\text{Coef}(A \Rightarrow B)}{\text{sup}(B)}$$

$\rightarrow \text{lift}(A, B) < 1$

strong negative correlation

occurrence of A is negatively correlated with " " of B

$\rightarrow \text{lift}(A, B) > 1$ , " of A is positively " with " of B

$\rightarrow \text{lift}(A, B) = 1$ , ★ A and B are independent  
★ there is no correlation

$\chi^2$  testing:

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

where;

$E_{ij}$  = expected frequency of  $(A_i, B_j)$   
 $O_{ij}$  = observed " of  $(A_i, B_j)$

$$c_{ij} = \frac{\text{count}(A=a_i) \times \text{count}(B=b_j)}{n}$$

↑ number of data tuples

## Problem of Apriori

two nontrivial costs:

→ It may still need to generate a huge number of candidate sets

Ex  $10^4$  frequent 1-itemset →  $10^7$  candidate-2 itemsets

→ need to repeatedly scan the whole database

→ Check a large set of candidates

by pattern mining

Sol<sup>1</sup>:

→ Frequent pattern growth (FP-growth)

## FP growth:

→ it adopts a divide-and-conquer strategy

### Steps:

→ Scan DB once, find frequent-1 itemset

→ sort frequent itemset

in frequency decreasing order

→ scan DB again, construct FP-tree

→ construct conditional pattern base

→ conditional FP-tree

→ .. frequent pattern

## Example 6.5

Given that,

TID	List of item-IDs
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

$$\text{min-sup} = 2$$

Step-1

item	sup-count	item	sup-count
I1	6	I1	2
I2	7	I5	2
I3	6		

As min-sup = 2 . so all preceding table contain all frequent-1 itemset.

### Step-2

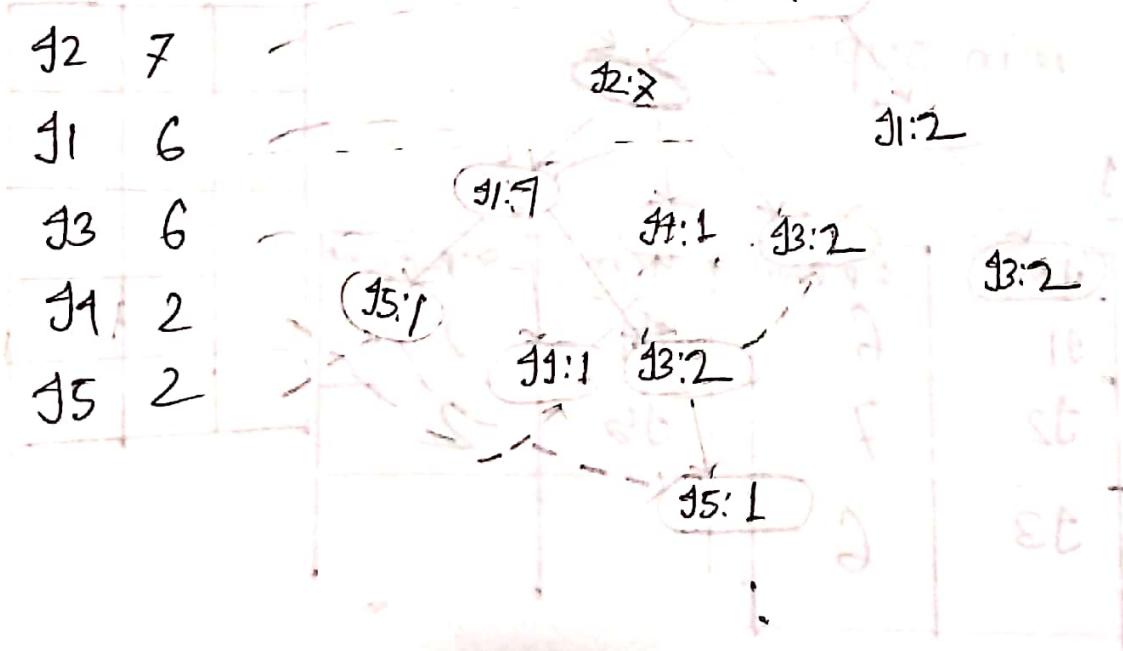
Sort frequent-1 itemset in decreasing order

item	Sup.count
g2	7
g1	6
g3	6
g4	2
g5	2

pattern :

g2, g1, g3, g4, g5

### Step-3



## Step-1

Not a complete tree

item	conditional pattern Base	Conditional FP-tree	Frequent patterns generated
g5	{g2, g1:1}, {g2, g1, g3:1}	{g2:2, g1:2}	{g2, g5:2}, {g1, g5:2}, {g5:2}
g1	{g2, g1:1}, {g2:1}	{g2:2}	{g2, g1:2}
g3	{g2, g1:2}, {g2:2}, {g1:2}	{g2:4, g1:2}, {g1:2}	{g2, g3:4}, {g1, g3:2}, {g2, g1, g3:2}
g1	{g2:4}	{g2:4}	{g2, g1:4}

2017-3(a)

Given that,

TID	Item - bought	
100	g6, g1, g3	g1, g3
200	g1, g2, g4, g5, g3	g1, g2, g5, g3
300	g3, g2, g5	g2, g5, g3
400	g8, g7	
500	g1, g3, g2, g4, g5	g1, g2, g5, g3
600	g1, g3, g6	g1, g3
700	g1, g2, g5, g7	g1, g2, g5
800	g2, g8, g5, g1	g1, g2, g5
900	g4, g6	
1000	g1, g2, g5	g1, g2, g5

min-support = 10%

Step-1

Step-1

min-support = 10%

$$\text{min-sup.-threshold} = \frac{10}{100} \times 10 \\ = 1$$

item	sup-count
j1.	7
j2.	6
j3	5
j4	3
j5.	6
j6	3
j7	2
j8	2

Q8 Ans

## Step-2

As min-supp-threshold = 4

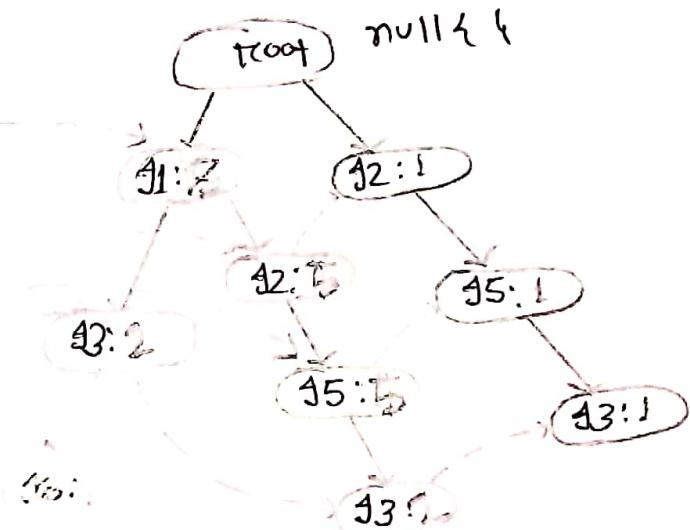
∴ frequent-1 itemset is given below:

item	min supcount
g1	7
g2	6
g5	6
g3	5

pattern: g1, g2, g5, g3

## Step-3:

item	supcount	node-link
g1	7	
g2	6	
g5	6	
g3	5	



## Step-4

## CPB: Conditional Pattern Base

CFPT: ..

## FP Tree

Item	CPB	CFPT	Frequent pattern
g3	$\{\{g1:2\}, \{g1,g2,g5:2\}, \{g2,g5:1\}\}$	$\{g1:1\}$	$\{g1, g3:1\}$
g5	$\{\{g1,g2:5\}, \{g2:1\}\}$	$\{\cancel{g1:5}, \{g1:5, g2:5\}\}$	$\{g1, g5:5\}$ $\{g2, g5:5\}$ $\{g1, g2, g5:5\}$
g2	$\{g1:5\}$	$\{g1:5\}$	$\{\cancel{g1,g5:5}\}$ $\{g1, g2:5\}$

Q 2017-1(b)

Given that:

TID	Items Bought
1	{Milk, Beerz, Diapers}
2	{Bread, Butter, Milk}
3	{Milk, Diapers, Cookies}
4	{Bread, Butter, Cookies}
5	{Beerz, Cookies, Diapers}
6	{Milk, Diapers, Bread, Butter}
7	{Bread}, Butter Diapers
8	{Beerz, Diapers}
9	{Milk, Diapers, Bread, Butter}
10	{Beerz, Cookies}

minsup > 0

Let,

Milk =  $m_i$

Beerz =  $Be$

Diapers =  $Di$

Butter =  $Bu$

Cookies =  $Co$

Bread =  $Br$

I length of transaction ID ~~is 6 and~~

2 9 is 4.

$\therefore$  As so maximum size of frequent itemsets is = 4  
Ans

ii Maximum number of size 3 itemsets.

$$\text{is } C_3 = \frac{6}{3} = 20 \quad \underline{\text{Ans}}$$

iii Let  $C_1$

item	sup_count
milk	5
Beer	4
Butter	5
Bread	5
Diapers	7
Cookies	4

As minsup > 0

$$\therefore C_1 = L_1$$

Now make pair

$O \rightarrow I$	$\text{conf}(\rightarrow) \text{sub}$	$\text{conf}(\rightarrow) \text{sup}$
Milk, Beer	4	5
Milk, Bottle ✓	5	5
Milk, Bread ✓	5	5
Milk, Diapers	7	5
Milk, Cookies	4	5
Beer, Butter	5	4
Beer, Bread	5	4
Beer, Diapers	7	4
Beer, Cookies ✓	9	4
Butter, Bread ✓	5	5
Butter, Diapers	7	5
Butter, Cookies	4	5
Bread, Diaper	7	5
Bread, Cookies	9	5
Diapers, Cookies	7	7

$\therefore$  (Milk, Butter) etc had some confidence  
in  $a \rightarrow b \Rightarrow$  Milk  $\rightarrow$  Butter  $\otimes$  and

$\text{Butter} \rightarrow \text{Milk}$  Ans

NB: Those pairs have some support  
count they have some confidence  
in  $b \rightarrow a$  and  $a \rightarrow b$  rules.

# # Finding support of size-2 itemsets

Itemset	Sup-count
Milk, Bread	3
Milk, Butter	3
Bitter, Milk	1
Milk, Cookies	1
Milk, Diaper	1
Bread, Butter	5 ✓
Bread, Beer	0
Bread, Cookies	1
Bread, Diaper	3
Butter, Beer	0
Butter, Cookies	1
Butter, Diaper	3
Beer, Cookies	2
Beer, Diaper	3
Cookies, Diaper	2

∴ Size-2 itemset with largest support  
is (Bread, Butter) Ans

## Exercise 6.6

Given that,

TID	items-bought
T100	{M, O, N, K, E, Y}
T200	{D, O, N, K, E, Y}
T300	{M, A, K, E}
T400	{M, V, C, K, Y}
T500	{C, O, N, F, E}

$$\text{min-SUP} = 60\%$$

$$\text{min-Cof} = 80\%$$

Step-1

item	sup-count
A	1
C	2
D	1
E	4
F	1
K	5
M	3
N	2
O	3
V	1
Y	3

Step-2

As  $\text{min-SUP} = 60\% = \frac{60}{100} \times S = 3$   
 $\therefore$  frequent-1 item(s) elis  
 in decreasing order

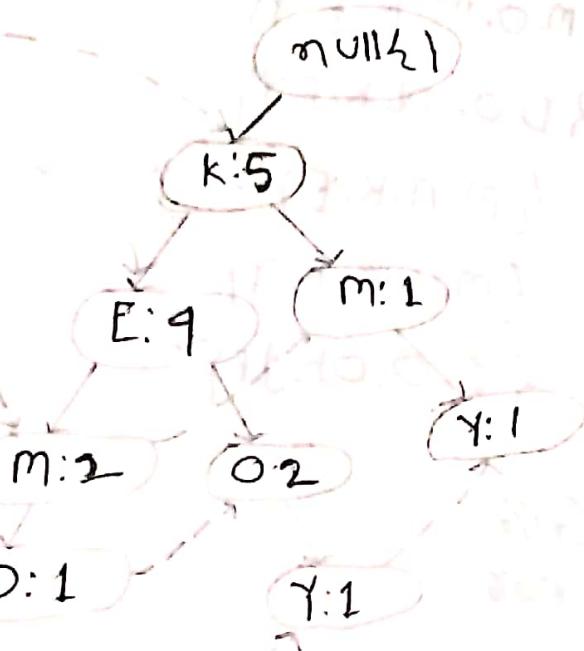
item	sup-count
K	5
E	4
M	3
O	3
Y	3

∴ pattern:  $\{K, E, M, O, Y\}$

### Step-3:

Item sup-count node link

K	5
E	4
M	3
O	3
Y	3



### Step-4

Item	CPB	CFPT	Frequent patterns
Y	$\{K, E, M, O, Y\}$ , $\{K, E, O, Y\}$ , $\{K, M, Y\}$	$\{K:3\}$	$\{K, Y:3\}$
O	$\{\{K, E, M:1\}, \{K, E:2\}\}$	$\{K:3, E:3\}$	$\{K:3\} \{E, O:3\}$ $\{K, E, O:3\}$
M	$\{\{K, E:2\}, \{K:1\}\}$	$\{K:3\}$	$\{K, M:3\}$
E	$\{K:4\}$	$\{K:4\}$	$\{K, E:4\}$

: complete set of frequent item-sets

$\{\{K:5\}, \{E:4\}, \{M:3\}, \{O:3\}, \{Y:3\},$   
 $\{K,Y:3\}, \{K,O:3\}, \{E,O:3\}, \{K,M:3\}, \{K,E:4\},$   
 $\{K,E,O:3\}\}$

Ans

### Apriori Vs FP-growth



Apriori	FP-growth
1) An <u>array based</u> algorithm	1) a <u>tree based</u> algorithm
2) It uses <u>join and prune</u> technique	2) It constructs <u>conditional frequent pattern tree</u> and <u>conditional pattern base</u> from database
3) It is a <u>breadth-first-search</u>	3) It uses <u>depth-first search</u>
4) It utilizes <u>level wide approach</u>	4) It utilizes a <u>pattern-growth approach</u>
5) Requires <u>large memory space</u>	5) Requires <u>small memory space</u>

Apriori

FP-growth

- \* 6) It scans database multiple time

- 6) It scans database only twice.

### Advantages of FP-growth over Apriori

→ FP-growth is more efficient

as

→ it able to mine the  
confidential pattern base

→ substantially reduce the  
size of the dataset

to be searched

## ■ Data format:

2 types

{TID: itemset}

→ horizontal data format

→ vertical data format

data format

→ {item: TID-set}

→ in vertical data format

Eclat

(Equivalence Class Transformation)

is used.

## ■ Eclat:

→ transform the horizontally formatted

data into the vertical format.

→ frequent k itemset can be used

to construct the candidate

(k+1) itemsets

based on the A priori property

→ Support count of an itemset

length of the TID-set.

### Advantages:

→ take the advantages of A priori property

→ no need to scan the database

to find the support of  $(k+1)$  itemset

### Disadvantages:

→ TID-sets can be quite long

→ taking substantial memory space  
→ computation time  
for intersecting long sets.

## Example - G.6

Given that

TID	List of item IDs
T100	g1, g2, g5
T200	g2, g4
T300	g2, g3
T400	g1, g2, g4
T500	g1, g3
T600	g2, g3
T700	g1, g3
T800	g1, g2, g3, g5
T900	g1, g2, g3

min-SUP = 2  
item → vertical data format

Item	TID-list
g1	{T100, T400, T500, T700, T800, T900}
g2	{T100, T200, T300, T400, T600, T800, T900}
g3	{T300, T500, T600, T800, T900}
g4	{T200, T400}
g5	{T100, T600}

As min-sup = 2

∴ All one frequent itemset

2-itemset vertical data format.

itemset	TD-set
{j1,j2}	{T100, T900, T800, T900}
{j1,j3}	{T500, T800, T900}
{j1,j4}	{T900}
{j1,j5}	{T100, T800}
{j2,j3}	{T300, T600, T800, T900}
{j2,j4}	{T200, T900}
{j2,j5}	{T100, T800}
{j3,j4}	∅
{j3,j5}	{T600}
{j4,j5}	∅

As min-sup = 2

∴ {j1, T4}, {j3, j4}, {j3, j5}, {j4, j5}

are not frequent.

3-itemset vertical dataset.

Applying "Apriori property" on  
2-itemset data:

3-itemset	TID-Det
{g1, g2, g3}	{T800, T900}
{g1, g2, g5}	{T100, T800}

Ans

# Chapter 3

## Data Preprocessing

### Data Quality:

Data have quality if they satisfy the requirement of intended use.

Many factors comprising data quality.

→ accuracy → The data are correct or not accurate

→ completeness

→ consistency

→ timeliness

→ believability

→ interpretability

absit proboscoz stob

Accuracy: The data are correct or not accurate or not.

Reasons for inaccurate data:

→ data collection instruments used

may be faulty

→ human or computer errors

occurring at data

→ errors in data transmission

★★ Disguised missing data:

→ User not wish to submit personal data

→ Purposely submit incorrect data for mandatory fields

## Completeness:

- Attributeness of interest may not always be available.
  - they were not considered important at the time of data entry
- Relevant data may not be recorded
  - due to equipment malfunction
- Inconsistent data may be deleted.

## Consistency:

## Timeliness:

Does the data are timely updated or not.

## Believability:

How trustable the data are.

## Interpretability:

How easily the data can be understood.

## Major tasks in data preprocessing:

→ Data cleaning

→ Data integration

→ Data reduction

    → dimensionality reduction  
    → numerosity

→ Data transformation and  
discretization

### Data cleaning:

Clean the data by

    → filling in missing value

    → smoothing noisy

    → value

    → identifying or  
removing outliers

    → resolving  
inconsistencies

## # Data integration

It involves integrating  
multiple databases

→ data cubes

→ files

Ex

[Customer-id, cust-id]

→ Redundant data may slow down / confuse

the knowledge discovery  
process

→ Data Cleaning and data integration

performed as a preprocessing  
step

## # Data reduction:

- Obtain a reduced representation of data set
- yet produce the same analytical result

### 2 strategies

- dimensionality reduction
  - numerosity
  - data compression

### Dimensionality reduction:

- obtain a compressed/reduced representation of original data.

Ex: → PCA

- Wavelet transformation
- attribute subset selection

## → attribute construction

→ small set of more useful  
~~information~~ attribute is derived  
from original set

## Numerosity reduction:

→ data are ~~represented~~ replaced  
by alternative, smaller representation

→ by using

→ parametric model

(regression,

long-linear models)

→ non-parametric

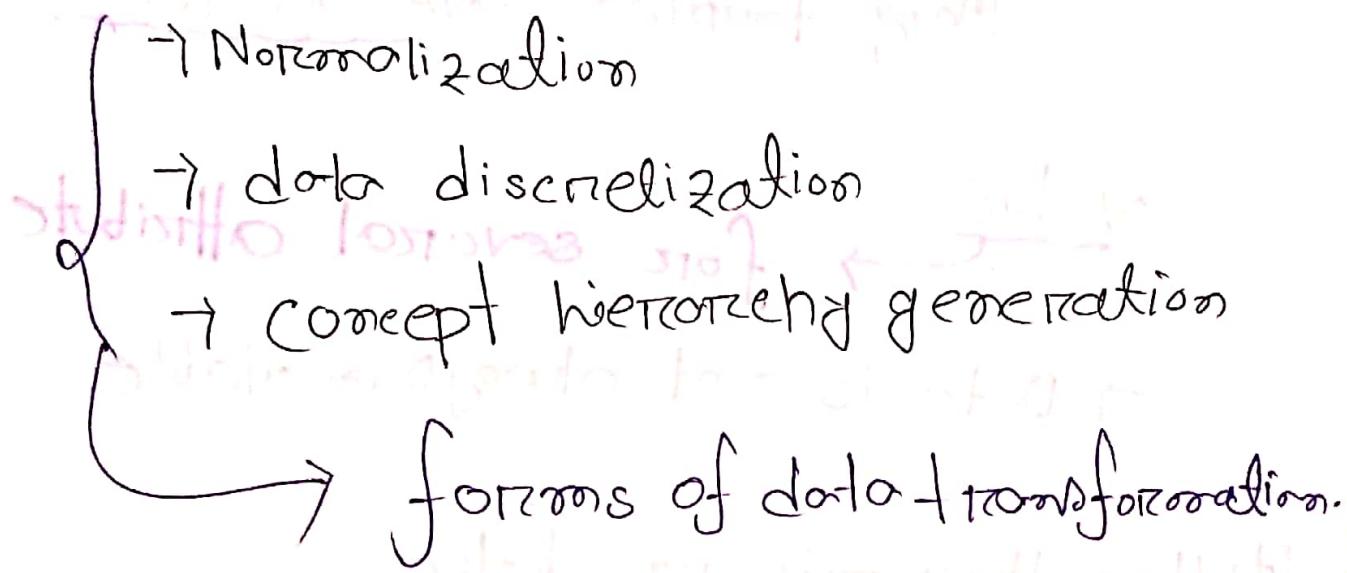
(histogram,

cluster,

sampling,

data aggregation)

## # Discretization and concept hierarchy:



## Reason of data preprocess

→ it can improve data quality

    ↑  
    helping to improve  
    ↑ accuracy  
    ↑ efficiency

→ important step in knowledge discovery process

    → quality decisions must be  
        based on quality data.

## Missing value:

Many tuples have no recorded

value

for several attribute

→ Data is not always available

## Handling missing data.



6 methods

→ Ignore the tuple

→ using when class label is  
missing

→ not very effective

→ Fill the missing value manually

→ time consuming

→ not feasible  
for large data

→ Use a global constant

→ Replace all missing attribute

values

by the same constant

Unknown

$\alpha$

→ Use a measure of central tendency for the attribute

→ for normal/symmetric data distribution

mean

→ " skewed "

median

→ Use the attribute mean or median

for all samples belonging to the same class

Class

→ Use the most probable value to fill

in missing value

→ determined with

• regression

• inference-based tools

→ Bayesian formulation

→ decision tree

## Reason of missing value:

→ equipment malfunction

→ inconsistent thus deleted

→ data not entered

    ↳ due to misunderstanding

→ certain data may not be considered important

    ↳ at the time of entry

## Noise:

A Random error or variance in a measured variable

## Reason of noisy/inaccurate data:

- ★ ★ → faulty data collection instrument
- Data entry problem
- Data transmission
- technology limitation
- inconsistency in mining convention

## Handling Noisy Data:

- 1) Binning method
- 2) Regression
- 3) Outlier analysis
- 4) Semi-automated method  
→ combined computer & human inspection

### Binning method:

- sort the data and partition into (equal frequency) bins
- then we can
- smooth each bin by
  - smoothing by mean
  - .. by median

- smoothing by bin boundary
  - maximum and minimum value of a bin is identified
  - each bin value is replaced by closest boundary value
- ★ Binning is also used as a discretization technique

Smoothing by mean:

a	b	c
---	---	---

$$\bar{x} = \frac{a+b+c}{3}$$

x	x	x
---	---	---

Smoothing by median

a	b	c
---	---	---

$$\bar{x} = b$$

6	6	6
---	---	---

Smoothing by bin boundary

a	b	c
---	---	---

$$\bar{x} = \min((|a - \text{bin}|), (|c - \text{bin}|))$$

if  $\bar{x} = |a - \text{bin}|$  then replaced by a

if  $x = |e - bin|$  then replaced by e

## Regression:

→ conforms data values form  
a function

Ex: Linear regression

$$\hat{y} = \theta_0 + \lambda \theta_1$$

finding the best line to fit two attributes

so that  
one attribute can be used

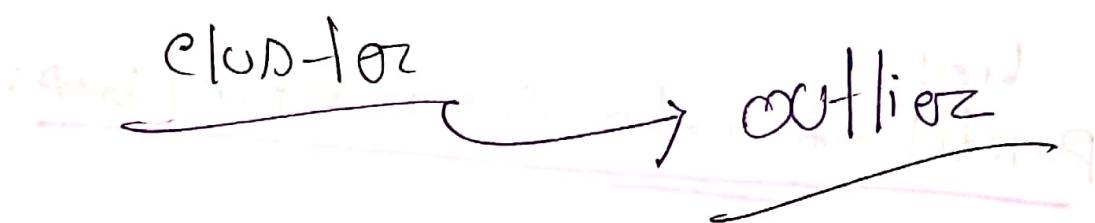
to predict another attribute

## Outlier analysis:

→ detected by clustering

→ similar values are organized  
into groups / clusters

→ values that fall outside the



## Example:

Sorted data for price

4, 8, 9, 15, 21, 21, 21, 25, 26, 28, 29, 31

bin size = 4

$$\therefore \text{number of bins} = \frac{12}{4} = 3$$

Bin 1

Partition into (equi-depth) bins:

→ Bin 1: 4, 8, 9, 15

→ Bin 2: 21, 21, 21, 25

→ Bin 3: 26, 28, 29, 31

Smoothing by bin mean:

→ Bin-1: 9, 9, 9, 9

→ Bin-2: 23, 23, 23, 23

→ Bin 3: 29, 29, 29, 29

## Smoothing by bin boundaries:

→ Bin-1: 4, 4, 4, 15

→ Bin-2: 21, 21, 25, 25

→ Bin 3: 26, 26, 26, ~~29~~ 34

## Data Integration

It is the merging of data from multiple data stores.

### Adv. of data integration:

→ help to reduce and avoid redundancies  
inconsistencies

→ help to improve the accuracy

→ help to " " " speed of  
  
Subsequent mining  
process

## Schema integration:

→ integrate metadata from  
different sources

e.g. A.cust-id ≡ B.cust-#

□ Considerable issue during data integration

2 issues

→ schema integration

→ object matching

□ Entity identification problem:

Equivalent real world entities

from multiple data sources

→ can be different.

Ex: Will Smith = W. Smith

{ customer\_id = cust\_number

→ how the data analyst / computer  
be pure

Sol: → attention must be paid to

structure of the data

→

## ■ Detecting & resolving value conflicts

→ For the same entry, attribute values from different sources are different

→ Possible reason:

→ different representation

→ " " scaled

E.g. metric vs. British units  
different scale.

## ■ Redundant attribute:

A attribute may be redundant if it can be derived from another attribute / set of attributes.

## Handling redundancy in data integration

→ by using correlation analysis

Chi-square correlation → Nominal data

Pearson's product moment coefficient / correlation coefficient

Nomeric data



brief soll für

(B.I.A) techn.



slight formalization

privat slight for u

Annot; Ocular

more privat u for u

slight id

## Chi-square correlation test / Pearson $\chi^2$ statistic

It can be computed as

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(O_{ij} - e_{ij})^2}{e_{ij}}$$

where,  $O_{ij}$  = observed frequency

$e_{ij}$  = expected

$$e_{ij} = \frac{\text{count}(A=a_i) \times \text{count}(B=b_j)}{n}$$

of the joint event  $(A_i, B_j)$

number of tuple

" of tuple having

" of " having value  
value  $a_i$  for A

" of " having value  
 $b_j$  for B

NB:

### Hypothesis:

→ A and B are independent, there is no correlation between them

→ test is based on a significant

level with

degree of freedom:  $(r-1)(c-1)$

### Correlation coefficient

Computed as

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

$$r_{A,B} = \frac{\sum_{i=1}^n (a_i b_i) - n \bar{A} \bar{B}}{n \sigma_A \sigma_B}$$

$$= \frac{\sum_{i=1}^n (a_i b_i) - n \bar{A} \bar{B}}{n \sigma_A \sigma_B}$$

where,

$n$  = number of tuples

$a_i$  = values of A in tuple  $i$

$b_i$  = " " of B in " i

$\sigma_A$  = Standard despeclive

$\sigma_A$  = Standard deviation of A

$\sigma_B$  = " " of B

$\bar{A}$  = mean values of A

$\bar{B}$  = " " of B

## Evaluation

$$-1 \leq \rho_{A,B} \leq 1$$

$\rightarrow \rho_{A,B} > 0$ , A and B are positively correlated

$\rightarrow \rho_{A,B} < 0$ , A and B " negatively"

$\rightarrow \rho_{A,B} = 0$ , A and B " independent"

## ↳ Covariance:

- Two numeric attributes A and B
- Set of  $n$  observations  $\{(a_1, b_1), \dots, (a_n, b_n)\}$

$$\begin{aligned}\text{Cov}(A, B) &= E((A - \bar{A})(B - \bar{B})) \\ &= \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}\end{aligned}$$

$$\boxed{\text{Cov}(A, B) = E(A \cdot B) - \bar{A} \bar{B}}$$

where,

$$E(A) = \bar{A}$$

$$E(B) = \bar{B}$$

## ↳ Relation between Correlation and Covariance

$$\star \star \quad R_{A,B} = \frac{\text{Cov}(A, B)}{\sigma_A \sigma_B}$$

## Variance:

Special case of covariance

where two attributes are identical

## Example - 3.1

Given that,

	male	female	Total
function	250(30)	200(36)	450
non-function	50(210)	1000(890)	1050
Total	300	1200	1500

$\chi^2 =$  Number in the parentheses are the expected frequencies

$$\begin{aligned}
 \chi^2 &= \frac{(250 - 50)^2}{90} + \frac{(50 - 210)^2}{210} + \frac{(1000 - 890)^2}{890} \\
 &\quad + \frac{(200 - 360)^2}{360} + \frac{(-160)^2}{360} + \frac{(-160)^2}{890} \\
 &= \frac{160^2}{90} + \frac{(-160)^2}{210} + \frac{(-160)^2}{360} + \frac{(160)^2}{890} \\
 &= 160^2 \left( \frac{1}{90} + \frac{1}{210} + \frac{1}{360} + \frac{1}{890} \right) \\
 &= 160^2 \times \frac{5}{252} \\
 &= 507.936 \\
 &\approx 507.91
 \end{aligned}$$

Ans

$\therefore \text{degree of freedom} = (2-1)(2-1)$   
 $= 1 \times 1 = 1$  Ans

$$F = F = \frac{8+8+1+3+3}{3+3} = (A) \quad : (A)$$

$$\bar{A} = \bar{B} = \frac{3+1+1+3+3}{2} = (B) \quad : (B)$$

### Example 3.2

Given that,

Time point	AllElectronics	HighTech
+1	6	20
+2	5	10
+3	4	19
+4	3	5
+5	2	5

(Q) & (O)

Let, A = AllElectronics

B = HighTech

$$\therefore E(A) = \frac{6+5+4+3+2}{5} = 4 = \bar{A}$$

$$E(B) = \frac{20+10+19+5+5}{5} = 10.8 = \bar{B}$$

$$\therefore \text{Cov}(A, B) = \frac{\sum (A - \bar{A})(B - \bar{B})}{n}$$

$$\begin{aligned} \text{Cov}(A, B) &= \frac{\{(6 \times 20) + (5 \times 10) + (4 \times 14) + (3 \times 5) \\ &\quad + (2 \times 5)\} - \bar{A} \cdot \bar{B}}{5} \\ &= \frac{251 - 4 \times 10.8}{5} \\ &= \frac{251 - 43.2}{5} \\ &= 50.2 - 43.2 \\ &= 7 \end{aligned}$$

As  $\text{Cov}(A, B) > 0$  ~~so A and B i.e.~~  
 So, positive covariance

So, Stock price for both companies will increase together.

$$P \cdot S_P =$$

$$(P \times 28) - P \cdot S_P = (8.87 \text{ ru})$$

$$P = P \cdot 28 - P \cdot S_P =$$

removes switching risk

## Example

Given that,

$$(A, B) = (2, 5); (3, 8); (5, 10); \\ (9, 11); (6, 19)$$

$$\text{Cov}(A, B) = ?$$

$$E(A) = \frac{2+3+5+7+6}{5} = \bar{A}$$

$$= 9$$

$$E(B) = \frac{5+3+5+8+10+11+19}{7} = \bar{B}$$

$$= \frac{58}{7} = 9.6$$

$$E(A \cdot B) = \frac{(2 \times 5) + (3 \times 8) + (5 \times 10) + (9 \times 11)}{5}$$

$$= \frac{10 + 24 + 50 + 99 + 81}{5}$$

$$= 42.4$$

$$\therefore \text{Cov}(A, B) = 42.4 - (9.6 \times 9)$$

$$= 42.4 - 86.4 = -44$$

$\therefore$  Positive covariance

Ans

## Wavelet transform:

A linear signal processing technique

that when applied to a data vector  $x$ ,

transforms it to a numerically different vector  $x'$

of wavelet coefficients

→ wavelet coefficients larger than some upper-specified threshold can be retained

→ All other coefficients are set to 0.

→ halves the data at each iteration

■ Hierarchical pyramid algorithm:

→ length  $L$  of the input data

vector must be an integer  
power of 2

→ Each transform involves applying two

functions

some data smoothing

sum

a weighted difference

weighted average

→ Functions are applied to pairs of data points in  $x$ .

→ pair of measurements  $(x_{2i}, x_{2i+1})$

→ results of two data sets of length  $L/2$

→ Functions are applied to the data sets  
until the resulting data sets obtained  
one of length 2

→ Selected values from the data sets  
obtained in the previous iterations

## Attribute subset selection

Reduces the data set size

by removing irrelevant or redundant attribute

goal: → to find a minimum set of

attribute → resulting probability distribution of data

so that

as close as possible to original distribution.

NB: For n attribute

there are  $2^n$  possible subsets

set of subset

subsets

## Greedy methods for attribute subset selection

Selection:

4 methods

→ Stepwise forward selection

→ " backward elimination

→ Combination of forward & backward selection

→ Decision tree induction

### Stepwise forward selection:

\* → Start with an empty set

→ best of the original attributes  
is determined

→ added to the reduced set.

for clustering still lossy  
binaries

## Chapter-10

### Cluster Analysis: Basic Concepts

#### and Methods

BBold

#### Clustering:

Business

It is the process of grouping  
a set of data objects  
into/multiple groups.

so

so that each group of data is known as cluster

→ objects within the a cluster

high similarity

→ very dissimilar to

objects in other clusters.

in individual analysis consists of  
the main concept of learning multimodal

## ★ Real life application of clustering

### Application of clustering:



→ taxonomy of living things

→ biology

→ security

→ business intelligence

→ web search

→ image pattern recognition

→ used to organize a large number of customers into groups

→ organize the search results into groups

→ present the result in a concise & easily accessible way

→ to discover clusters / subclust. in handwritten character recognition sys.

Chapters 6  
Cluster Analysis  
and Data Mining  
with Applications  
in R

- Clustering as a disjointed function:
- Application as a data mining function:

→ as a stand alone tool

→ as a preprocessing step

→ outlier detection

→ to gain insight  
into the distinction  
of data

to observe the  
characteristic of each  
cluster.

- Contributing areas of research:

→ data mining

→ biology

→ statistics

→ marketing

→ machine learning

→ spatial database

→ technology

→ web search

→ information retrieval

clustering মুলে কৃষ্ণ বিভিন্ন cluster  
cluster shape, clustering result বিষয়ে  
আলোচনা

Clustering is also called

→ automatic classification

→ data segmentation

→ outlier detection

→ unsupervised learning

learning by observation

as cloud label information  
the  
is not present

as it partitions large data sets  
into groups according to their  
similarity

a cluster of data object can be  
treated as an implicit cloud.

★ input data এর ফর্ম কিনা কোনো  
— যানেকা

## ■ Requirements of clustering:

↳ Fault tolerance

→ Scalability:

clustering all the data instead of sample

→ Ability to deal with different types of attribute:

→ typically designed to

cluster numeric data

→ applications may require

clustering other data types

↳ binary

→ nominal

→ ordinal

→ mixture of these data type

→ constraint-based clustering

→ user may give inputs on constraints

→ where domain knowledge is

used

→ यहाँ कौनी क्षमता देती है जो हमें इसका विवरण दे सकती है

→ Interpretability and usability:

→ User wants clusters

→ Results to be interpretable,

→ comprehensible

→ usable

→ of meaningful clusters

→ Discovery of clusters with arbitrary

space

→ Ability to deal with noisy data

→ Capability of clustering high-dimensional data

→ Incremental clustering & insensitive to input order

high dimensionality  
quality of clustering / Requirements for domain knowledge  
to determine input parameters

## □ Considerations for clustering Analysis

★ ★ Orthogonal aspects of comparing Clustering methods:

→ The partitioning criteria

All the objects are partitioned

partition data object

hierarchically

→ separation of clusters:

Exclusive

data object may belong to only one cluster

" to more than one cluster

" non-exclusive

→ Similarity measure:

→ Distance based method

↳ Euclidean space

↳ road network

↳ vector space

→ connectivity

→ cluster of arbitrary space

→ take advantage of optimization

→ Clustering space

→ search for clusters within the entire given data space

→ subspace clustering

↳ useful for high-dimensional

## Basic clustering methods:

14

~~syntax~~ → partitioning methods

Biostatistics → Find mutually exclusive clusters

- Distance based
- Use mean or medoid etc to represent cluster center
- Effective for small-to medium size data sets

Ex: K-means, K-medoids,

CLARA, CLARANS

→ Hierarchical methods

↳ Clustering is a hierarchical decomposition.

2 types  $\rightarrow$  agglomerative (bottom-up)  
 $\rightarrow$  divisive (top-down)

→ distance based or  
density and continuity based

→ may incorporate other techniques

Problem: like microclustering

→ if one step is done,

it can never be undone

→ can't correct erroneous decisions

Density based method

Continue to grow cluster  
as long as density in neighbourhood  
exceeds a threshold

→ can find arbitrary shaped

cluster

→ used to detect outliers

Clustering: A  $k$ -means  $\rightarrow$   $k$

density:

- Each point must have a minimum number of points within its "neighborhood":

Ex:  $\rightarrow$  DBSCAN

- OPTICS
- Dense

Grid-based methods:

quantize the object space  
into a finite number of cells

that form a grid structure

Advantages:

Fast processing time

Ex

- STING
- WaveCluster

\* CLIQUE

## K-Means: A Centroid-Based Technique

•

→ uses the centroid of a cluster

initial

→ uses the centroid of a cluster  
c: to represent that cluster.

→ centroid is defined by

→ mean  
→ median  
of the objects assigned to the cluster

→ difference between an object

$P \in C_i$  and  $c_i$  other

→ measured by

$\text{dist}(P, c_i)$

→ quality of cluster  $C_i$

→ measured by

sum of squared errors

$$E = \sum_{i=1}^k \sum_{P \in c_i} \text{dist}(P, c_i)^2$$

→ centroid of cluster  
 point in space

standard for reducing  $E_{\text{tot}} = 0$

★ ★  
K-means procedure:

Step-1: Randomly selects  $k$  of the objects in  $D$

Step-2: Assign each object to the cluster

↳  $c_i$  → object assigned to which the object is the most similar

Step-3: update the cluster centroid

Step-4: Repeatedly perform step 2 and step 3 until the assignment is stable / no change

↳ until the assignment is stable / no change

**Complexity of k-means:**  $\underbrace{3 \cdot 3}_{\text{steps}} = 9$

→ time complexity is  $O(nkt)$

where,

$n$  = total number of objects

$k$  = number of clusters

$t$  = " of iterations

\* As  $K < n$  and

$t < n$

so method is → relatively scalable

→ efficient

**Problem:**

\* it is not guaranteed to converge to the global optimum

→ terminates at a local optimum

Soln:

run the k-means algorithm

multiple times

with different initial  
cluster center

→ can't be applied to nominal data

Soln: Use (k-modes)  
How can we make the k-means  
more scalable?

★  
★  
Soln

→ uses a good-sized sample  
set of samples in clustering

→ employ a filtering approach:

→ uses a spatial hierarchical  
data index

→ to save cost  
when computing means

~~→ unclustered soft way~~

for this we will use euclidean distance

### Example 10.2



Given:

Data = {1, 2, 3, 8, 9, 10, 25}

$k=2$  then

$$C_1 = \{1, 2, 3\}$$

$$C_2 = \{8, 9, 10, 25\}$$

$$\therefore \text{mean}(C_1) = 2$$

$$\text{mean}(C_2) = 13$$

$$\therefore E = (1-2)^2 + (2-2)^2 + (3-2)^2$$

$$(8-13)^2 + (9-13)^2 + (10-13)^2$$

$$+ (25-13)^2$$

$$= 1 + 0 + 1 + 25 + 16 + 9 + 144$$

$$= 196$$

But comparing with this

$$C_1 = \{1, 2, 3, 8\}$$

$$C_2 = \{9, 10, 25\}$$

$$\text{mean}(C_1) = 4.5$$

$$\text{mean}(C_2) = 14.67$$

$$\therefore E = (1-4.5)^2 + (2-4.5)^2 + (3-4.5)^2 + (8-4.5)^2$$

$$+ (9-14.67)^2 + (10-14.67)^2 + (25-14.67)^2$$

$$= (6.25 + 6.25 + 6.25 + 20.25 + 32.1989 + 21.8089 + 106.7089)$$

$$= 189.67$$

So K-mean is less sensitive to outliers.

due to

## Strength of k-means:

→ Efficient ( $O(n^2k)$ )

Comparing → PAM

→ CLARA

## Weakness of k-means:

→ sensitive to noisy data &

→ outliers

→ Need to specify K

→ good for numerical data.

K-modes: for categorical

data

K-medoids: applied to

a wide range  
of data

→ terminate at global optimum

→ not suitable to discover

clusters → with non-convex  
SP shapes

global optimum

of clusters from

several local optima

global optimum  $\rightarrow$  (x)

(global)

## What are K-Medoids?

Instead of taking mean value of the object in a cluster, as a reference point, medoids can be used.

### Medoids:

→ most centrally located object in a cluster.

Ex: → PAM (Partitioning Around Medoids)

which is

K-mean vs K-medoid

$\rightarrow$  K-medoid is more robust than

K-mean

$\rightarrow$  in the presence of

noise

outlier

$\rightarrow$  complexity of K-mean  $O(kn)$

of K-medoid  $O(k(n-k))$

if  $k=1$   $O(Kn^2)$  otherwise if it is a NP hard problem

How can we modify the K-mean algorithm to diminish sensitivity

to outliers?

Soln:  $\rightarrow$  pick actual objects to represent

the cluster.

→ partitioning method is then performed based on ~~on~~  
~~absolute error criterion~~

the principle of ~~minimizing~~  
the sum of the dissimilarities between  
each object and

its two corresponding representatives

Absolute-error criterion:

$$E = \sum_{i=1}^K \sum_{P \in C_i} d(p, o_i)$$

where,

$o_i$  = representative object  
of cluster  $C_i$

$P$  = data point.

$K$  = number of clusters

Q: How can we scale up the k-medoids method?

Sol: ~~divide and conqu~~

Answer:

→ use a sample based method

Ex: CLARA (Clustering LArge Application)

CLARA:

→ uses a random sample of data points to find medoids.

→ PAM is applied then to

compute

best medoids from sample

to avoid getting stuck

$\rightarrow$  sample should closely represent the original data.

$\rightarrow$  CLARA builds clustering

based on significant from multiple random sample

Complexity:  $O(Ks^2 + K(n-k))$

where,

$s$  = size of the sample

$K$  = number of clusters

$n$  = number of objects

Weakness:

$\rightarrow$  efficiency depends on sample size

$\rightarrow$  sample may be biased

## PAM

Step-1: Select k representative objects

arbitrarily

Step-2:

For each pair of non-selected

object h and selected object i

calculate total swapping cost  $\delta$

$TC_{ih}$

Step-3

For each pair of  $i$  and  $h$

→ if  $TC_{ih} < 0$ ,  $i$  is replaced

by  $h$

→ assign each non-selected

object to the most similar representative object.

## Step 4:

Repeat steps 2-3 until  
there is no change.

### Efficiency improvement on PAM

→ CLARA / PAPP on sampled

→ CLARANS

Biggest / total / Randomized re-sampling

## Hierarchical clustering

- This method works by grouping data objects into a hierarchy or "tree" of clusters.
- Useful for data summarization for "visualisation"
- Does not require numbers of clusters as input
- needs a termination condition

## Types of hierarchical method

Can be categorized into

→ algorithmic method

→ probabilistic

→ Bayesian

algo → method

the algorithmic method

\* Agglomerative method

↳ joining of objects

\* divisive

\* multiphasic

→ considered data object

↳ determine

→ compute clusters according to the

## deterministic distance

### # Probabilistic method

the applies → probabilistic model to  
estimate capture clusters.

Model of adopted → measure quality of clusters  
always no with by the features of  
model.

### # Bayesian method

→ compute or distribution of  
possible clustering

it returns a group of  
clustering structure  
their probabilities  
conditional on given data

## Aggregation vs Divide and Conquer

## Agglomerative vs. Divisive

~~Star~~ ~~Circle~~ ~~Triangle~~ ~~Diamond~~ ~~Cross~~ ~~Star~~ ~~Circle~~ ~~Triangle~~ ~~Diamond~~ ~~Cross~~ ~~Star~~ ~~Circle~~ ~~Triangle~~ ~~Diamond~~ ~~Cross~~

- starts by letting → starts by placing each object from all objects in one its own cluster cluster. (hi branch'd root)
  - iteratively merge → divides the root cluster into clusters into smaller and larger clusters. subclusters.
  - recursively partition until all the those cluster into single objects one.
  - single cluster → until each cluster at the lowest level

□ clustering opproaches

2 types

→ single-linkage

→ dendrogram

\* Single-linkage:

→ each cluster is represented

by all the objects in the cluster

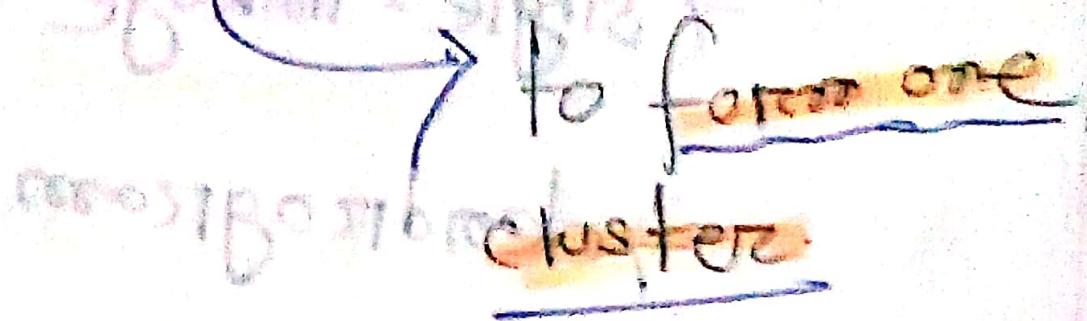
→ similarity between two clusters

is measured by the similarity

of the closest points of data

points belonging two different clusters.

→ merging process repeats until all the objects are eventually merged



\* dendrogram:

→ it shows how objects are grouped together or positioned step-by-step.

\* → it fuse similarity scale & Leveling

## E) Distance Measure in Algorithms

Method:

Four widely used method

minimum distance:

$$\text{dist}_{\min}(C_i, e_j) = \min_{P \in C_i, P' \in e_j} \{ |P - P'| \}$$

maximum distance

$$\text{dist}_{\max}(C_i, e_j) = \max_{P \in C_i, P' \in e_j} \{ |P - P'| \}$$

mean distance

$$\text{dist}_{\text{mean}}(C_i, e_j) = \frac{1}{m_i m_j} \sum_{P \in C_i, P' \in e_j} |P - P'|$$

Average distance

$$\text{dist}_{\text{avg}}(C_i, e_j) = \frac{1}{m_i m_j} \sum_{P \in C_i, P' \in e_j} |P - P'|$$

whereas, single linkage

$|P - p'| = \text{distance between two points}$

$m_i = \text{mean for clusters } C_i$

$n_i = \text{number of object in } C_i$

nearest-neighbour clustering:

When an algorithm uses the  
minimum distance  $d_{min}(C_i, C_j)$   
to measure the distance between  
clustering

single-linkage algorithm:

minimum distance between any  
two cluster objects of two cluster  
exceeds a user defined threshold

Complete

furthest neighbor clustering:

when algorithm uses maximum distance  $d_{max}(c_i, c_j)$  to measure the distance between clusters

complete-linkage algorithm:

if maximum distance  $d_{max}(c_i, c_j)$  exceeds a user defined threshold

minimum spanning tree algorithm:

if An agglomerative hierarchical clustering algorithm

uses minimum distance measure

for similarity

Spanning tree

A tree that connects all vertices

Minimal spanning tree:

One with the least sum of edge weight.

Minimum Spanning Tree

• Average and minimum measure is sensitive to outliers or missing and noisy data.

• Median will be unaffected by outliers or missing data because it is the middle value.

→ Outlier sensitivity problem

→ average distance can handle  
→ categorical data  
→ numeric "

BR

BREH

→ designed for clustering a large amount of numerical data

by integrating hierarchical clustering & other clustering method

(hierarchical clustering + other clustering methods)

overcome two difficulties



1) Scalability

2) inability to undo what was done in the previous step

done in the

previous step

1 2 3

→ it uses the ~~method~~<sup>→</sup> of clustering feature ~~of~~<sup>→</sup> clustering feature  
clustering feature ~~of~~<sup>→</sup> clustering feature  
clustering feature (CF) tree

### Clustering Feature (CF) tree

→ Clustering feature is a 3-D vector summarizing information about clusters of objects.

$$CF = (n, LS, SS)$$

Where,

LS = linear sum of the n points  
$$\sum_{i=1}^n x_i$$

SS = square sum of data points

$$\sum_{j=1}^n x_i^2$$

→ A CF is essentially a summary of the statistics

cluster's centroid,  $x_0$

"radius", R

"diameter", D

$$x_0 = \frac{\sum n_i}{n}$$

$$R = \sqrt{\frac{\sum_{i=1}^n (x_i - x_0)^2}{n}} = \sqrt{\frac{\sum_{i=1}^n x_i^2 - 2\sum_{i=1}^n x_i x_0 + \sum_{i=1}^n x_0^2}{n}}$$

$$= \sqrt{\frac{nSS - 2LS^2 + nL^2}{n^2}} = \sqrt{\frac{nSS - 2LS^2 + nL^2}{n^2}} = \sqrt{\frac{15 - 2(1)(1) + 1}{9}} = \sqrt{\frac{13}{9}} = \sqrt{1.44} = 1.2$$

$$D = \sqrt{\frac{\sum_{i=1}^n \sum_{j=1}^n (x_i - u_j)^2}{n(n-1)}}$$

$$= \sqrt{\frac{2nSS - 2LS^2}{n(n-1)}}$$

-  $R$  is average distance

from member object to the centroid

-  $D$  is the average pairwise distance within cluster.

\* Clustering features are additive.

If  $c_1$  and  $c_2$  are two distinct clusters,  $CF_1 = (n_1, LS_1, SS_1)$   
 $CF_2 = (n_2, LS_2, SS_2)$

then,  $EF_1 + EF_2 = (n_1 + n_2, LS_1 + LS_2, SS_1 + SS_2)$

Ex. Example 10.5:

Given that 3 points of  $C_1$  is

~~points of~~  $(2, 5), (3, 2)$  and  $(4, 3)$

~~points of~~  $CF_2 = (3, (35, 36), (417, 418))$

$\therefore EF_1 = (n_1, LS_1, SS_1)$

~~points of~~  $n_1 = (3, (9, 10), (29, 38))$

$\therefore EF_3 = EF_1 + EF_2$

$= (6, (49, 46), (996, 978))$

Ans

## OF-tree (m) (n=10) and

→ highest-balanced tree

→ Same the clustering factor

→ for a hierarchical clustering

→ nonleaf nodes store

leaf index

sum of leafs of their

children

has two parameters:

→ branching factor, B

→ threshold, T

leaf index, leaf node, leaf

leaf

browsing cost) to a search of  
browsing factor,  $B$ :

maximum number of children  
per nonleaf node

of each cell need

threshold,  $T$ :

maximum diameter of

subclusters

conditions on how  
to store at the  
leaf nodes.

# Multiphase clustering technique:

→ single scan of the data set  
yields basic good clustering

→ additional scans can optionally

be used to further improve  
the quality

## □ Phases of CF-tree

\*\* two phases

### Phase 1:

→ Scan the database to

build an initial in-memory CF tree

shift to bottom up clustering

viewed as multilevel composition  
of data

### Phase 2:

→ applied on (selected) cluster

algorithm

→ to cluster the leaf nodes

## Chameleon

→ uses dynamical modeling

to determine the similarity between pairs of clusters.

Similarity is assessed based on

- 1) how well connected objects one within a cluster
- 2) the proximity of clusters

★ Two clusters are merged if their interconnectivity is high and they are close together

Similarity measure between pair of  
cluster  $e_i$  and  $e_j$

→ relative interconnectedness

$$RJ(e_i, e_j)$$

→ relative closeness

$$RC(e_i, e_j)$$

relative interconnectedness  $RJ(e_i, e_j)$ :

Defined on the absolute

interconnectedness between  $e_i$  and  $e_j$ .

normalized with respect to

internal interconnectedness

of two clusters.

(largest) and (smallest)

$$+ (|E_{cl}| + |E_{c1}|)$$

minimum cut

Where

$|E_{cl}|$  = Edge cut for a cluster

$|E_{c1}|$  = Edge cut for a part

② Considering both  $E_{cl}$  &  $E_{c1}$

③ To find minimum sum of the

cut edge that

$\frac{2d}{2} = 69 \text{ of } 90$

→ that partitions G into  
two equal parts

Algorithm for finding minimum sum of the cut edges

Algorithm for finding minimum sum of the cut edges

## Relative closeness:

Defined as abs. closeness

Closeness between  $C_i$  and  $C_j$

normalized

internal closeness

of the two clusters.

$$RC(C_i, C_j) = \frac{\bar{S}_{EC\{C_i, C_j\}}}{\bar{S}_{EC\{C_i\}} + \bar{S}_{EC\{C_j\}}}$$

$$\text{where } \bar{S}_{EC\{C_i, C_j\}} = \frac{|C_i|}{|C_i| + |C_j|} \bar{S}_{EC\{C_i\}} + \frac{|C_j|}{|C_i| + |C_j|} \bar{S}_{EC\{C_j\}}$$

where,

$\bar{S}_{EC\{C_i\}}$  - average weight

of the edge

absolute weight

of the edge

## B Chameleon problem

types of nearest neighbor graph

graph  $\rightarrow$  connects a

spoke graph  $\rightarrow$  many

short paths between nodes

$\rightarrow$  uses agglomerative clustering

algorithm to partition the

$k$ -nearest-neighbor graph

longer number of small subgraphs

$\rightarrow$  use agglomerative hierarchical

clustering that iteratively  
merges subclusters

$\rightarrow$  based on their  
similarity

## ■ Density-Based Methods

If we model clusters as a dense regions in the data space, separated by sparse regions, to find clusters of arbitrary shape, called density based model.

## ■ Types of density-based method

- \* DBSCAN
- \* OPTICS
- \* DBNCLUE

## ■ DBSCAN's features

→ Discover clusters of arbitrary shape

→ Handle noise

→ Core point

→ Need density parameters

• Set  $\epsilon$  as termination condition

■ Q: How can we find dense regions in density-based clustering?

Answer:

→ density of an object  $Q$  can be measured by

→ number of objects close to  $Q$ :

→ H finds core objects

Definition for objects that have dense neighborhoods

■ Re → maximum radius of  
■  $\epsilon$ -neighborhood:  $\text{Neigh}(o)$  neighbourhood

The  $\epsilon$ -neighborhood

of an object  $o$  is the space within a radius  $\epsilon$

space within a radius  $\epsilon$

centred at  $o$  as well.

■ MinPts: Number of neighbors

minimum number of points in an  $\epsilon$ -neighborhood of that point.

→ specifies the density threshold of dense regions.

### Candidate set

subset -  $N_{\epsilon \text{PS}}(P) \subseteq \{q \in D \mid \text{dist}(P, q) \leq \epsilon\}$

### Core object

An object is core object if the  $\epsilon$ -neighbourhood of the object contains at least  $M_{\text{min}} \text{pts}$  objects.

### Directly density-reachable:

For a core object  $q$  and another object  $P$  we say that  $P$  is directly density-reachable from  $q$  if  $P$  is within the  $\epsilon$ -neighbourhood of  $q$ .

## Density reachable

A point  $P$  is density-reachable from  $q$  if there is a chain of objects  $P_1, \dots, P_n$ , such that

$P_1 = q$ ,  $P_n = P$ , and  $P_{i+1}$  is

directly-density-reachable

from  $P_i$

## Density-connected:

A point  $p$  is density-connected to a point  $q$  w.r.t  $\epsilon$  and minpts if there is a point  $o$  such that both  $p$  and  $q$  are density-reachable from  $o$ .

OR,

Two points  $P_1, P_2 \in D$  are density-connected w.r.t  $\epsilon$  and minpts if there is an object  $q \in D$  such that both  $P_1, P_2$  are density-reachable from  $q$ .

**Q** Question: How does DBSCAN find clusters?

Answer:

- All the objects in a given data set D one marked as "unvisited"
- randomly selects an unvisited object  $p$  for further processing
- marks  $p$  as visited
- checks whether  $p$  is core object or not
- If not  $p$  is marked as noise point
- Otherwise → a new cluster  $C$  is created using  $\delta$  for  $P$ . later
- All the objects in the neighbourhood of  $P$  added to  $N_P$

- If an object is not visited
- If  $p^*$  is a core object, then
  - all objects within  $\epsilon$ -neighborhood of  $p^*$  are added to  $N$ .
- Add  $p^*$  to  $C$ .
- iteratively adds those until  $N$  is empty.
- To find the next cluster
  - randomly selects an unvisited object from remaining ones
- continues the process until all of the points have been processed.

## Computational complexity of DBSCAN:

- if a spatial index is used  
complexity is  $O(n \log n)$
- otherwise complexity is  $O(n^2)$
- n → number of datasets of  
database objects

## Evaluation of Clustering

The major task of clustering evolution include the following

→ Assessing clustering tendency

For a given data set, we assess whether a nonrandom structure exists in the data

→ Determining the number of clusters in a data set

It is desirable to estimate the number of cluster before a clustering algorithm is used

to drive detailed clusters

## → Measuring clustering quality

↳ How good the resulting  
group of clusters are

↳ measure how well the  
clusters fit the ground truth

↳ the stub rows fit the data set  
→ measure how well the  
cluster match

→ the ground truth  
↳ the stubs of marsh grass

↳ number of marsh grass in a

→ number of pixels to measure  
↳ multiple groups

→ number of pixels of

## QUESTION

### QUESTION

Ques. How can we assess the clustering tendency of a dataset?

Ans. From 1 to 10

Answer:

\* By using 'Hopkins Statistics'

### Calculation

1. Sample  $n$  points  $p_1, \dots, p_n$  uniformly from  $D$

For each point  $p_i$  we find the nearest neighbour of  $p_i$  in  $D$

$$x_i = \min \{ \text{dist}(p_i, v) \}$$

( $v \in D$ )

Conclusion

## 2. 1 Sample - points having neighbors

For each  $x_i$  we find the nearest

neighbor of  $x_i$  in  $D_{\text{test}}$

to be nearest neighbor

$$d_i = \min \{ \text{dist}(x_i, v) \mid v \in D_{\text{test}}, v \neq x_i \}$$

## 3. Hopkins Statistic

$$H = \frac{\sum d_i}{\sum x_i + \sum d_i}$$

$$\sum x_i = \sum d_i$$

If  $H=0.5$  (normal)  $D$  were uniform

distributed

If  $H < 0.5$  (abnormal)  $D$  is highly

skewed

If  $H > 0.5$   $D$  were nonuniformly distributed

homogeneous hypothesis

D is uniformly distributed

and thus contains no meaningful  
clusters

nonhomogeneous hypothesis:

D is nonuniformly distributed

and thus contains significant

clusters

(e.g.)

bottom cluster

## Determining the number of clusters

### Empirical method

\* number of clusters in a data

$$\text{set is } \sqrt{n}/2$$

\* number of data points in each

cluster

$$\text{is approximately } \sqrt{2m} \text{ points}$$

A data set of  $n$  points

(ex)

### Elbow method

→ based on observation, that  
increasing the number of clusters  
can help to reduce the  
sum of within cluster variance

of each cluster and step by step

→ having more clusters allows one to capture finer groups of data objects

→ that are more similar to each other.

Procedure:

→ compute the clustering algorithm (like k-means) for different values of  $k$  ( $k > 0$ )

→ calculate the sum of within-cluster variance ( $\text{var}(k)$ )

→ plot the curve of  $\text{var}$  w.r.t  $k$

→ the first turning point of the curve

suggests the "right" number.

### Cross-validation

divide the given data set D  
into m parts

→ use m-1 parts to build a  
clustering model

→ use the remaining part  
to test the

quality of the clustering

→ for any k > 0 repeat the  
process m times

→ average of the quality measure

is taken as the overall quality

quality measure

→ compute the overall quality measure

w.r.t different values of  $k$ .

→ find the number of clusters

that best fit the data

measures

clustering algorithm

starts with many

individuals to group

bottom up and each

bottom up merging from

bottom up

## Measuring Clustering quality

### Question:

How good is the clustering generated by method and how can we compare the clusterings generated by different methods?

### Answer:

Two methods to choose from for measuring the quality of a clustering.

two types of method:

- Extrinsie methods
- Intrinsie

## Clustering

## Agglomerative

- 1) The ground truth is available
- 1) The ground truth is unavailable
- 2) Also known as supervised method
- 2) Also known as unsupervised method
- 3) Compare the clustering against goodness of a cluster  
the ground truth ring
- 3) Evaluate the clustering quality measure

using certain clustering quality measure

clustering techniques

the ground truth +

clustering techniques

good +

clustering techniques

## 4) External Methods

It measures clustering quality  $\text{Q}(\text{c}, \text{S})$  for

clustering  $C \rightarrow$  clustering  $C$   
of samples  $\rightarrow$  giving the group  
threshold  $\epsilon$

$\text{Q}$  is good if it satisfies the  
following 4 essential criteria

$\rightarrow$  cluster homogeneity

$\rightarrow$  cluster completeness

$\rightarrow$  Raw bag

$\rightarrow$  small cluster preservation

## Clustering homogeneity

The more pure the clusters in a clustering one, the better the clustering.

## Cluster completeness

If only two objects belong to the same category, they should be assigned to some cluster.

## Raw bag:

Putting a heterogeneous object into a pure cluster should be penalized more than putting it into a raw bag.

## Small clusters preservation

Splitting a small category

into pieces is more harmful

than splitting a large category

into pieces instead

split above out profit

plant - prepared since salt

young tubers put

also some of berries add

best and

are segregated to different

parts plants are to other

and from each other so

and more difficult to find

## B Guttman's Method:

1. (1)

Evaluate a clustering by examining

→ how well the clusters are separated

→ how compact the clusters are

## Silhouette Coefficient:

Let  $D$  = a set data set of  $n$  objects

$c = k$  number of clusters

of  $D$  =  $(c)$

also

$O \in D$

$O$  = object

$\sigma(O)$  = average distance between

$O$  and all other

objects in the same cluster.

$\alpha(0) = \sum_{i \in C_0} \text{dist}(a_i)$

$$\{C_1\} \cup \dots$$

3.  $\beta$  minimum distance

between minimum average

distance from  $a_i$  to  
all clusters  $C_j$

$$\beta(0) = \min_{i \in C_0} \sum_{j \neq 0} \text{dist}(a_i)$$

if  $\beta(0) < \alpha(0)$  then  $C_0 \leftarrow C_0 \cup \{a_i\}$

$$\delta(0) = \min_{i \in C_0} \frac{\sum_{j \neq 0} \text{dist}(a_i)}{\alpha(0)}$$

e.g. 1536.371 {51}

4. silhouette coefficient of  $\alpha$

$$\text{silhouette}(\alpha) = \frac{\alpha(0) - \beta(0)}{\max\{\alpha(0), \beta(0)\}}$$

if  $\text{silhouette}(\alpha) > 0$  then  $C_0 \leftarrow C_0 \cup \{a_i\}$

- Silhouette coefficient is between -1 and 1.
- $a(o)$  reflects the compactness of cluster
- $b(o)$  captures the degree to which  $o$  is separated from other cluster.
  - \* positive: good
  - negative: bad (exception case)

## Chapter-12

### Outlier Detection

#### Outlier detection:

It is the process of finding data objects with behaviors that are very different from expectation.

\* also known as anomaly detection.

#### Outlier:

An outlier is an object that deviates significantly from the rest of the objects.

## Outliers & noise

- Noise is anything that is not the true signal
- It may have values close to the true signal
- An outlier is something that is much different than the other values
- \* noise is a random error or variation in the measured variable
- Noise should be removed before before outlier detection

## Types of outliers

→ three categories

\* global outliers

\* contextual (conditional) outliers

\* collective outliers

global outliers:

A data object is a global outlier if it deviates significantly from the rest of the data set.

data set

\* also called **point anomalies**

\* simplest type of outliers

## determination

Critical issue is to find  
an appropriate measurement  
(of deviation) w.r.t  
the application in question

### Contextual outliers:

A data object is a contextual outlier if it deviates significantly w.r.t a specific context of the object.

\* also called conditional outliers

as they are conditional  
on the selected context.

## 3 types of attributes of data object

→ Contextual attributes

→ Behavioral

### Contextual attribute

define the object

in context

### Behavioral attribute

define the object

### Behavioral characteristics

- \* it is a generalization of local outliers

### Global outliers

An object is a data set

is a local outlier if its density significantly deviates from local area in which it occurs

### ■ Collective outliers:

Given a data set, a subset of data objects from a collective outliers, if the objects as a whole deviate significantly from the entire data set.

- \* It requires background information to model the relative or skip among objects to find group of outliers.
- \* Consider not only behavior of individual objects, but also that of groups of objects.

ND: → A data set may have multiple types of outliers  
→ One object may belong to more than one type of outlier

## IV Application of decision detection

→ Credit card ~~fraud detection~~

→ ~~Telecom~~ → ~~customer segmentation~~

→ ~~Medical analysis~~

→ ~~Marketing~~

→ ~~Finance~~

→ ~~Anti-money laundering~~

→ ~~Banking~~

→ ~~Healthcare~~

→ ~~Food industry~~

→ ~~Transport~~

## Outlier detection methods:

Based on user-labeled

Example of outliers can be  
obtained as:

- \* Supervised
- \* Unsupervised
- \* Semi supervised

Based on assumptions about  
normal data

- \* Statistical methods
- \* proximity-based methods
- \* clustering-based methods

## Supervised methods

→ can be approached by model

→ outlier detection can be modeled  
as a classification problem

→ sample is used for training  
and some not in training for testing

\* Model normal objects & report  
those not matching the model  
as outliers

\* Model the outliers and  
test object not matching the  
model is outlier.

Challenge: 2 challenges

→ imbalanced classes: i.e.  
→ outliers are rare.

→ method of handling outlier

imbalanced classes is oversampling

→ catch as many outliers as possible

→ sensitivity and recall of a classifier

outlier is far more important

than from accuracy

about 50% of data are outliers

are outliers

and outliers are about 5%

outliers for t-test

outliers about 10%

approx 10% & 100%  
are outliers

→ needs bootstrap

then one outlier

## ■ Unsupervised methods

### Assumptions:

→ Assume that normal objects are somewhat clustered into multiple groups.

→ Normal objects follow a pattern more frequently than outliers.

→ An outlier is expected to occur far away in feature space from any those groups of normal objects.

## Weaknesses in Unsupervised

Cannot detect outliers

collective

outliers

~~→ to some case~~

~~→ to some case normal object~~

\* normal objects do not always

~~distributed~~. Some may follow  
many such objects do not follow

~~strong pattern~~

Ex: ~~method in~~

for some intrusion detection and

computer virus detection problem,

normal activities are diverse

~~so~~ ~~in~~ ~~some~~ ~~cases~~

~~→ unsupervised may have high false~~  
positive rate

The supervised method is more effective for this

Process: ~~start with one cluster~~

→ find cluster's first

→ next pick object ~~not belonging to~~  
→ data objects not belonging to  
~~one cluster~~

~~accidentally~~ → ~~more detected as~~  
~~outliers~~ ~~being outliers~~

Problem: ~~cluster has~~

→ Hard to distinguish it from noise from  
outliers on the side

→ It is often costly to find  
clusters first and then  
find outliers at ~~stab~~

## • Semi-supervised method

- small set of the labeled data for other objects are labeled
- most of the data are unlabeled
- regarded as application of semi-supervised learning not methods are used
- \* If some labeled, unlabeled objects are available
  - needs the labeled sample to promote unlabeled object to train the model for new data

→ those ~~not~~ fitting if the model  
of normal objects

→ detected as  
outlier

\* If only some labeled outliers  
one available.

A small number of labeled  
outliers may not cover the  
possible outliers well.

→ to improve the quality of  
outlier detection, we can get  
help from normal objects  
→ learned from unsupervised  
method

## Chapter 7 Advanced Pattern Mining

### E) Roadmap on pattern mining

three pattern mining aspects

- kind of patterns
- mining methodologies
- applications

### F) Kind of patterns and their

3 types

\* Basic patterns

\* Multi-level & multi-dimensional patterns

\* Extended patterns

## One pattern

Simple frequent pattern

Closed or max pattern, flone  
and negative pattern.

## Frequent pattern

Pattern that satisfies

minimum support and threshold.

## Closed pattern

P is closed pattern if there  
is no super-pattern  $p'$  with the  
same support or  $p$ .

## Max pattern

P is max pattern if there  
is no frequent super-pattern of P.

## 2. Multilevel and multidimensional Pattern:

- \* Based on the association levels involved in a pattern
  - multilevel association rule
  - single-level "
- \* Based on the number of dimensions involved in the rule or pattern
  - single-dimensional association rule
  - multidimensional "
- \* Based on the type of value handled in the rule or pattern
  - Boolean association rule
  - Qualitative "

Example → 2 or more items

Multiple ref. composition rule

bug<sub>1</sub> (e., "computer")  $\Rightarrow$  bug<sub>1</sub><sup>2</sup> (e., "problem")

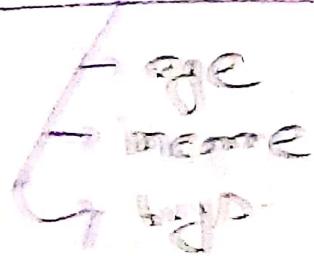
bug<sub>2</sub> (e., "laptop computer")  $\Rightarrow$

bug<sub>2</sub> (e., "color laser printer")

Computer  $\rightarrow$  higher level of abstraction.

Multilevel:

• role/pattern references two or more dimensions.



get(e., "20...27") ^ income(e., "52k...58k")  $\Rightarrow$  hyp (e., "Pd")

## • Extended pattern

Based on the constraints or  
criteria used to mine selective  
patterns:

- constraint
- oppositional

~~K~~ → uncertain  
satisfying a set of user defined constraints → constrained  
near match

→ top-k

→ redundancy-aware top-k

$k$  most frequent items for a user-specified value  $k$

top-k patterns with redundant pattern excluded

→

## Method for mining multi-level multidimensional Space

→ multilevel association rule

→ multidimensional " "

→ quantitative " "

→ tone patterns and negative pattern

→ involves numeric attributes that have an implicit ordering among values.

show negative correlations between them.

multilevel association rule:

concept hierarchy:

It defines a set of high-level concepts from a set of low-level concepts to a higher level more general concept set.

- \* last level is the most specific abstraction level of hierarchy.

Association rules generated from mining data at multiple abstraction levels are called multiple-level association rules.

NB: can be mined effectively using concept hierarchies

with a puppet from

## Approach to mine multi-level association rule

rule

using

→ uniform minimum support

→ relaxed support

→ group-based mixed minimum support

uniform

uniform

uniform minimum support

\* The same minimum support

threshold is used when mining at each abstraction level

Level 1

min sup = 5%

Computer  
[Support = 100%]

Level 2

Laptop Computer  
[Support = 64%]

min sup = 5%

Desktop Computer  
[Support = 42%]



## Adv:

- search procedure is simplified
- requires only one minimum threshold

## drawback:

- \* if min-sup threshold is too high
  - miss some meaningful associations occurring at low abstraction level
- \* min-sup threshold is too low:
  - generate many uninteresting association occurring at high abstraction level.

## B Reduced support

Each abstraction level had its own minimum support threshold  
→ The deeper the abstraction level

4  
the smaller the corresponding threshold.

Level 1  
min\_sup = 5%

Computer [Supp<sub>10=12</sub>]

Level 2  
min\_sup = 3% [Support = 6%]

laptop computer  
[support = 9.5%]

Effect on group-based minimum support

a) User specific group based minimum support:

Example:

High price company may have lower threshold.

B side effect of multi-level organization

problem

generation of many redundant rules

occurred multiple obstruction

levels

→ due to complex relationships

## Ancestors rule:

A rule  $R_1$  is an ancestor of rule  $R_2$ , if  $R_1$  can be obtained by replacing the items in  $R_2$  by their ancestors in a concept hierarchy.

### Example:

$\text{buys}(X, \text{"laptop computer"}) \Rightarrow$

$\text{buys}(X, \text{"HP printer"})$

[support = 8%, confidence = 70%]

$\text{buys buys}(X, \text{"Dell laptop computer})$

$\Rightarrow \text{buys}(X, \text{"HP printer"})$

[support = 2%, confidence = 72%]

## Multidimensional association rule

Rules that involve two or more dimensions on predicates

## 1-dimensional association rule

Multidimensional association rule with no repeated predicates.

Ex:  $\text{Age} \in [20, \dots, 29] \wedge \text{Occupation} \in \{\text{L}, \dots, \text{H}\}$   
 $\Rightarrow \text{bug} \in G, \{\text{top}, \text{P}\}$ .  
1-dimensional association

## 2-dimensional association rule

Ex: Multidimensional association rule with repeated predicates

Ex:  $\text{Age} \in [\text{*}, \dots, 29] \wedge \text{bug} \in (\text{*}, \text{"topo})$   
 $\Rightarrow \text{bug} \in (\text{*}, \text{"HP private")}$

## Handling Quantitative Association Rules:

two approaches

dynamic

\* Static discretization:

Quantitative attributes are  
discretized using → predefined concept  
hierarchies

\* dynamic discretization:

Quantitative attributes are  
discretized or clustered into  
"bins" → based on the data distribution

Simple discretization:

- discretize quantitative attributes into multiple intervals
- need them on numerical data in association mining

Problem of simple discretization:

generates enormous number of rules.

Solution:  
three methods to overcome this difficulty

1) a data cube method

2) a clustering-based method

3) a statistical analysis method

↑  
Uncover exponential behaviors

## Data Cube

Quantitative  $\rightarrow$  nominal

$\rightarrow$  by discretization

Transformed multi-dimensional data  
may be used to construct a data  
cube

### Example

o-P (age) cuboid

(age)  
(age, income)

1-D cuboid  
(bigP)

2-D cuboid  
(income, bigP)

3-D cuboid  
(age, income, bigP)

\* cells of an ordimensional cuboid  
can be used to store the  
support counts of the corresponding  
n-predicate fact

## B Mining Clustering Based Quantitative Association

\* Interesting frequent patterns or association rules are in general found at relatively dense clusters of quantitative attributes.

two approaches

\* top-down approach

\* bottom-up

topdown approach:

standard clustering algorithm can be applied to find clusters in this dimension that satisfy the minimum support threshold.

↳ K-means

↳ density based clustering

- bottom-up approach
- first clustering in high-dimensional space to form clusters with support-min-sup
  - projecting and merging those clusters in space containing fewer dimensional combinations

## Quantitative association based

### on Statistics:

It is possible to disclose exponential behavior:

(sex = female)  $\Rightarrow$  (mean wage = \$7.29/h)

(overall mean wage = \$9.29/h)

Involves applying statistical tests to confirm the validity of our rules.  $\rightarrow z$ -test

To different negative patterns:

Pattern with a freezing appeal  
and no belief (you know) a user  
specified number support - Handshaking

### Example

long before writing

### To negative pattern:

if hands x and y are both  
frozen but rarely come together  
(i.e.  $\text{sup}(x \cup y) < \text{sup}(x \cap y)$ )  
then hands x and y are  
negatively correlated. And the  
pattern  $x \cap y \times y$  is negatively

Correlated pattern

Strongly correlated pattern of debt

If  $\text{w}_t(\text{C}_t)$  is significantly positive  
then  $\text{z}_t(\text{D}_t)$  are strongly negatively  
correlated and the pattern  
 $y_{t+1}$  is a strongly negatively  
correlated.

Problem

full investment problem

full participation

transaction that does not exclude  
one of the agents being crowded

### Example 7.4

\* Given that,

A store sold two needle  
A and B 100 packages out of  
200 transaction

Only one transaction contained

A and B.

$$S(A \cup B) = \frac{1}{200} = 0.005$$

$$S(A) = \frac{2100}{200} / : S(A) \times S(B) = \frac{100}{200} + \frac{100}{200}$$

$$S(B) = \frac{100}{200} / = 0.5 \times 0.5 = 0.25$$

$$S(A \cup B) < S(A) \times S(B)$$

Suppose we have  $10^6$  null transaction

$$\therefore S(A \cup B) = \frac{1}{10^6} = 10^{-6} \quad S(A) \times S(B)$$

$$\therefore S(A \cup B) & S(A) \times S(B) = \frac{10^6}{10^{12}} = \frac{1}{10^6} = \frac{100}{10^6} \times \frac{100}{10^6}$$

Def-2

If  $x$  and  $y$  one strongly negating correlated then

$$\text{sup}(x \cup y) \times \text{sup}(x \cup y) \gg$$

$$\text{sup}(x \cup y) \times \text{sup}(y)$$

this definition also can't remove the problem of null formation.

Example 7.5

Q.  $\text{sup}(A \cup B) \times \text{sup}(A \cup B) = \frac{99}{200} \times \frac{99}{200}$   
 $= 0.295$

$$\text{sup}(A \cup B) \times \text{sup}(A \cup B) = \frac{1}{200} \times \frac{199}{200}$$

$$= 0.005$$

$$\sup(C \cup B) > \sup(A \cup B) \gg \sup(C) \cup B \\ \times \sup(A \cup B)$$

$\therefore A \text{ and } B$  are negatively correlated

Now for 1e<sup>6</sup> null fraction.

$$\sup(C \cup B) \times \sup(A \cup B) = \frac{99}{10^6} \times \frac{99}{10^6}$$

$$= 9.8 \times 10^{-9}$$

Again,  $\sup(C \cup B) \approx \sup(A \cup B)$

$$\approx \frac{99}{10^6} \times \frac{(10^6 - 99)}{10^6}$$

$$\approx 1.99 \times 10^{-9}$$

$\therefore$  No null fraction creates null invariant problem.

Def 3

$\beta_1, \gamma$ , and  $\delta$  are both frequent  
 $\text{avg}(\text{PC11}) + \text{PC112}) / 2 < 0$   
where  $\text{C11}$  is negative pattern  
threshold. Then pattern  $\gamma^{\text{PC11}}$   
is a negatively correlated pattern

Europe-T

@ Greenland

## Constraint Based Frequent Pattern

Mining:

Constraint are added to specific mining interest pattern.

It helps to confine the search space.

The constraints can include:

→ Knowledge type constraints:

specify the type of knowledge to be mined.

- association
- correlation
- classification
- clustering

→ Data constraints:

specify the set of task-relevant data.

→ Dimension / level constraints:

Specify the desired

→ dimension of the data

→ the abstraction levels

→ the level of the concept hierarchy

used in the mining

→ interestingness constraints:

specify thresholds on statistical measure of rule interestingness

support  
confidence  
correlation

## → Rule constraints

→ specify the form of, or  
conditions on, the rules to be  
mined

→ constraints may be expressed  
as metarules / rule template  
maximum or minimum  
number of prediction

## Q1 Metarule Guided Mining

Metarules allow us to specify the syntactic form of

rules interested in mining.

Q2: How can metarules be used to guide the mining process?

Answer: A metarule is a rule template of the form

$$P_1 \wedge P_2 \wedge \dots \wedge P_k \Rightarrow Q_1 \wedge Q_2 \wedge \dots \wedge Q_l$$

where  $P_i$  ( $i=1 \dots k$ ) and  $Q_j$  ( $j=1 \dots l$ ) are either instantiated predicates or predicates variable.

## Method to find Meta Rule:

1. Find frequent C1+H1 predicates based on min threshold.
2. Push constraints deeply when possible into the mining process.
3. Use confidence, correlation and other measure when possible.

## B) Constraint Based Pattern Generation

1. Pruning Pattern Space
2. Pruning Data space

## Pattern Pruning constraints:

Constraint that facilitate pattern space prun. is called pattern pruning constraints.

## Data Pruning constraints:

Constraints that can be used for data space pruning is called data pruning constraints

## Types of pattern space pruning constraints:

5 categories

→ antimonotonic → convertible

→ monotonic → inconvertible

→ succinct

### \* Antimonotonic:

If constraint  $c$  is violated  
its further mining can be  
terminated.

### \* Monotonic:

If constraint  $c$  satisfied,  
no need to check  $c$  again

### \* Succinct:

Constraint  $c$  must be  
satisfied, so one constraint  
with the data set satisfying  
 $c$ .

\* Convertible:

$\Leftrightarrow$  e is not monotonic note  
antimonotonic, but it can be  
converted into it, if items in the  
transaction can be properly  
ordered.

Example of Antimonotonic:

$$\text{Count}(S) \leq \vee$$

$$\text{max}(S) \leq \vee$$

$$\text{min}(S) \geq \vee$$

$$\text{sum}(S) \leq \vee \quad (\forall a \in S, a > 0)$$

$$\text{range}(S) \leq \vee$$

## Example of Monotonic

$\min(S) \leq v$

$\max(S) \geq v$

$\text{count}(S) \geq v$

$\text{sum}(S) \geq v$

$\text{range}(S) \geq v$

## ■ Data Space Pruning Constraint:

two properties/types

→ data succinctness

→ data cardinality

\* Data succinctness:

Data space can be pruned

at the initial pattern mining  
process.

Ex: mined pattern must contain  
digital camera

\* Data anti-monotone:

If a transaction  $t$  does not satisfy c.  $t$  can be pruned from its further mining.

## Mining High-Dimensional Data

Can't mine by using conventional mining approach.

two approaches

method-1

- \* data set with a large number of dimensions but small number of rows
- \* pattern of very long length

method-1

Extends or pattern growth approach by further exploring the vertical data format

to handle data of first to

Ex:

useful for the analysis  
of gene expressions in bioinformatics  
or others.

\* Column enumeration/horizontal  
data format:

D is viewed as a set of row IDs.

consists of an itemset

\* Row enumeration/vertical data  
format:

D is viewed as a set of  
row IDs where the item appears  
in the traditional view of D.

## Process:

→ The original data set, D can easily be transformed into a transposed data set T.

$$D \xrightarrow{\text{transpose}} T$$

→ efficient pattern growth method can be developed

## Global pattern:

Pattern of very long length.

### method:

(Pattern-Fusion) method