

# Association Rule mining

OR

## Association Rule Learning

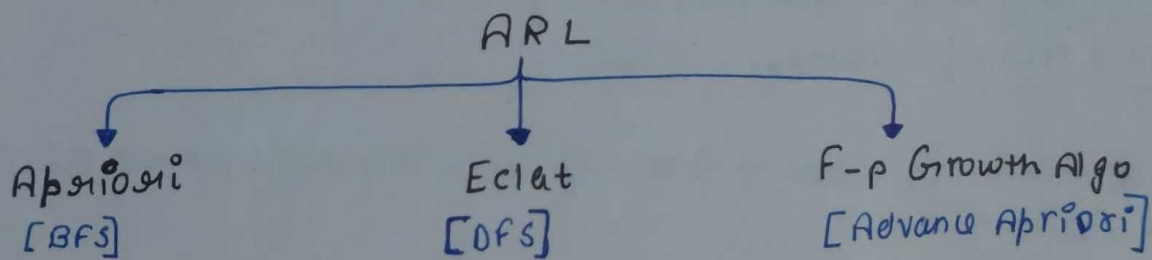
- Type of unsupervised learning that checks for dependency of one data item on another data item and maps accordingly so that it can be more profitable

If A then B  $[A \Rightarrow B]$  Consequent  
 Product  $\uparrow$  Antecedent

For ex  $\rightarrow$  Buying Bread [A]  $\left\{ \begin{array}{l} \text{milk} \\ \text{egg} \\ \text{butter} \end{array} \right\}$  [B] 3 options a person might buy after buying bread maps A to possible B's To make most profit

[If customer buys Bread, he's likely 70% to buy milk]

$\uparrow$   
THE WHOLE IDEA



Working  $\rightarrow$  Support  $\rightarrow$  Frequency of A  $\rightarrow \text{Supp}(x) = \frac{\text{Freq}(x)}{T}$

OR  $\rightarrow$  Relative Frequency that the items in the rule appear together in the dataset.

[0.5] = A and B occur together in half of the Transaction

$$\text{Support}(A \Rightarrow B) = \frac{(T \text{ containing Both A and B})}{\text{Total Transaction}} \rightarrow \boxed{P(A \cap B)} \quad \boxed{\frac{\text{Freq}(A, B)}{N}}$$

**Confidence**  $\rightarrow$  Cond<sup>n</sup> Probab of an itemset occurring given that another itemset has already occurred.  $\rightarrow$  Ratio of the no. of Transaction contain Both itemsets to the total no. of Transactions containing the first one  $\rightarrow \frac{P(A \cap B)}{P(A)} * [A \Rightarrow B = \frac{\text{Supp}(AB)}{\text{Support}(A)}]$

[no. of times  $A \Rightarrow B$  were true] have found to be TRUE

- Lift → Strength of Association B/w two items in a dataset.
- ↳ Dividing the probability of finding the two items together in a Transaction by the probability of finding each item independently in a Transaction.
- ↳  $Lift(A \Rightarrow B) = \frac{P(A \cap B)}{P(A) * P(B)}$  |  $> 1 \rightarrow$  Strong association | To see if we need to spend time for that Rule
- $\downarrow$   
 $\frac{\text{Support}}{P(A) * P(B)}$
- $\leq 1 \rightarrow$  weak association  
less

## Apriori Algorithm

- Uses frequent item sets to generate 'association Rule'.
  - Concept → subset of a frequent itemset must also be a frequent itemset
  - \* Frequent itemset is an itemset whose support value > Threshold value
  - Bottom-up Approach → Identify individual items that appear in at least a minimum no. of T, and then extend these items into larger and larger itemsets, using Join and Prune strategy.
- 1] Set a minimum support and confidence threshold for the desired association Rules.
  - 2] Identify all individual items that appear in at least minimum no. of T  
 ↳ Frequent itemsets  $[S_1, S_2, S_3, S_4]$   $\uparrow$  [ $>$  Threshold value]
  - 3] Use the frequent itemsets to generate candidate Rules.  
 $[\text{subsets}] \uparrow$   $[X \Rightarrow Y]$  [ $X, Y$  both of sets of items]
  - 4] For each candidate rule → calculate the support and confidence of the Rule.  
 $\uparrow$  [verifying if those subsets on the basis of confidence level]
  - 5] Select the Rules that meet the min support and confidence Threshold and [output these as the Resulting Association Rules].  
 $\downarrow$
  - 6] Repeat the process, using the Resulting Association Rules as the new set of frequent itemsets and generate new candidate rules from these.  
 ↳ [Process continue until no more association rules can be found]



in a  
density

\* confidence Threshold = 60%

TID	Items
T1	1 3 4
T2	2 3 5
T3	1 2 3 5
T4	2 5
T5	1 3 5

\* Support threshold = 20%

Items	Freq	Support
1	3	60% → <u>S<sub>1</sub></u>
2	3	60%
3	4	80%
<u>4</u>	1	20% — <u>X</u>
5	4	80%

Rem elements → 1, 2, 3, 5

Set of 2 <u>S<sub>2</sub></u>	→ Items	Freq	Support
	<u>1, 2</u>	1	20% — <u>X</u>
	1, 3	3	60%
	1, 5	2	40%
	2, 3	2	40%
	2, 5	3	60%
	3, 5	3	60%

Rem elements → 1, ~~2~~, 3, 5 → using S<sub>2</sub> Create S<sub>3</sub>

Set of 3 <u>S<sub>3</sub></u>	→ Items	Freq	Support
	1, 2, 3 X		
	1, 2, 5 X		
	1, 3, 5 →	2	40%
	2, 3, 5 →	2	40%

[At S<sub>3</sub> we can check for Pruning → If subset are not frequent itemsets so that set is not frequent itemset]

$[1, 2, 3 \rightarrow \{1\}, \{2\}, \{3\}, \{1, 2\}, \{1, 3\}, \{2, 3\}]$  from S<sub>2</sub> so  $(1, 2, 3) \rightarrow \underline{X}$   
 $[1, 2, 5 \rightarrow \{1\}, \{2\}, \{5\}, \{1, 2\}, \{2, 5\}, \{1, 5\}]$  so  $(1, 2, 5) \rightarrow \underline{X}$

Set of 4 <u>S<sub>4</sub></u>	Items	Freq	Support
<u>S<sub>4</sub></u> →	(1, 2, 3, 5) →	1	20% <u>X</u> Revert Back to <u>S<sub>3</sub></u>

Ans 1, 3, 5 and 2, 3, 5

Frequency itemsets are identified on the basis of support Now find rules from freq itemsets on the basis of confidence.

→ (1, 3, 5) → (1, 3), (3, 5), (1, 5), (1), (3), (5)  
 → (2, 3, 5) → (2, 3), (2, 5), (3, 5), (2), (3), (5)

verify these subsets  
using confidence

$$\boxed{1, 3, 5} \rightarrow (1, 3) \mid (1, 5) \mid (3, 5) \mid (1) \mid (3) \mid (5)$$

$\boxed{1 \Rightarrow 3}$  (1,3) Rule 1  $\rightarrow$

$$S \rightarrow (I-S) \quad \begin{matrix} \text{proper} \\ \text{subset} \end{matrix} \quad (S \text{ recommends } I-S)$$

$$\text{Conf-level} \rightarrow \frac{\text{support}(I)}{\text{support}(S)}$$

R1  
 $(1, 3) \rightarrow \frac{2/5}{3/5} = \frac{2}{3} > 60\% \text{ - } \checkmark$

R2  
 $(1, 5) \rightarrow \frac{2/5}{2/5} = 100\% > 60\% \checkmark$

R3  
 $(3, 5) \rightarrow \frac{2/5}{2/5} > 60\% \checkmark$

R4  
 $(1) \Rightarrow \frac{2/5}{3/5} > 60\% \checkmark$

R5  
 $(3) = \frac{2/5}{4/5} = 50\% < 60\% \text{ X}$

R6  
 $(5) = \frac{2/5}{4/5} = 50\% < 60\% \text{ X}$

$\boxed{R}$

library(arules) //apriori

library(arulesviz) //plot(rules)

data(Groceries)

data\_lst <- as(Groceries, 'list')

data\_transaction <- as(Groceries, 'transactions')

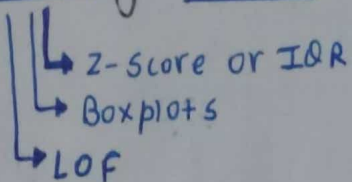
$\boxed{\text{rules}} \leftarrow \text{apriori}(\text{data} = \text{Groceries}, \text{parameters} = \text{list}(\text{support} = 0.2, \text{confidence} = 0.15))$

inspect(rules) OR inspect(sort(rules, by = 'lift'))

$\boxed{\text{plot(rules)}}$

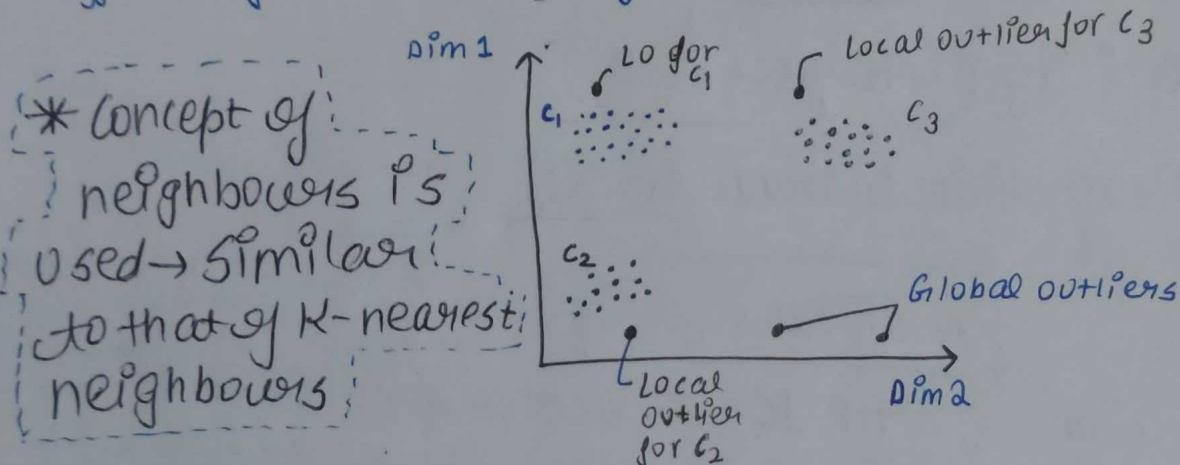


## • Local Outlier Factor Anomaly Detection → LOF



Global Outliers → Data points which are 'significantly' 'different' from the rest of the dataset.

Local Outliers → Data points which are significantly different from their [neighbours in the dataset]



LOF → • A score (scalar value) = LOF is the Deciding factor

For a given Dataset

$$D_n = \{ (x_i, y_i) \mid x_i \in \mathbb{R}^2, y_i \in \{xyz\} \}$$

$$* LOF(x_i) = \frac{\sum_{x_j \in N(x_i)} Lrd(x_j)}{|N(x_i)|} \times \frac{1}{Lrd(x_i)}$$

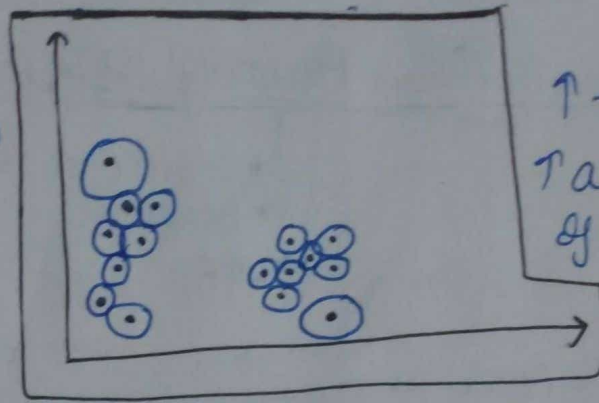
$|N(x_i)|$  = No. of elements in the neighbourhood of  $x_i$

$Lrd(x_i)$  = Local Reachability Density of  $x_i$

• WORKING → ① LOF is assigned to each data points. These assigned LOF scores of each data points are compared to find outliers.

② ↑ value of LOF of any DP → ↑ chances of that DP being an outlier

Radius of circle  $\propto (=)$  to the LOF scores of each Data Point.



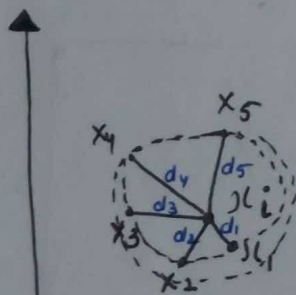
↑ the Radius,  
There are the chances  
of that DP being  
an outlier.

## 2] Parameters for calculating LOF

- i)  $k$ -distance ( $x_i$ )
- ii] Nearest Neighbor  $N_k(x_i)$
- iii] Reachability Distance  $(x_i, x_j)$
- iv] Local Reachability Distance  $lrd(x_i)$

$k$ -distance ( $x_i$ )

The  $k$ -Distance of a DP  $x_i$  in a dataset is the Distance of the  $k^{\text{th}}$  nearest neighbour of  $x_i$  from  $x_i$ .



→  $x_1, x_2, x_3, x_4, x_5$  are nearest neighbours of  $x_i$

→  $d_i$  is the distance b/w  $x_i$  and  $x_j$  where  $x_j \in (x_1, x_2, x_3, x_4, x_5)$

→  $d_1 < d_2 < d_3 < d_4 < d_5$

• 5-distance ( $x_i$ )

↳ Distance b/w  $x_i$  and the 5th nearest neighbour of  $x_i$  i.e.  $d_5$

• 4-distance ( $x_i$ )

↳  $d_4$

• 1-distance ( $x_i$ )

↳  $d_1$

$x$ -data

`sklearn.neighbors`  
`(import) LocalOutlierFactor`

`lof = LocalOutlierFactor(n_neighbors=20, contamination=0.1)`

`outlier_scores = lof.fit_predict(x)`

↳ 1 → Inlier (+ve)  
↳ -1 → Outlier (-ve)

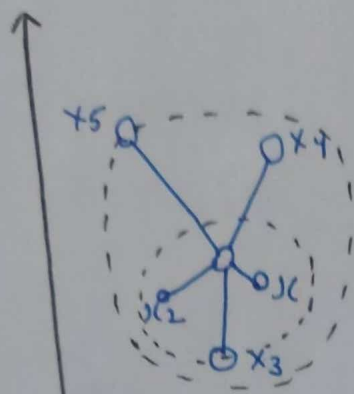
`outlier_indices = np.where(outlier_scores == -1)`  
↳ `np.delete(x, outlier_indices, axis=0)`



chances  
of being  
missed,

# Nearest Neighbour $N_k(x_i)$

- Just like K-Nearest Neighbours of D.P  $x_i$  denoted by  $N_k(x_i)$  is a SET of all points that belong to the  $k^{th}$  nearest neighbour of  $x_i$  (i.e. points in the neighborhood of  $x_i$ )



$$N_5(x_i) = \{x_1, x_2, x_3, x_4, x_5\}$$

$$N_3(x_i) = \{x_1, x_2, x_3\}$$

$|N(x_i)|$  = Total no. of data points present in the neighborhood of  $x_i$  (5)

actual meaning

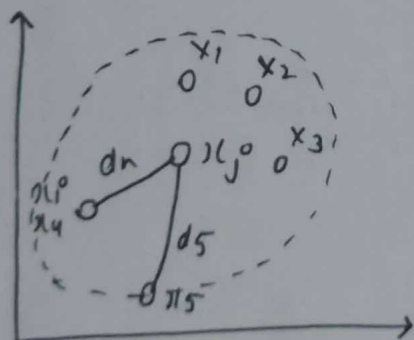
2 points ke Beech ka actual distance  
ya source wala ka  
K-distance (Total neighbours)

## Reachability - Distance $(x_i, x_j)$

→ The R-D of a DP  $x_j$  from  $x_i$  is the max of K-distance of  $x_j$  and the actual distance b/w  $x_i$  and  $x_j$

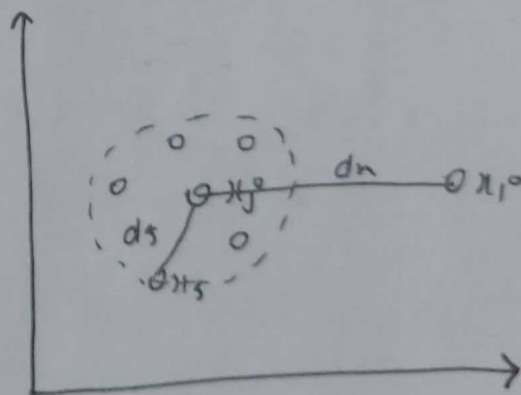
$$R-D(x_i, x_j) = \max(K\text{-distance}(x_j), \text{dist}(x_i, x_j))$$

Case 1 ( $x_i$  lies within the cluster)



$$\max(d_s, d_n) = d_s$$

Case 2 ( $x_i$  lies outside the cluster)



$$\max(d_s, d_n) = d_n$$

# Local Reachability Density $lrd(x_i)$

- $lrd$  of a data point  $x_i$  is the Inverse of the Avg RD of  $x_i$  from its neighborhood.

\* In Short  $\rightarrow$  It measures How CLOSE the neighborhood of a points of  $x_i$  are from it.

\* If  $lrd(x_i)$  is HIGH  $\rightarrow x_i$  is in dense neighborhood

\* If  $lrd(x_i)$  is LOW  $\rightarrow x_i$  is in sparse neighborhood

$$lrd(x_i) = \frac{1}{\sum_{x_j \in N(x_i)} \left\{ \frac{RD(x_i, x_j)}{|N(x_i)|} \right\}}$$

L Set                      L Total Number

$$\underline{LOF}(x_i) \rightarrow \frac{\sum_{x_j \in N(x_i)} lrd(x_j)}{|N(x_i)|} \times \frac{1}{lrd(x_i)} \quad \begin{matrix} A & B \end{matrix}$$

a)  $LOF(x_i)$  is large when A is large and B is small  
Outlier

b)  $LOF(x_i)$  is small when A will be small and B will be large  $\rightarrow$  Not outlier

- $\sim 1$  Similar Density as neighbors,
- $< 1$  Higher Density than neighbors (Inliers)
- $> 1$  Lower Density than neighbors (outliers)

Imp  $\rightarrow$  we can either visualize the circles to show outliers ~~or~~ we can show red dots.