

Big Data: Large amounts of data

→ Data that is so large and complex that it is difficult to process using traditional methods

→ Includes structured, unstructured, semi structured data

→ act of storing and analysing

↓  
cloud

↓

for insights

Types of data

Structured: Table, pdfs, word [CSV, relational data]

↳ stored, accessed with fixed format

Unstructured: Google search results

↳ might have missing / wrong values, unknown form

Semistructured: XML, HTML, <tag>

↳ combination

why is big data important?

- i) cost reduction due to storing data in hadoop etc.
- ii) Time reduction, speed in memory analysis
- iii) New product development & optimization
- iv) Smart decision making
- v) Multiple key technologies like ML, AI, mining etc.

Characteristics

speed & data generation

↳ heterogenous data

3V: Volume, Velocity, Variety

extra 2V: Veracity, Value

Validating Big Data: Find out feasibility, reasonability, value, integrity and sustainability to validate data and accept a technique and tool.

class label is in  
supervised learning and it is the supervisor

Uses: Healthcare, Govt, Education, Telecommunication, Manufacture, Financial Sector.

Different sources of data

→ sensor, social media, E-commerce, transaction, log, web-related data.

Veracity: verify truthfulness of data.

Value: enhances business profit by providing decision making systems.

### Traditional Data

- manual data generation
- standard format
- able to store and process with the help of traditional database management tools

### Big Data

- machine generates data
- messy and ugly format
- unable to capture & store & process using traditional DB management tools within elapsed time

### Then vs Now

Healthcare → <sup>Paper</sup> ~~EHR~~ (Electronic Health Record)

↙ ↘  
MRI Lab ...

Telecommunication : Bandwidth, call rate, N/W capacity, number of users, customer feedback, demographic location

Government : aadhar, job, income, tax

Education : students count, interest, age,

Financial sector : Bank account, credit score, loan required, income, profession, age, address, phone  
→ targeted marketing

To process large data, one can replace a single high end computer with many many computers with resources.

erce, transaction, log.

data.

providing decision

Big Data

- machine generates data
- messy and ugly format
- unable to capture all
- store & process using traditional DB management tools within elapsed time

(Health Record)

h, call rate, N/W capacity

Feedback, demographic

ne, tax

rest, age,

credit score, loan

age, address, phone

can replace a  
many many computers

Distributed computation

writing distributed pgms require high level of expertise}

Issues and challenges

Points to be considered when designing a system

- Scalable
- Cost effective
- Flexibility
- Ease of use

HDFS : Hadoop Distributed File System

YARN : yet another Resource Navigator

Map Reduce : Programming Model to handle big data

Hbase	Flink	Giraph	Zookeeper
Cassandra	spark	Hive	
MongoDB	storm	Pig	

Hadoop : A scalable, fault tolerant, cost effective open source software that handles large data in a distributed system

Hadoop

HDFS      Map Reduce (MR)

Map and Reduce

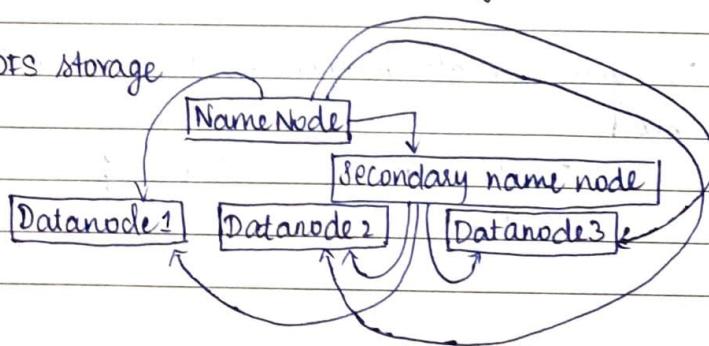
→ 2 phases : mapper phase and reducer phase

→ in between these phases there is shuffle and sort phase

→ programmer needs to identify those phases.

→ output of phases is key-value pair

HDFS storage



Namenode:

one of the three main components in HDFS storage

It is the master.

① Master Daemon

② Metadata information

→ Size of file

→ Permission of file

→ Location of a file

③ Receives heartbeat from the datanode

Replication factor (Default)  
in 3 different places.

Rack Awareness Algorithm

→ Blocks should

be replicated in different racks.

RACK 1

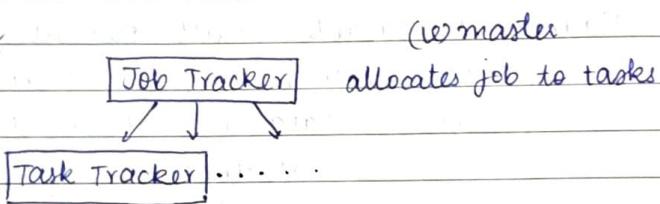
RACK 2

Datanode

① Serving read requests and write requests from the client

② Actual data is stored

Datanode 1	Datanode 2
Block 1	Block 2
Datanode 3	Datanode 4
Block 3 (1)	Block 3 (2)
Datanode 5	Datanode 6
Block 3 (3)	Block 3 (4)

Map Reduce

hdfs-site.xml

Hadoop V1

Hadoop V2

→ Original and

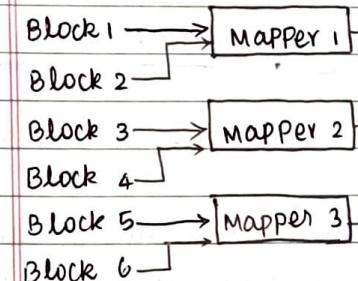
racks.

Large I/P file  
is divided into  
blocks

128 mb

Map Reduce Programm  
→ Divide

QN 1: How many tasks will be created if file of size 514 mb  
is copied to hdfs?



If you allocate hadoop V2 then we will  
need 5 blocks

Block 1: 128 ..... Block 4: 128,

Block 5: 2 mb

Divide and  
conquer

Replication factor (Default) is 3. Blocks can be replicated in 3 different places.

### Rack Awareness Algorithm

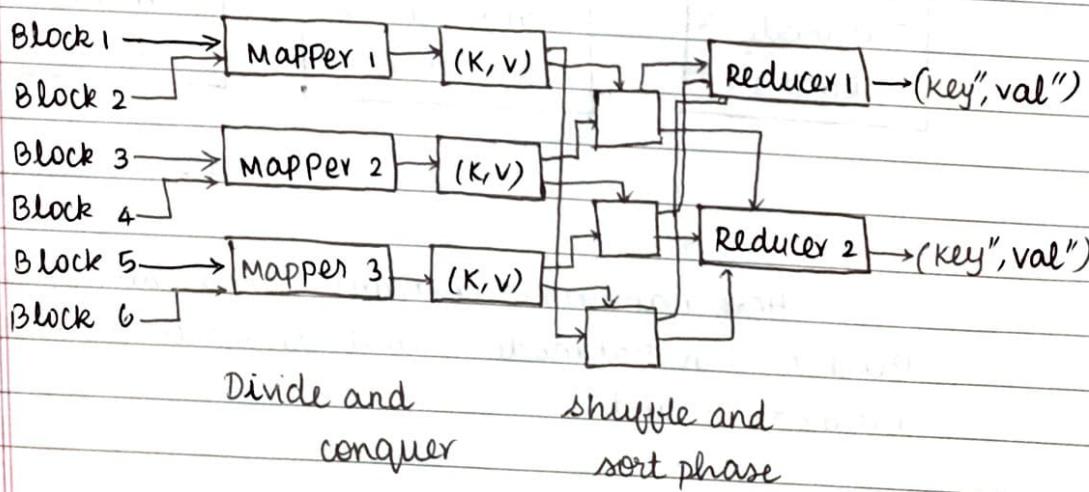
→ Blocks should come under same rack if they are replication data. For example,

RACK 1	RACK 2	RACK 3	RACK 4
Datanode 1 Block 1	Datanode 2 Block 2	Datanode 7 Block 1 (rep)	Datanode 8 Block 3
Datanode 3	Datanode 4 Block 3 (rep)	Datanode 9 Block 1 (rep)	Datanode 10 Block 2 (rep)
Datanode 5	Datanode 6 Block 3 (rep)	Datanode 11	Datanode 12 Block 2 (rep)

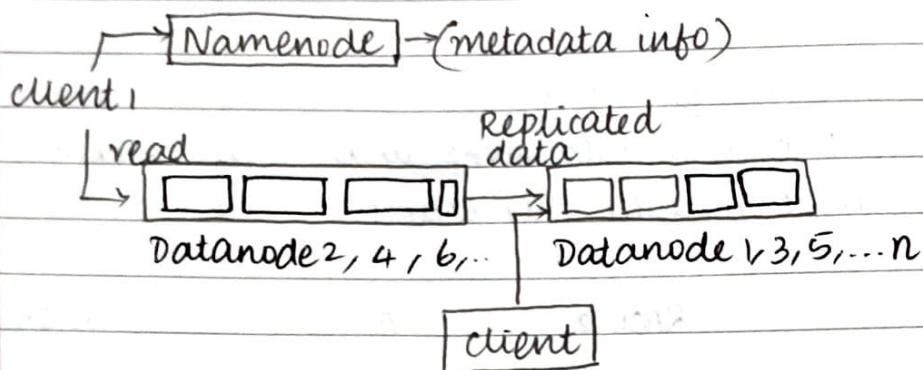
→ original and replication data must be in diff racks.

### Map Reduce Programming Paradigm

→ Divide & Conquer Approach



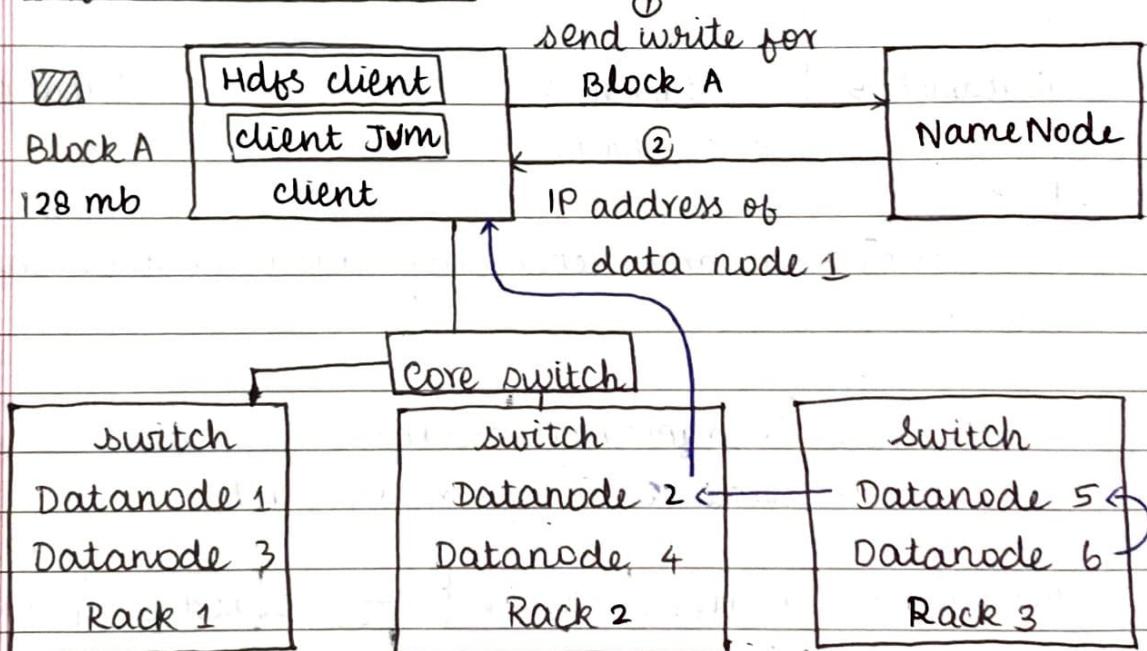
## HDFS architecture



## HDFS read/write mechanism

- HDFS write mechanism
- HDFS write mechanism with ACK
- HDFS multi write pipeline : has many blocks
- HDFS read mechanism : sends read request @ ①

## HDFS write mechanism with ACK



→ Acknowledgement

Here name node contains size, location, permission  
Block A is in Datanode 2 and replications are in  
datanode 5 and 6.

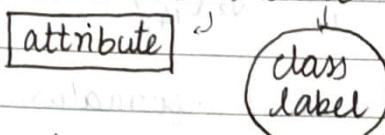
## Classification (Type of supervised learning)

uses classification algorithm/model

- ↳ decision tree algorithm
- ↳ naive based Bayes algorithm
- ↳ support vector machine
- ↳ neural network

## Decision Trees (Greedy approach)

consists of Internal nodes, leaves, branches



Root node: Attribute that provides maximum information

Example (IV or not)



## Attribute Selection

Gini Index    Info Gain : Entropy formula



Gain(D) - Gain<sub>A</sub>(D)

$$\sum_{i=1}^n -P_i \log_2 P_i$$

Example

Day	outlook	Temperature	Humidity	wind	Play Tennis
1	sunny	Hot	High	weak	No
2	sunny	Hot	High	Strong	No
3	Overcast	Hot	High	weak	Yes
4	Rainy	Mild	High	weak	Yes
5	Rainy	Cool	Normal	weak	Yes
6	Rainy	Cool	Normal	strong	No
7	overcast	Cool	Normal	strong	Yes
8	sunny	High-Mild	High	weak	No
9	sunny	Cool	Normal	weak	Yes
10	Rainy	Mild	Normal	weak	Yes

11	Sunny	High	Mid	Normal	Strong	Yes
12	overcast	Mild		High	Strong	Yes
13	Overcast	Hot		Normal	weak	Yes
14	Rainy	High	Mid	High	Strong	No

SolutionStep 1: Find Entropy

$$\text{Entropy} = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.94$$

Step 2: Find 1st highest information gain

## ① Outlook

sunny overcast Rainy  
 Y N N Y N N  
 2 3 4 0 3 2

$$\begin{aligned} \text{Entropy} (\text{outlook} = \text{sunny}) &= -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) \\ &= 0.971 \end{aligned}$$

$$\text{Entropy} (\text{outlook} = \text{overcast}) = -\frac{4}{4} \log_2\left(\frac{4}{4}\right) = 0$$

$$\begin{aligned} \text{Entropy} (\text{outlook} = \text{Rainy}) &= -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) \\ &= 0.971 \end{aligned}$$

$$\begin{aligned} \text{Gain}_{\text{outlook}} &= \frac{5}{14} \times 0.971 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.971 \\ &= \frac{5}{14} \times 0.971 \times 2 = \frac{5}{7} \times 0.971 = 0.693 \end{aligned}$$

$$\text{Info Gain} (\text{outlook}) = 0.94 - 0.693 = 0.247$$

## ② Temperature

Hot High Mild Cold  
 2 2 4 2 3 1

$$\text{Entropy}(\text{Temp} = \text{Hot}) = 0.1$$

$$\text{Entropy}(\text{Temp} = \text{Mild}) = -\frac{4}{6} \log_2\left(\frac{4}{6}\right) - \frac{2}{6} \log_2\left(\frac{2}{6}\right) = 0.910$$

$$\text{Entropy}(\text{Temp} = \text{Cold}) = -\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) = 0.811$$

$$\begin{aligned} \text{Gain}_{\text{Temp}} &= \frac{4}{14} \times 0.1 + \frac{6}{14} \times 0.910 + \frac{4}{14} \times 0.811 \\ &= 0.911 \end{aligned}$$

$$\text{Info Gain}(\text{temp}) = 0.94 - 0.911 = 0.0292$$

### (3) Humidity

/ \

	Normal	High
\	1	3
6	1	4

$$\begin{aligned} \text{Entropy}(\text{Humidity} = \text{Normal}) &= -\frac{6}{7} \log_2\left(\frac{6}{7}\right) - \frac{1}{7} \log_2\left(\frac{1}{7}\right) \\ &= 0.591 \end{aligned}$$

$$\begin{aligned} \text{Entropy}(\text{Humidity} = \text{High}) &= -\frac{3}{7} \log_2\left(\frac{3}{7}\right) - \frac{4}{7} \log_2\left(\frac{4}{7}\right) \\ &= 0.983 \end{aligned}$$

$$\begin{aligned} \text{Gain}_{\text{Humidity}} &= \frac{7}{14} \times 0.591 + \frac{7}{14} \times 0.983 \\ &= 0.787 \end{aligned}$$

$$\text{Info Gain}(\text{Humidity}) = 0.94 - 0.787 = 0.153$$

### (4) wind

/ \

	Strong	weak
\	3	3
6	2	2

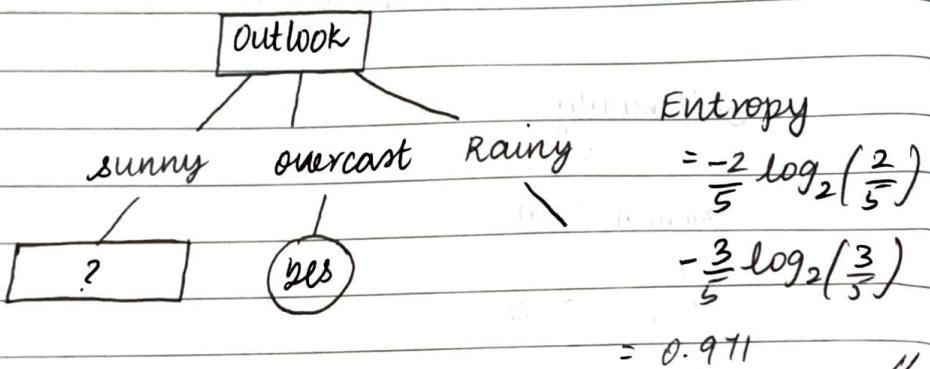
$$\begin{aligned} \text{Entropy}(\text{wind} = \text{weak}) &= -\frac{6}{8} \log_2\left(\frac{6}{8}\right) - \frac{2}{8} \log_2\left(\frac{2}{8}\right) \\ &= 0.811 \end{aligned}$$

$$\text{Entropy}(\text{wind} = \text{Strong}) = 1$$

$$\text{Gain}_{\text{wind}} = \frac{6}{14} \times 1 + \frac{8}{14} \times 0.811 \\ = 0.892$$

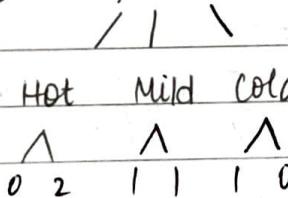
$$\text{Info gain}_{(\text{wind})} = 0.94 - 0.892 \\ = 0.048$$

Step 3: Find root



Repeat Step 2:

① Temperature



$$\text{Entropy} (\text{Temp} = \text{Hot}) = 0$$

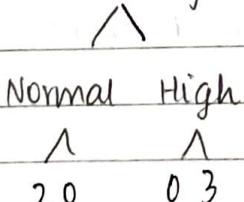
$$\text{Entropy} (\text{Temp} = \text{Mild}) = 1$$

$$\text{Entropy} (\text{Temp} = \text{Cold}) = 0$$

$$\text{Gain} (\text{Temp}) = 1 \times \frac{2}{5} = \frac{2}{5} = 0.4$$

$$\text{Info Gain} (\text{Temp}) = 0.971 - 0.4 = 0.571$$

② Humidity

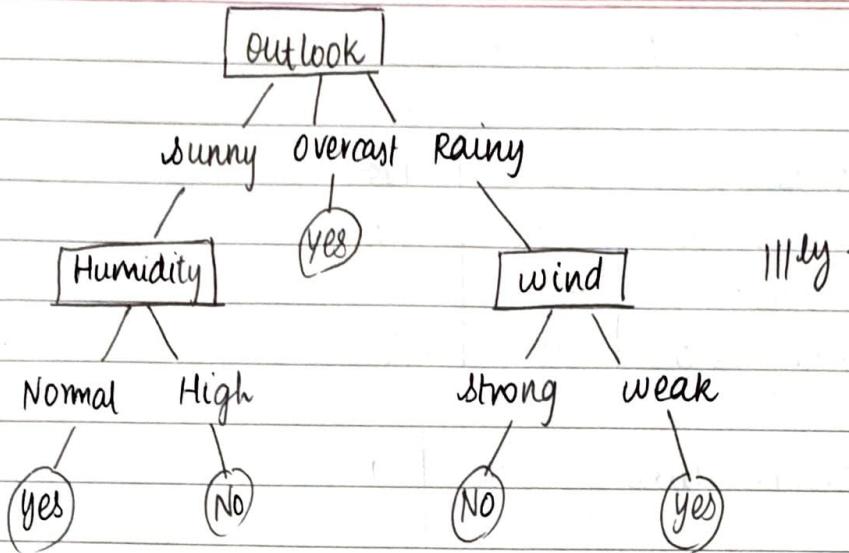


$$\text{Entropy} (\text{Humid} = \text{Normal}) = 0$$

$$\text{Entropy} (\text{Humid} = \text{High}) = 0$$

$$\text{Gain} (\text{Humidity}) = 0$$

$$\text{Info Gain} (\text{Humidity}) = 0.971$$



## ② Bayesian Classifier (Supervised learning)

- i) Based on class independence : the effect of values of one attribute is independent of values of another attribute
- (ii) Age, Income, student/not are independent.
- ii) Each feature equally contributes / has equal importance in predicting the class label

Baye's Theorem :

$$P(\text{Hypothesis} | (x_1, x_2, \dots, x_n)) = \frac{P(x_i | H) \cdot P(H)}{P(x_i)}$$

↓                      ↓                      ↓  
 given              sample tuple            evidence  
 ↓                      ↓                      ↓  
 Posterior probability      likelihood      prior probability

$$P(C/x_i) = P(C/x_1) * P(C/x_2) * P(C/x_3) * \dots * P(\text{class})$$

$$P(C/x_i) = \sum_{i=1}^n P(C/x_i)$$

$$P(C/x_i) \propto \sum_{i=1}^n P(C/x_i)$$

$$P(C/x_i) = \text{argmax } P(C/x_i)$$

Let us assume Play Tennis dataset

Test data : outlook      Temperature      Humidity      wind  
 sunny                  Hot                    Normal            weak

$$P(\text{Yes}) = \frac{9}{14} \quad P(\text{No}) = \frac{5}{14}$$

Outlook	Yes	No
Sunny = 5	2/9	3/5
Overcast = 4	4/9	0/5
Rainy = 5	3/9	2/5

Temperature	Yes	No
Hot = 4	2/9	2/5
Mild = 6	4/9	2/5
Cold = 4	3/9	1/5

Humidity	Yes	No
Normal = 7	6/9	1/5
High = 7	3/9	4/5

Wind	Yes	No
Strong = 6	3/9	3/5
Weak = 8	6/9	2/5

$P(\text{Yes} / \text{outlook} = \text{sunny}, \text{Temp} = \text{hot}, \text{humidity} = \text{normal}, \text{wind} = \text{weak})$

= likelihood  $\times$  Prior Probability

Evidence

$$= \frac{P(\text{outlook} = \text{sunny})}{P(\text{Yes})} \times \frac{P(\text{Temp} = \text{hot})}{P(\text{Yes})} \times \frac{P(\text{Humid} = \text{N})}{P(\text{Yes})} \\ \times \frac{P(\text{wind} = \text{weak})}{P(\text{Yes})} \times \frac{9}{14}$$

$$= \frac{2}{9} \times \frac{2}{9} \times \frac{6}{9} \times \frac{6}{9} \times \frac{9}{14}$$

$$= 8/567 \approx 0.014$$

$P(\text{No} / \text{outlook} = \text{bunny}, \text{Temp} = \text{hot}, \text{humidity} = \text{normal}, \text{wind} = \text{weak})$

$$= \frac{P(O = \text{Sunny})}{P(\text{No})} \times \frac{P(T = \text{hot})}{P(\text{no})} \times \frac{P(\text{Humid} = \text{Normal})}{P(\text{no})} \times \frac{P(\text{wind} = w)}{P(\text{no})} \\ \times P(\text{no}) \\ = \frac{3}{5} \times \frac{2}{5} \times \frac{1}{5} \times \frac{2}{5} \times \frac{5}{14} = \frac{6}{875} \\ \approx 0.006$$

$\therefore$  For the given data, Bott PlayTennis = YES

For BuysComputer Dataset

1/4/22

K-Means clustering (Unsupervised learning)

Within a cluster, similarities are high

Application of clustering

- i) outlier detection Eg: fraud detection/credit card
- ii) grouping for google search

Clustering can also be called cluster analysis/ data segmentation

Challenges

- i) Scalability
- ii) Need to consider dealing with high dimensional data
- iii) Ability to deal with data of arbitrary shape (types)
- iv) Sensitive to noise
- v) Constraint based
- vi) Different set of attributes
- vii) choosing for input parameters
- viii) Incremental updates

- ix) Discovery of clusters with arbitrary shape

- iv) Requirement for domain knowledge to determine input parameters
- v) Ability to deal with noisy data
- vi) Incremental clustering as insensitive to input
- vii) Capability of clustering high dimensional data
- viii) Interpretability and usability
- ix) constraint based clustering

Clustering algo → Partitioning clustering Eg: K-mean  
Types

- Hierarchical clustering
- Density based clustering
- Grid based clustering

These algos should cover the following

- i) Operation of clusters    ii) Similarity measure
- iii) Clustering method                  [Euclidean, Manhattan, Jaccard distance etc.]
- iv) partitioning criteria

Intracluster distance: within a cluster (must be low)

Intercluster distance: datapoints b/w two clusters

similarity high

similarity is low

Euclidean distance

$$= \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

K-means clustering

- i) Assume a centroid value for some cluster
- ii) Find distance b/w centroid and datapoint.
- iii) Put datapoint under cluster where distance is small.

Qn: Cluster the following 8 points [with (x, y) representing location] into three clusters.

A <sub>1</sub> (2, 10)	B <sub>1</sub> (5, 8)	C <sub>1</sub> (1, 2)
A <sub>2</sub> (2, 5)	B <sub>2</sub> (7, 5)	C <sub>2</sub> (4, 9)
A <sub>3</sub> (8, 4)	B <sub>3</sub> (6, 4)	

Distance function is Euclidean distance. Suppose initially we assign A<sub>1</sub>, B<sub>1</sub>, and C<sub>1</sub> as the center of each cluster respectively, use K-means algorithm to show

a) 3 cluster centers after first round of execution

b) Final 3 clusters

FIRST ROUND

K=3; Centroids are A<sub>1</sub>, B<sub>1</sub>, C<sub>1</sub>

	cluster	1	2	3
A <sub>1</sub>	A <sub>3</sub>	A <sub>2</sub>		
B <sub>1</sub>		B <sub>1</sub>	C <sub>1</sub>	
B <sub>2</sub>			B <sub>2</sub>	
B <sub>3</sub>			B <sub>3</sub>	
C <sub>2</sub>				C <sub>2</sub>

A<sub>2</sub>(2, 5)

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

$$\text{Distance} = \sqrt{(2-2)^2 + (10-5)^2} = \sqrt{25} = 5$$

ii) B<sub>1</sub>(5, 8)

$$\text{Distance} = \sqrt{(5-2)^2 + (8-5)^2} = \sqrt{9+9} = 4.2$$

iii) C<sub>1</sub>(1, 2)

$$\text{Distance} = \sqrt{(2-1)^2 + (5-2)^2} = \sqrt{1+9} = \sqrt{10} = 3.1$$

⇒ A<sub>2</sub> belongs to cluster 3

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

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

$$\text{Distance} = \sqrt{(8-2)^2 + (4-10)^2} = \sqrt{36+36} = 10.09$$

ii) B<sub>1</sub>(5, 8)

$$\text{Distance} = \sqrt{(8-5)^2 + (4-8)^2} = \sqrt{9+16} = 5$$

iii) C<sub>1</sub>(1, 2)

$$\text{Distance} = \sqrt{(8-1)^2 + (4-2)^2} = \sqrt{49+4} = 7.2$$

⇒ Belongs to cluster 2

B<sub>2</sub>(7, 5)

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

$$\text{Distance} = \sqrt{(7-2)^2 + (5-10)^2} = \sqrt{25+25} = 7.07$$

i)  $B_1(5, 8)$ 

$$D = \sqrt{(7-5)^2 + (5-8)^2} = \sqrt{4+9} = 3.6$$

iii)  $C_1(1, 2)$ 

$$D = \sqrt{(7-1)^2 + (5-2)^2} = \sqrt{36+9} = 6.7$$

 $\Rightarrow$  Cluster 2 $B_3(6, 4)$ i)  $A_1(2, 10)$ 

$$D = \sqrt{(6-2)^2 + (4-10)^2} = \sqrt{16+36} = 7.2$$

ii)  $B_1(5, 8)$ 

$$D = \sqrt{(6-5)^2 + (4-8)^2} = \sqrt{1+16} = 4.12$$

iii)  $C_1(1, 2)$ 

$$D = \sqrt{(6-1)^2 + (4-2)^2} = \sqrt{25+4} = 5.3$$

 $\Rightarrow$  Cluster 2 $C_2(4, 9)$ i)  $A_1(2, 10)$ 

$$D = \sqrt{(4-2)^2 + (9-10)^2} = \sqrt{4+1} = \sqrt{5} = 2.2$$

ii)  $B_1(5, 8)$ 

$$D = \sqrt{(5-4)^2 + (8-9)^2} = \sqrt{2} = 1.4$$

iii)  $C_1(1, 2)$ 

$$D = \sqrt{(4-1)^2 + (9-2)^2} = \sqrt{9+49} = 7.6$$

 $\Rightarrow$  Cluster 2Centroid of cluster 1:  $A_1(2, 10)$ 2:  $(5-4, 6) (6, 6)$ 3:  $(1.5, 3.5)$ 

$$\left( \frac{5+8+7+6+4}{5}, \frac{8+4+5+4+9}{5} \right)$$

$$\left( \frac{1+2}{2}, \frac{2+5}{2} \right)$$

$$(5.4, 6) (6, 6)$$

$$\left( \frac{3}{2}, \frac{7}{2} \right)$$

=

SECOND ROUNDcentroid  $(2, 10) (6, 6) (1.5, 3.5)$ 

clusters 1 2 3

A1 A3 A2

C2 B1 C1

B2

B3

 $A_2(2, 5)$ i)  $A_1(2, 10)$ 

$$D = \sqrt{0+5^2} = 5$$

ii)  $B_1(6, 6)$   $D = \sqrt{(6-2)^2 + (6-5)^2} = \sqrt{16+1} = 4.12$ iii)  $(1.5, 3.5)$ 

$$D = \sqrt{(2-1.5)^2 + (5-3.5)^2} = 1.58 //$$

 $A_3(8, 4)$ 

$$i) D = \sqrt{(8-2)^2 + (4-10)^2} = \sqrt{36+36} = 10.09$$

$$ii) D = \sqrt{(8-6)^2 + (4-6)^2} = \sqrt{4+4} = 2.8 //$$

$$iii) D = \sqrt{(8-1.5)^2 + (4-3.5)^2} = 6.51$$

 $B_1(5, 8)$ 

$$i) D = \sqrt{(5-2)^2 + (10-8)^2} = \sqrt{9+4} = 3.6$$

$$ii) D = \sqrt{(5-6)^2 + (8-6)^2} = \sqrt{1+4} = 2.23 //$$

$$iii) D = \sqrt{(5-1.5)^2 + (8-3.5)^2} = 5.7$$

 $B_2(7, 5)$ 

$$i) D = \sqrt{(7-2)^2 + (5-10)^2} = \sqrt{25+25} = 7.07$$

$$ii) D = \sqrt{(7-6)^2 + (5-6)^2} = \sqrt{2} = 1.4 //$$

$$iii) D = \sqrt{(7-1.5)^2 + (5-3.5)^2} = 5.70$$

 $C_1(1, 2)$ 

$$i) D = \sqrt{(2-1)^2 + (10-2)^2} = \sqrt{1+64} = 8.06$$

$$ii) D = \sqrt{(6-1)^2 + (6-2)^2} = \sqrt{25+16} = 6.4 //$$

$$iii) D = \sqrt{(1.5-1)^2 + (3.5-2)^2} = 1.5 //$$

 $C_2(4, 9)$

$$\text{i) } D = \sqrt{(4-2)^2 + (9-10)^2} = 2.2 //$$

$$\text{ii) } D = \sqrt{(6-4)^2 + (6-9)^2} = \sqrt{4+9} = 3.6$$

$$\text{iii) } D = \sqrt{(4-1.5)^2 + (9-3.5)^2} = 6.04$$

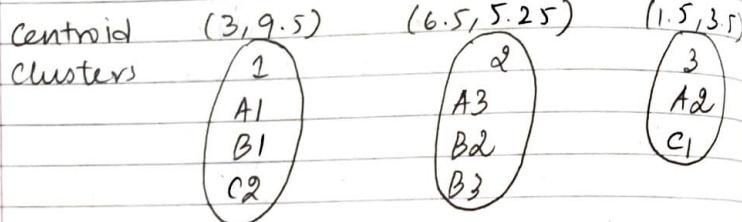
clusters

$$\text{i) } \left( \frac{2+4}{2}, \frac{10+9}{2} \right) = (3, 9.5)$$

$$\text{ii) } \left( \frac{5+8+7+6}{4}, \frac{8+4+5+4}{4} \right) = (6.5, 5.25)$$

$$\text{iii) } (1.5, 3.5)$$

### THIRD ROUND



A1(2, 10)

$$\text{i) } \sqrt{(3-2)^2 + (9.5-10)^2} = \sqrt{1.25} = 1.1 //$$

$$\text{ii) } \sqrt{(6.5-2)^2 + (5.25-10)^2} = \sqrt{42.8} =$$

$$\text{iii) } \sqrt{(2-1.5)^2 + (10-3.5)^2} = 6.51$$

A2(2, 5)

$$\text{i) } \sqrt{(3-2)^2 + (9.5-5)^2} = 4.60$$

$$\text{ii) } \sqrt{(6.5-2)^2 + (5.25-5)^2} = 4.50$$

$$\text{iii) } \sqrt{(2-1.5)^2 + (5-3.5)^2} = 1.58 //$$

A3(8, 4)

$$\text{i) } \sqrt{(8-3)^2 + (4-9.5)^2} = 7.43$$

$$\text{ii) } \sqrt{(8-6.5)^2 + (4-5.25)^2} = 1.95 //$$

$$\text{iii) } \sqrt{(8-1.5)^2 + (4-3.5)^2} = 6.51$$

B1(5, 8)

$$\text{i) } \sqrt{(5-3)^2 + (8-9.5)^2} = 2.06 //$$

$$\text{ii) } \sqrt{(6.5-5)^2 + (5.25-8)^2} = 3.13$$

$$\text{iii) } \sqrt{(5-1.5)^2 + (8-3.5)^2} = 5.70$$

B2(7, 5)

$$\text{i) } \sqrt{(7-3)^2 + (5-9.5)^2} = 6.02$$

$$\text{ii) } \sqrt{(7-6.5)^2 + (5-5.25)^2} = 0.55 //$$

$$\text{iii) } \sqrt{(7-1.5)^2 + (5-3.5)^2} = 5.70$$

B3(6, 4)

$$\text{i) } \sqrt{(6-3)^2 + (4-9.5)^2} = 6.26$$

$$\text{ii) } \sqrt{(6.5-6)^2 + (5.25-4)^2} = 1.34 //$$

$$\text{iii) } \sqrt{(6-1.5)^2 + (4-3.5)^2} = 5.53$$

C1(1, 2)

$$\text{i) } \sqrt{(3-1)^2 + (9.5-2)^2} = 7.76$$

$$\text{ii) } \sqrt{(6.5-1)^2 + (5.25-2)^2} = 6.38$$

$$\text{iii) } \sqrt{(1.5-1)^2 + (3.5-2)^2} = 1.51 //$$

C2(4, 9)

$$\text{i) } \sqrt{(4-3)^2 + (9-9.5)^2} = 1.11 //$$

$$\text{ii) } \sqrt{(6.5-4)^2 + (5.25-9)^2} = 4.5$$

$$\text{iii) } \sqrt{(4-1.5)^2 + (9-3.5)^2} = 6.04$$

## R-Programming

- ↳ statistical Programming language
- ↳ graphical visualization
- ↳ used to develop Machine learning algo
- ↳ open source software
- ↳ platform independent
- ↳ packages

## Assign variables

a <- 5

subject = "bda"

## Datatypes

- \* Number → int eg: x <- 10L  
→ numeric eg: m = 7.0
  - \* complex 7i+9
  - \* character " "
- single quotes automatically becomes double
- \* Boolean/Logical True or false

## Data structures

- i) vector
- ii) List
- iii) Matrix
- iv) Array
- v) Dataframe
- vi) Factors

- i) c() → c function stands for concatenate or combine
- ii) supports different datatype

- iii) matrix (c(elements, nrow=2, ncol=3))
- iv) array (c(elements, dim = (row, column, dimensions)))
- v) data.frame () → create tabular data
- vi) Shows levels which are distinct values in the factor

## Name 1 Association Rule mining algorithm /

- helps for product suggestion or recommendation
- finding frequent item set which leads to the discovery of association rules and correlations
- ↳ support value
- confidence value

Name 2

Eg: Market-Basket Analysis

no. of times item is present in db % of interest in buying things together

- ↳ low: large number of items can be retrieved
- ↳ high: nothing will be there

## Two step process

- ↳ find the frequent itemset
- ↳ association rules from the frequent itemset

## Name 3 Apriori algorithm : Iterative approach

scan the entire database and generate C1: candidate frequent itemset and L1: first frequent itemset.

The second frequent itemset L2 can be created by joining operation

C<sub>k</sub> can be generated by joining L<sub>k-1</sub>

↳ Join property and prune property

## Example:

Transaction Id	List of items
100	I <sub>1</sub> I <sub>2</sub> I <sub>5</sub>
101	I <sub>2</sub> I <sub>4</sub>
103	I <sub>2</sub> I <sub>3</sub>
104	I <sub>1</sub> I <sub>2</sub> I <sub>4</sub>
105	I <sub>1</sub> I <sub>3</sub>
106	I <sub>2</sub> I <sub>3</sub>
107	I <sub>1</sub> I <sub>3</sub>
108	I <sub>1</sub> I <sub>2</sub> I <sub>3</sub> I <sub>5</sub>
109	I <sub>1</sub> I <sub>2</sub> I <sub>3</sub>

Suppose

min support count = 2

min confidence

value = 70%

## Steps

i) generate first frequent itemset

$$I_1 \rightarrow 6$$

$$I_1 \rightarrow 6$$

$$I_2 \rightarrow 7$$

$$I_2 \rightarrow 7$$

$$I_3 \rightarrow 6$$

$$\Rightarrow I_3 \rightarrow 6$$

$$I_4 \rightarrow 2$$

$$\text{Support} = 2$$

$$I_4 \rightarrow 2$$

$$I_5 \rightarrow 2$$

$$I_5 \rightarrow 2$$

$$C_1$$

$$L_1$$

first frequent itemset

ii) generate second frequent itemset

(Join operation)

$$I_1, I_2 \rightarrow 4$$

$$I_1, I_2 \rightarrow 4$$

$$I_1, I_3 \rightarrow 4$$

$$I_1, I_3 \rightarrow 4$$

$$I_1, I_4 \rightarrow 1$$

$$I_1, I_5 \rightarrow 2$$

$$I_1, I_5 \rightarrow 2$$

$$I_2, I_3 \rightarrow 4$$

$$I_2, I_3 \rightarrow 4$$

$$I_2, I_4 \rightarrow 2$$

$$I_2, I_4 \rightarrow 2$$

$$I_2, I_5 \rightarrow 2$$

$$I_2, I_5 \rightarrow 2$$

$$I_3, I_4 \rightarrow 0$$

$$I_3, I_5 \rightarrow 1$$

$$I_4, I_5 \rightarrow 0$$

$$L_2$$

$$C_2$$

candidate pair

iii) generate third frequent itemset

$$I_1, I_2, I_3 \rightarrow 2$$

$$I_1, I_2, I_3 \rightarrow 2$$

$$I_1, I_2, I_5 \rightarrow 2$$

$$I_1, I_2, I_5 \rightarrow 2$$

$$I_1, I_2, I_3 \rightarrow 1$$

$$\text{support} = 2$$

$$I_1, I_2, I_5 \rightarrow 2$$

not all subsets are in  $L_2$ , so delete

$$I_1, I_2, I_3, I_5 \rightarrow 1$$

$$I_1, I_3, I_2, I_4 \rightarrow 0$$

$$I_1, I_3, I_2, I_5 \rightarrow 1$$

$$F_3$$

$$C_3$$

iv) generate fourth frequent itemset

$$I_1, I_2, I_3, I_5 \nrightarrow 1 \times \xrightarrow{\text{support} = 2} C_4 = \emptyset$$

$(I_2, I_3, I_5)$  is not there (e) prune property

v) generate association rules from frequent item set

Frequent item set  $\{I_1\}, \{I_2\}, \{I_3\}, \{I_4\}, \{I_5\}$

$\{I_1, I_2\}, \{I_1, I_3\}, \{I_1, I_5\}, \{I_2, I_3\}, \{I_2, I_4\}$

$\{I_2, I_5\}, \{I_1, I_2, I_3\}, \{I_1, I_2, I_5\}$

Let us consider  $I_1, I_2, I_5$

subsets are  $\{I_1, I_2\}, \{I_2, I_5\}, \{I_1, I_5\}, \{I_1\}, \{I_2\}, \{I_5\}$

Rule 1:  $I_1, I_2 \rightarrow I_5$

$$\frac{\text{support}(I_1, I_2, I_5)}{\text{support}(I_1, I_2)} = \frac{2}{4} = \frac{1}{2} = 50\% \quad X$$

Rule 2:  $I_1, I_5 \rightarrow I_2$

$$\frac{\text{support}(I_1, I_2, I_5)}{\text{support}(I_1, I_5)} = \frac{2}{2} = 1 = 100\% \quad \checkmark$$

Rule 3:  $I_2, I_5 \rightarrow I_1$

$$\frac{\text{support}(I_1, I_2, I_5)}{\text{support}(I_2, I_5)} = \frac{2}{2} = 1 = 100\% \quad \checkmark$$

Rule 4:  $I_1 \rightarrow I_2, I_5$

$$\frac{\text{support}(I_1, I_2, I_5)}{\text{support}(I_1)} = \frac{2}{6} = \frac{1}{3} = 33.3\% \quad X$$

Rule 5:  $I_2 \rightarrow I_1, I_5$

$$\frac{\text{support}(I_1, I_2, I_5)}{\text{support}(I_2)} = \frac{2}{7} = 28.5\% \quad X$$

Rule 6:  $I_5 \rightarrow I_1, I_2$

$$\frac{\text{support}(I_1, I_2, I_5)}{\text{support}(I_5)} = \frac{2}{8} = \frac{1}{4} = 25\% = 33.3\% + 100\% \quad \checkmark$$

⇒ Association Rules

$$I_1, I_5 \rightarrow I_2$$

$$I_2, I_5 \rightarrow I_1$$

$$I_5 \rightarrow I_1, I_2$$

If support value = 5% then do  $\frac{5}{100}$  and work.

On

suppose the association rule hot dogs → hamburgers is mined. give a min support threshold of 25% and a minimum confidence of 50%. does the association mining hold?

	hotdog	hotdog	Total
hamburger	2000	500	2500
hamburger	1000	1500	2500
Total	3000	2000	5000

hotdogs → hamburgers

$$\text{support}(\text{hotdog}, \text{hamb}) = \frac{2000}{5000} = \frac{2}{3} = 66.6\%$$

Min confidence = 50% ✓

⇒ Association holds

Packages (Decision Tree)

caret : classification & regression

rpart, e1071

21/4

in RStudio

Qn): A database has five transactions. Let minimum support = 50% and min confidence = 80%.

Transaction Id	items/- bought	$\frac{60}{100} \times 5$
T100	{M, O, N, K, E, Y}	= 3
T200	{D, O, N, K, E, Y}	
T300	{M, A, K, E, Y}	→ 27
T400	{M, V, C, K, Y}	
T500	{C, O, O, K, I, E}	

- Find all frequent item sets using apriori algorithm.
- List all the strong association rules with support s and confidence c matching the following meta-rule where x is a variable representing customers, item<sub>i</sub> denotes variables representing items  
 $x \in \text{transaction, buys}(x, \text{item}_1) \wedge \text{buys}(x, \text{item}_2) \Rightarrow \text{buys}(x, \text{item}_3)$

Ans: First frequent items

$$M \rightarrow 3$$

$$O \rightarrow 4$$

$$N \rightarrow 2$$

$$K \rightarrow 5$$

$$E \rightarrow 4$$

$$Y \rightarrow 3$$

Support = 3

$$D \rightarrow 1$$

$$A \rightarrow 1$$

$$V \rightarrow 1$$

$$C \rightarrow 2$$

$$I \rightarrow 1$$

Second iteration

$$M, O \rightarrow 1 \quad O, K \rightarrow 3 \quad K, Y \rightarrow 3 \quad M, K \rightarrow 3$$

$$M, K \rightarrow 3 \quad O, E \rightarrow 3 \quad E, Y \rightarrow 2 \quad \overrightarrow{S=3} \quad O, K \rightarrow 3$$

$$M, E \rightarrow 2 \quad O, Y \rightarrow 2 \quad O, E \rightarrow 3$$

$$M, Y \rightarrow 2 \quad K, E \rightarrow 4 \quad K, E \rightarrow 4$$

$$K, Y \rightarrow 3 \quad K, Y \rightarrow 3$$

Third iteration

$$M \rightarrow K \rightarrow O$$

$$M \rightarrow K \rightarrow O \rightarrow E$$

$$M \rightarrow K \rightarrow E$$

$$M \rightarrow K \rightarrow Y$$

$$O \rightarrow K \rightarrow E \rightarrow Y \rightarrow 3$$

$$S = 3$$

Third frequent itemset

$$O \rightarrow K \rightarrow Y$$

$$K \rightarrow E \rightarrow Y$$

$$O \rightarrow E \rightarrow K \rightarrow Y$$

OKE by prune property  
generate association rules

Frequent itemset  $\{OY, KY, EY, OKY, OKE, OEY\}$   
for  $\{OKE\}$

Let us consider

$$i) O \rightarrow K \rightarrow E$$

$$\frac{\text{Support}(OKE)}{\text{support}(OK)} = \frac{3}{3} = 100\%$$

$$ii) O, E \rightarrow K$$

$$\frac{\text{Support}(OKE)}{\text{support}(OE)} = \frac{3}{3} = 100\%$$

$$iii) O, K, E \rightarrow O$$

$$\frac{\text{Support}(OKE)}{\text{support}(KE)} = \frac{3}{4} = 75\%$$

iv) confidence value 80%.

Association rules

$O, K \rightarrow E$
$O, E \rightarrow K$

Ques: Database has four transactions. Let min support = 60%. min confidence = 80%.

CartId	TId	Items_bought
01	T100	$\{A, B, C, D\} \cap \{KC, SM, DC, BB\}$
02	T200	$\{BC, DM, GA, TP, WB\}$
03	T300	$\{WA, DM, WB, TP\}$
04	T400	$\{WB, SM, DC\}$

 $K^C$  - King's crab $SM$  - Sunset milk $DC$  - Dairyland cheese $BB$  - Best Bread

BC - Best cheese

DM - Dairyland Milk

GA - Goldenfarm Apple

TP - Tasty Pie

WB - Wonder Bread

WA - Westcoast Apple

- i) At the granularity of item category (item could be milk) for the following rule template

$\forall x \in \text{transaction}, \text{buys}(x, \text{item}_1) \wedge \text{buys}(x, \text{item}_2) \Rightarrow \text{buys}(x, \text{item}_3)$   
Find and list k frequent itemset and all the strong association rules with a support s and confidence c.

- ii) At the granularity of brand item category (item could be sunset milk following rule template

$\forall x \in \text{customer}, \text{buys}(x, \text{item}_1) \wedge \text{buys}(x, \text{item}_2) \Rightarrow \text{buys}(x, \text{item}_3)$

### ① Item category granularity

		Support val
01	T100	$\{G, M, C, B\}$
02	T200	$\{C, M, A, D, B\}$
01	T300	$\{A, M, B, P\}$
03	T400	$\{B, M, C\}$

$$\Rightarrow \frac{60}{100} \times 4$$

$$= 0.4$$

First freq.

$$C \rightarrow 1$$

$$M \rightarrow 4$$

$$C \rightarrow 3$$

$$B \rightarrow 4$$

$$A \rightarrow 2$$

$$P \rightarrow 2$$

$$\text{Support} = 2.4$$

$$M \rightarrow 4$$

$$C \rightarrow 3$$

$$B \rightarrow 4$$

Second freq.

$$M, C \rightarrow 2$$

$$M, B \rightarrow 3$$

$$C, B \rightarrow 3$$

$$M, B \rightarrow 3$$

$$C, B \rightarrow 3$$

$$M, C \rightarrow 3$$

$$\text{Support} = 2.4$$

Third freq.

$$M, B, C \rightarrow 2$$

$$\text{Support} = 2.4$$

$$\phi$$

Freq. item set  $\{\{M\}, \{B\}, \{C\}\}$  $\{\{M, B\}, \{M, C\}, \{B, C\}\}$  $M, B \Rightarrow C$ 

$$\frac{\text{Support}(M, B, C)}{\text{Support}(M, B)} = \frac{2}{4} = 50\% \quad 75\%$$

 $M, C \Rightarrow B$ 

$$\frac{\text{Support}(M, B, C)}{\text{Support}(M, C)} = \frac{2}{3} = 66.6\% \quad 100\%$$

 $B, C \Rightarrow M$ 

$$\frac{\text{Support}(M, B, C)}{\text{Support}(B, C)} = \frac{2}{3} = 66.6\% \quad 100\%$$

Association Rule

$$MC \Rightarrow B$$

$$BC \Rightarrow M$$

## (2) Brand item granularity

First freq.

$$KC \rightarrow 1$$

$$SM \rightarrow 2$$

$$DC \rightarrow 2$$

$$BB \rightarrow 1$$

$$BC \rightarrow 1$$

$$DM \rightarrow 2$$

$$GA \rightarrow 1$$

$$TP \rightarrow 2$$

$$WB \rightarrow 3$$

$$WA \rightarrow 1$$

$$\text{Support} = 2.4$$

$$WB \rightarrow 3$$

Antecedent  $\rightarrow$  consequentEg: Milk  $\rightarrow$  Bread

$$\text{Lift} = \frac{\text{support}(\text{Milk, Bread})}{\text{support}(\text{Milk}) * \text{support}(\text{Bread})}$$

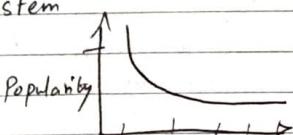
$$\text{Confidence} = \frac{\text{support}(\text{Milk, Bread})}{\text{support}(\text{Milk})}$$

5/5/22

Recommendation system / Algo

Search vs recommendation system

→ long tail phenomenon



Types of recommendation systems

→ collaborative based rec.

→ content-based rec.

→ knowledge-based

find similarity b/w 2

based on past historical info  
 Eg: based on movie rating given by user features of that movie are extracted for recommendation

Collaborative  $\rightarrow$  memory based  
 $\rightarrow$  model based

construct utility function / utility matrix  
 $R: X \times S$

$X \rightarrow$  set of customers

$S \rightarrow$  set of items

	Item 1	Item 2	Item 3	Item 4
User 1	5		3	
User 2		6	4	2
User 3		7		
User 4			5	
User 5	3			

$\rightarrow$  Data gathering

$\hookrightarrow$  collected in 2 ways

$\hookrightarrow$  explicit data collection

$\hookrightarrow$  implicit data collection

$\rightarrow$  Extrapolating the unknown values

$\rightarrow$  Evaluate the values

2 major challenges in recommendation system

1. Cold start problem

2. data sparsity problem

$\hookrightarrow$  new user will not have historical data so finding similarity in collaborative technique is not possible

Aggregate recommendation system

Recent videos in youtube, top 10, most popular.

Personalized recommendation system

Similarity

$$\text{Sim}(u_i, u_j) = \cos(u_i, u_j) = \frac{u_i \cdot u_j}{\|u_i\| \|u_j\|}$$

Pearson Correlation Coefficient

$\text{sim}(u_i, u_j)$

$$= \frac{\sum_k (r_{ik} - \bar{r}_i)(r_{jk} - \bar{r}_j)}{\sqrt{\sum_k (r_{ik} - \bar{r}_i)^2} \sqrt{\sum_k (r_{jk} - \bar{r}_j)^2}}$$

$$\bar{r}_{ik} \rightarrow \text{rating of user } i \text{ for the product } k$$

$r_{jk}$  - rating of user  $j$  for the product  $k$

$$r_{ui} = \bar{r}_u + \sum_{v \in N(u)} \text{sim}(u, v) (r_{vi} - \bar{r}_v)$$

$$\sum_{v \in N(u)} \text{sim}(u, v)$$

So finding unknown rating in matrix

$r_{vi}$   $\rightarrow$  observed rating of user  $v$  for the product  $i$

$\bar{r}_v$   $\rightarrow$  user  $v$  mean rating

Memory based collaborative  $\xrightarrow{\text{User based}}$  user-based

$\xrightarrow{\text{Item based}}$

User based

1) Calculate average rating

2) Similarity

3) Neighbourhood size

4) Predict the unknown rating value

	Item 1	Item 2	Item 3	Item 4	Item 5
Anu	5	3	4	4	?
Priya	3	1	2	3	3
Gam	4	3	4	3	5
Ravi	3	3	1	5	4
Raj	1	5	5	2	1

Step 1: Calculate Unknown Avg rating

$$\bar{Y}_{\text{anu}} = \frac{5+3+4+4}{4} = 4$$

$$\bar{Y}_{\text{priya}} = \frac{3+1+2+3+3}{5} = 2.4$$

$$\bar{Y}_{\text{sam}} = \frac{4+3+4+3+5}{5} = 3.8$$

$$\bar{Y}_{\text{ravi}} = \frac{3+3+1+5+4}{5} = 3.2$$

$$\bar{Y}_{\text{Raj}} = \frac{1+5+5+2+1}{5} = 2.8$$

Step 2: Similarity

$$\text{Sim}(\text{anu}, \text{priya}) = \frac{(5 \times 3) + (3 \times 1) + (4 \times 2) + (4 \times 3)}{\sqrt{(5)^2 + (3)^2 + (4)^2 + (4)^2} \cdot \sqrt{(3)^2 + (1)^2 + (2)^2 + (3)^2}}$$

$$= \frac{15 + 3 + 8 + 12}{\sqrt{66} \cdot \sqrt{23}}$$

$$= \frac{38}{\sqrt{66} \cdot \sqrt{23}} = 0.975$$

$$\text{Sim}(\text{anu}, \text{sam}) = \frac{(5 \times 4) + (3 \times 3) + (4 \times 4) + (4 \times 3)}{\sqrt{5^2 + 3^2 + 4^2 + 4^2} \cdot \sqrt{4^2 + 3^2 + 4^2 + 3^2}}$$

$$= \frac{57}{\sqrt{66} \cdot \sqrt{50}} = \frac{57}{\sqrt{66} \cdot \sqrt{44}} = 0.992$$

$$\text{Sim}(\text{anu}, \text{ravi}) = \frac{(5 \times 3) + (3 \times 3) + (4 \times 1) + (4 \times 5)}{\sqrt{66} \cdot \sqrt{3^2 + 3^2 + 1^2 + 5^2}}$$

$$= \frac{5 \cdot 9.08}{\sqrt{66} \cdot \sqrt{44}} = 0.890$$

$$\text{Sim}(\text{anu}, \text{Raj}) = \frac{(5 \times 1) + (3 \times 5) + (4 \times 5) + (4 \times 2)}{\sqrt{66} \cdot \sqrt{1^2 + 5^2 + 5^2 + 2^2}}$$

$$= \frac{48}{\sqrt{66} \cdot \sqrt{55}} = 0.796$$

Step 3: Neighbourhood size

Assume neighbourhood size = 2

⇒ Top 2 are

- i) Anu + Sam
- ii) Anu + Priya

Step 4: Predict unknown rating value

$$\bar{Y}_{\text{anu, items}} = 4 + \frac{0.992(5 - 3.8) + 0.975(3 - 2.4)}{0.992 + 0.975}$$

$$\textcircled{i} = \bar{Y}_{\text{anu}} + \frac{\text{sim}(\text{anu}, \text{sam})(\bar{Y}_{\text{sam}} - \bar{Y}_{\text{anu}}) + \text{sim}(\text{anu}, \text{priya})(\bar{Y}_{\text{priya}} - \bar{Y}_{\text{anu}})}{\text{sim}(\text{anu}, \text{sam}) + \text{sim}(\text{anu}, \text{priya})}$$

$$= 4 + \frac{1.188 + 0.582}{1.96}$$

$$\bar{Y}_{\text{anu, items}} = 4.903$$

Item based (Consider same data)

① Avg rating

$$\bar{X}_{i1} = \frac{5+3+4+3+1}{5} = 3.2$$

$$\bar{X}_{i2} = \frac{3+1+3+3+5}{5} = 3$$

$$\bar{X}_{i3} = \frac{4+2+4+1+5}{5} = 3.2$$

$$\bar{X}_{i4} = \frac{4+3+3+5+2}{5} = 3.4$$

$$\bar{X}_{i5} = \frac{3+5+4+1}{4} = 3.25$$

② Similarity

$$\text{sim}(i1, i5) = \frac{(3 \times 3) + (4 \times 5) + (3 \times 4) + (1 \times 1)}{\sqrt{3^2 + 5^2 + 4^2 + 1^2} \cdot \sqrt{3^2 + 5^2 + 4^2 + 1^2}}$$

$$= \frac{42}{\sqrt{51} \cdot \sqrt{35}} = 0.994$$

$$\text{sim}(i2, i5) = \frac{(3 \times 1) + (5 \times 3) + (4 \times 3) + (1 \times 5)}{\sqrt{3^2 + 5^2} \cdot \sqrt{1^2 + 3^2 + 3^2 + 5^2}}$$

$$= \frac{35}{\sqrt{51} \cdot \sqrt{44}} = 0.738$$

$$\text{sim}(i3, i5) = \frac{(3 \times 2) + (5 \times 4) + (4 \times 1) + (1 \times 5)}{\sqrt{51} \cdot \sqrt{2^2 + 4^2 + 1^2 + 5^2}}$$

$$= \frac{35}{\sqrt{51} \cdot \sqrt{46}} = 0.722$$

$$\text{sim}(i4, i5) = \frac{(3 \times 3) + (5 \times 3) + (4 \times 5) + (1 \times 2)}{\sqrt{51} \cdot \sqrt{3^2 + 3^2 + 5^2 + 2^2}}$$

$$= \frac{46}{\sqrt{51} \cdot \sqrt{47}} = 0.939$$

③ Neighbourhood zone = 2

i) Item1, Item5

ii) Item4, Item5

④ Find unknown rating

$$\bar{X}_{\text{items}} = \frac{3.25 + 0.994(5 - 3.2) + 0.939(4 - 3.4)}{0.994 + 0.939}$$

$$= 3.25 + \frac{1.7892 + 0.558}{1.92}$$

$$= 4.46$$

Content based filtering

I/P : User i's profile information, item description for each items  $j \in \{1, 2, 3, \dots, n\}$  keywords,  $\delta$  number of recommendation

Return  $n$  number of recommended items  
 $U = \{u_1, u_2, u_3, \dots, u_n\}$  user profile vector

$$(I_j)_{j=1}^n = \{I_{j,1}, I_{j,2}, I_{j,3}, \dots, I_{j,k}\}$$

↳ Item description vector

$$S_{i,j} = \text{sim}(U_i, I_j)$$

Return top  $\delta$  with maximum similarity

On: Data: Movie Item Feature Matrix

	Fantasy	Action	Cartoon	Drama	Comedy
M1	0	0	1	0	1
M2	0	0	0	1	0
M3	1	0	0	0	1
M4	1	1	0	0	1
M5	1	1	0	1	0
M6	0	1	1	0	1
M7	0	0	0	1	0

On: User Item Matrix

$$\text{User1: } M1 = 7 \quad M3 = 10 \quad M5 = \quad M7 = \\ M2 = 4 \quad M4 = \quad M6 =$$

Similarity b/w User and Item

$$= \begin{bmatrix} 7 \\ 4 \\ 10 \end{bmatrix} \times \begin{bmatrix} 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 & 7 & 0 & 7 \\ 0 & 0 & 0 & 4 & 0 \\ 10 & 0 & 0 & 0 & 10 \end{bmatrix} = \begin{bmatrix} 10 & 0 & 7 & 4 & 17 \end{bmatrix} \xrightarrow{\text{SUM}=38}$$

$$= [0.26 \ 0 \ 0.18 \ 0.10 \ 0.44]$$

Preference

$$\downarrow \quad \frac{10}{38} \quad \frac{0}{38} \quad \frac{7}{38} \quad \frac{4}{38} \quad \frac{17}{38}$$

- i) Comedy
- ii) Fantasy
- iii) Cartoon
- iv) Drama

Future recommendation

$$\begin{bmatrix} 1 & 1 & 0 & 0 & 1 \\ 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix} \times [0.26 \ 0 \ 0.18 \ 0.10 \ 0.44]$$

$$\begin{array}{cccccc} \text{Fantasy} & \text{Action} & \text{Cartoon} & \text{Drama} & \text{Comedy} \\ \hline 0.26 & 0 & 0 & 0 & 0.44 & M4 \\ 0.26 & 0 & 0 & 0.10 & 0 & M5 \\ 0 & 0 & 0.18 & 0 & 0.44 & M6 \\ 0 & 0 & 0 & 0.10 & 0 & M7 \end{array}$$

$$= \begin{bmatrix} 0.7 \\ 0.36 \\ 0.62 \\ 0.10 \end{bmatrix} \begin{matrix} M4 \\ M5 \\ M6 \\ M7 \end{matrix}$$

Top two recommended: M4, M6 by content based filtering

{ Implement content & collab based filtering in R / Python using inbuilt movie dataset (IMDB) }