

## Unit - 4

### Business Intelligence

- ✓ Linear Regression
- ✓ Apriori Algorithm Qs
- ✓ Clustering Theory
- ✓ Data mining
- ✓ Data mining outliers
- ✓ K-Medoids
- ✓ K-Means clustering

### Apriori Algorithm

Q1. Suppose there are 4000 customer transactions. Calculate the Support, Confidence and lift for two products, Biscuits & Chocolate.

Out of the 4000 transactions, 400 contain biscuits, 600 contain chocolate, and these 600 transactions include 200 with both biscuits and chocolate.

$$\text{Support (Biscuits)} = \frac{\text{Transactions w/ biscuits}}{\text{Total transactions}}$$

$$= \frac{400}{4000} = 10\%$$

$$\text{Confidence} = \frac{\text{Transactions w/ biscuits \& chocolate}}{\text{Transactions w/ biscuits}}$$

$$= \frac{200}{400} = 50\%$$

$$\text{lift} = \frac{\text{confidence}}{\text{support}} = \frac{50}{10} = 5$$

Q) Apply the Apriori Algorithm on the following product set.

Transaction ID	Rice	Pulse	Oil	Milk	Apple
t <sub>1</sub>	1	1	1	0	0
t <sub>2</sub>	0	1	1	1	0
t <sub>3</sub>	0	0	0	1	0
t <sub>4</sub>	1	1	0	1	1
t <sub>5</sub>	1	1	1	0	0
t <sub>6</sub>	1	1	1	1	1

Let min-sup = ~~2~~ > 50%.

Ans

Product	Frequency
R	4
P	5
O	4
M	4
A	3



prune the apples

apply  
self join  
→

Product	Frequency
R	4
P	5
O	4
M	4

Product	Frequency
RP	4
RO	3
RM	2
PO	4
PM	3
OM	2

prune RO, PM, RM & OM

Product	Frequency
RP	4
PD	4
RO	3
PM	3

$\geq 3$

$R_{CU}$  is the customer's set of product

=>

? S

Example Apply the Apriori algorithm on

TID	items
100	1 3 4 6
200	2 3 5
300	1 2 3 5
400	1 5 6

$$\text{min-sup} = 0.5$$

Ans

TID	1	2	3	4	5	6
100	1	0	1	1	0	1
200	0	1	1	0	1	0
300	1	1	1	0	1	0
400	1	0	0	0	1	1

$$0.5 \times \text{num}(TID) = 2$$

Itr 1

Item	Freq
1	3
2	2
3	3
4	1
5	3
6	2

prune

Item	Freq
1	3
2	2
3	3
5	3
6	2

join

↓

Item	Freq
$\{1, 2\}$	1
$\{1, 3\}$	2
$\{1, 5\}$	2
$\{1, 6\}$	2
$\{2, 3\}$	2
$\{2, 5\}$	2
$\{2, 6\}$	0
$\{3, 5\}$	2
$\{3, 6\}$	1
$\{5, 6\}$	1

prune

→

Item	Freq
$\{1, 3\}$	2
$\{1, 5\}$	2
$\{1, 6\}$	2
$\{2, 3\}$	2
$\{2, 5\}$	2
$\{3, 5\}$	2

join

↓

Item	Freq
$\{1, 3, 5\}$	
$\{1, 3, 6\}$	X
$\{1, 5, 6\}$	X
$\{2, 3, 5\}$	

prune by  
⇒ checking  
subsets

Item	Freq
$\{1, 3, 5\}$	1
$\{2, 3, 5\}$	2

prune by  
⇒ missup

Item	Freq
$\{2, 3, 5\}$	2

↓

cannot join anymore

$\{2, 3, 5\}$  is the Frequent Item set

Example 2 Use the apriori algorithm to find frequent itemsets.

for min sup = 2

TID	Items
T <sub>1</sub>	1, 2, 5
T <sub>2</sub>	2, 4
T <sub>3</sub>	3, 2
T <sub>4</sub>	1, 2, 4
T <sub>5</sub>	1, 3
T <sub>6</sub>	2, 3
T <sub>7</sub>	1, 3
T <sub>8</sub>	1, 2, 3, 5
T <sub>9</sub>	1, 2, 3

TID	1	2	3	4	5
T <sub>1</sub>	1	1	0	0	1
T <sub>2</sub>	0	1	0	1	0
T <sub>3</sub>	0	1	1	0	0
T <sub>4</sub>	1	1	0	1	0
T <sub>5</sub>	1	0	1	0	0
T <sub>6</sub>	0	1	1	0	0
T <sub>7</sub>	1	0	1	0	0
T <sub>8</sub>	1	1	1	0	1
T <sub>9</sub>	1	1	1	0	0

Item	Frequency
I <sub>1</sub>	6
I <sub>2</sub>	7
I <sub>3</sub>	6
I <sub>4</sub>	2
I <sub>5</sub>	2

$$\text{minsup} = 2$$

↓ no pruning needed

join

Item	Frequency
{I <sub>1</sub> , I <sub>2</sub> }	4
{I <sub>1</sub> , I <sub>3</sub> }	4
{I <sub>1</sub> , I <sub>4</sub> }	1
{I <sub>1</sub> , I <sub>5</sub> }	2
{I <sub>2</sub> , I <sub>3</sub> }	4
{I <sub>2</sub> , I <sub>4</sub> }	2
{I <sub>2</sub> , I <sub>5</sub> }	2
{I <sub>3</sub> , I <sub>4</sub> }	0
{I <sub>3</sub> , I <sub>5</sub> }	1
{I <sub>4</sub> , I <sub>5</sub> }	0

prune  
⇒

Item	Frequency
{I <sub>1</sub> , I <sub>2</sub> }	4
{I <sub>1</sub> , I <sub>3</sub> }	4
{I <sub>1</sub> , I <sub>5</sub> }	2
{I <sub>2</sub> , I <sub>3</sub> }	4
{I <sub>2</sub> , I <sub>4</sub> }	2
{I <sub>2</sub> , I <sub>5</sub> }	2

join

Item	
{I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> }	
{I <sub>1</sub> , I <sub>2</sub> , I <sub>5</sub> }	
{I <sub>1</sub> , I <sub>3</sub> , I <sub>5</sub> }	X
{I <sub>2</sub> , I <sub>3</sub> , I <sub>4</sub> }	X
{I <sub>2</sub> , I <sub>3</sub> , I <sub>5</sub> }	X
{I <sub>2</sub> , I <sub>4</sub> , I <sub>5</sub> }	X

⇒

Item	Frequency
I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub>	2
I <sub>1</sub> , I <sub>2</sub> , I <sub>5</sub>	2

↓  
join

{I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub>, I<sub>5</sub>} X

The frequent itemsets are {I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub>}  
 {I<sub>1</sub>, I<sub>2</sub>, I<sub>5</sub>} X

## \* Clustering

Cluster = a collection of data objects which are similar/dissimilar

→ cluster analysis involves - given a set of data points, partition them into a set of groups, which are as similar as possible

→ a type of unsupervised learning.

→ Good clustering gives:

- (i) High intra-class similarity: cohesive within clusters
- (ii) Low Inter-class similarity: distinctive between clusters

## \* Distance Measures

→ used to measure how the similarity of 2 elements is calculated  
it influences the shape of the clusters.

These include:

- (i) Euclidean distance
- (ii) Manhattan distance
- (iii) Max Norm  $d(x, u) = \max |x_i - u_i|$
- (iv) Inner product space
- (v) Hamming Distance

## \* Applications of Cluster Analysis

- (i) intermediate step for other data mining tasks
- (ii) data summarization, compression and reduction
- (iii) dynamic trend prediction

## (iv) multimedia data analysis

### \* Clustering Methods

#### A. Partitioning Approach

- construct various  $k$  partitions & then evaluate them by some criterion

eg. k-means, k-medoids, CLARANS

#### B. Hierarchical Approach

- create a hierarchical decomposition of the set of data using some criterion
- classified using agglomerative (bottom-up) or divisive (top-down)

eg. Diana, Agnes, BIRCH, CAMELO

#### C. Density Based methods

- based on connectivity and density functions
- used to filter out noise & outliers

eg. DBSCAN, OPTICS

#### D. Grid-based Methods

- based on multiple level granularity structure
- quantize object space into finite no. of cells that form a grid structure

eg. STING, CLIQUE

## \* $k$ -Means Clustering Algorithm

① Cluster the following 8 points into 3 clusters

A<sub>1</sub>(2,10) A<sub>4</sub>(5,8) A<sub>7</sub>(1,2)

A<sub>2</sub>(2,5) A<sub>5</sub>(7,5) A<sub>8</sub>(4,9)

A<sub>3</sub>(8,4) A<sub>6</sub>(6,4)

Consider the cluster centers to be:

A<sub>1</sub>(2,10)

A<sub>4</sub>(5,8)

A<sub>7</sub>(1,2)

and the distance metric is defined as:  $D(a,b) = |x_0 - x_1| + |y_0 - y_1|$

Distance from

	C <sub>1</sub> A <sub>1</sub> (2,10)	C <sub>2</sub> A <sub>4</sub> (5,8)	C <sub>3</sub> A <sub>7</sub> (1,2)	Cluster
A <sub>1</sub> (2,10)	0+0 = 0	3+2 = 5	1+8 = 9	C <sub>1</sub>
A <sub>2</sub> (2,5)	0+5 = 5	3+3 = 6	1+3 = 4	C <sub>3</sub>
A <sub>3</sub> (8,4)	6+6 = 12	3+4 = 7	7+2 = 9	C <sub>2</sub>
A <sub>4</sub> (5,8)	3+2 = 5	0+0 = 0	4+6 = 10	C <sub>2</sub>
A <sub>5</sub> (7,5)	5+5 = 10	2+3 = 5	6+3 = 9	C <sub>2</sub>
A <sub>6</sub> (6,4)	4+6 = 10	1+4 = 5	5+2 = 7	C <sub>2</sub>
A <sub>7</sub> (1,2)	1+8 = 9	4+6 = 10	0+0 = 0	C <sub>3</sub>
A <sub>8</sub> (4,9)	2+1 = 3	1+1 = 2	3+7 = 10	C <sub>2</sub>

$c_1$   
 $A_1(2,10)$  $c_2$  $A_3(8,4)$  $c_3$  $A_2(2,5)$  $A_4(5,8)$  $A_7(1,2)$  $A_5(7,5)$  $A_6(6,4)$  $A_8(4,9)$ 

Find new cluster centers

 $c_1 = (2,10)$  $c_2 = (6,6)$  $c_3 = (1.5, 3.5)$ 

Distance from

Point	$c_1$ (2,10)	$c_2$ (6,6)	$c_3$ (1.5,3.5)	belongs to cluster
$A_1(2,10)$	0+0 0	4+4 8	0.5+6.5 7	$c_1$
$A_2(2,5)$	0+5 5	4+1 5	0.5+1.5 2	$c_3$
$A_3(8,4)$	6+6 12	8+2 4	6.5+0.5 7	$c_2$
$A_4(5,8)$	3+2 5	1+2 3	3.5+4.5 8	$c_2$
$A_5(7,5)$	5+5 10	1+1 2	5.5+1.5 7	$c_2$
$A_6(6,4)$	4+6 10	6+2 2	4.5 0.5 5	$c_2$
$A_7(1,2)$	1+8 9	5+4 9	0.5+1.5 2	$c_3$
$A_8(4,9)$	2+1 3	2+3 5	2.5+5.5 8	$c_1$

The clusters are:

$C_1$	$C_2$	$C_3$
$A_1(2,16)$	$A_3(8,4)$	$A_8(2,5)$
$A_8(4,9)$	$A_4(5,8)$	$A_7(1,2)$
	$A_5(7,5)$	
	$A_6(6,4)$	

Repeat this process

## K-Medoids Clustering

→ also called PAM

### Partitioning Around Medoid

A medoid is a point in the cluster whose dissimilarities with all the other points in the cluster are minimum

The cost in the K-Medoids algorithm is given as:

$$C = \sum_{c_i} \sum_{p_i \in c_i} |p_i - c_i|$$

Algorithm: Initialize  $k$  random points out of the  $n$  data points as medoids

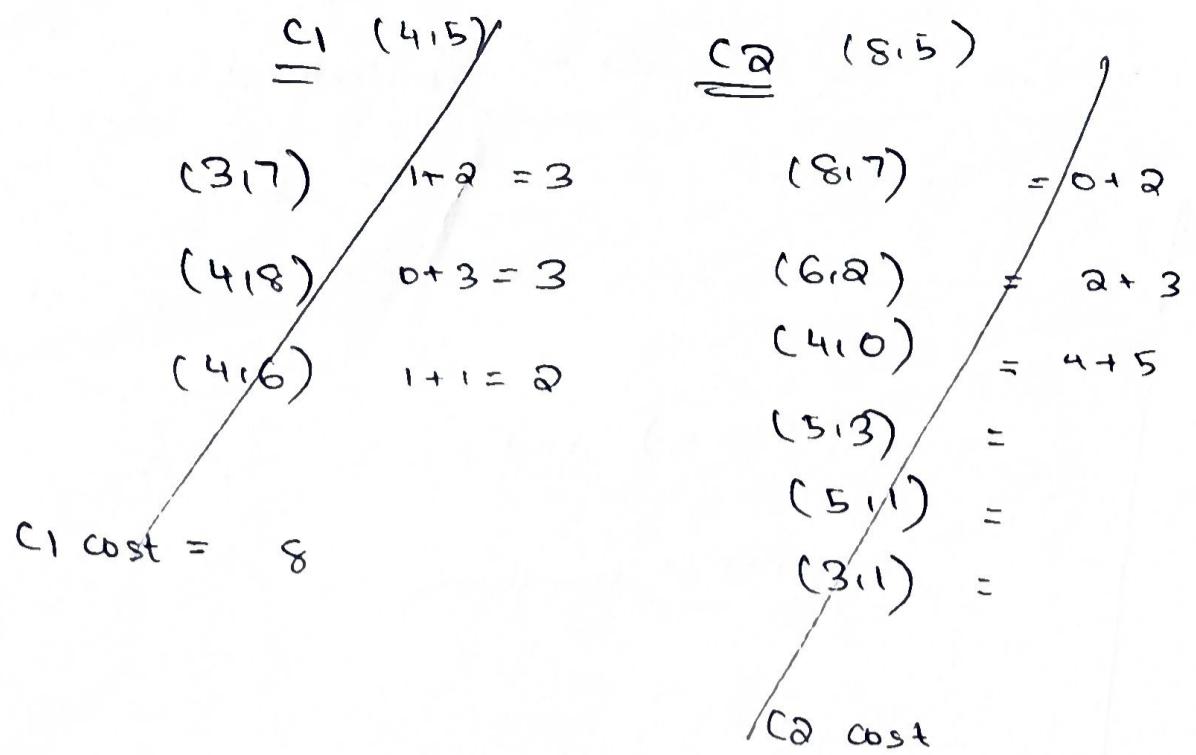
2. Associate each data point to the closest medoid by using any common distance metric method
3. While the cost decrease:
  - (i) swap  $m$  and  $o$ , associate each data point to the closest medoid and recompute the cost
  - (ii) If the total cost is more than that in the previous step, undo the swap

\* Apply K-medoids clustering on the following dataset

X	4
8	7
3	7
4	9
9	6
8	5
5	8
7	3
8	4
7	5
4	5

Let  $(8, 5)$  and  $(4, 5)$  be the initial medoids

Point	Distance from $(8, 5)$	Distance from $(8, 5)$	Cluster
$(8, 7)$	$4+2$ 6	$0+2$ 2	$C_2$
$(3, 7)$	$1+2$ 3	$5+2$ 7	$C_1$
$(4, 9)$	$0+4$ 4	$4+4$ 8	$C_1$
$(9, 6)$	$5+1$ 6	$1+1$ 2	$C_2$
$(8, 5)$	$4+0$ 4	0	$C_2$
$(5, 8)$	$1+3$ 4	$3+3$ 6	$C_1$
$(7, 3)$	$3+2$ 5	$1+2$ 3	$C_2$
$(8, 4)$	$4+1$ 5	1	$C_2$
$(7, 5)$	$3+0$ 3	1	$C_2$
$(4, 5)$	0	-	$C_1$



$$c_1 = (4,5)$$

$$c_2 = (8,5)$$

$$(3,7) = 1+2 = 3$$

$$(9,6) = 1+1 = 2$$

$$(4,9) = 0+4 = 4$$

$$\cancel{(6,6)} = 0+0 = 0$$

$$(5,8) = 1+3 = 4$$

$$(7,3) = 1+2 = 3$$

$$\text{cost } c_1 = 11$$

$$(8,4) = 0+1 = 1$$

$$(7,5) = 1+0 = 1$$

$$(8,7) = 0+2 = 2$$

$$= 9$$

$$\underline{\text{Total cost} = 20}$$

choose some random point find distance for (4,5) and 0  
 random point

compute cost again.

## Outliers

- An outlier is a point that lies abnormally far away from the other values in a dataset.
- A data point is an outlier if it is 1.5 times the interquartile range greater than the third quartile or 1.5 times the interquartile range less than the first quartile.

Eq1 25<sup>th</sup> percentile (Q1) income = 15,000

75<sup>th</sup> percentile (Q3) income = 120,000

Find the range of outliers

$$IQR = 105,000$$

$$\text{Lower boundary} = Q_1 - 1.5 \times IQR$$

$$= 15,000 - 1.5 \times 105,000$$

$$= -142,500$$

$$\text{Upper boundary} = Q_3 + 1.5 \times IQR$$

$$= 120,000 + 1.5 \times 105,000$$

$$= 277,500$$

Eq2 25<sup>th</sup> percentile (Q1) for how long individuals can hold their

breath = 15 seconds

75<sup>th</sup> percentile (Q3) = 75 seconds

$$IQR = 60 \text{ sec}$$

$$\text{Lower boundary} = Q_1 - 1.5 \text{ IQR}$$
$$= 15 - 1.5(60)$$
$$= \underline{\underline{-75}}$$

$$\text{Upper boundary} = Q_3 + 1.5 \text{ IQR}$$
$$= 15 + 1.5(60)$$
$$= 165$$

eq3 - 25<sup>th</sup> percentile ticket sales = \$2 million

75<sup>th</sup> percentile ticket sales = \$15 million

$$\text{IQR} = \$13 \text{ million}$$

$$\text{Lower boundary} = Q_1 - 1.5 \text{ IQR}$$
$$= 2 - 1.5(13)$$
$$= -17.5$$

$$\text{Upper boundary} = Q_3 + 1.5 \text{ IQR}$$
$$= 15 + 19.5$$
$$= \underline{\underline{34.5}}$$