

Unit - 3

→ Mining Data Streams :- process of extracting knowledge structures from continuous rapid data.

→ Introduction to Stream Concepts

- ↳ Temporally ordered, fast changing, massive
- ↳ Huge volume of continuous data
- ↳ fast changing

Data Stream mining is the process of extracting knowledge & from continuous rapid data records which comes to the system in a stream.

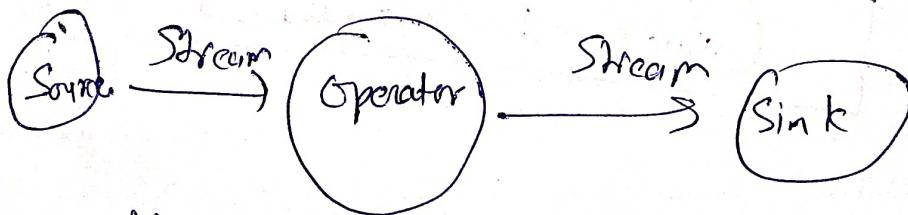
- A data stream is an ordered sequence of instances in time.
- Data Stream transmit from a Source & receives at processing end.
- Also refers to communication of bytes

Ex:-

- Monitoring & detection
- analysis of social media
- Watching online video lecture (Rewind, forward)

Data Stream Model.

- Stream is data in model
- Three approaches for updating the end points

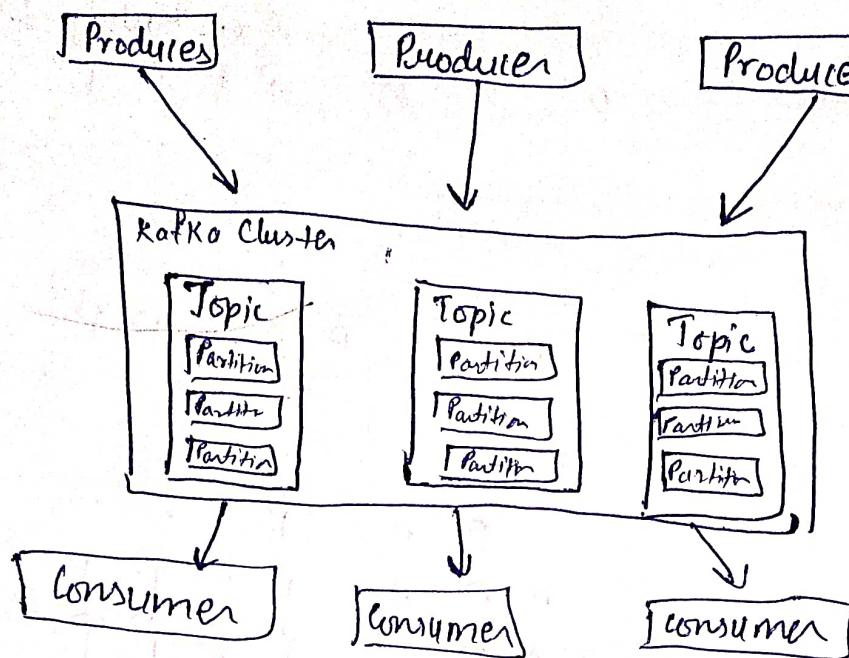


Types of Data Stream

- Data Stream :- Sequence of tuples
- Transaction data Stream: Credit Card, Web
- Measurement data Stream: Sensor data, IP Network, Earth Climate

Model & Architecture

Architecture



Continuous flow
Stream management
can be done
externally.

Tweets

Stream Processing

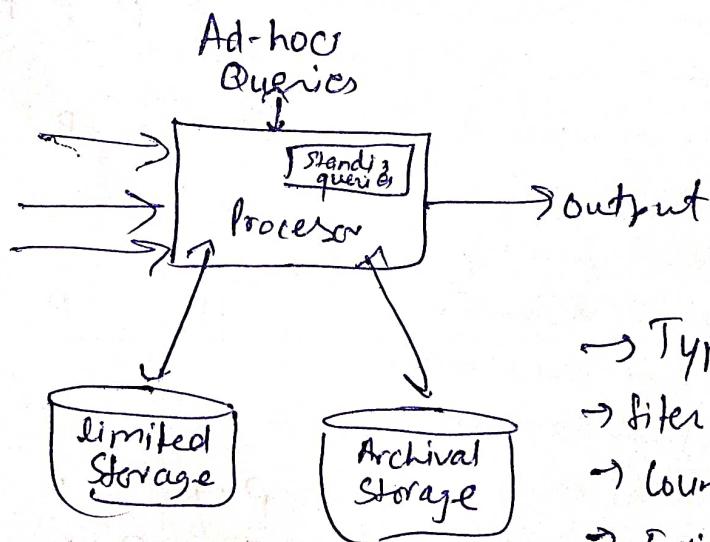
- Data scope
- Data size
- Performance

Model:-

- The input element enter at a rapid rate at one or more input ports.
- The system cannot store the entire stream accessibly.

1, 5, 2, 7, 0, 9, 3
9, n, v, t, y, k, b
0, 0, 1, 0, 1, 1, 6
← time

Stream Entering
Each Stream
is Composed
of element D
tuples.



- Types of query's
- Filtering a data stream
- Counting distinct elements
- Estimating moment
- finding frequent element

Sec-A

(63)

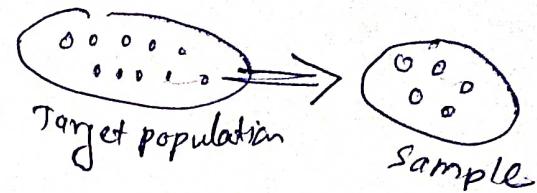
7S - B

~~Architecture~~ Sampling data in a stream

- Since we can not store the entire stream, one obvious approach is to store a sample random selection
- Two different problems? Probability, Non Probability
 - (1) Sample a fixed proportion of element in stream. every member has equal chance to be selected
 - (2) maintain a random sample of fixed size over a potentially infinite stream

① Sampling fixed proportion Search engine query

- Stream of tuple (user, query, time)
- Answer question such as: How often did a user run the same query in a single day?
- Have space to store $\frac{1}{10^m}$ of query stream
- Random integers



Covid \rightarrow Vaccine \Rightarrow
different by people

Ex:- Say $s = 2$

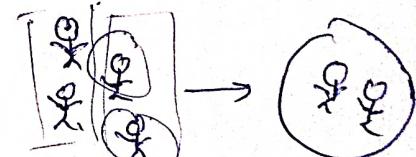
Stream: a x c y z, k c deg ...

At $n=5$, each of the first 5 tuples is included in the sample S with equal prob.

At $n=7$, each of the first 7 tuples is included in the sample S with equal prob

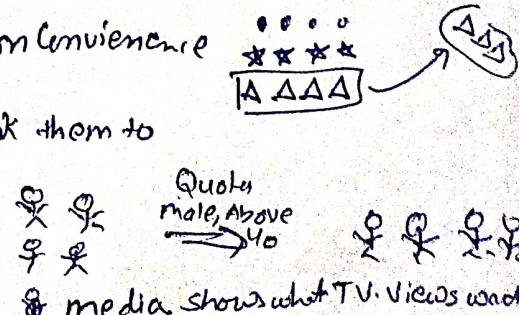
Probability
Sampling

- ↳ Simple Random: - equal chance of selection
- ↳ Systematic Sampling: - every n th element selected
- ↳ Cluster Sampling: - randomly selecting cluster
- ↳ Stratified Sampling: - dividing the population in groups



Non Probability Sampling

- ↳ Convenience Sampling: - sample based on convenience
- ↳ Snowball Sampling: - select sample & ask them to refer them
- ↳ Quota Sampling: - some characteristics
- ↳ Judgemental Sampling: - based on judgment



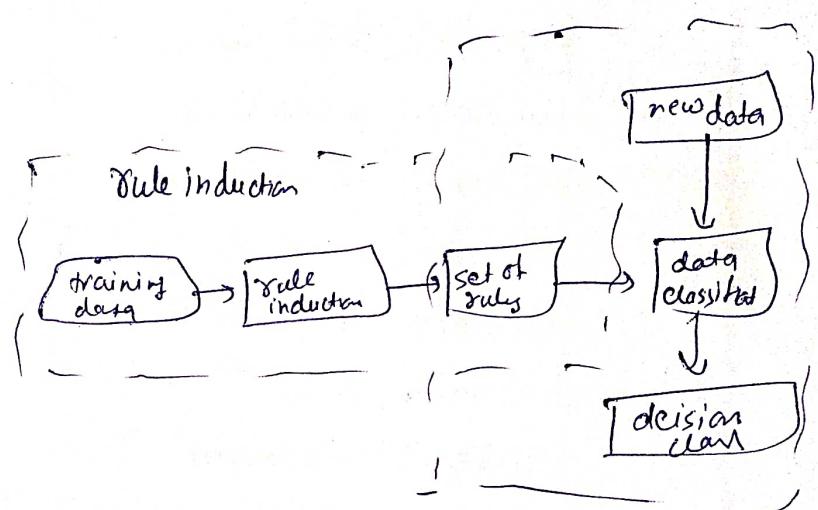
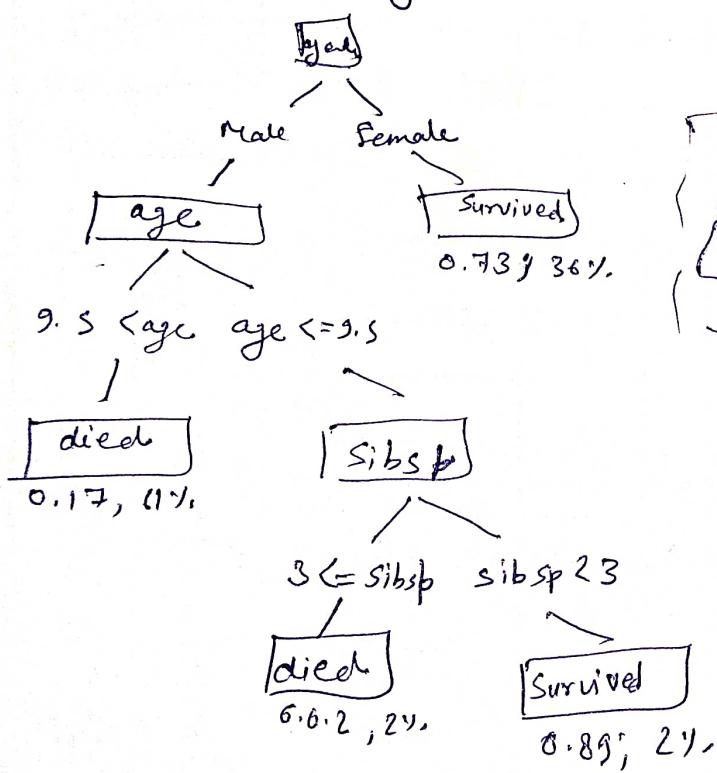
* Rule Induction

⑥

- actions are given and we have to discover the rules.
- ex:- Do not give discount on 2 items. At 8 a.m on weekday, traffic is heavy, At 8 p.m on Sunday.
- This scheme make use of If -Then rules.
- This is the reverse of a rule-based agt, where the rules are given and the agent must act.

- Prevalence = Probability that LHS & RHS occur together
- Predictability = probability of RHS given LHS.
- Builds a Decision tree representing rule for distinguish between different cases.

Survival of Passengers in Titanic



Neural Networks

- Collection of neurons that take input and, in conjunction with information from other nodes, develop output without programmed rules.
- A simple neural network has an input layer, output layer and hidden layer.
- These neurons are computational units, These networks transform data like the pixels in an image or the word in a document until they can classify it as an output, such as naming an object in an image or tagging unstructured text data.

→

Learning and Generalizations

Neural Networks learn by adjusting the strength of their connections to better convey input signals through multiple layers of neurons associated with the eight general concepts.

- Learning :- The network must learn decision surfaces from a set of training patterns so that these training patterns are classified correctly.
- Generalization :- After training, the network must also be able to generalize, i.e. correctly classify test patterns it has never seen before.

Competitive learning rule

→ In CL, neurons compete among themselves to be activated.

Learning

- Supervised :- Applications :- Speech Recognition, Bioinformatics, Spam Detection, object Recognition
- Unsupervised, Amazon using clustering,

Supervised

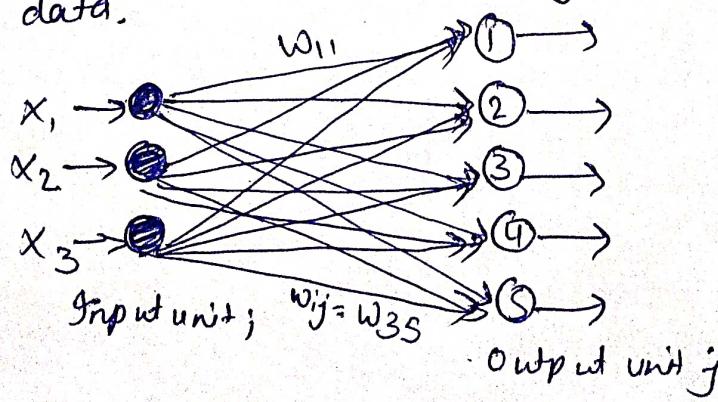
- are trained using labeled data
- takes direct feedback to check if the output is correct or not.
- simpler method
- predicts the Output
- More Accurate
- Categorized in Classification & Regression problems
- ex:- House prices, checking Weather

Unsupervised

- using unlabeled data
- does not take any feedback
- computationally complex
- hidden pattern in data
- less Accurate
- Clustering & Association problems
- finding customer segment in Marketing data (gender, age, location)

Competitive learning

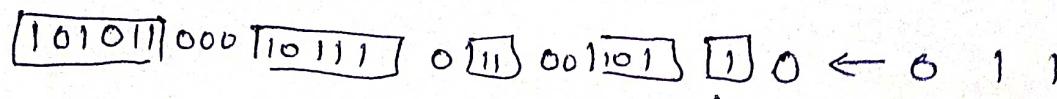
- form of unsupervised learning,
- In which nodes compete for the right to respond to a subset of input data.



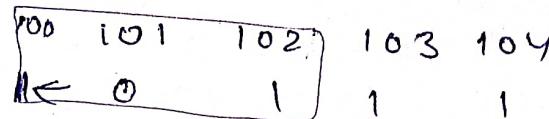
Given the give data stream, let us assume the new bit arrives from the (2) sig bit

If new bit = 0 \Rightarrow No change in bucket,

Timestamp \rightarrow 87 ... 92 ... 93 ... 98 ... 10 ... 101 102 103



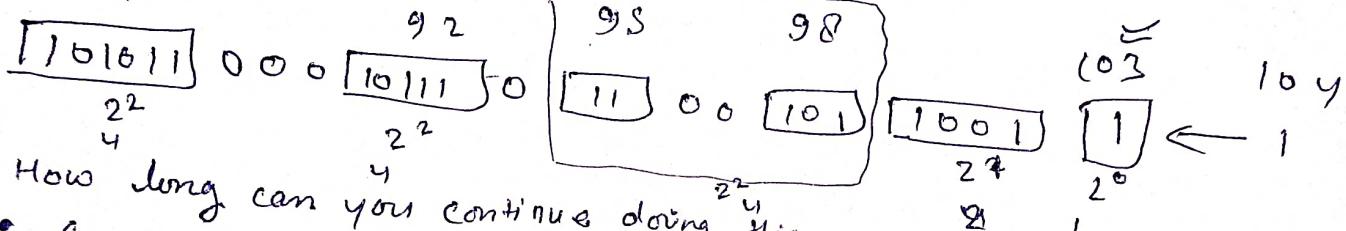
when new bit = 0 enters
it does not effect bucket



when new bit = 1 enters

Combine

when next bit = 1 enters



How long can you continue doing this...

- Current time stamp - left most bucket timestamp $< N$ (window size = 24 bits)
- Eg : $103 - 87 = 16 < 24$ so I continue

finally answer to query

How many no. of 1's in last 20 bits?

≈ 20 bit

\rightarrow No. of 1's in last 20 bits ≈ 11

X X

Real time Analytics Platform (RTAP)

\rightarrow Allows Business to get insights.

\rightarrow ex :- Continuously updated Customer activity like page views and shopping cart.

• Components of RTAP:

\rightarrow Aggregator \rightarrow Compiles real time streaming data

\rightarrow Broker \rightarrow Makes data in Real time available for use.

\rightarrow Analytics Engine \rightarrow Correlates values & blends data streams together during analysis

\rightarrow Stream Processor \rightarrow Executes Real time app analytics & logic by receiving & sending data stream.

Real Time Analytic ~~System~~^{uni} Platform

-> Amazon uses
Amazon kinesis

- Refers to finding meaningful patterns in data at the actual time of receiving
- Analyse the data, correlate
- Predicts the outcome in real time
- Manages & processes data
- evaluate Business intelligence
- Apache Spark Streaming - a Big Data platform for data stream analytics in real time
- Cisco Connected Streaming Analytics (CSA) - delivers insights from high-velocity streams of live data from multiple sources
- Oracle Stream Analytics (OSA) - a platform provides a graphical interface to fast data
- SAP HANA → streaming analytics tools
- IBM Stream Computing → a data streaming tool that analyze a broad range of streaming data - unstructured text, audio, video, sensor - to spot the risk & opportunities

Applications

- 1) frauds detection systems for online transactions.
- 2) Log analysis for understanding usage pattern.
- 3) Click analysis for online Recommendation
- 4) Social Media Analytics
- 5) Push notification to the customers for location based platforms

disadvantage

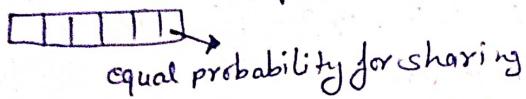
→ we don't know the value of n

→ we need to define window size (by storage capacity) → Count the element in every window of size k

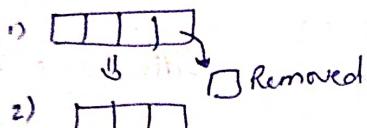
Sampling :-

Unit	1	2	3	4	5	6
x	1	1	2	2	2	3

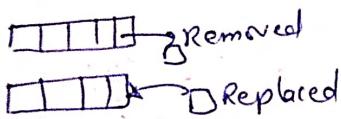
- Simple Random Sample



- Simple Random Sample Without Replacement



- SRS With Replacement



- Stratified Sampling



Counting distinct element (Flajolet Martin Algo.)

→ Count the element in every window of size k

$$a = [1, 2, \dots, 2, 1, 3, 1, 1, 3]$$

$$k = 4$$

Output

$$\text{Total windows} = ((n-k+1) \times k^2)$$

2
3
3
2
2

$$\begin{array}{ccccccc} & 2 & & x & & x & \\ \text{Stream} : & - 1, 4, 2, 1, 2, 4, 4, 4, 1, 2, 4, 1, 7 - & x \\ h(a) = (3x+1) \bmod 8 & & & & & & \end{array}$$

$$\therefore \text{After } h(a) = 4, 3, 2, 4, 2, 3, 3, 3, 4, 2, 3, 4,$$

$$\therefore \text{Convert it to binary} = 100, 011, 010, 100, 010, 011, 011, 100, 010, 011, 100, 010$$

$\stackrel{3}{=}$ Count no. of 0s (Trailing zero \Rightarrow end zero's)

$$r(a) = 2, 0, 1, 2, 1, 0, 0, 0, 2, 1, 0, 2, 1$$

$$\stackrel{4}{=} R = \max[r(a)] = 2$$

$$\stackrel{5}{=} E = 2^R = 2^2 = 4$$

There are 4 distinct elements in Stream

$$1, 2, 4, 7$$

- Estimating moments [Mean Mode Median]
- measuring data uniformity
- moments for stream: 1, 3, 2, 1, 2, 3, 4, 3, 1, 2, 3, 1
- occurrence = $(1 \times 4), 2 \times 3, 3 \times 1, 4 \times 1$
- 0th moment = $4^0 + 3^0 + 2^0 + 1^0 = 4$ (Distinct elements)
- 1st moment = $4^1 + 3^1 + 2^1 + 1^1 = 12$ (Stream length)
- 2nd moment = $4^2 + 3^2 + 2^2 + 1^2 = 42$ (Surprise Number)

Moments :-

Stream distinct elements (frequency, how they are distributed)

→ A is ordered set

→ $m_i \rightarrow$ no. of times value i occurs

→ kth moment is

$$\sum_{i \in A} (m_i)^k$$

→ 0th moment = no. of distinct element.

→ 1st moment = Length of the stream.

→ 2nd moment = Surprise no. S = a measure of how uneven the distribution is

Ex:- Stream of length 100 1 3 4 5 ... 100
11 distinct values
find 2nd moment

→ Item count: 10, 9, 9, 9, 9, 9, 9, 9, 9, 9
Surprise S = 910

$$= 10 \times (1)^2 + 9 \times (10)^2$$

→ Item count: 90, 1, 1, 1, 1, 1, 1, 1, 1, 1,
= 8, 110

Ams method
→ Alon - matias - Szegedy

→ To calculate second moment.

→ estimated moments provided, gives an unbiased estimate

→ X.val X.el

$$Q_x = \underline{a}, b, \underline{c}, \underline{b}, d, a, c, \underline{d}, \underline{a}, b, d, c, \underline{a}, \underline{a}, b$$

$$a = 3, b = 4, c = 3, d = 3$$

$$(3)^2 + (4)^2 + (3)^2 + (3)^2 = 59$$

Alc to Ams algorithm

$n = 15$ bits (lets assume no. of stream)

$$f(x) = n(2c - 1)$$

we will calculate variables on different time interval.

$$x_1 = \{c, 3\}$$

$$x_2 = \{d, 2\}$$

$$x_3 = \{a, 2\}$$

$$\text{Calculate estimate} = f(x_1) = n(2c - 1) \\ = 15(2 \times 3 - 1) \\ = 15 \times 5 = 75$$

$$f(x_2) = 15(4 - 1) \\ = 15 \times 3 = 45$$

$$f(x_3) = 15(3)$$

$$\text{Final estimate} = \frac{75 + 45 + 45}{3} = \frac{165}{3} = \cancel{\cancel{55}}$$

Unit- 4

Frequent itemset:- Set of words that occurs together in some minimum fraction of document in a dataset.

Mining frequent itemsets (Association Rule Mining)

- Searches for frequent items in the data-set.
- Frequent mining shows which items appear together in a transaction or relation.

Need of Association Mining:

$$\rightarrow [\text{milk}] \wedge [\text{bread}] \Rightarrow [\text{butter}]$$

Important Definitions:-

1. Support: It is one of the measure of interestingness

Unit - 4

Frequent Itemsets & Clustering.

⇒ Mining frequent itemsets.

- Association Mining →
- If there are two items X & Y purchased frequently then it's good to put them together in stores or provide some discount. It increase Sale
 $[milk] \wedge [bread] \Rightarrow [butter]$
- Support : measure of interestingness, tells about usefulness & certainty of rule.
 $\text{Support}(A \rightarrow B) = \text{Support_Count}(A \cup B)$

- Confidence: A confidence of 60% means that 60% of customer who purchased a milk & bread also bought butter.

$$\text{Confidence}(A \rightarrow B) = \frac{\text{Support_Count}(A \cup B)}{\text{Support_Count}(A)}$$

- Maximal Itemset: An itemset is maximal frequent if none of its supersets are frequent
- Closed Itemset: if none of immediate superset have same support.
- K-itemset: contains k-itemset.

Ex :-

TransactionId	Items
1	{A, C, D}
2	{B, C, D}
3	{A, B, C, D}
4	{B, D}
5	{A, B, C, D}

- Let say minimum support count is 3
- If all maximal frequent ⇒ closed ⇒ frequent

→ Market basket Model

Associate Rule Mining (ARM)

- Market Basket Analysis
- Affinity Analysis
- Set of items in a transaction is called Market Basket
- Mostly used in Retail
- if 'A' then 'B' $\{ A \Rightarrow B \}$

- Support (s): % of transaction of both 'A' & 'B'
measure frequency $(A \rightarrow B) = P(A \cap B)$ } measures frequency of association
- Confidence (c): transaction set 'T' if 'C' is true,
Strength of association of item 'B' is present in all the transaction in
 $C = P(B|A) = \frac{P(A \cap B)}{P(A)}$ }

Parameters:-

- (i) finding all items that appear frequently in transaction } support
- (ii) finding strong association among frequent items } confidence value.

Apriori Algorithm:-

generate candidate item set

- (i) All singleton item are candidate
- (ii) This algo is given by R. Agrawal

By R. Srikant in 1994 for finding frequent itemset in a data set.

Steps:-

- 1 Create a table containing support count of each item.
- 2 Compare candidate set item's support count
- 3 Generate candidate set Using Condition
- 4 Check all subset

Tid	item
1	2, 3
2	1, 3, 5
3	1, 2, 4
4	2, 3

Eg:- min support = 2

item	Support	itemsets	Support
1 → 2		{1, 2} → 1	
2 → 3		{1, 3} → 1	
3 → 3		{2, 3} → 2	elim.
4 → 1			
5 → 1			eliminated

Apriori Algorithm (To generate association rule)

- Min Support = 50%.
- Threshold confidence = 70%.

eg:-

Tid	Items			
100	1	3	4	
200	2	3	5	
300	1	2	3	5
400	2	5		

I

Item set	Support
1	2/4 → 50%.
2	3/4 → 75%.
3	3/4 → 75%.
4	1/4 → 25%.
5	3/4 → 75%.

{ item set → 1, 2, 3, 5 }

II

Itemset	Support
{1, 2, 3}	1/4 → 25%
{1, 3, 5}	2/4 → 50%
{1, 2, 5}	1/4 → 25%
{2, 3, 5}	2/4 → 50%
{2, 5}	3/4 → 75%
{3, 5}	2/4 → 50%

III

Itemset	Support
{1, 2, 3, 5}	1/4 = 25%
{2, 3, 5}	2/4 = 50%
{1, 2, 3}	1/4 = 25%

{ 2, 3, 5 } = This holds the minimum support

Now generate rules :-

Rules	Support	Confidence
(2^3) → 5	2	2/2 = 100%.
(3^5) → 2	2	2/2 = 100%.
(2^3) → 3	2	2/3 = 66%.
2 → (3^5)	2	2/3 = 66%.
5 → (2^3)	2	2/3 = 66%.
3 → (2^5)	2	2/3 = 66%.

final Association
rules are

$$\begin{aligned} & (2^3) \rightarrow 5 \\ & (3^5) \rightarrow 2 \end{aligned}$$

$$\text{Confidence} = S(A \cup B) / S(A)$$

$$= \frac{2^3 \rightarrow 5}{A} = \frac{S((2^3) \cup 5)}{S(2^3)} = 2/2 = 100\%.$$

→ Handling large data sets in main m/m.

- Allocate more Memory
- Work with a Smaller Sample
- Use a Computer with more memory
- Change the data format
- Stream data or use progressive loading
- Use a Relational database
- Use big data platform

→ Limited Pass Algorithm :- finding frequent itemsets

- Algorithms used to compute exact collection of frequent itemsets of size k in k passes

→ A Priori

→ PLY

→ multistage

→ multihash

- Algo. that find all frequent itemsets

→ Sampling

→ SON

→ Toivonen's Algo.

→ SON algorithm (Sarwar, Omiecinski & Navathe)

- known as Partition algorithm
- MapReduce

→ first map function

→ first Reduce function

→ second map

→ second Reduce

Limited pass Algo. (PCY Algorithm) frequent item set
 Park Chen Run Counting

Min support = 2

Tid	Items
1	1, 3, 4
2	2, 3, 5
3	1, 2, 3, 5
4	2, 5

P
A
S
S
1

T_1	$\{1, 3\}_2 \{1, 4\}_1 / \{3, 4\}_1$
T_2	$\{2, 3\}_2 \{2, 5\}_3 - \{3, 5\}_5 + 2$
T_3	$\{1, 2\} \{1, 3\} \{1, 5\} \{2, 3\} \{2, 5\} \{3, 5\}$
T_4	$(2, 5)$

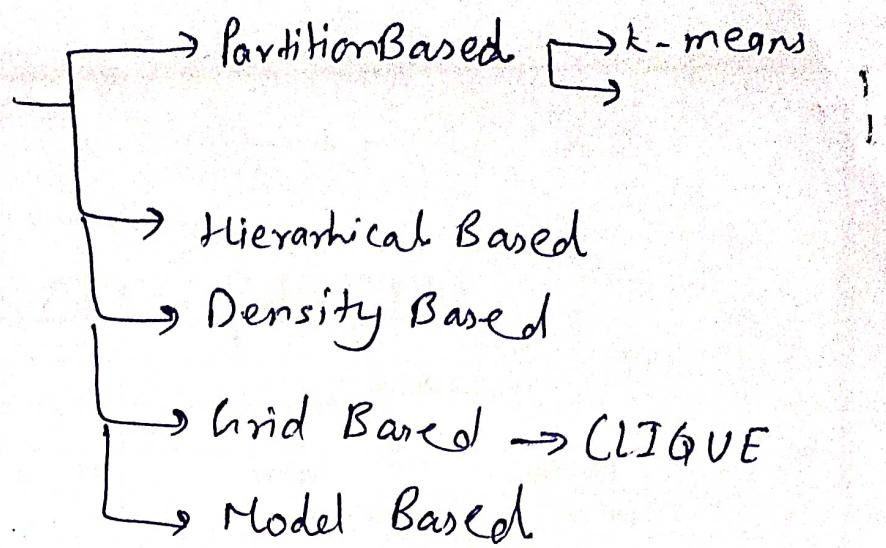
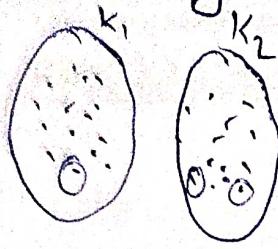
Pass 2 Apply hash fn,

Bucket -	1	2	3	4	5
Count -	4	6	3	0	0

Construct Transaction Id table

Counting frequent itemset in a stream
 → Count Sketch Algorithm.

Clustering Techniques.

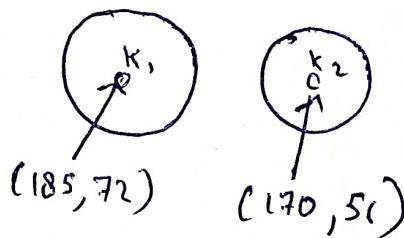


k-means clustering:

	Height	Weight
①	185	72
②	170	56
③	168	60
④	179	68
⑤	182	72
⑥	188	77
⑦	180	71
⑧	180	70
⑨	183	84
⑩	180	88
⑪	180	67
⑫	177	76

Euclidean distance

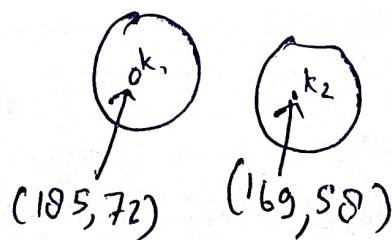
$$\text{Euclidean Distance} = \sqrt{(x_0 - x_c)^2 + (y_0 - y_c)^2}$$



$$\begin{aligned}
 \text{ED for } ③ \rightarrow k_1 &\rightarrow \sqrt{(168-185)^2 + (60-72)^2} \\
 &= 20.80 \\
 \rightarrow k_2 &\rightarrow \sqrt{(168-170)^2 + (60-56)^2} \\
 &= 4.48
 \end{aligned}$$

New Centroid Calculation

$$\text{for } k_2 = \left(\frac{170+168}{2}, \frac{60+56}{2} \right) = (169, 58)$$



$$\begin{aligned}
 \text{ED for } ④ \rightarrow k_1 &\rightarrow \sqrt{(75-185)^2} \\
 &+ (68-72)^2 \\
 \rightarrow k_2 &\rightarrow \sqrt{(75-169)^2} \\
 &+ (68-58)^2 \\
 &= 14.14
 \end{aligned}$$

$$k_1 \rightarrow \{1, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$$

$$k_2 \rightarrow \{2, 3\}$$

Hierarchical Clustering

→ Agglomerative

→ Divisive

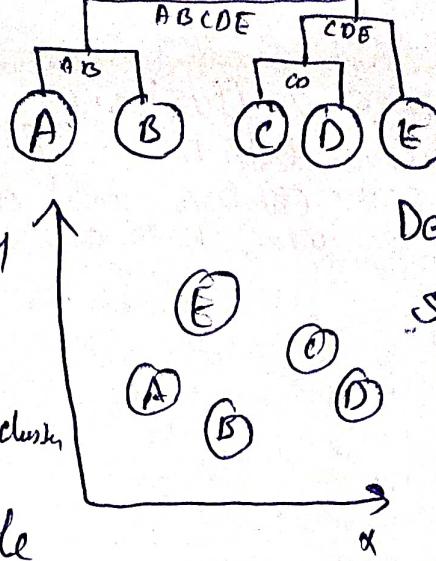
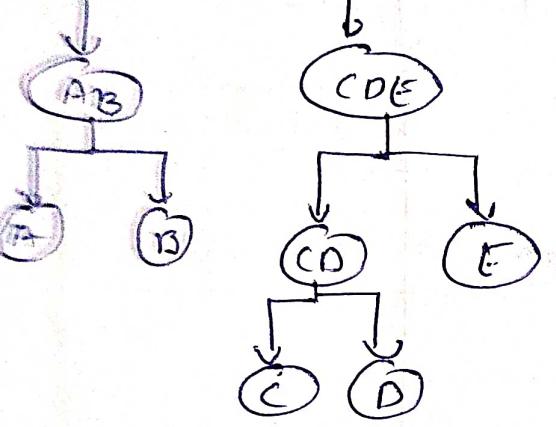
Start with individual data item
bottom to up

top to down

Put all the data points in one cluster
divide them

ABCDE

Starts with whole data points



Dendrogram :-

Structure which represent the hierarchical clustering

→ Clustering high-dimensional data.

- Subspace clustering :- Search for clusters already existing
 - Such as :- CLIQUE, Proclus
- Dimensionality Reduction approaches :- Construct much lower dimensional space
 - Spectral Clustering
- Clustering should not only consider dimension

CLIQUE Cluster

→ Grid based Cluster & Density Based

→ estimate the density by points

Clique & Percolation Method

use this method (PM)

find the communities for $k=3$

Step 1:- Divide the graph into Small Clusters (Cliques)

Step 2:- Write down all the cliques and group them into different categories

Step 3:- 6 cliques found here

$$a = (1, 2, 3)$$

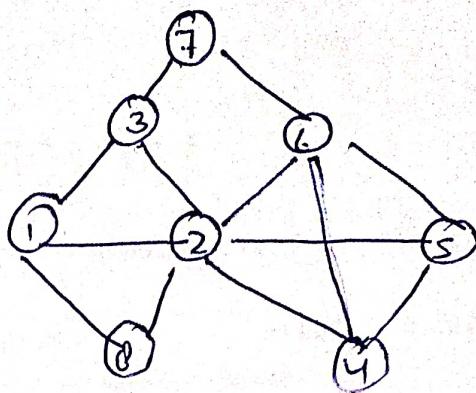
$$b = (1, 2, 8)$$

$$c = (2, 4, 6)$$

$$d = (2, 4, 5)$$

$$e = (4, 6, 5)$$

$$f = (2, 6, 5)$$



Clique

k = no. of nodes



$N=2$

Step 4:- Community 1 = $(1, 2, 3, 8)$

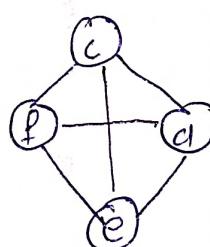
Community is group of clique.

$$C_2 = (2, 4, 5, 6)$$

To find community ($k-1$ node should be common)

for $k=3$

Step 5:- Now Create New Graph



$$C_1 = \{a, b\}$$

$$C_2 = \{c, d, e, f\}$$

Step 6:- Formulae is $k-1$
 $3-1=2$

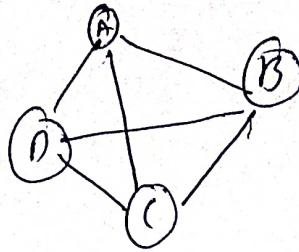
Step 7:- Find the Communities

$$C = \{C_1, C_2\}$$

Diagram for Communities

(i) Community = {a, b} A → B

(ii) Community ~ {a, b, c, d}



Ques: ~~Path Clusters~~

Frequent Pattern Based Clustering Methods.

To remove drawback of apriori algorithm.
• It is expensive operation for calculation

Ex:-

Tid	item	min-sup=2
T100	11, 12, 15	12, 11
T200	12, 14	12, 14
T300	12, 13	12, 13
T400	11, 12, 14	12, 11, 14
T500	11, 11, 13	11, 13

Step1

item	sup.
11	3
12	4
13	2
14	2
15	1

derending
→
Check min
Support
8 write in decreasing order

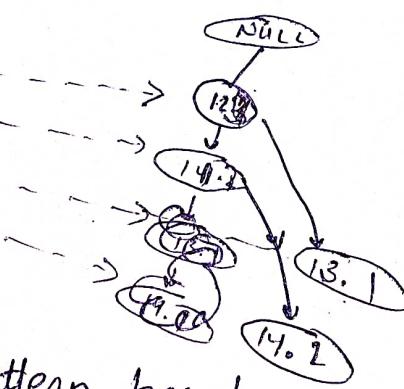
item	sup.
I ₁	15
I ₂	12
I ₃	13
I ₄	11
I ₅	14

N ₁	12
I ₁	4
I ₂	3
I ₃	2
I ₄	2

↑
frequent pattern

Step3 Construct FP Tree : ordered item set

ITEM ID	SUP.	N.L.
I ₂	4	-
I ₄	3	-
I ₃	2	-
I ₁	2	-



Step4 Conditional Pattern based

Steps:- finding frequent patterns

14	{12, 11, 14} {12, 14}
13	{11, 13} {12, 13}
11	{11, 13} {12, 11, 14} {12, 14}
12	{ }

Frequent Pattern Growth Algorithm

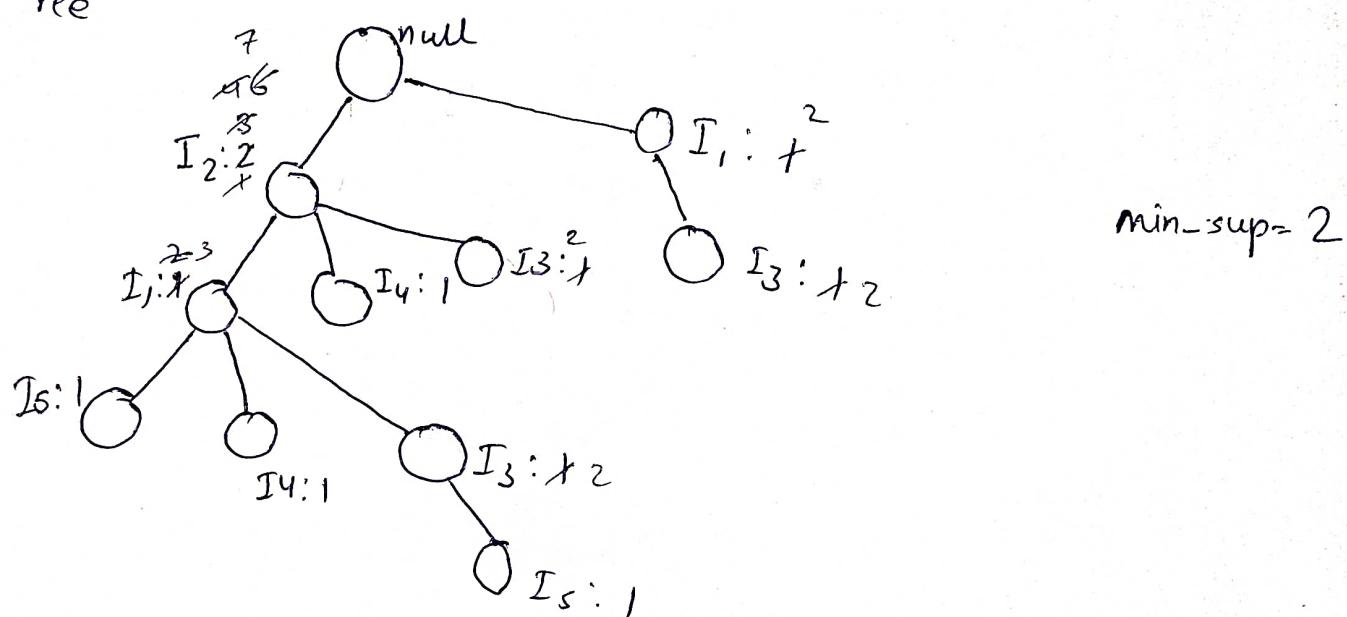
Trans ID	Item	L ₁
T100	I ₁ I ₂ I ₅	I ₂ I ₁ I ₅
T200	I ₂ I ₄	I ₂ I ₄
T300	I ₂ I ₃	I ₂ I ₃
T400	I ₁ I ₂ I ₄	I ₂ I ₃ I₄
T500	I ₁ I ₃	I ₂ I ₁ I ₄
T600	I ₂ I ₃	
T700	I ₁ I ₃	
T800	I ₁ I ₂ I ₃ I ₅	I ₂ I ₁ I ₃ I ₅
T900	I ₁ I ₂ I ₃	I ₂ I ₁ I ₃

Item set	Sup. Count
I ₁	6
I ₂	7
I ₃	6
I ₄	2
I ₅	2

Items	Sup. Count
I ₂	7
I ₁	6
I ₃	6
I ₄	2
I ₅	2

⇒

FP Tree



Item	Conditional Pattern Base	Conditional FP-tree	Frequent Patterns Generated
I ₅	{I ₂ I ₁ : 1}	{I ₂ , I ₁ , I ₃ : 1}	{I ₂ ; I ₅ : 2} {I ₁ , I ₅ : 2} {I ₂ , I ₁ , I ₅ : 2}
I ₄	{I ₂ , I ₁ : 1} {I ₂ : 1}		{I ₂ , I ₄ : 2}
I ₃	{I ₂ , I ₁ : 2} {I ₂ : 2}, {I ₁ : 2}	(I ₂ : 2, I ₁ : 2) (I ₂ : 2)	{I ₂ , I ₃ : 4} {I ₁ , I ₃ : 4} {I ₂ , I ₃ : 2}
I ₁	{I ₂ : 4}	(I ₂ : 4) (I ₂ : 4, I ₁ : 2), (I ₁ : 2)	{I ₂ , I ₁ : 4}

Fuzzy decision Tree

ID3 Algorithm

Age	Competition	Type	Profit
old	Yes	S/W	Down
old	No	S/W	Down
old	No	S/W	Down
old	Yes	H/W	Down
mid	Yes	S/W	Down
mid	Yes	H/W	Down
mid	No	H/W	Up
mid	No	S/W	Up
new	Yes	S/W	Up
new	No	H/W	Up
new	No	S/W	Up

I_G

$$\text{Info. Gain} = \frac{P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

$$\text{Entropy}(A) = \sum_{\text{Attribute}} \frac{P_i + N_i}{P+N} I(P_i N_i)$$

$$\text{Gain} = I_G - E(A)$$

$$\log_2 x = \frac{\log_{10} x}{\log_{10} 2}$$

Target attribute

Age :

Age	Down	Up
old	3	0
mid	2	2
new	0	3

$$I(\text{old}) = - \left[\frac{3}{3} \log_2 \left(\frac{3}{3} \right) + \frac{0}{3} \log_2 \left(\frac{0}{3} \right) \right] \\ = 0 \times \frac{3}{10} = 0$$

$$I(\text{mid}) = - \left[\frac{2}{4} \log_2 \left(\frac{2}{4} \right) + \frac{2}{4} \log_2 \left(\frac{2}{4} \right) \right] = 0.$$

$$I(\text{new}) = 0$$

$$E(A_g) = 0 + 0.4 + 0 = 0.4$$

$$I_G = - \left[\frac{5}{10} \log_2 \left(\frac{5}{10} \right) + \frac{5}{10} \log_2 \left(\frac{5}{10} \right) \right]$$

$$= - [0.5 \times \log_2 2^{-1} + 0.5 \log_2 2^{-1}]$$

$$= - [0.5 \times (-1 \log_2 2) + 0.5 \times (-1 \log_2 2)]$$

$$= [-0.5 - 0.5] = [-1]$$

$$I.G. = 1$$

$$\text{Gain} = 1 - 0.4 \\ = 0.6$$

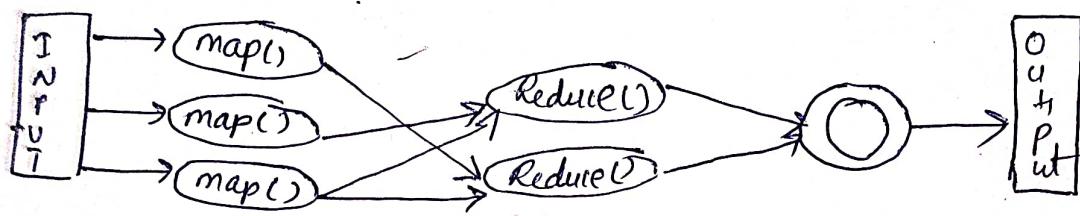
~~Clustering in non-euclidean Space~~

Unit-5

MapReduce

- Word Count program :- In Hadoop architecture
- Key value pair → Parallel process

Map Reduce Component



→ key / Value → Combines
key values
and make
Set of tuples

(one)

- ① master Job tracker
 - (1) managing Resources
 - (2) Resource management
 - (3) Scheduling task
 - (4) Monitoring task

(many)

- ② Slave Job tracker
 - ① Executes the task
 - ② Provide task status