



Data Warehousing and Mining Lab (DWML)-

Case study on OLAP & ETL process tools

Data Warehousing - OLAP

Online Analytical Processing Server (OLAP) is based on the multidimensional data model. It allows managers, and analysts to get an insight of the information through fast, consistent, and interactive access to information. This chapter covers the types of OLAP, operations on OLAP, differences between OLAP, and statistical databases and OLTP.

>>Types of OLAP Servers

We have four types of OLAP servers –

- Relational OLAP (ROLAP)
- Multidimensional OLAP (MOLAP)
- Hybrid OLAP (HOLAP)
- Specialized SQL Servers

>Relational OLAP

ROLAP servers are placed between relational back-end server and client front-end tools. To store and manage warehouse data, ROLAP uses relational or extended-relational DBMS.

- ROLAP includes the following –
- Implementation of aggregation navigation logic.
- Optimization for each DBMS back end.
- Additional tools and services.

>Multidimensional OLAP

MOLAP uses array-based multidimensional storage engines for multidimensional views of data. With multidimensional data stores, the storage utilization may be low if the data set is sparse. Therefore, many MOLAP server use two levels of data storage representation to handle dense and sparse data sets.



>Hybrid OLAP

Hybrid OLAP is a combination of both ROLAP and MOLAP. It offers higher scalability of ROLAP and faster computation of MOLAP. HOLAP servers allow large data volumes of detailed information. The aggregations are stored separately in the MOLAP store.

>Specialized SQL Servers

Specialized SQL servers provide advanced query language and query processing support for SQL queries over star and snowflake schemas in a read-only environment.

OLAP Operations Since OLAP servers are based on a multidimensional view of data, we will discuss OLAP operations in multidimensional data.

Here is the list of OLAP operations -

- Roll-up
- Drill-down
- Slice and dice
- Pivot (rotate)

>>Top OLAP Marketing Tools

You can use OLAP tools to analyze large volumes of multidimensional data from different perspectives. They make it easy to filter, analyze, and visualize key data insights. These tools are often part of a Business Intelligence Suite.

OLAP marketing tools should have the following features:

1. The ability to analyze large volumes of (big) data
2. The ability to perform analytical operations
3. A high degree of interactivity
4. Fast response times
5. Different types of **data visualizations**
6. The ability to analyze why things happen



Some OLAP tools used in marketing include:

>**IBM Cognos**

IBM Cognos is a web-based reporting and analytical tool to help you understand your organizational data. It's used to view or create detailed business reports, analyze data, and help you make effective business decisions.

>**MicroStrategy**

MicroStrategy is a business analytics platform that helps enterprises build and deploy analytics and mobile apps to transform their business. The MicroStrategy platform provides interactive **dashboards**, highly formatted reports, **ad hoc queries**, and automated report distribution. The software's ROLAP architecture is a key differentiator from other vendors who offer full-featured solutions.

>**Palo OLAP Server**

Palo is a MOLAP (Multidimensional Online Analytical Processing) server typically used as a BI tool for controlling and budgeting. It is a Jedox AG product. Palo enables multiple users to share one centralized data storage. It works with real-time data. Data can then be consolidated or written back with the help of multidimensional queries. Palo stores run-time data in its memory to give faster data access to users.

>**Sisense**

Sisense is an agile business intelligence (BI) solution that provides advanced tools to manage **big data in marketing analytics**. It helps you simplify complex data and transform it into powerful analytic apps to give you a more comprehensive understanding of your data.

>**icCube**

icCube owns a business intelligence software that offers an end-to-end BI solution. This is great for software companies looking to embed data analytics, visualization, and reporting into their product. icCube sells an online analytical processing server that is implemented in Java as per J2EE standards. It's an in-memory OLAP server and is compatible with any data source that holds its data in tabular form.



>**SAP NetWeaver Business Warehouse**

SAP NetWeaver Business Warehouse provides a high-performance infrastructure that helps you evaluate and interpret data. It provides reporting, analysis, and interpretation of business data quickly and in line with market needs.

>**Oracle Business Intelligence Enterprise Edition (OBIEE)**

Oracle Business Intelligence Enterprise Edition helps customers discover new data insights and make faster business decisions by offering interactive dashboards, powerful operational reporting, and real-time alerts. It reduces the total cost of ownership and increases return on investment for the entire organization.

>**Apache Kylin**

Apache Kylin is an open-source, distributed Analytical Data Warehouse for Big Data. It provides an SQL interface and MOLAP combined with Hadoop and Spark to support large data. In addition, Kylin reduces query processing time and quickly filters billions of data rows.

>**Final Thoughts**

Businesses continuously need to plan, analyze, and report on sales and marketing activities to maximize efficiency. OLAP applications can help increase the productivity of business managers, developers, marketing analysts, and whole organizations. In addition, they can also help you **transform data into actionable insights**.



What is ETL

ETL stands for Extract Transform and Load. ETL combines all the three database functions into one tool to fetch data from one database and place it into another database.

>**Extract:** Extract is the process of fetching (reading) the information from the database. At this stage, data is collected from multiple or different types of sources.

>**Transform:** Transform is the process of converting the extracted data from its previous form into the required form. Data can be placed into another database. Transformation can occur by using rules or lookup tables or by combining the data with other data.

>**Load:** Load is the process of writing the data into the target database.

>>**ETL** is used to integrate the data with the help of three steps Extract, Transform, and Load, and it is used to blend the data from multiple sources. It is often used to build a **data warehouse**. In the ETL process, data is extracted from the source system and converted into a format that can be examined and stored into a **data warehouse** or any other system. ETL is an alternative but a related approach which is designed to push processing down to the database to improve the performance.

>>Types of ETL Tools

ETL tools can be categorized into the following main types:

>Batch ETL Tools

In this type of ETL tool, batch processing is used to acquire data from the source systems. The data is extracted, transformed, and loaded into the repository in batches of ETL jobs. It's a cost-effective method because it uses limited resources in a time-bound way.



>**Real-Time ETL Tools**

Data is extracted, cleansed, enriched, and loaded to the target system in real-time ETL tools. These tools offer you faster access to information and improve time to insights.

As the need to gather and analyze the data in the shortest possible time has augmented, these ETL tools are becoming more popular among businesses.

>**On-Premise ETL Tools**

Many companies operate legacy systems that have both the data and the repository configured on-premise. The main reason behind such an implementation is data security. That's why companies prefer having an ETL tool deployed on-site.

>**Cloud ETL Tools**

As the name suggests, these tools are deployed on the cloud as various cloud-based applications form an essential part of enterprise architecture. Companies opt for cloud ETL tools to manage data transfer from these applications. Cloud-based ETL tools let businesses leverage flexibility and agility in the ETL process.

>**Oracle**

Oracle is the industry-leading database. It offers a wide range of choice of Data Warehouse solutions for both on-premises and in the cloud. It helps to optimize customer experiences by increasing operational efficiency.

STEPS:

1. Installing WEKA

Open Terminal

> sudo apt update

> sudo apt -y install weka

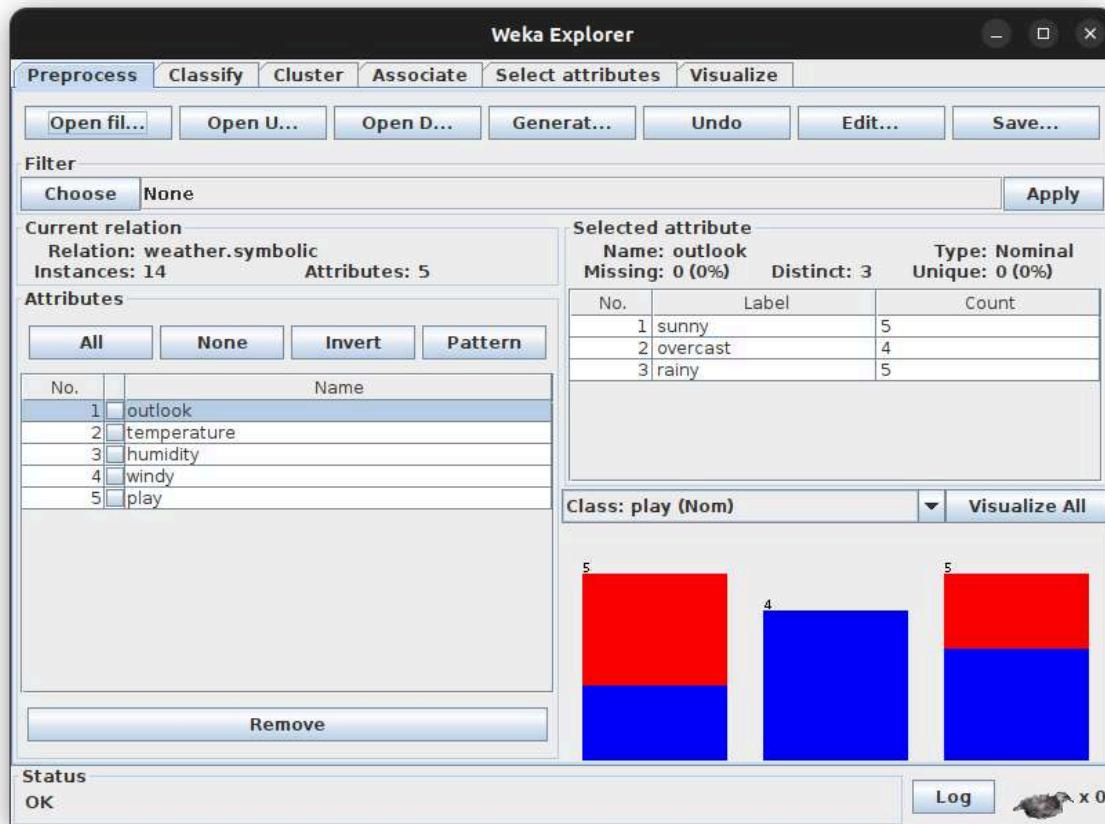
```
(base) computer@computer-ThinkCentre:~$ sudo apt update
Hit:1 http://archive.ubuntu.com/ubuntu jammy InRelease
Hit:2 http://archive.ubuntu.com/ubuntu jammy-updates InRelease
Hit:3 https://ppa.launchpadcontent.net/gns3/ppa/ubuntu jammy InRelease
Hit:4 http://archive.ubuntu.com/ubuntu jammy-backports InRelease
Hit:5 http://archive.ubuntu.com/ubuntu jammy-security InRelease
Hit:6 http://packages.microsoft.com/repos/code stable InRelease
Hit:7 https://deb.nodesource.com/node_18.x jammy InRelease
Hit:8 https://dl.google.com/linux/chrome/deb stable InRelease
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
250 packages can be upgraded. Run 'apt list --upgradable' to see them.
(base) computer@computer-ThinkCentre:~$ sudo apt -y install weka
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
weka is already the newest version (3.6.14-3).
0 upgraded, 0 newly installed, 0 to remove and 250 not upgraded.
(base) computer@computer-ThinkCentre:~$ weka
[warning] /usr/bin/weka: Unable to locate libsvm.jar in /usr/share/java
---Registering Weka Editors---
Trying to add database driver (JDBC): jdbc.idbDriver - Error, not in CLASSPATH?
```

2. Open WEKA.

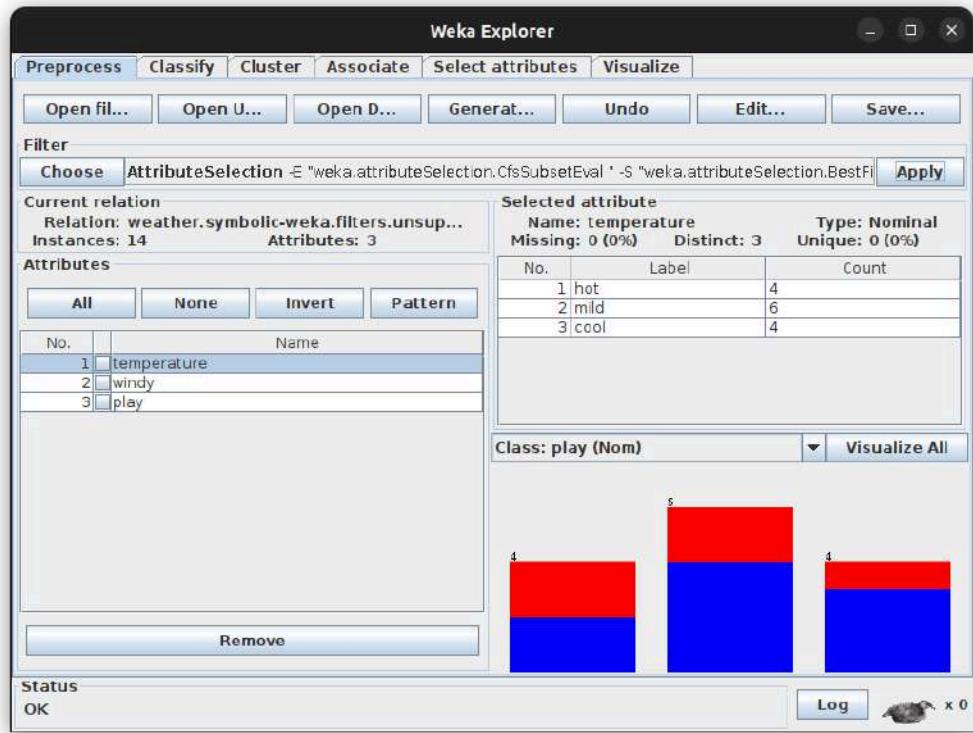
> weka



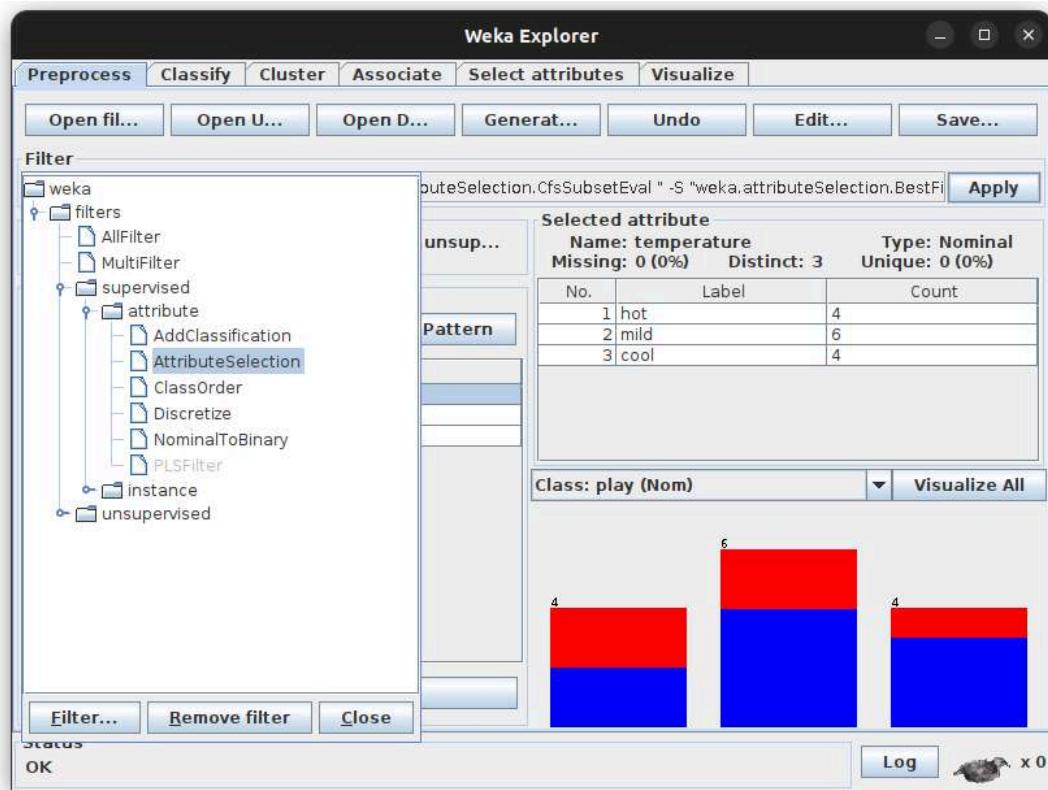
3. Open file: /->usr->share>doc-> weka->examples->weather.nominal.arff



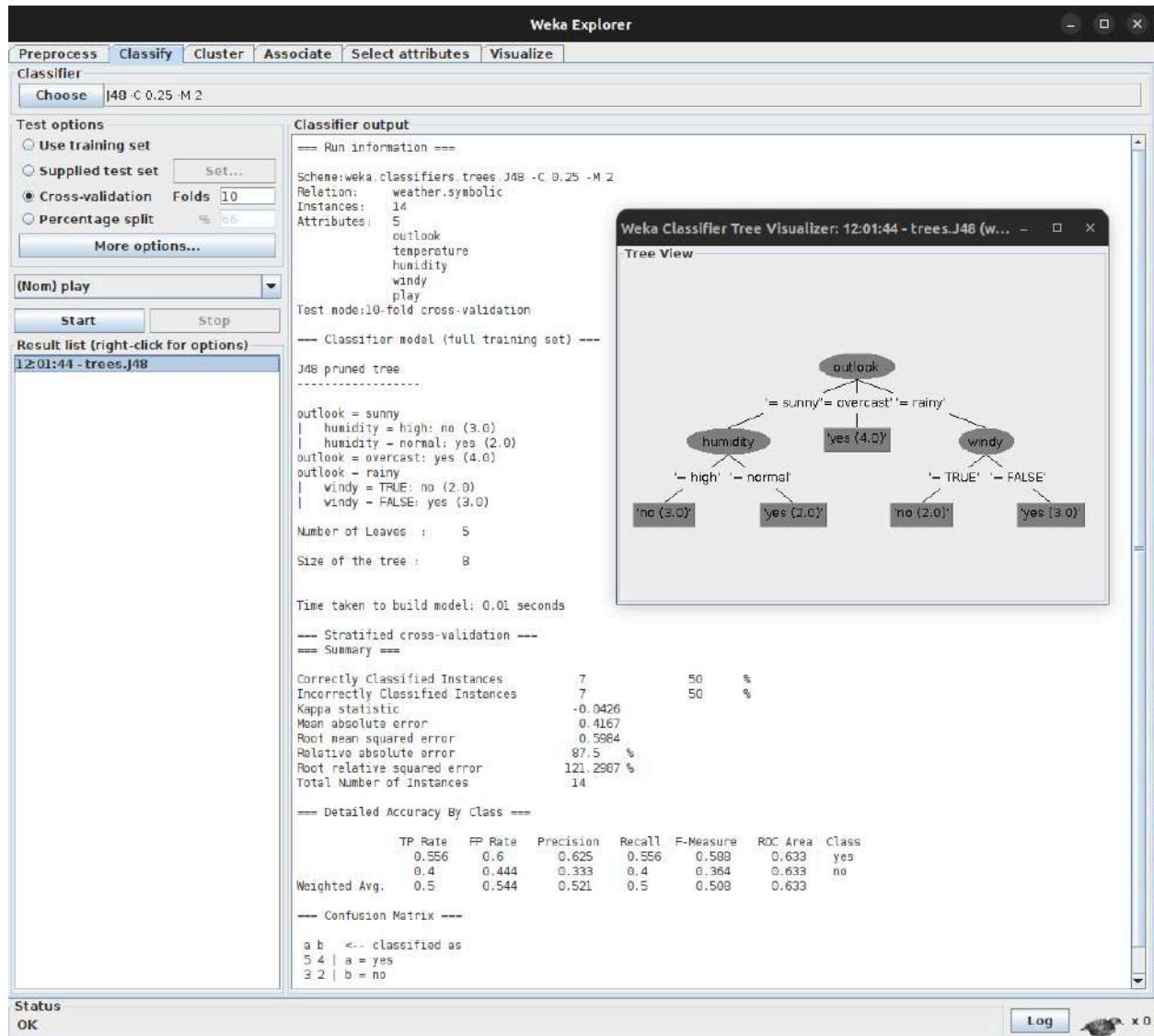
4. Removing attributes: To remove Attribute/s select them and click on the Remove button at the bottom.



5. Applying Filters: Click on the Choose button in the Filter subwindow and select the following filter >weka>filters>supervised>attribute>AttributeSelection

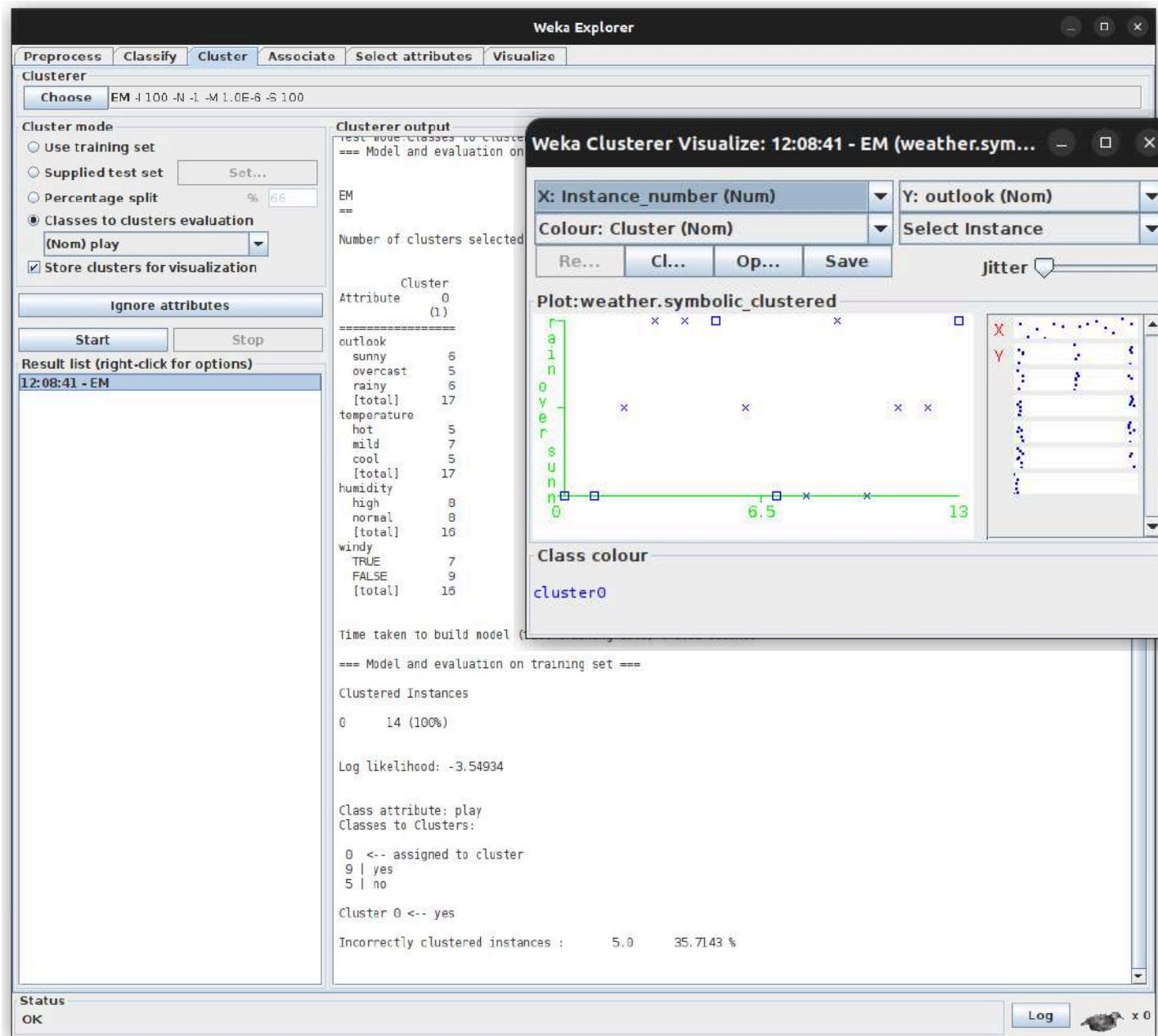


6. Selecting Classifier: Click on the Choose button and select the following classifier
 -weka→classifiers>trees>J48



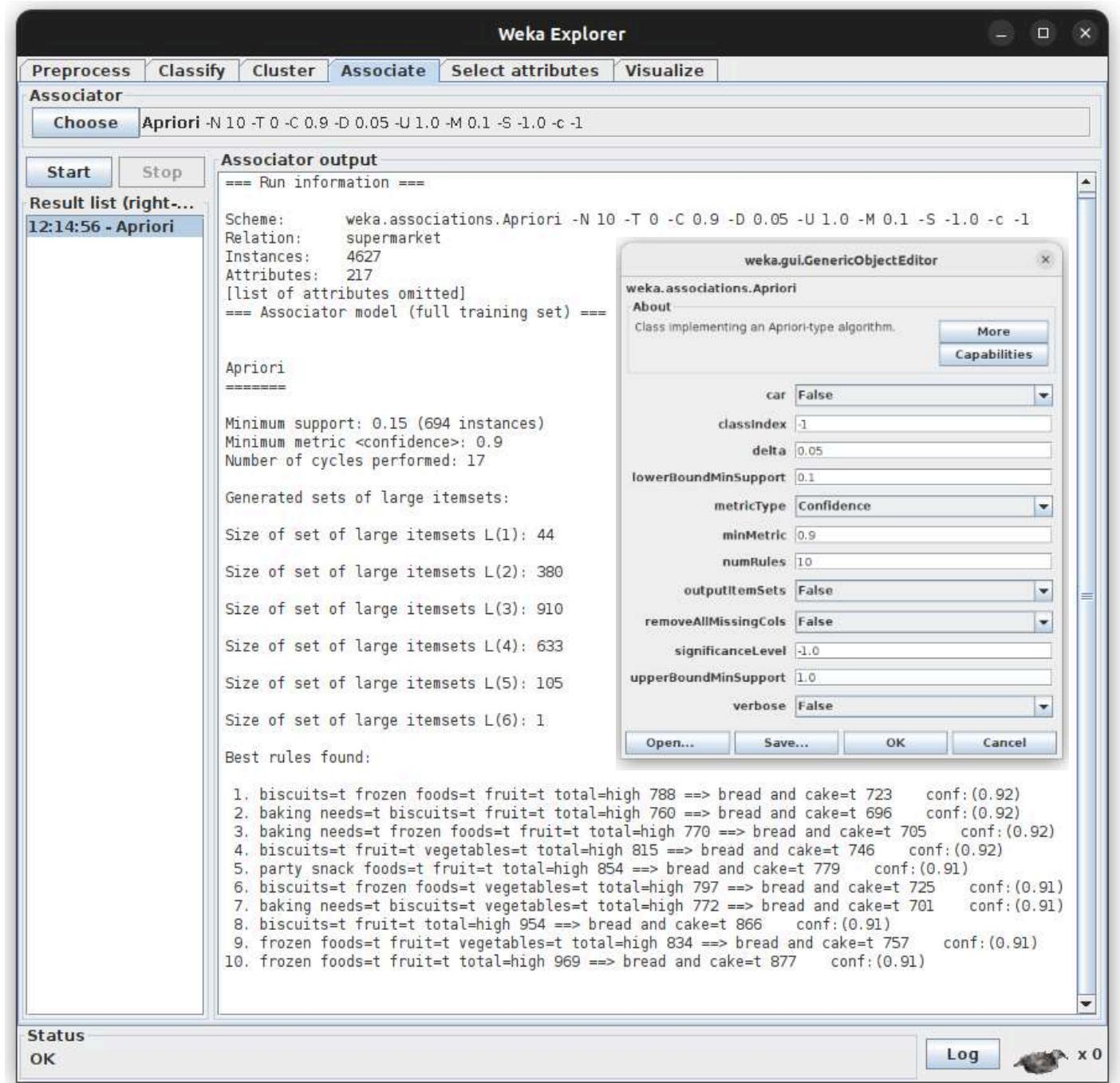
STEPS:

1. Open file - weather.arff
2. CLUSTER tab
3. Choose
4. Select EM
5. Start



STEPS:

1. Open supermarket.arff
2. Open Associate Tab
3. Choose
4. Select Apriori association
5. Start



CODE:

```
data = [
    ['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']]
]
init = []
for i in data:
    for q in i[1]:
        if(q not in init):
            init.append(q)
init = sorted(init)
print(init)
sp = 0.4
s = int(sp*len(init))

from collections import Counter
c = Counter()
for i in init:
    for d in data:
        if(i in d[1]):
            c[i]+=1
print("C1:")
for i in c:
    print(str([i])+": "+str(c[i]))
print()
l = Counter()
for i in c:
    if(c[i] >= s):
        l[frozenset([i])]+=c[i]
print("L1:")
for i in l:
    print(str(list(i))+": "+str(l[i]))
print()
pl = l
pos = 1
for count in range (2,1000):
    nc = set()
```

```

temp = list(l)
for i in range(0,len(temp)):
    for j in range(i+1,len(temp)):
        t = temp[i].union(temp[j])
        if(len(t) == count):
            nc.add(temp[i].union(temp[j]))
nc = list(nc)
c = Counter()
for i in nc:
    c[i] = 0
    for q in data:
        temp = set(q[1])
        if(i.issubset(temp)):
            c[i]+=1
print("C"+str(count)+":")
for i in c:
    print(str(list(i))+": "+str(c[i]))
print()
l = Counter()
for i in c:
    if(c[i] >= s):
        l[i]+=c[i]
print("L"+str(count)+":")
for i in l:
    print(str(list(i))+": "+str(l[i]))
print()
if(len(l) == 0):
    break
pl = l
pos = count
print("Result: ")
print("L"+str(pos)+":")
for i in pl:
    print(str(list(i))+": "+str(pl[i]))
print()
from itertools import combinations
for l in pl:
    c = [frozenset(q) for q in combinations(l,len(l)-1)]
    mmax = 0
    for a in c:
        b = l-a
        ab = l
        sab = 0
        sa = 0

```

```

sb = 0
for q in data:
    temp = set(q[1])
    if(a.issubset(temp)):
        sa+=1
    if(b.issubset(temp)):
        sb+=1
    if(ab.issubset(temp)):
        sab+=1
temp = sab/sa*100
if(temp > mmax):
    mmax = temp
temp = sab/sb*100
if(temp > mmax):
    mmax = temp
print(str(list(a))+" -> "+str(list(b))+" = "+str(sab/sa*100)+"%")
print(str(list(b))+" -> "+str(list(a))+" = "+str(sab/sb*100)+"%")
curr = 1
print("choosing:", end=' ')
for a in c:
    b = l-a
    ab = l
    sab = 0
    sa = 0
    sb = 0
    for q in data:
        temp = set(q[1])
        if(a.issubset(temp)):
            sa+=1
        if(b.issubset(temp)):
            sb+=1
        if(ab.issubset(temp)):
            sab+=1
    temp = sab/sa*100
    if(temp == mmax):
        print(curr, end = ' ')
    curr += 1
    temp = sab/sb*100
    if(temp == mmax):
        print(curr, end = ' ')
    curr += 1
print()
print()

```

OUTPUT:

(base) computer@computer:~/Desktop\$ python apriori.py

['I1', 'I2', 'I3', 'I4', 'I5']

C1:

['I1']: 6

['I2']: 7

['I3']: 6

['I4']: 2

['I5']: 2

L1:

['I1']: 6

['I2']: 7

['I3']: 6

['I4']: 2

['I5']: 2

C2:

['I2', 'I4']: 2

['I2', 'I3']: 4

['I3', 'I1']: 4

['I5', 'I4']: 0

['I4', 'I1']: 1

['I5', 'I1']: 2

['I3', 'I4']: 0

['I2', 'I5']: 2

['I5', 'I3']: 1

['I2', 'I1']: 4

L2:

['I2', 'I4']: 2

['I2', 'I3']: 4

['I3', 'I1']: 4

['I5', 'I1']: 2

['I2', 'I5']: 2

['I2', 'I1']: 4

C3:

['I3', 'I2', 'I4']: 0

['I2', 'I5', 'I1']: 2

['I2', 'I4', 'I1']: 1

['I5', 'I3', 'I1']: 1

['I2', 'I5', 'I4']: 0

['I2', 'I3', 'I1']: 2

['I2', 'I5', 'I3']: 1

L3:

['I2', 'I5', 'I1']: 2

['I2', 'I3', 'I1']: 2

C4:

['I5', 'I1', 'I2', 'I3']: 1

L4:

Result:

L3:

['I2', 'I5', 'I1']: 2

['I2', 'I3', 'I1']: 2

['I2', 'I5'] -> ['I1'] = 100.0%

['I1'] -> ['I2', 'I5'] = 33.3333333333333%

['I2', 'I1'] -> ['I5'] = 50.0%

['I5'] -> ['I2', 'I1'] = 100.0%

['I5', 'I1'] -> ['I2'] = 100.0%

['I2'] -> ['I5', 'I1'] = 28.57142857142857%

choosing: 1 4 5

['I2', 'I3'] -> ['I1'] = 50.0%

['I1'] -> ['I2', 'I3'] = 33.3333333333333%

['I2', 'I1'] -> ['I3'] = 50.0%

['I3'] -> ['I2', 'I1'] = 33.3333333333333%

['I3', 'I1'] -> ['I2'] = 50.0%

['I2'] -> ['I3', 'I1'] = 28.57142857142857%

choosing: 1 3 5

**CODE:**

```
dataset = [
    [0,0,1,0,0],
    [0,0,1,1,0],
    [1,0,1,0,1],
    [2,1,1,0,1],
    [2,2,0,0,1],
    [2,2,0,1,0],
    [1,2,0,1,1],
    [0,1,1,0,0],
    [0,2,0,0,1],
    [2,1,0,0,1],
    [0,1,0,1,1],
    [1,1,1,1,1],
    [1,0,0,0,1],
    [2,1,1,1,0]
]

mp = dict()
for i in range(len(dataset)):
    row = dataset[i]
    y = row[-1]
    if (y not in mp):
        mp[y] = list()
    mp[y].append(row)

for label in mp:
    print(label)
    for row in mp[label]:
        print(row)

test = [2,1,0,1]

probYes = 1

count = 0
total = 0
for row in dataset:
    if(row[-1] == 1):
        count+=1
    total+=1
```



```
print("Total yes: "+str(count)+" / "+str(total))
probYes *= count/total
for i in range(len(test)):
    count = 0
    total = 0
    for row in mp[1]:
        if(test[i] == row[i]):
            count += 1
            total += 1
    print('for feature '+str(i+1))
    print(str(count)+" / "+str(total))
    probYes *= count/total

probNo = 1
count = 0
total = 0
for row in dataset:
    if(row[-1] == 0):
        count+=1
        total+=1
probNo *= count/total
print("Total no: "+str(count)+" / "+str(total))
for i in range(len(test)):
    count = 0
    total = 0
    for row in mp[0]:
        if(test[i] == row[i]):
            count += 1
            total += 1
    print('for feature '+str(i+1))
    print(str(count)+" / "+str(total))
    probNo *= count/total

print(probYes)
print(probNo)

prob = probYes/(probYes+probNo)
print("Probability of playing golf: "+str(prob*100)+"%")
```



OUTPUT:

```
(base) computer@computer-ThinkCentre:~$ python NaivesBayes.py
0
[0, 0, 1, 0, 0]
[0, 0, 1, 1, 0]
[2, 2, 0, 1, 0]
[0, 1, 1, 0, 0]
[2, 1, 1, 1, 0]
1
[1, 0, 1, 0, 1]
[2, 1, 1, 0, 1]
[2, 2, 0, 0, 1]
[1, 2, 0, 1, 1]
[0, 2, 0, 0, 1]
[2, 1, 0, 0, 1]
[0, 1, 0, 1, 1]
[1, 1, 1, 1, 1]
[1, 0, 0, 0, 1]
Total yes: 9 / 14
for feature 1
3 / 9
for feature 2
4 / 9
for feature 3
6 / 9
for feature 4
3 / 9
Total no: 5 / 14
for feature 1
2 / 5
for feature 2
2 / 5
for feature 3
1 / 5
for feature 4
3 / 5
0.021164021164021163
0.006857142857142859
Probability of playing golf: 75.5287009063444%
```



CODE:

```
data = [
[5,2],[2,4],[9,5],[4,6],[5,2],[1,5],[6,7],[4,2],[6,4],[9,2],[4,5],[1,6],[4,7],[3,6],[1,1],[8,4],[8,7],
[7,2],[2,2],[2,1],[1,2],[1,4],[2,6],[7,7],[7,4],[3,4],[1,4]
]
x = [i[0] for i in data]
y = [i[1] for i in data]
import math
def dist(center, point):
    d = 0.0
    for i in range(0,len(point)):
        d += (center[i]-point[i])**2
    return math.sqrt(d)

def assignCenters(centers, dataset):
    clusters = []
    for i in range(len(dataset)):
        distances = []
        for center in centers:
            distances.append(dist(center, dataset[i]))
        temp = [z for z, val in enumerate(distances) if val==min(distances)]
        clusters.append(temp[0])
    return clusters

def mean_center(k, dataset, clusters):
    nCenters = []
    for i in range(k):
        x = 0.0
        y = 0.0
        count = 0
        for j in range(len(clusters)):
            if(i == clusters[j]):
                x += dataset[j][0]
                y += dataset[j][1]
                count += 1
        x = x/count
        y = y/count
        nCenters.append([x,y])
    return nCenters
print("enter k")
```



```
k = int(input())
centers = []
for i in range(k):
    print("enter center "+str(i))
    temp = [int(x) for x in input().split()]
    centers.append(temp)
print("Initial centers: ")
print(centers)
print("Initial clusters: ")
clusters = assignCenters(centers, data)
for i in range(k):
    print("cluster "+str(i))
    for j in range(len(clusters)):
        if(i == clusters[j]):
            print(data[j],end=' ')
print()
print()
for itr in range(10):
    print("Iteration "+str(itr))
    centers = mean_center(k,data,clusters)
    print("Updated centers: ")
    print(centers)
    clusters = assignCenters(centers, data)
    print("Updated clusters: ")
for i in range(k):
    print("cluster "+str(i))
    for j in range(len(clusters)):
        if(i == clusters[j]):
            print(data[j],end=' ')
print()
print()
```

**OUTPUT:**

```
computer@computer-ThinkCentre:~/Documents/CSE-AIML/TE/AIML11$ python
-u "/home/computer/Documents/CSE-AIML/TE/AIML11/kMeans.py"
enter k
2
enter center 0
6 4
Initial centers:
[[6, 4]]
Initial clusters:
enter center 1
9 2
Initial centers:
[[6, 4], [9, 2]]
Initial clusters:
cluster 0
[5, 2] [2, 4] [4, 6] [5, 2] [1, 5] [6, 7] [4, 2] [6, 4] [4, 5] [1, 6]
[4, 7] [3, 6] [1, 1] [8, 4] [8, 7] [2, 2] [2, 1] [1, 2] [1, 4] [2, 6]
[7, 7] [7, 4] [3, 4] [1, 4] cluster 1
[9, 5] [9, 2] [7, 2]

Iteration 0
Updated centers:
[[3.6666666666666665, 4.25], [8.333333333333334, 3.0]]
Updated clusters:
Iteration 1
Updated centers:
[[2.9, 4.0], [7.857142857142857, 4.428571428571429]]
Updated clusters:
Iteration 2
Updated centers:
[[2.555555555555554, 3.833333333333335], [7.444444444444445,
4.666666666666667]]
Updated clusters:
Iteration 3
Updated centers:
```



```
[[2.555555555555554, 3.833333333333335], [7.444444444444445,  
4.666666666666667]]  
Updated clusters:  
Iteration 4  
Updated centers:  
[[2.555555555555554, 3.833333333333335], [7.444444444444445,  
4.666666666666667]]  
Updated clusters:  
Iteration 5  
Updated centers:  
[[2.555555555555554, 3.833333333333335], [7.444444444444445,  
4.666666666666667]]  
Updated clusters:  
Iteration 6  
Updated centers:  
[[2.555555555555554, 3.833333333333335], [7.444444444444445,  
4.666666666666667]]  
Updated clusters:  
Iteration 7  
Updated centers:  
[[2.555555555555554, 3.833333333333335], [7.444444444444445,  
4.666666666666667]]  
Updated clusters:  
Iteration 8  
Updated centers:  
[[2.555555555555554, 3.833333333333335], [7.444444444444445,  
4.666666666666667]]  
Updated clusters:  
Iteration 9  
Updated centers:  
[[2.555555555555554, 3.833333333333335], [7.444444444444445,  
4.666666666666667]]  
Updated clusters:  
cluster 0  
[5, 2] [2, 4] [4, 6] [5, 2] [1, 5] [4, 2] [4, 5] [1, 6] [4, 7] [3, 6]  
[1, 1] [2, 2] [2, 1] [1, 2] [1, 4] [2, 6] [3, 4] [1, 4] cluster 1  
[9, 5] [6, 7] [6, 4] [9, 2] [8, 4] [8, 7] [7, 2] [7, 7] [7, 4]
```



CODE:

```
import networkx as nx                                     hits.txt
import numpy as np
from numpy import array
import matplotlib.pyplot as plt
with open('hits.txt') as f:
    lines = f.readlines()

G = nx.DiGraph()
for line in lines:
    t = tuple(line.strip().split(','))
    G.add_edge(*t)

h, a = nx.hits(G, max_iter=100)
h = dict(sorted(h.items(), key=lambda x: x[0]))
a = dict(sorted(a.items(), key=lambda x: x[0]))

print(np.round(list(a.values()), 3))
print(np.round(list(h.values()), 3))

pr = nx.pagerank(G)
pr = dict(sorted(pr.items(), key=lambda x: x[0]))
print(np.round(list(pr.values()), 3))

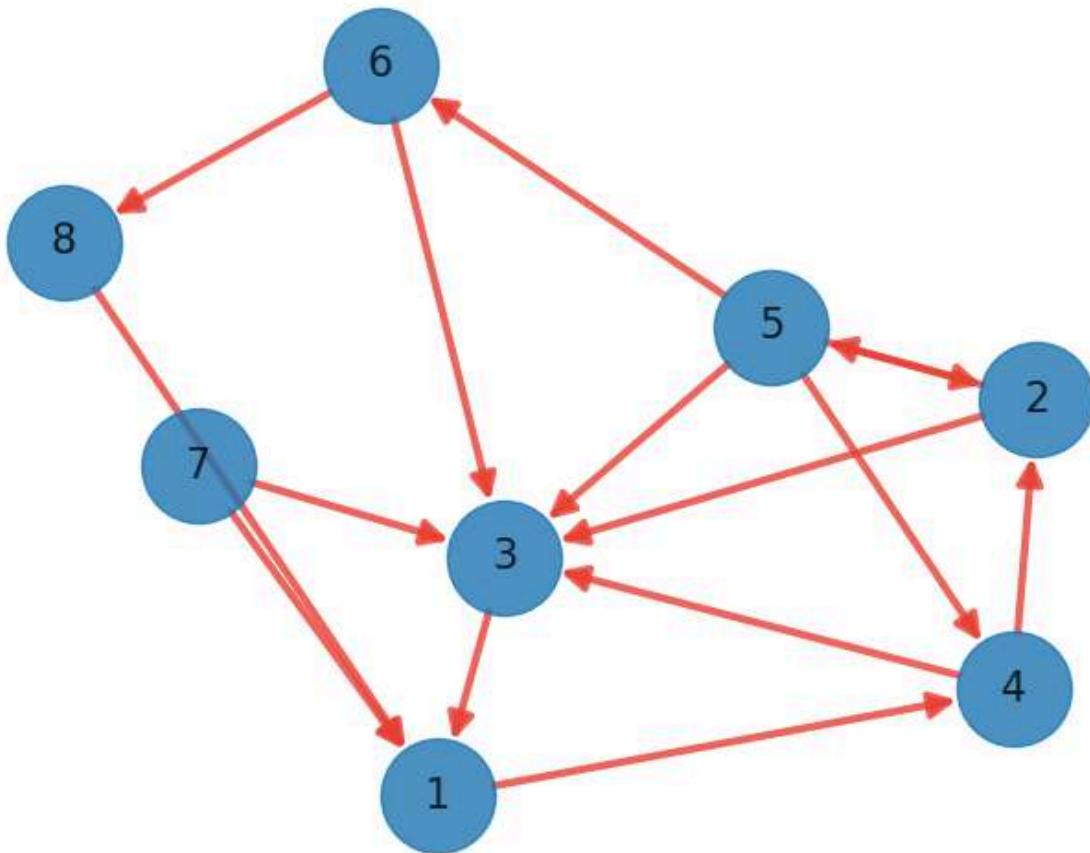
sim = nx.simrank_similarity(G)
lol = [[sim[u][v] for v in sorted(sim[u])] for u in sorted(sim)]
sim_array = np.round(array(lol), 3)
print(sim_array)

nx.draw(G, with_labels=True, node_size=2000, edge_color='#eb4034', width=3,
font_size=16, font_weight=500, arrowsize=20, alpha=0.8)
plt.savefig("graph.png")
```

OUTPUT:

```
computer@computer-ThinkCentre:~/Documents/CSE-AIML/TE/AIML11$ python -u  
"/home/computer/Documents/CSE-AIML/TE/AIML11/pageRank.py"  
/home/computer/anaconda3/lib/python3.9/site-packages/networkx/algorithms/l  
ink_analysis/hits_alg.py:78: FutureWarning: adjacency_matrix will return a  
scipy.sparse array instead of a matrix in Networkx 3.0.  
    A = nx.adjacency_matrix(G, nodelist=list(G), dtype=float)  
[0.088 0.187 0.369 0.128 0.059 0.11  0.     0.059]  
[0.043 0.144 0.03  0.187 0.268 0.144 0.154 0.03 ]  
[0.241 0.137 0.218 0.24  0.077 0.035 0.019 0.034]  
[[1.    0.207 0.221 0.193 0.217 0.269 0.     0.171]  
 [0.207 1.    0.355 0.369 0.302 0.553 0.     0.369]  
 [0.221 0.355 1.    0.242 0.4   0.324 0.     0.427]  
 [0.193 0.369 0.242 1.    0.229 0.548 0.     0.243]  
 [0.217 0.302 0.4   0.229 1.    0.271 0.     0.498]  
 [0.269 0.553 0.324 0.548 0.271 1.    0.     0.244]  
 [0.     0.     0.     0.     0.     1.    0.     0.    ]  
 [0.171 0.369 0.427 0.243 0.498 0.244 0.     1.    ]]
```

graph.png





CODE(pageHit.py):

```
import networkx as nx
import matplotlib.pyplot as plt
G = nx.DiGraph()
G.add_edges_from([('A', 'D'), ('B', 'C'), ('B', 'E'), ('C', 'A'), ('D', 'C'), ('E', 'D'), ('E', 'B'), ('E', 'F'), ('E', 'C'), ('F', 'C'), ('F', 'H'), ('G', 'A'), ('G', 'C'), ('H', 'A')])
plt.figure(figsize =(10, 10))
nx.draw_networkx(G, with_labels = True)
hubs, authorities = nx.hits(G, max_iter = 50, normalized = True)
print('Hub Scores: ')
for i in hubs:
    print("{}: {}".format(i, hubs[i]))
print("\nAuthority Scores: ")
for i in authorities:
    print("{}: {}".format(i, authorities[i]))
```

OUTPUT:

```
computer@computer-ThinkCentre:~/Documents/CSE-AIML/TE/AIML11$ python -u
"/home/computer/Documents/CSE-AIML/TE/AIML11/pageHit.py"
    A = nx.adjacency_matrix(G, nodelist=list(G), dtype=float)

Hub Scores:
A: 0.04642540403219996,
D: 0.13366037526115382,
B: 0.15763599442967322,
C: 0.03738913224642651,
E: 0.2588144598468665,
F: 0.15763599442967322,
H: 0.03738913224642651,
G: 0.17104950750758036,


Authority Scores:
A: 0.10864044011724336,
D: 0.13489685434358004,
B: 0.11437974073336449,
C: 0.388372800387618,
E: 0.06966521184241475,
F: 0.11437974073336449,
H: 0.06966521184241477,
G: 0.0,
```

Name :- Prathamesh b. CHIKONKAR

ROLL NO. :- AIML11

BRANCH :- CSE - (AI & ML)

Year :- TE Subject :- DW & M [Data warehousing and mining]

Topic :- Assignment No. 01

Sign :- Prathamesh

Date :- August' 22

1. Explain Data warehouse Architecture in detail.

Architecture is the proper management / arrangement of the components.

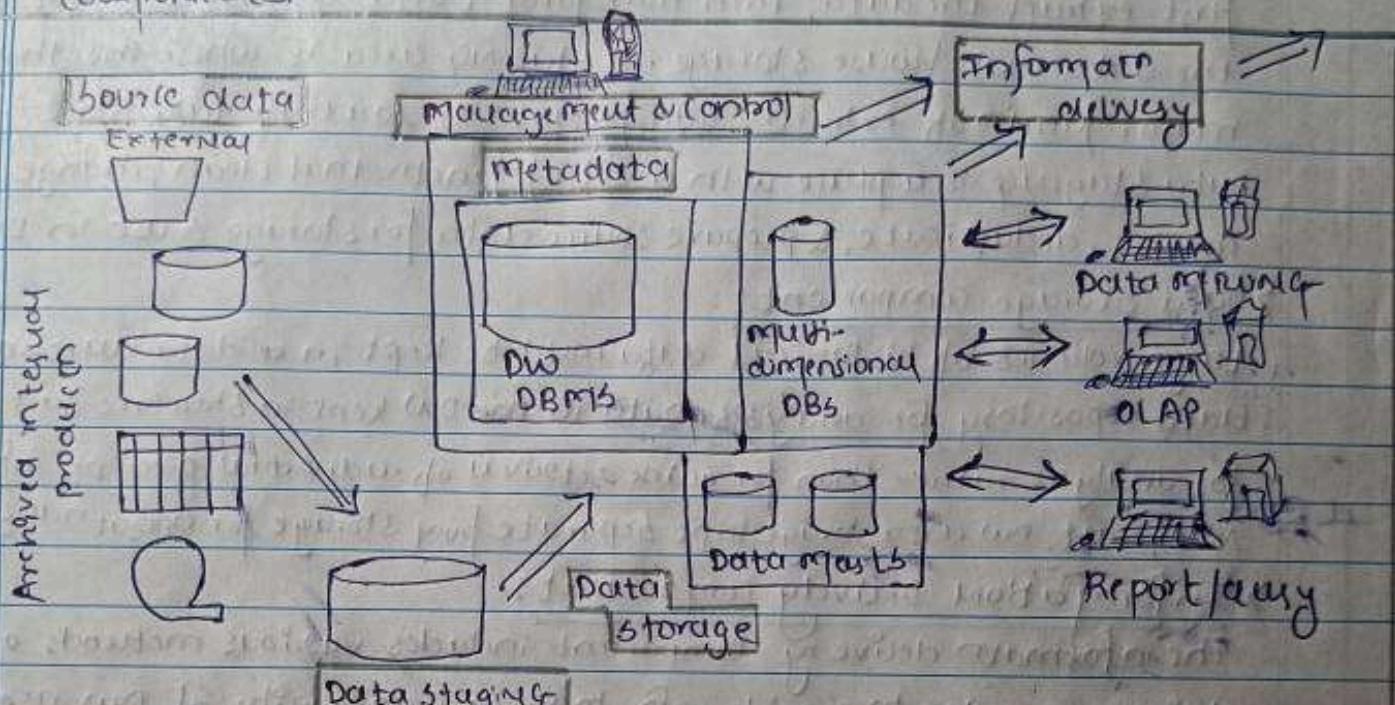


Fig. Data warehouse Architecture.

> Source data component

Source data coming into the data warehouse may be grouped in four broad categories, as discussed here:

- Product data: This type of information originates from the company's numerous operational systems. (data from various operational systems)
- Internal data: Users preserve their "private" spreadsheets, documents, customer profiles, & sometimes even departmental DB in every company.
- External data: most executives rely on data from outside sources for a large part of their information they use. Rely on competitor market share data. Use standard values of financial indicators.
- Archived data: operational systems are largely designed to run a company's existing operation. Data kept in operational system DB.

> Data Staging Component : (ETL process)

To get the data ready, three major functions must be completed. You must first extract the data, then transform it and at last load it into the data warehouse storage. A staging area is where the three major functions of extraction, transformation & loading take place.

Data Staging is a place with a set of functions that clean, change, combine, convert, deduplicate, & prepare source data for storage & use in DW.

> Data Storage Component :

Large volumes of historical data must be kept in a data warehouse's data repository for analysis. Data in the DW kept in structure suitable for analysis rather than for quick retrieval of individual piece of information. As a result, DW data storage kept separate from storage for operational systems.

> Information Delivery Component :

The information delivery component includes various methods of info. delivery in order to provide info. to large community of DW users. This information can be delivered through online, intranet, internet or email mode.

> Metadata Components :

Metadata is "data that describes other data". It is similar to data dictionary or data catalog in DBMS. Keeps the information about logical data structure, info. about files & addresses, info. about indices & so on. Types → Operational metadata, Extract and Transformation metadata and End-user metadata. e.g.: - web page, document, etc.

> Management and Control Component :

This component of DW sits on the top of all other components. The management & control component coordinates services & activities within the DW. Controls data transformation, mediates info. delivery.

The management & control component integrates with the metadata component to perform the management & control function.



2. Explain Data warehouse schemas in detail.

Schema is a logical description of entire database. A database uses a relational model, while DW uses Star and Snowflake Schema.

> Star Schema :

- A star schema is the elementary form of a dimensional model, in which data are organised into facts and dimensions.
- A fact is an event that is counted or measured, such as sale or dealer.
- It is a relational schema (represents multidimensional data model)
- It is known as star Schema because the entity-relationship diagram of this schema simulates a star, with points, diverge from central table.
- centre consists of a large fact table & points of star are dimension Table.

Fact Tables :

- It is a table in a star schema which contains facts and connected to dimensions. A fact has two types of column: those that include fact and those that are foreign key to the dimension table. primary key of the fact table is generally a composite key that is made up of all its foreign key.

Dimension Tables :

- A dimension is an architecture usually composed of one or more hierarchies that categorize data. If a dimension has not got hierarchies and levels, it is called a flat dimension or list. small in size than fact table.

Characteristics :

It creates denormalized database that can quickly provide query response.
It provides a flexible design that can be changed easily.

It reduces the complexity of metadata for both developers & end users.

Keys : primary keys, surrogate keys, foreign key.

Advantages :

- Query performance & query as well, - Built-in referential integrity

Disadvantages :

- Data integrity is not enforced well since in a highly denormalized schema.

> Snowflake Schema:

- Snowflake schema is variant of star schema.
- Here the centralised fact table is connected to multiple dimensions.
- In the snowflake schema, dimensions are present in a normalized form in multiple related tables.
- Snowflake schema affects only the dimension tables & not fact tables.

Characteristics of snowflake schema:

- uses small disk space.
- Easy to implement dimension that is added to schema.
- Multiple tables, so performance is reduced.

Advantages:

- It provides structured data which reduces the problem of integrity.
- uses small disk space because data are highly structured.

Disadvantages:

- Reduces space consumed by dimension tables, but compared with the entire DW the saving is usually insignificant.
- Avoid snowflaking or normalization of dimensions tables, unless required and appropriate. - complex structure
- Do not snowflake hierarchies of one-dimensional table into separate table. - query performance degrades.

3. why is entity-relationship modelling technique not suitable for the data warehouse?

How is dimensional modelling different.

- Dimensional Modelling (DM) is a DB technique optimized for data storage in a DW. The purpose of dimensional modelling is to optimize the DB for faster retrieval of data. The concept of DM was developed by Ralph Kimball and consists of "fact" & "dimension" tables. Dimension table records information on each



dimension, and fact records all the "fact", or measures.

- A dimensional model in data warehouse is designed to read, summarize, analyze numeric information like values, balances, counts, weights, etc. in a DW. In contrast, fact models are optimized for add, update and delete of data in real-time online transaction system.
- These dimensional and relational models have their unique way of data storage that has specific advantages.
 - for instance, in relational model, normalized and ER models reduce redundancy in data. On the contrary, dimensional model in DW arranges data in such a way that is easier to retrieve information and generates reports.
 - Hence, dimensional models are used in the DW systems and not a good fit for relational systems.

4. what is metadata? Explain different types of metadata.

Metadata is "data that describes other data" (data dictionary).
Metadata in a data warehouse is similar to the data dictionary or the data catalog in a DBMS.

In data dictionary, you keep the information about the logical data structures, the information about files, addresses, info. about indexes, and so on, so data dictionary contains data about data in DB. Similarly, metadata component is data about data in the DW.

Types of metadata:

> Operational metadata.

Data for the DW comes from various operational systems of enterprise. These source systems include different data structures. The data elements selected for the DW have various field length & data types. When we deliver information from the source system to the end-users,

We must be able to tie that back to the source data sets. Operational metadata contains all this info. about the operational data sources.

> Extract and Transformation Metadata -

Extract & transform metadata include data about the removal of data from the source systems, namely, the extract frequencies, extract methods, and business rules for the data extract.

Contains info. about all data transformation takes place in data staging.

> End-User Metadata -

The end-user metadata is the navigational map of the DW. It enables the end-users to find data from the DW. The end-user metadata allows the end-users to use their business terminology and look for the info. in those ways in which they usually work of the business.

5. A dimension table is wide, the fact table is deep. Explain.

- fact table is a table in star schema which contains facts and connected to dimensions, has two types of column: those that include fact & those that use foreign keys to dimension table.

- whereas on other hand dimension is an architecture usually composed of one or more hierarchies that categorize data. If a dimension has not good got hierarchies and levels, it is called a flat dimension or list.

> fact tables are deep whereas dimension tables are wide as fact tables will have a higher no. of rows and a lesser number of columns.

A primary key defined in the fact table is primary key to identify each row separately. The primary key is also called a composite key in fact table. A dimension table contains a higher granular information so have less no. of records and it needs to have all necessary details (more column(s)) related to the grain of the table. On the other side,

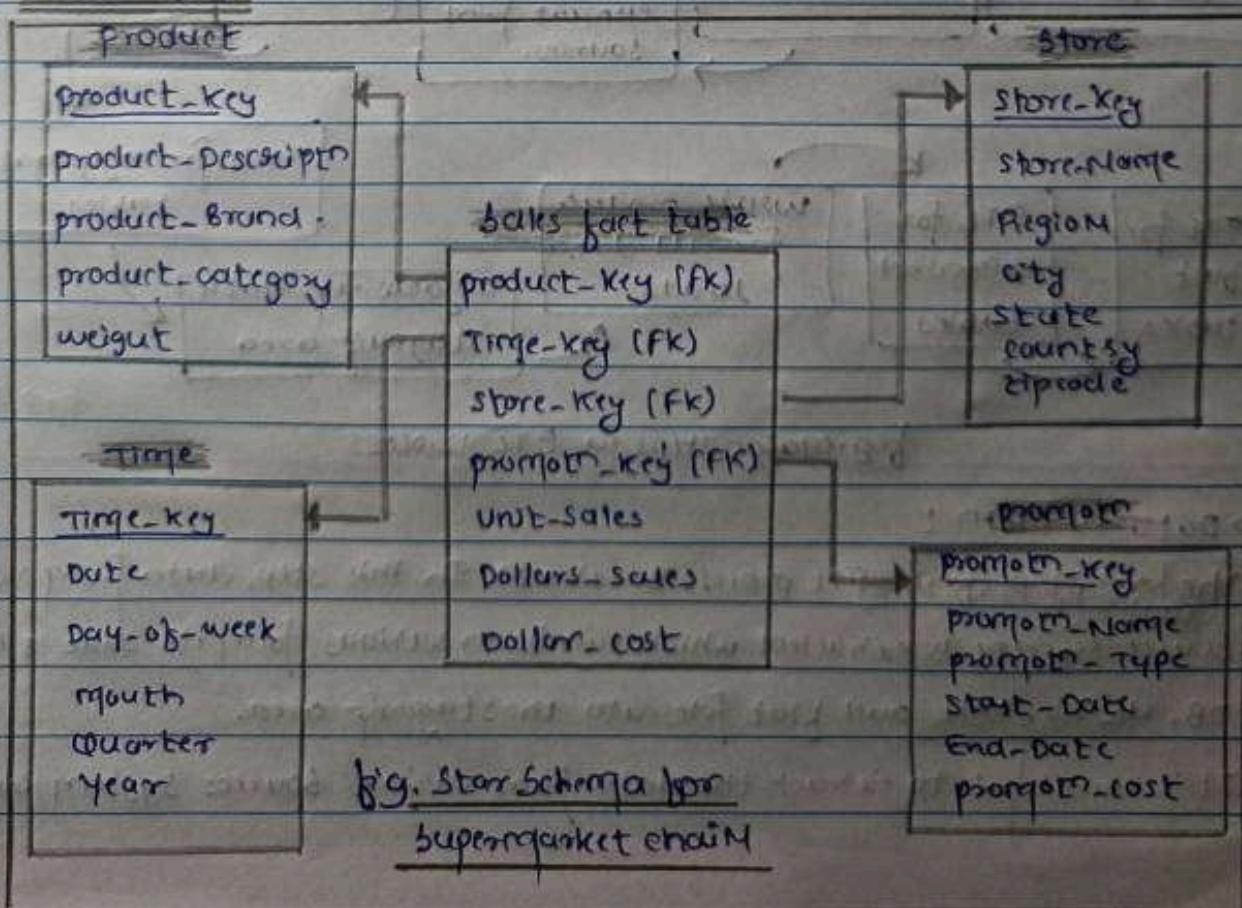
a fact table has the lowest level grain of a subject area. Lower grain causes more number of rows in the fact table.

6. For a Supermarket chain, consider the following dimensions namely product, store, time and promotion. The schema contains a central fact table for sales with three measures unit-sales, dollars sales and dollar-cost.

i. Draw star schema

- ii. Calculate the maximum no. of base fact table records for warehouse with the following values given below:
- Time period - 5 years
 - Store - 300 stores reporting daily sales
 - Product - 40,000 products in each store (about 4000 sell in each store daily)
 - Promotion - a sold item may be in only one promotion in a store on given day.

Star Schema :



> fact table records

→ time period = 5 years = 5×365 days = 1825 days

→ No. of stores = 300

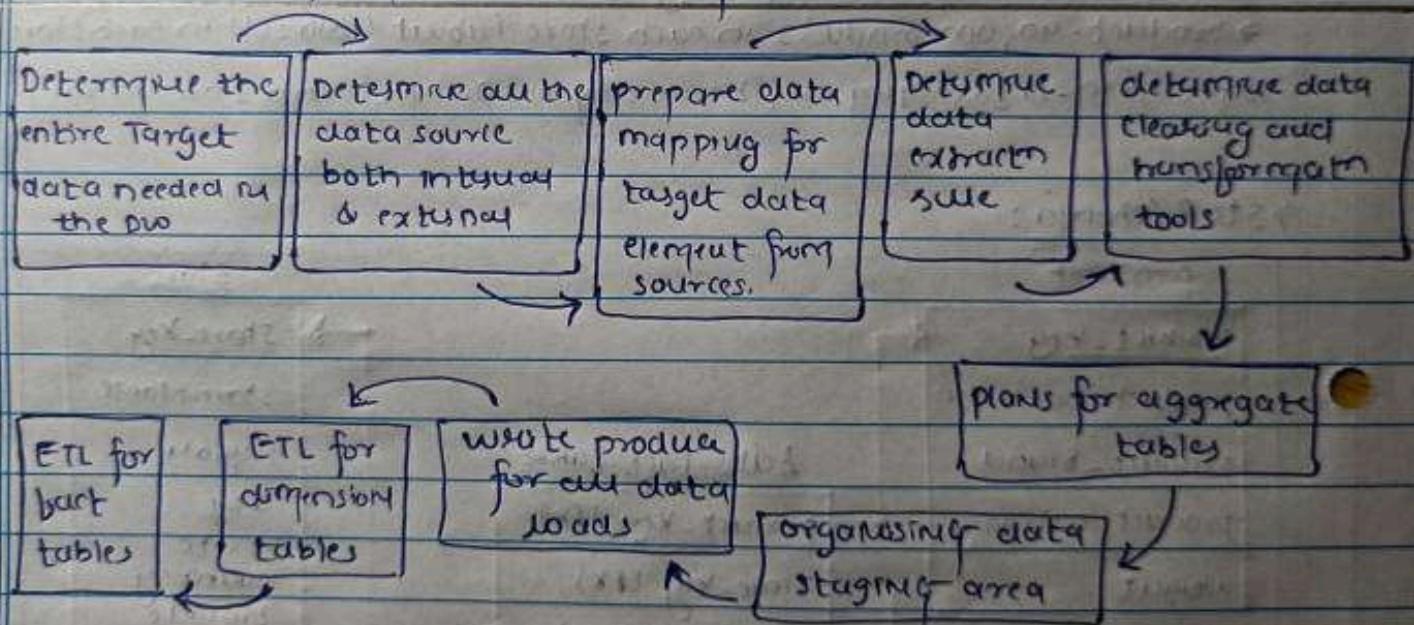
→ product sell in each store daily = 4000

→ promod? = 1

∴ maximum no. of fact table records = $1825 \times 300 \times 4000 \times 1$
 $= 2190000000$ records.

7. Explain the major steps in ETL process:

The process of extracting data from source systems and bringing it into the data warehouse is commonly called ETL, which stands for extract, transform, and loading.



Major steps in ETL process

> Data extract :

- The first step of the ETL process is extract; In this step, data from various source system is extracted which can be in various formats like relational DB, XML, CSV and flat files into the staging area.
- It is important to extract the data from various source system and



Store it into the staging area first and not directly into the DW because the extracted data is in various formats and can be corrupted also.

Hence, loading it directly into the data warehouse may damage it & rollback will be much more difficult. Therefore, this is one of the most important steps of ETL process.

→ identifying sources → deciding methods → Extract frequency → defining the time window → Exceptional condition (which will be handled) → handling of data sources.

Techniques:

⇒ Immediate data Extract Technique: (Real time extract)

→ capture through transaction logs → database triggers → source application

⇒ delayed data Extract

→ capture based on date and timestamp → comparing files → without transaction logs for data & time setup.

Data transformation:

- The second step of ETL process is transformation. In this step, a set of rules or functions are applied on the extracted data to convert it into single standardized format.

⇒ Task involved in Data Transformation

→ format conversion (data types, data field lengths)

→ decoding of files → splitting of fields → merging information

→ character set conversion → date & time conversion

→ summarization → key restructuring → de duplication

Techniques

⇒ Data smoothing: used for removing the noise from a dataset.

⇒ Data aggregation: collecting data from variety of sources & storing it in a single format. This method helps in collecting vast data.

⇒ Discretization: process of converting continuous data into set of data intervals.

⇒ Generalization: low-level data attributes pass into high-level attributes.

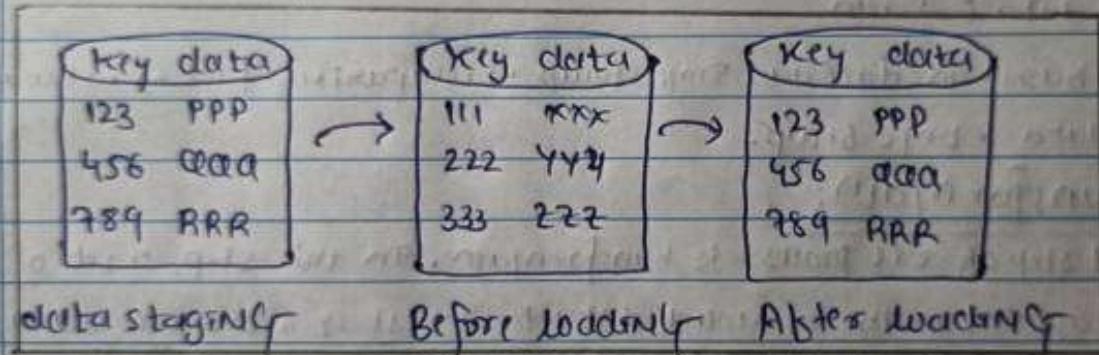
⇒ Normalization: data transformation in data mining.

7 Data loading:

- The third and final step of the ETL process is loading. In this step the transformed data is finally loaded into the data warehouse.
- sometimes the data is updated by loading it to the data warehouse very frequently & sometimes it's done after long but regular intervals.
- The rate & period of loading solely depends on the requirements and various from system to system.
 ⇒ Total load, Incremental load, Full refresh

Techniques:

data staging



⇒ load mode

⇒ append mode of data loading

⇒ descriptive merge (w.r.t. f.k)

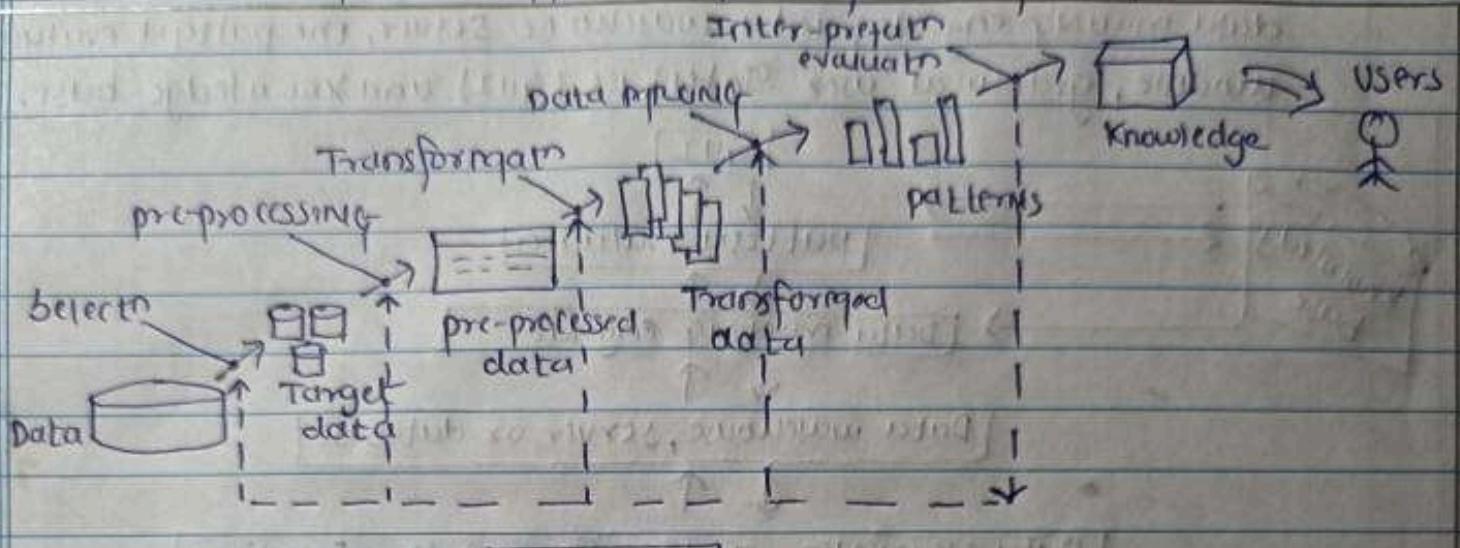
⇒ constructive merge

8. Describe the steps involved in Data Mining when viewed as a process of Knowledge Discovery.

- Knowledge Discovery in the database (KDD) is the process of searching for hidden knowledge in the massive amounts of data that we are technically capable of generating and storing.
- The basic task of KDD is to extract knowledge from information from a lower level data (database).



- The overall process of finding and interpreting patterns from data involves the repeated application of the following steps:



- > Data Cleaning: Removal of noise, inconsistency data & outliers.
- > Data Integration: Data from various sources such as databases, data warehouse, and transactional data are integrated.
- > Data Selection: Data relevant to the analysis task is retrieved from DB. collecting only necessary information to the model.
- > Data Transformation: Transformed & consolidated into forms appropriate for mining by performing summary aggregation operations.
- > Data Mining: Essential process where intelligent methods are applied to extract data patterns. (which model parameters may be appropriate)
- > Pattern Evaluation: To identify truly interesting patterns representing knowledge based on interesting measures.
- > Knowledge presentation: visualization and knowledge representation techniques are used to present mined knowledge to users. visualization can be in form of graphs, charts or tables. This all are valuable or enlightenment in its discovery.

q. Explain Data mining architecture in detail.

The significant components of data mining systems are data sources, data mining engine, data warehouse server, the pattern evaluation module, displayed user interface (GUI) and knowledge base.

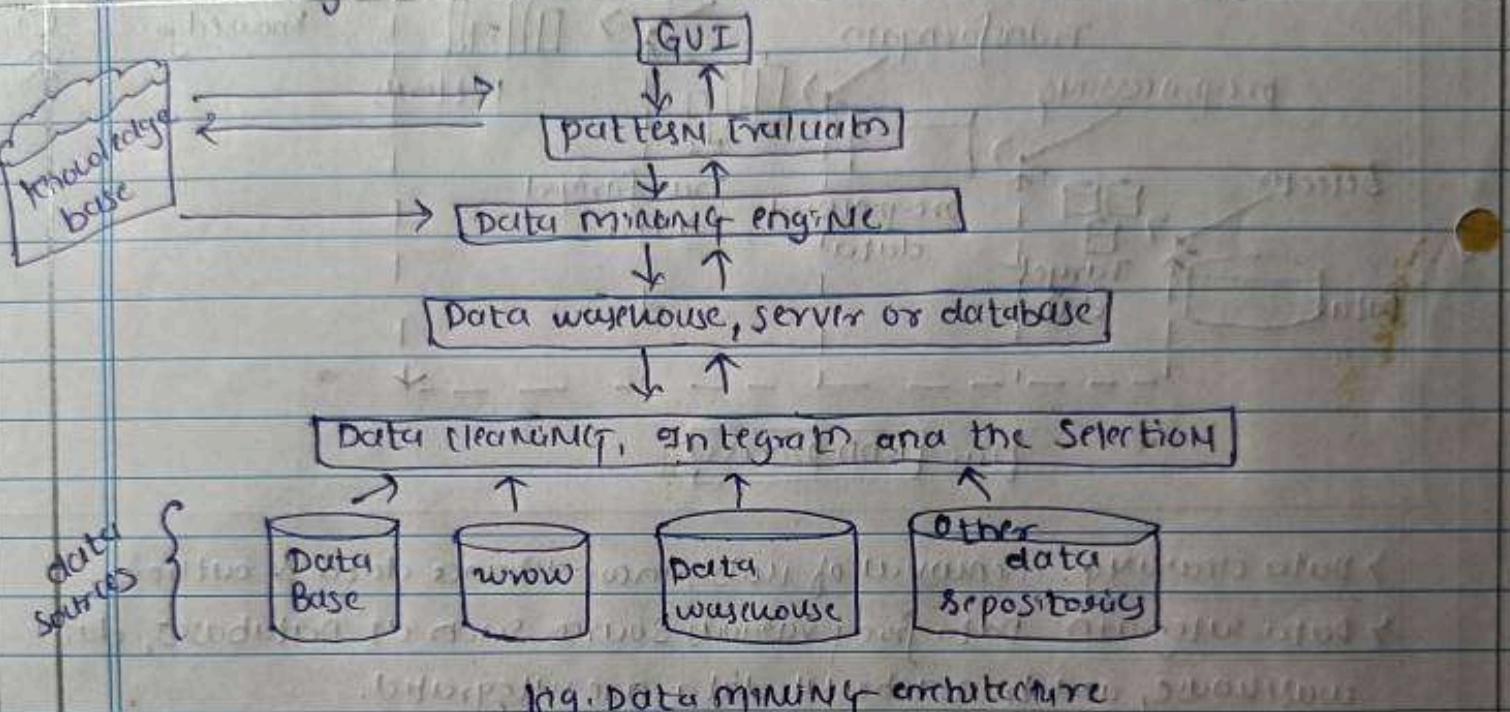


Fig. Data mining architecture

- Data Sources: Actual sources of data is the Database, data warehouse, www, txt files, and other documents. You need a huge amount of historical data for data mining. Successful.
- Data cleaning, Integration and selection: Before passing data to the database or DW server, data must clean, integrated & selected.
- Database or Data warehouse behaviour: The DB or DW server consists of the original data that is ready to be processed. Hence, the server is used for retrieving the relevant data that is based on data mining as per the user's request.
- Data mining engine: Major component of any data mining system. Contains several modules for executing data mining tasks, including association, correlation, classification, clustering, prediction, large-series,

analysis, etc.

> Pattern Evaluation module: The pattern evaluation module is primarily responsible for the measure of investigation of the patterns by using a threshold value. It collaborates with the data mining engine to focus the search on existing patterns.

It might utilize a static threshold to filter out discovered patterns.

> Graphical User Interface: The graphical user interface (GUI) module communicate b/w the data mining system and the users. This module helps the user to easily and efficiently use the system without knowing the complexity of the process. Cooperates with data mining system when user specifies a query or a task & displays the results.

> Knowledge Base: Helpful in the entire process of data mining.

Contain user view and data from user experiences that might be helpful in the sensible, data mining process. The pattern assessment module regularly interacts with knowledge base to get inputs & also update it.

10.

At last.

11.

Discuss the types of attributes & basic statistical descriptor of data in detail.

Types of attributes:

> Data objects: comprise to form data sets. Represents entity.

Typically described by attributes. If stored in DB, the data objects are referred to as data tuples.

> Attribute types: represent characteristic or feature of data object.

- Nominal attribute: quantitative attribute associated to names.

Nominal attributes are referred as categorical attributes and there is no order (rank position) among values of nominal attribute.

eg:-	Own House	:	1. Yes	2. No
	Marital status	:	1. Unmarried	2. Married



- Binary Attribute: quantitative attribute. Data has only 2 values / states
for e.g.: - Yes or no, true or false.

Symmetric Binary attribute: Both values are equally important
e.g.: Gender

Assymmetric Binary attribute: not equally important (e.g. Result)

Example:	Gender	Male, Female
	Cancer Detected:	Yes, No
	Result	Pass, Fail

- Ordinal attribute: Quantitative attribute

meaningful sequence or ranking (order) b/w them, but the magnitude b/w values is not actually known.

Example:	Grade:	A, B, C, D, E, F, O
	Income:	Low, medium, High
	Product rating:	0, 1, 2, 3, 4, 5

- Numerical attribute:

quantitative because, it is measurable quantity, represent in integers or real values. Types → real interval or ratio scale.

→ interval - data can be added or subtracted but cannot be multiplied or divided. (Don't have correct reference point / zero point)

→ ratio-scaled attribute - fix zero-point (can be multiplied), compare difference b/w values, mean, median, mode, quantile-range, five no.

Summary can be given.

- Discrete attribute: quantitative attribute, can be numerical and also be in categorical form; have finite values.

Example:	Profession:	Principal, Teacher, Clerk
	Zip code:	400050, 400051, 400052

- Continuous attribute: quantitative, can take values b/w two specified values

Example:	Height:	5.2, 5.4, 5.6, ..
	Weight:	50.33, ..



> Basic statistical descriptors of data

essential to have an overall picture of the data, if data processing to be made successful.

following are the basic statistical descriptors of data.

⇒ measures of central Tendency: used to represent centre of set data.

- mean (\bar{x}): mean of n no. sum of no. divided by n.

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

for weighted arithmetic mean

$$\bar{x} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i} = \frac{w_1 x_1 + w_2 x_2 + \dots + w_n x_n}{w_1 + w_2 + \dots + w_n}$$

- median: median of n no. is the middle no. when no. are written in order. If n is even, median is avg/mean of two middle no.

$$\text{median} = L_1 + \left[\frac{n/2 - (\sum \text{freq})_1}{\text{freq}_\text{median}} \right] \text{width}$$

$L_1 \rightarrow$ lower boundary of median interval; $n \rightarrow$ no. of values in entire set

$(\sum \text{freq})_1 \rightarrow$ sum of freq.; $\text{freq}_\text{median} \rightarrow$ median interval width

- mode: mode of n no. is the no. or nos. that occurs most frequently. may be one mode, no mode, or more than one mode.

for unimode, numeric data are g. asymmetric.

$$\text{mean} - \text{median} \approx 3(\text{mean} - \text{mode})$$

- midrange: average of smallest & largest value in the set.

⇒ Dispersion of data: statistic that tells you how dispersed, or spread out, data values are.

- quartiles: values that divides your values into quarters

As quantiles divides no. up according to where their position is on no. line, you have to put the no. in order before you can fig. out which quantile it is.

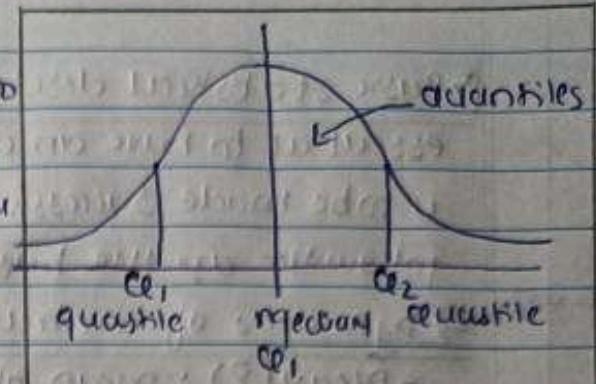


Fig. Dispersion of data.

- Range: difference b/w the upper and lower quartile value in set of data, commonly referred as IQR

$$\text{IQR} = Q_3 - Q_1$$

- Five no. summary: Rough data about your data set looks like. includes five items: minimum value, Q_1 , (first quartile), median, third quartile (Q_3), and the maximum value.
- Boxplot: (or whisker plot) defined as graphical method of displaying variation in set of data. Incorporates the five-summary as follows.

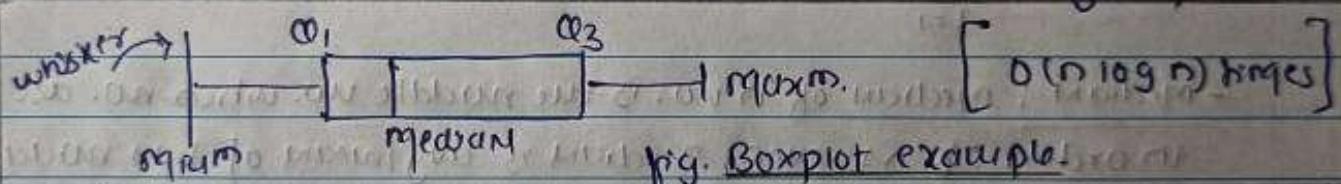


Fig. Boxplot example.

- Outliers:

value higher or lower than $1.5 \times \text{IQR}$ (Inter-Quartile range)

- Variance and standard deviation: measures of data dispersion

For data set x_1, x_2, \dots, x_n variance is calculated as

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \left(\frac{1}{n} \sum_{i=1}^n x_i^2 \right) - (\bar{x})^2$$

The standard deviation is square root of variance σ^2 .

when all observation have same values, $\sigma=0$; otherwise $\sigma>0$.

\Rightarrow Graphic displays of basic statistical descriptors of data:

Helpful for visual inspection of data, useful for data preprocessing.

Includes quantile, quantile-quantile plots, histograms & scatter plots.

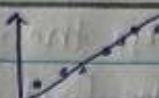
- Quantile plot: graphical way of checking whether your data are normally distributed i.e. unimodal distribution.



we compute for each data value, the normal quantile value as

$$z = \frac{x_i - \bar{x}}{s}$$

rough
Fig:-



→ Normal quantiles

- quantile-quantile plot: plots of two quantiles against each other.

part where certain values fall below that quantile.

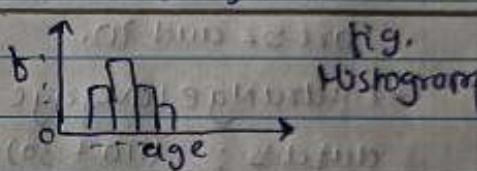
purpose of Q-Q plots is to find out if two sets of data comes from the same distribution. [comparison]

If two data sets come from common distribution, points will fall on single reference line.

- Histograms: usually shows the distribution of values of single variable divides values into bins & shows a bar plot of no. of object in each bin.

height of each bar indicates no. of objects.

shape depend on no. of bins.



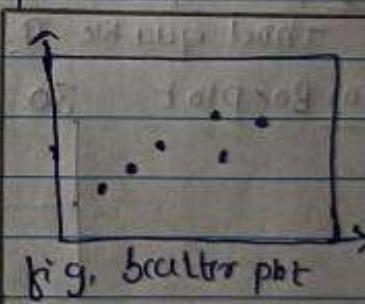
- Scatter plot: determining if there is relationship, pattern, or trend existing b/w two numeric attributes.

Bivariate distribution of data

Comparison of two attributes

(correlation), amount of data in cluster

correlation can be, +ve, -ve or zero.



10. college wants to record the marks for the courses completed by students using the dimensions. i) course, ii) student, iii) time & a measure Aggregate marks. Create a cube and performing following OLAP operation : i) Rollup ii) Drill down, iii) slice, iv) dice, v) pivot.
pg.no.

12. suppose that the data for analysis includes the attributes salary (in thousand of dollars), shown in increasing order.

30, 36, 47, 50, 52, 56, 60, 63, 70, 70, 110.

i. what are the mean, median, mode & midrange of data?

ii. find the first quartile (Q_1) & third quartile (Q_3) of data.

iii. Show the boxplot of the data.

i) Mean = $\bar{x} = \frac{\sum x_i}{n} = \frac{696}{12} = 58$

- median (mean of two middle values of ordered set, as the no. of values in the set is even) of the data is 54.

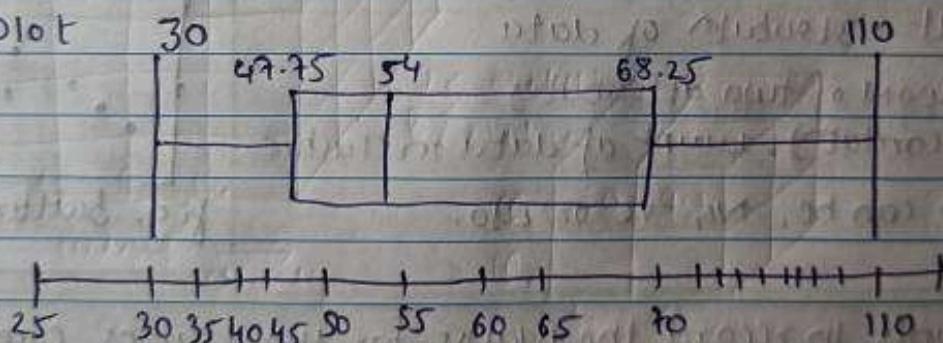
- This value set has two values that occurs with the same highest frequency and is, therefore, bimodal.

Other modes (values occurring with greatest frequency) of the data are 52 and 70.

- Midrange (average of largest & smallest value in dataset) of the data is: $(110 + 30)/2 \Rightarrow 70$.

ii) First quartile (Q_1 , corresponding to 25th percentile) of data is: 47.75
Third quartile (Q_3 , — 1 - 75th percentile) of data is: 68.25

iii) Boxplot



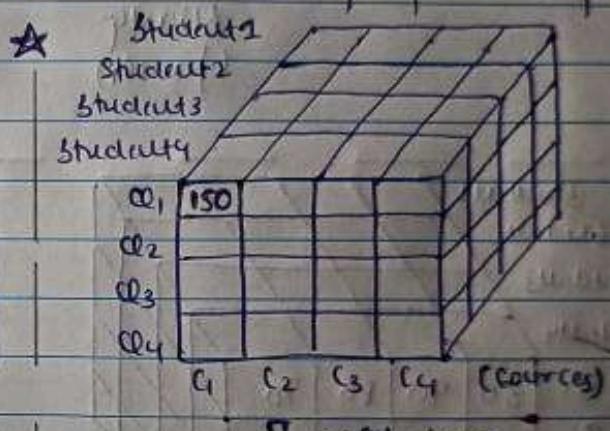
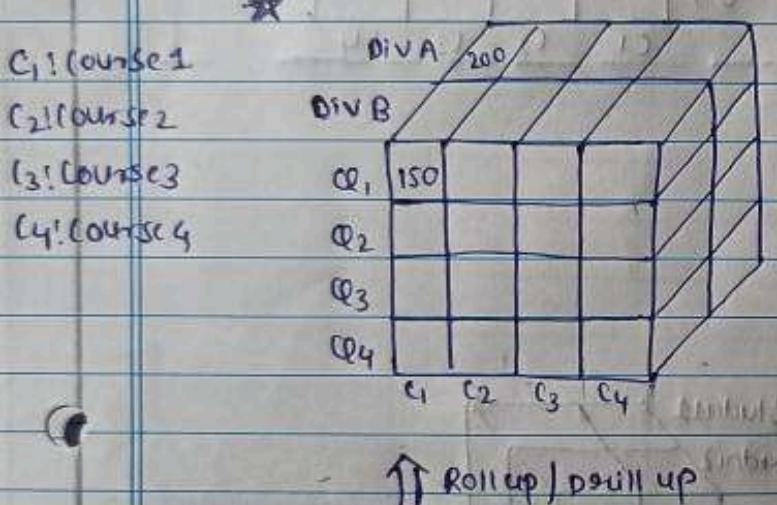
Boxplot for salary attribute.

OLAP operators:

1. slice: for single student get the marks for all courses and all semester is slice operatⁿ as cube is generated with respect to one dimension Student-Dim
2. Dice: Get the marks of single (or multiple) students for one (or more)

Semester but for all courses in dice operation an cube is generated with respect to two dimensions i.e. Student-Dim and Time-Dim.

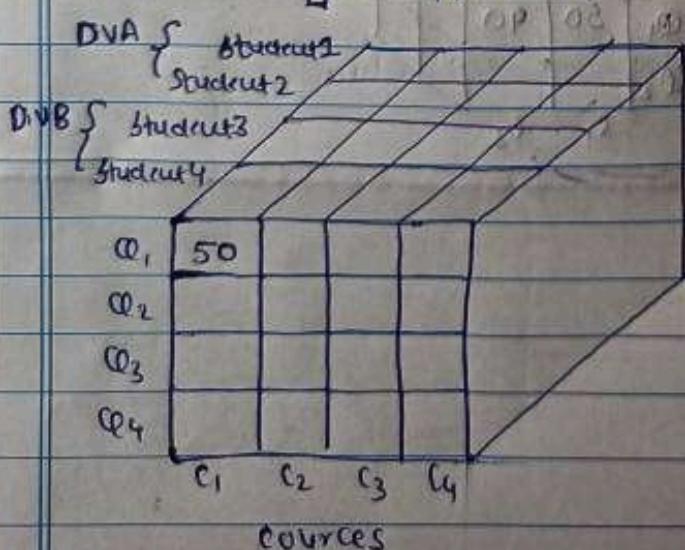
3. Roll up: Get the year wise marks for all students is roll up as semester wise marks are added to get year wise marks and its one level up hierarchy of dimension : Time-Dim.
4. Drill Down: Get the marks based on second_id is drill down operation as its one level down hierarchy for Course-Dim, discussion.
5. Pivot: Getting the month wise reports or course wise reports is pivot operation. In pivot one can view or rotate the information as per requirement.



↓ Drill down



↑ Roll up / Drill up





Star

student1
student2
student3
student4

	C1	C2	C3	C4	
Q1	50	90	80	70	b1
Q2					b2
Q3					b3
Q4					b4
courses					
SPLITTING (for Q1)					
	C1	C2	C3	C4	

Star

student1
student2
student3
student4

	C1	C2	C3	C4	
Q1	50	90	80	70	
Q2					
Q3					
Q4					
courses					
CHUNK					
	C1	C2			

DIVISIONS

A) student1
B) student2

Q1	50	90		
Q2				
	C1	C2		



Name :- Prathamesh S. Chikankar

ROLL NO. :- AIML11 BRANCH :- CSE-(AI & ML)

YEAR :- TE SUBJECT :- DWM

Topic :- Assignment No. 02

Sign :- (Prathy)

Date :- October '22



1. Explain Market Basket Analysis with an example.

Market Basket Analysis is a data mining technique that is used to uncover patterns in any retail setting.

- The goal of Market Basket analysis is to understand consumer behaviour by identifying relationships b/w the items that people buy.
- A co-occurrence is when two or more things take place together.
- Creates If-then scenario rules, for e.g., if item A is purchased then item B is likely to be purchased.
- The rules could be written as : If $\{A\}$ then $\{B\}$
- A association rule has three measures that express degree of confidence in the rule. They are support, confidence and lift.

> The support is the no. of transact that include items in the $\{A\}$ and $\{B\}$ pairs of the rule as percentage of total no. of transact.
 Also called as occurrence frequency, frequency, support ratio or count.

$$\text{support } (A \rightarrow B) = \frac{\text{No. of transact containing both A and B}}{\text{Total no. of transact.}}$$

> The confidence of the rule is the ratio of no. of transact that include all items in $\{B\}$ as well as the $\{BA\}$ to the no. of $\{A\}$.

$$\text{confidence } (A \rightarrow B) = \frac{\text{no. of transact containing both A and B}}{\text{No. of transact containing A}}$$

> The lift or lift ratio, is the ratio of confidence to expected confidence. Expected confidence is the confidence divided by the frequency of B. Greater lift values indicate stronger association.

$$\text{lift } (A \rightarrow B) = \frac{\text{no. of transact containing both A & B}}{\frac{\text{no. of transact containing B}}{\text{total no. of transact.}}}$$

Example: Consider there are nine laptops containing varying combinations of Asus, Acer, MSI, and Microsoft.



Laptop	Product 1	Product 2	Product 3
1	Acer	ASUS	
2	Acer	MSI	ASUS
3	MSI	MICROSOFT	
4	Acer	ASUS	
5	MSI	MICROSOFT	
6	Acer	ASUS	MICROSOFT
7	Acer	ASUS	
8	ASUS	MICROSOFT	
9	ASUS	Acer	

> support ($A \rightarrow B$)

$$= \frac{\text{no. of laptops containing both Acer \& ASUS}}{\text{Total no. of laptops}}$$

$$= \frac{6}{9} \Rightarrow 0.67$$

> confidence ($A \rightarrow B$)

$$= \frac{\text{no. of laptops containing both Acer and ASUS}}{\text{no. of laptops containing Acer}}$$

$$= \frac{6}{6} \Rightarrow 1.00$$

> lift ($A \rightarrow B$) =

$$= \frac{\text{no. of laptops containing both Acer and ASUS}}{\text{no. of laptops containing Acer}} \times \frac{\text{no. of laptops containing ASUS}}{\text{no. of laptops containing both Acer and ASUS}}$$

$$\times \frac{\text{no. of laptops containing both Acer and ASUS}}{\text{no. of laptops containing both Acer and ASUS}}$$

$$\times \frac{\text{no. of laptops containing both Acer and ASUS}}{\text{no. of laptops containing both Acer and ASUS}}$$

$$= \frac{6/6}{7/9} = 1.29$$

Applications of market basket analysis:

- Retail, • Telecommunications, • Banks, • Insurance & • medical.



Q. 2.

Discuss Association Rule mining and Apriori Algorithm

Association rule mining finds interesting association and support among large sets of data items.

This rule shows how frequently a itemset occurs in a transact.

A typical example is market Basket Analysis.

An implict expression of the form $x \rightarrow y$, where x and y are any 2 itemsets.

Example : {Milk, Diaper} \rightarrow {Beer}

This rule must satisfy the confidence value.

Rule Evaluation Metrics -

γ support (s) - γ confidence (c)

γ support = $\frac{s}{(x+y)} + \text{total}$ γ lift (l)

γ conf ($x \Rightarrow y$) = supp ($x \cup y$) \div supp (x)

γ lift ($x \Rightarrow y$) = conf ($x \Rightarrow y$) \div supp (y)

The association rule is very useful in analyzing datasets.

Apriori algorithm was the first algorithm that was proposed for frequent itemset mining.

It is an iterative approach to discover the most frequent itemsets.

This algorithm uses two steps "join" and "prune" to reduce the search space.

γ Join step: this step generates $(k+1)$ itemset from k -itemsets by joining each item with itself.

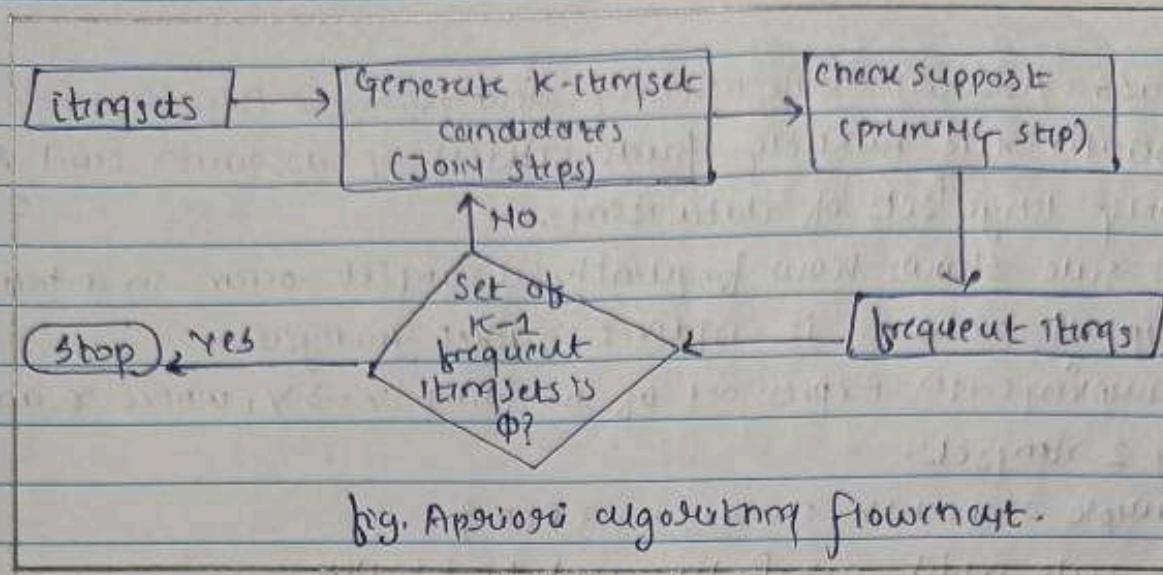
γ Prune step: this step scans the count of each item in the database.

- Advantages: Easy to understand.

Join & prune steps are easy to implement on large itemsets in large DB.

- Disadvantages: Entire database needs to be scanned.

Requires high computation if the itemsets are very large and minimum support is kept very low.



3. A database has four transactions. Let minimum support = 50% and minimum confidence = 50%.

Transaction_id	Items bought
T100	A, B, C
T200	A, C
T300	A, D
T400	B, E, F

Find all frequent item sets using apriori algorithm.

list strong- association rules

The given minimum support = 50%.

The no. of transactions = 4

$$\therefore \text{minimum support} = 0.5 \times 4 = 2$$

Step 1: Calculate the minimum count of each item.

∴ Items	minimum count
A	3
B	2
C	2
D	1
E	1
F	1

Step 2 - Delete the items that do not have a minimum support count of 2.

∴ Items D, E and F are deleted.

Items	minimum cost
A	3
B	2
C	2

Step 3 - Combine 2-items and find out the minimum cost count of the occurrences of the 2-items.

Items	minimum cost
A, B	1
A, C	2
B, C	1

Now, only got the item-set {A, C} that is frequent.

Step 4 - Generate Association Rules from the frequent itemset discovered in the above Step 3.

Rule 1 - $\{A\} \Rightarrow \{C\}$

$$\text{Confidence} = \frac{\text{Support } \{A, C\}}{\text{Support } \{A\}} = \frac{2}{3} \times 100 = 66.67\%$$

Rule 2 - $\{C\} \Rightarrow \{A\}$

$$\text{Confidence} = \frac{\text{Support } \{A, C\}}{\text{Support } \{C\}} = \frac{2}{2} \times 100 = 100\%$$

The minimum given confidence = 50%.

∴ This shows that all the above association rules are strong.

Apriori is an algorithm for frequent item set mining and association rule learning over the given dataset.

It works by identifying the frequent individual items in the

dataset and executing them to large and large item sets as long as those item sets appear sufficiently often in the dataset.

Algorithm: Apriori. Find frequent itemsets using an iterative level-wise approach based on candidate generation.

Input:

D: a database of transaction;

min-sup: the minimum support count threshold.

Output: L: frequent itemsets in D.

Procedure:

1) $L_1 = \text{find-frequent-1-itemset}(D)$;

2) for ($k=2$) $L_{k-1} \neq \emptyset$; $k++$ {

3) $C_k = \text{apriori-gen}(L_{k-1})$;

4) for each transaction $t \in D$ {

5) $C_t = \text{subset}(C_k, t)$;

6) for each candidate $c \in C_t$

7) $c.\text{count}++$;

8) }

9) $L_k = \{c \in C_k \mid c.\text{count} \geq \text{min-sup}\}$

10) }

11) return $L = \cup_k L_k$;

procedure apriori-gen (L_{k-1} : frequent $(k-1)$ -itemsets)

for each itemset $l_1 \in L_{k-1}$

for each itemset $l_2 \in L_{k-1}$

if $((l_1[1] = l_2[1]) \wedge (l_1[2] = l_2[2]))$

$\wedge \dots \wedge (l_1[k-2] = l_2[k-2]) \wedge (l_1[k-1] \subset l_2[k-1])$

then {

$c = l_1 \wedge l_2$

if not frequent subset (c, L_{k-1}) then

delete c

```

else add c to Cx
}
return Cx;
procedure has-frequent-subset (c : scanable k-itemset;
                               Lk-1 : frequent (k-1)-itemsets);
for each (k-1)-subset s of c
  if s ∈ Lk-1 then
    return TRUE;
return FALSE;

```

- Q.4. A database has five transactn. Let num-sup = 60% & num-cost = 80%. Find all the frequent itemsets using F-P growth.

TID	items_bought	ordered-itemsets
T100	{M, O, N, K, E, Y}	K, E, M, O, Y
T200	{D, O, N, K, E, Y}	K, E, O, Y
T300	{M, A, K, E}	K, E, M
T400	{M, U, K, K, Y}	K, M, Y
T500	{C, O, O, K, I, E}	K, E, O

Given : num-Sup = 60%

$$\therefore \text{Sup-count to be satisfied} = 5 \times 0.6 = 3$$

Step 1: Scan the databases for count of each itemset.

Itemset	Sup-count	Itemset	Sup-count
{A}	1	{M}	3
{C}	2	{N}	2
{D}	1	{O}	4
{E}	4	{U}	1
{I}	1	{Y}	3
{K}	5		



Step 2: Sort the set of frequent itemsets in the order of descending support count and denote that list as L.

L:	Itemset	Sup-Count
	{K}	5
	{E}	4
	{O}	4
	{M}	3
	{Y}	3

Step 3: Scan the database for second hinge and sort items in each transaction according to descending support count

i.e. ordered-itemset

TID	List of Items
T100	{K, E, M, O, Y}
T200	{K, E, O, Y}
T300	{K, E, M}
T400	{K, M, Y}
T500	{Y, E, O}

Step 4: construct the F-P tree

4.1. Create a node with label 'NULL':

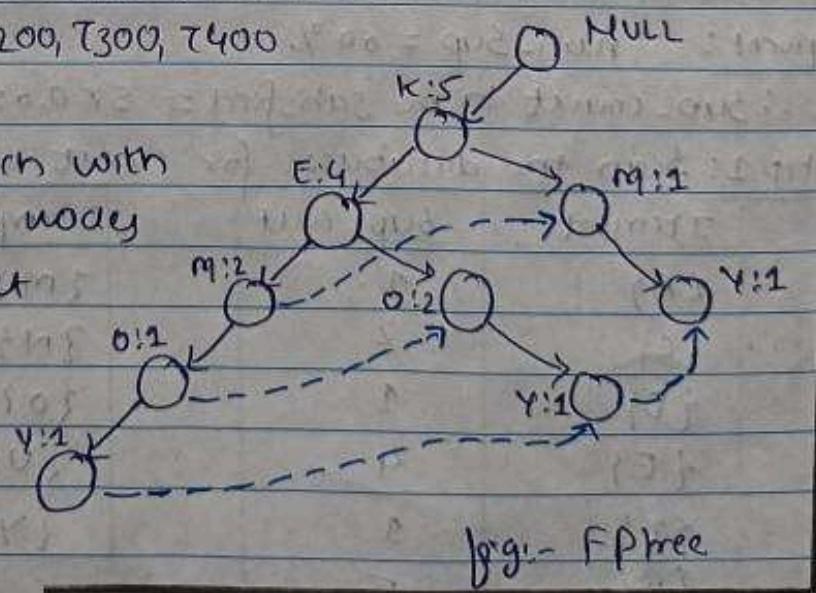
4.2. Scan T100, T200, T300, T400

T500.

4.3. Construct Branch with their respective nodes

4.4. Also increment its count by 1

4.5. connect all similar nodes.



Step 5: mining FP-tree

itemset	conditional pattern base	conditional FP-tree	frequent patterns generated
{Y}	{K, E, M, O:1}, {K, M:1}	{K:3}	{K, Y:3}
{}	{K, E, O:1}		
{O}	{K, E, M:1}, {K, E:2}	{K:3, E:3}	{K, O:3}, {E, O:3}, {K, E, O:3}
{M}	{K, E:2}, {X:1}	{X:3}	{K, M:3}
{E}	{X:4}	{K:4}	{K, E:4}
{X}	-	-	-

FP tree

The frequent pattern tree is a tree-like structure that is made with the frequent itemsets of the dataset.

- The purpose of the FP tree is to mine the most frequent data patterns.
- Each node of the FP tree represents an item of the itemset.

Algorithm!

FP growth: mine frequent itemsets using FP tree by pattern fragmentation growth.

Input:

D, a transaction DB;

min-sup, the minimum support count threshold.

Output:

complete set of frequent patterns.

Method:

Scan

insert-tree ([P1P], T)

insert-tree (P, N)

FP-growth (FP tree, min)

procedure FP-growth (tree, d)



If tree contains a single path P then

for each combination

else for each ai in the header of tree Z

:

If Tree B = \emptyset then

call FP-growth (tree, B)) }

Q.5. Demonstrate multidimensional & multilevel association rule mining with suitable examples.

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

> Support and confidence of multilevel association rules.

Support of rules increases specialized and decreases generalized itemsets.

Confidence - Support is below threshold value, rule becomes invalid.

> Approaches of multilevel association rules:

-using uniform support level for all levels.

If the support threshold is too low, it generates too many high level association.

		Priority (support = 10%)	fig. with minimum support	
Level 1				
min-sup = 5%				
Level 2	Laserjet priority		Dot matrix priority	
min-sup = 5%		(support = 6%)		(support = 10%)

-using reduced minimum support at lower level.

Level 1		Priority (support = 10%)	fig. with reduced support	
min-sup = 5%				
Level 2		Laserjet priority (support = 8%)	Dot matrix priority (support = 10%)	
min-sup = 3%				

at every abstraction level, there is its own minimum support threshold; so minimum support at lower levels reduces.

-using item or group-based minimum support.

Multidimensional association rules. (clustering)

> Single-dimensional rule : buys(x,"milk") \rightarrow buy(y, "bread")

> multi-dimensional rule : contain more than one predicate.

> categorical attribute : finite no. of values, eg:- color.

> quantitative attribute : numeric values, eg:- age

Techniques : can be categorical or quantitative

- Static Discretization of quantitative attributes.

We refer to this as strong multidimensional association rules using static discretization of quantitative attributes.

The discretized numerical

attribute, with their

interval labels, can then

be treated as categorical

attributes, where each

interval is category consider-

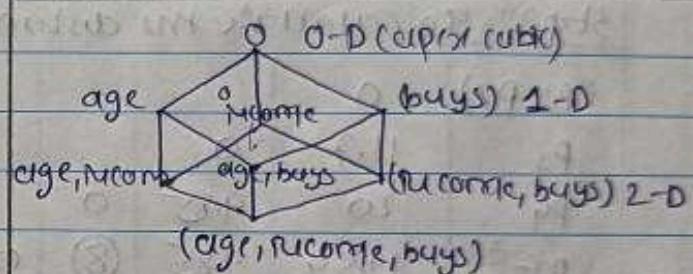
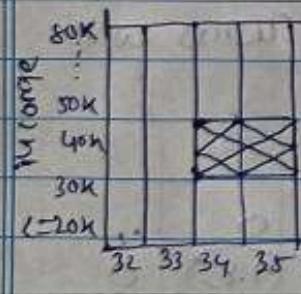


Fig. Lattice of cuboids to form 3-D data cube

- Dynamic quantitative association rules



The strong association rules

Fig. A 2-D obtained by mapped a 2-D good for good as shown.

tuples representing this b&m further combining customers purchasing SUV

Q. 6. Differentiate b/w simple linkage, average linkage and complete linkage algorithms. Use a complete linkage algorithm to find the clusters from the following dataset.

X	4	8	15	24	324
Y	4	4	8	4	12

Step 1: Calculate distance from each object (point) using Euclidean distance



$$\text{Distance } [(x, y), (a, b)] = \sqrt{(x-a)^2 + (y-b)^2}$$

Distance matrix:

P1	0			
P2	4	0		
P3	11.7	8.06	0	
P4	20	16	9.85	0
P5	21.54	17.89	9.85	8
	P1	P2	P3	P4
				P5

P1 and P2 are two clusters with shortest distance 4
merge P1 and P2

Step 2: To calculate the distance of P3 from (P1, P2):

(P1, P2)	0			dist((P1, P2), P3) = \max(\text{dist}(P1, P3),
P3	11.7	0		\text{dist}(P2, P3))
P4	20	9.85	0	= \max(11.7, 8.06)
P5	21.54	9.85	8	max 11.7
(P1, P2)	P3	P4	P5	

Step 3: In above matrix, P4 and P5 are two clusters with shortest distance 8, merge P4 and P5

Recompute the

distance matrix

(P1, P2)	0		
P3	11.7	0	
(P4, P5)	21.54	9.85	0
(P1, P2)	(P3)	(P4, P5)	

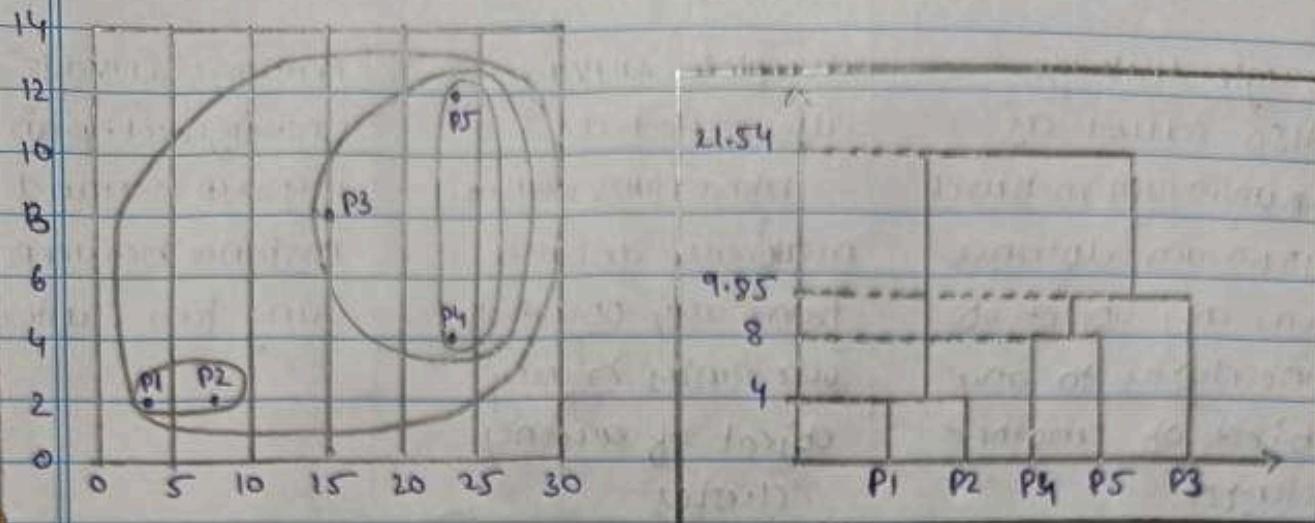
Step 4:

P3 and (P4, P5) are two clusters with shortest distance 9.85
merging

Recompute the

distance matrix

(P1, P2)	0		
(P3, P4, P5)	21.54	0	
(P1, P2)	(P3, P4, P5)		

Fig. Dendrogram

Step 5: Looking at the above distance matrix in Step 4, we see that (P1, P2) and (P3, P4, P5) have the smallest distance 21.54 (only one left) so, we merge those two single vertex clusters. No need to re-compute. As there are no more clusters to merge.

Hierarchical clustering

↳ produces a set of nested clusters organized as hierarchical tree
↳ can be visualized as dendrogram - a tree-like diagram that records the sequence of merges or splits.

- agglomerative
- divisive

Algorithm:

1. compute the distance matrix b/w input data points
2. let each data point be a cluster
3. Repeat
4. merge the two closest clusters
5. update the distance matrix
6. until only a single cluster remains.

Advantage: no. prior info. required about clusters.

Disadvantage: can never undo what was done previously.



	<u>single linkage</u>	<u>complete-linkage</u>	<u>average-linkage</u>
>	Also called as minimum method	also called as maximum method	unweighted-pair group method
>	minimum distance from any object of one cluster to any object of another cluster.	maximum distance from any object of one cluster to any object of another cluster	distance between any two clusters
>	D(A,B) computed as $D(A,B) = \min \{ d(i,j) : i \in \text{cluster A}, j \in \text{cluster B} \}$	$D(A,B) = \max \{ d(i,j) : i \in \text{cluster A}, j \in \text{cluster B} \}$	$D(A,B) = \text{mean} \{ d(i,j) : i \in \text{cluster A}, j \in \text{cluster B} \}$
>	'loose' clusters	"tight" clusters	"tight" clusters
>	can handle non- elliptical shape	less susceptible to noise and outliers	less susceptible to noise & outliers
>	sensitive to noise and outliers	tends to break large clusters	biased towards globular clusters
>			

Q7. Explain Naïves Bayesian Classification in detail.

Naïves Bayes classifiers are a collection of classification algorithms based on 'Bayes' theorem. It is not a single algorithm but a family of algorithms, where all of them share a common principle, i.e. every pair of features being classified is independent of each other. To start with, let us consider a dataset. Consider a fictitious dataset that describes the weather, condition for playing a game of golf. Given the weather condition, each tuple classifies that condition as



fit ("yes") or unfit ("no") for playing golf.

classify using Naïve Bayes classifier

outlook	sunny	Humidity	windy	play
Rainy	cool	high	true	?

Bayes theorem: find probability of an event occurring given the probability of another event that has already occurred.

Stated mathematically as:

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

... event B is true i.e. evidence

... event A prior

regards to our dataset

$$P(y|x) = \frac{P(x|y) \cdot P(y)}{P(x)}$$

... $y \rightarrow$ class variable

... $x \rightarrow$ dependent feature

Naïve assumption

$$P(y|x_1, \dots, x_n) = \frac{P(x_1|y) \cdot P(x_2|y) \cdots P(x_n|y) \cdot P(y)}{P(x_1) \cdot P(x_2) \cdots P(x_n)}$$

(x_1, x_2, \dots, x_n)

also mathematically

$$y = \arg \max_y P(y) \prod_{i=1}^n P(x_i|y)$$

to calculate $P(y)$ & $P(x_i|y)$... class prob. & conditional prob.

Steps:-

importing library

importing datasets

Reading data

splitting data

creating model

fitting model

making prediction

Evaluating result

Apply the Naïve Bayes classifier algorithm to classify an unknown sample X (outlook=sunny, temperature=cool, humidity=high,



windy = false) the sample dataset by follows:

outlook	Temperature	Humidity	windy	class
sunny	Hot	High	false	N
sunny	Hot	High	true	N
overcast	Hot	High	false	P
Rain	Mild	High	false	P
Rain	Cool	Normal	false	P
Rain	Cool	Normal	true	N
overcast	Cool	Normal	true	P
sunny	Mild	High	false	N
sunny	Cool	Normal	false	P
Rain	Mild	Normal	false	P
sunny	Mild	Normal	true	P
overcast	Mild	High	true	P
overcast	Hot	Normal	false	P
Rain	Mild	High	true	P

~~$P_c(\text{class} = \text{Yes}) = 10/14 = 0.714$~~

~~$P_c(\text{class} = \text{No}) = 5/14 = 0.357$~~

outlook	Y	N	Humidity	Y	N
sunny	2/9	3/5	High	3/9	4/5
overcast	4/9	0	Normal	6/9	1/5
rain	3/9	2/5			

Temperature	windy
Hot	Strong
Mild	Weak
Cool	

(outlook = sunny, Temperature = cool, humidity = High)

wind = true)

wind = false)

Soln:-

Consider the tuple X to classify.

$$x = [\text{outlook} = \text{sunny}, \text{temperature} = \text{cool}, \text{humidity} = \text{high}, \text{windy} = \text{false}]$$

The data tuple describes by the attribute look(outlook), temperature, humidity, windy

The class label class which has 2 distinct values P and N

c_1 corresponds to the class P & c_2 corresponds to the class N

$$x = [\text{outlook} = \text{sunny}, \text{temperature} = \text{cool}, \text{humidity} = \text{high}, \text{windy} = \text{false}]$$

We need to calculate posterior probability

$$P(X|c_i) P(c_i) \text{ for } i=1, 2$$

$$P(\text{class} = P) = 10/14 = 0.174$$

$$P(\text{class} = N) = 4/14 = 0.285$$

To compute $P(X|c_i)$ for $i=1, 2$, we compute the following conditional probability.

$$P(\text{outlook} = \text{sunny} | \text{class} = P) = 2/10 = 0.2$$

$$P(\text{outlook} = \text{sunny} | \text{class} = N) = 3/4 = 0.75$$

$$P(\text{temperature} = \text{cool} | \text{class} = P) = 3/10 = 0.3$$

$$P(\text{temperature} = \text{cool} | \text{class} = N) = 1/4 = 0.25$$

$$P(\text{humidity} = \text{high} | \text{class} = P) = 4/10 = 0.4$$

$$P(\text{humidity} = \text{high} | \text{class} = N) = 3/4 = 0.75$$

$$P(\text{windy} = \text{false} | \text{class} = P) = 6/10 = 0.6$$

$$P(\text{windy} = \text{false} | \text{class} = N) = 2/4 = 0.5$$

$$P(X | \text{class } P) = P(\text{outlook} = \text{sunny} | \text{class} = P) \times P(\text{temperature} = \text{cool} | \text{class} = P) \times$$

$$P(\text{humidity} = \text{high} | \text{class} = P) \times P(\text{windy} = \text{false} | \text{class} = P)$$

$$= 0.2 \times 0.3 \times 0.4 \times 0.6 \Rightarrow 0.0144$$



$$\begin{aligned}
 P(X | \text{class} = N) &= p(\text{outlook} = \text{sunny} | \text{class} = N) \times p(\text{temperature} = \text{cool} | \text{class} = N) \\
 &\quad p(\text{humidity} = \text{high} | \text{class} = N) \times p(\text{windy} = \text{false} | \text{class} = N) \\
 &= 0.25 \times 0.25 \times 0.75 \times 0.5 \Rightarrow 0.09375
 \end{aligned}$$

To find the class with the maximum $P(X|C)P(C)$ we compute

$$\begin{aligned}
 P(X | \text{class} = P)P(\text{class} = P) &= 0.0144 \times 0.0714 \Rightarrow 0.01028 \\
 P(X | \text{class} = N)P(\text{class} = N) &= 0.09375 \times 0.285 \Rightarrow 0.02600
 \end{aligned}$$

The Naive Bayes' classifier algorithm predicts class N for tuple X.

Q. 8. Explain K-means clustering algorithm.

- > **K-Means:** Each cluster is represented by the centre of cluster.
- > K-means clustering is simple unsupervised learning algorithm
- > K-means tries to put points or data points into the set of K clusters where each data point is assigned to its closest cluster.
- > This method is defined by the objective function which tries to minimize the sum of all squared distances within a cluster, for all the clusters.
- > Objective function: $\text{args min } \sum_{i=1}^n \sum_{j \in S_i} \|x_j - \mu_i\|^2$ $\because x_i = \text{data point}$
 $S_i = \text{cluster}$
 $\mu_i = \text{cluster mean}$

K-means clustering Algorithm

Algorithm: K-means. \rightarrow Each cluster is represented by mean value.

Input: K: the no. of clusters.

D: dataset consisting n objects

Output: A set of K clusters.

Method:

- ① Arbitrarily choose K objects from D as initial clusters' centres,
- ② Repeat
- ③ (re) assign each object to the cluster to which the object is most



function, based on the mean value for the objects in the clusters.

- (4) update the cluster means, calculate mean value of object
- (5) until no changes.

Explain K-Medoids clustering algorithm.

- > K-medoids (also called as partitional based method) algorithm.
 - > A medoid can be defined as the point in the cluster, whose dissimilarity with all the other points in the cluster is minimum.
 - > The dissimilarity of the medoid (c_i) and object (p_i) is calculated by using $E = |p_i - c_i|$
 - > The cost in k-medoids algorithm given by
- $$C = \sum_{c_i} \sum_{p_i \in c_i} |p_i - c_i|$$

- > Algorithm : K-medoids. PAM .. based on medoid or central, objects.

Input: K: no. of clusters

D: dataset containing n objects

Output: A set of K clusters.

Method:

- ① arbitrary choose K objects in D as initial or preputative seeds.

- ② repeat

- ③ assign each object to nearest representative object.

- ④ randomly select non-representative object, Orandom.

- ⑤ compute totalcost, s, object oj with Orandom;

- ⑥ If $s < 0$ then swap o_j with Orandom to form new set of k

- ⑦ until no change;

- > Time complexity $O(K \times (n-k)^2)$

- > Advantage: fast K-medoid & less sensitive PAM

- > disadvantage: may obtain different results.

Find clusters using k-means clustering algorithm, if we have several objects (4 types of medicines) and each object have two attribute.

as features as shown in table below. The mission is to group these objects into $K=2$ group of medicine based on the two features (pH and weight index).

Object	Attribute 1 (x) weight index	Attribute 2 (y) pH
Medicine A	1	1
Medicine B	2	1
Medicine C	4	3
Medicine D	5	4

Soln:-

No. of clusters $K=2$

Initial cluster centres be $C_1 = (1,1)$ and $C_2 = (2,1)$

We will check distance b/w data points and all cluster centres.

We will use Euclidean distance formula for finding distance.

$$\text{Distance } (x, a) = \sqrt{(x-a)^2} \quad \text{OR}$$

$$\text{Distance } [(x,y), (a,b)] = \sqrt{(x-a)^2 + (y-b)^2}$$

As given data in pic, we will use second formula of Euclidean dist.

ITEM 1:

Following notation: D_1 = distance from cluster $C_1(1,1)$

D_2 = distance from cluster $C_2(2,1)$

- Data point $(1,1)$: $D_1 = 0$ and $D_2 = 1$

$$D_1 = [(1,1), (1,1)] = \sqrt{(1-1)^2 + (1-1)^2} = 0$$

$$D_2 = [(1,1), (2,1)] = \sqrt{(1-2)^2 + (1-1)^2} = 1$$

∴ smallest distance, so point $(1,1)$ belongs to cluster C_1

- Data point $(2,1)$: $D_1 = 1$ and $D_2 = 0$ ∴ $(2,1)$ belongs to C_2

- Data point $(4,3)$: $D_1 = 3.6$ and $D_2 = 2.83$ ∴ $(4,3)$ belongs to C_2

- Data point $(5,4)$: $D_1 = 5$ and $D_2 = 4.24$ ∴ $(5,4)$ belongs to C_2

Hence clusters are $C_1: \{A(1,1)\}$

$C_2: \{B(2,1), C(4,3), D(5,4)\}$

Accurating centre of clusty C_2 as $= \frac{(2+4+5)}{3}, \frac{(1+3+4)}{3} \Rightarrow (3.67, 2.67)$

7 Iteration 2:

D_1 = Distance from clusty $C_1(1,1)$

D_2 = Distance from clusty $C_2(3.67, 2.67)$

- data point $(1,1)$: $D_1 = 0$ & $D_2 = 3.14 \therefore (1,1)$ belongs to C_1
- data point $(2,1)$: $D_1 = 1$ & $D_2 = 2.75 \therefore (2,1)$ belongs to C_1
- data point $(4,3)$: $D_1 = 3.6$ & $D_2 = 0.47 \therefore (4,3)$ belongs to C_2
- data point $(5,4)$: $D_1 = 6$ & $D_2 = 1.88 \therefore (5,4)$ belongs to C_2

Accurating centre of clusty C_1 & C_2 i.e. $C_1 = (1.5, 1)$ & $C_2 = (4.5, 3.5)$

7 Iteration 3:

D_1 = Distance from clusty $C_1(1.5, 1)$

D_2 = Distance from clusty $C_2(4.5, 3.5)$

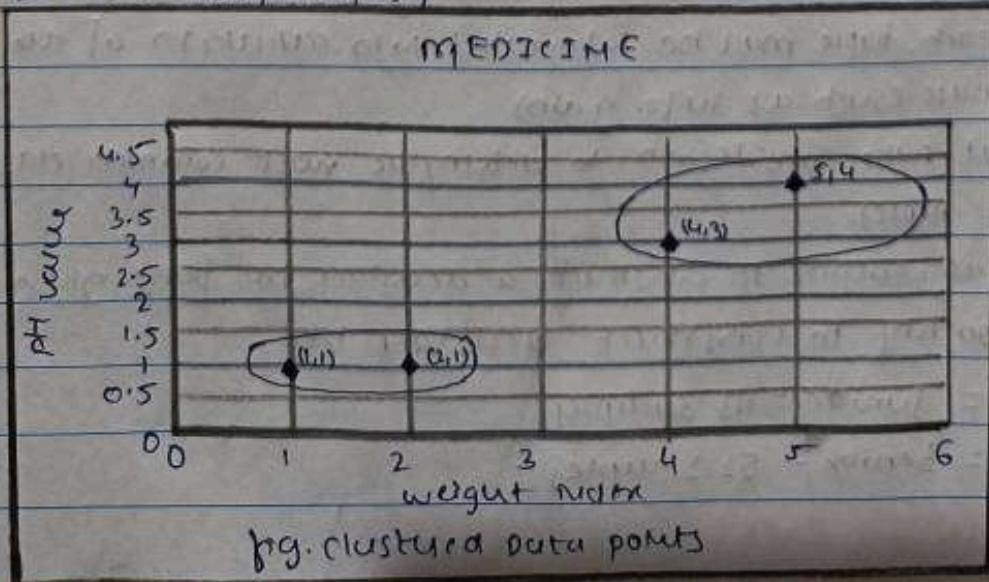
- data point $(1,1)$: $D_1 = 0.5$ & $D_2 = 4.3 \therefore (1,1)$ belongs to C_1
- data point $(2,1)$: $D_1 = 0.5$ & $D_2 = 3.54 \therefore (2,1)$ belongs to C_1
- data point $(4,3)$: $D_1 = 3.2$ & $D_2 = 0.71 \therefore (4,3)$ belongs to C_2
- data point $(5,4)$: $D_1 = 4.6$ & $D_2 = 0.71 \therefore (5,4)$ belongs to C_2

Accurating centre of clusty C_1 & C_2 $C_1 = C_2$:

Comparing iteration 2 & iteration 3

$C_1: \{A(1,1), B(2,1)\}$

$C_2: \{C(4,3), D(5,4)\}$



- q. The following table consists of training data from the employee DB. The data have been generalized. For example, "31...35" for age represents the age range of 31 to 35. For a given new entry, count implements the no. of data tuples having values for department, study, age and salary given in that row.

department	status	age	salary(brk)	count
sales	senior	21..35	40K..50K	30
sales	junior	26..30	26K..30K	40
sales	junior	31..35	31K..35K	40
systems	junior	26..35	40K..50K	20
systems	senior	31..35	60K..70K	5
systems	junior	26..30	40K..50K	3
systems	senior	41..45	60K..70K	3
marketing	junior	36..40	40K..50K	10
marketing	junior	31..35	41K..45K	7
secretary	senior	46..50	36K..40K	4
secretary	junior	26..30	26K..30K	6

a) How would you modify the basic decision tree algorithm to take into consideration the count of each generalized data tuple (new entry)?

The basic decision tree algorithm should be modified as follows to take into consideration the count of each generalized data tuple.

- count of each tuple must be integrated into calculation of the attribute selection measure (such as info. gain)
- Take count into consideration to determine most common class among the tuples.

b) Use your algorithm to construct a decision tree for given data.

ID3 algorithm to construct decision tree

C1 = study = junior = 13 samples

C2 = study = senior = 52 samples



$$\therefore P(C_1) = 113/165 \text{ and } P(C_2) = 52/165$$

(1) Entropy before split for guru database D:

$$H(D) = \sum_{i=1}^5 p_i \log_2 (p_i) \quad \therefore H(D) = \frac{113}{165} \log_2 \frac{165}{113} + \frac{52}{165} \log_2 \frac{165}{52}$$

$$\therefore H(D) = 0.3740 + 0.5250 \Rightarrow 0.899$$

(2) splitting attribute

> department as C1: status=junior & C2: status=senior with H.

$$\therefore H(\text{department}) = \frac{110}{165} \times H(\text{issues}) + \frac{31}{165} \times H(\text{systems}) + \\ \frac{145}{165} \times H(\text{marketing}) + \frac{10}{165} \times H(\text{salary})$$

$$= \frac{110}{165} \times 0.8454 + \frac{31}{165} \times 0.8238 + \frac{14}{165} \times 0.8631 + \frac{10}{165} \times 0.9709 = 0.8504$$

$$\therefore H(\text{department}) = 0.8504$$

$$\text{Gain}(\text{department}) = H(D) - H(\text{department}) = 0.899 - 0.8504 \Rightarrow 0.0486$$

> age as C1: status=junior & C2: status=senior with entropy H

$$\therefore H(\text{age}) = \frac{20}{165} \times 0 + \frac{49}{165} \times 0 + \frac{79}{165} \times 0.9906 + \frac{10}{165} \times 0 + \frac{3}{165} \times 0 + \frac{4}{165} \times 0 = 0.4743$$

$$\text{Gain}(\text{age}) = H(D) - H(\text{age}) = 0.899 - 0.4743 \Rightarrow 0.4247$$

> salary

$$\therefore H(\text{salary}) = \frac{46}{165} \times 0 + \frac{40}{165} \times 0 + \frac{4}{165} \times 0 + \frac{4}{165} \times 0 + \frac{63}{165} \times 0.9468 + \frac{8}{165} \times 0$$

$$H(\text{salary}) = 0.3615$$

$$\text{Gain}(\text{salary}) = H(D) - H(\text{salary}) = 0.899 - 0.3615 \Rightarrow 0.5375$$

Summary: Gain(department) = 0.0486

$$\text{Gain}(\text{age}) = 0.4247 \text{ and } \text{Gain}(\text{salary}) = 0.5375$$

salary has highest attribute gain

∴ it is used as decision tree having 6 possible values.

For the dataset D, the root node has 6 branches

All the branches are gain except only one branch,

i.e. 46K...50K is not giving a unique class label attribute

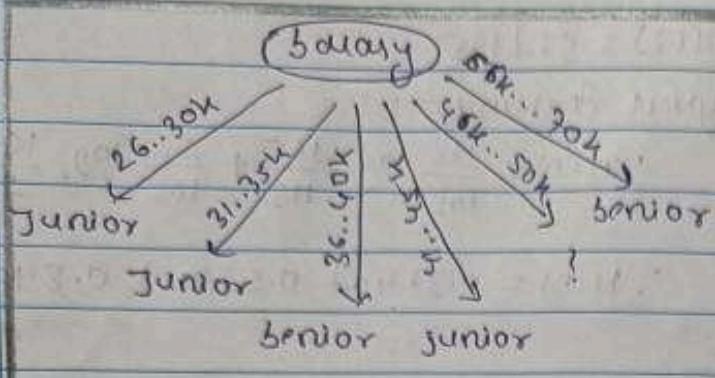


Fig. P.

Now consider salary 46k..50k and count the no. of samples from the original dataset D. let us denote it as D1.

Let the class label attribute be as follows:

$$C_1 = \text{status} = \text{junior} | \text{salary} = 46k..50k = 23 \text{ samples}$$

$$C_2 = \text{status} = \text{senior} | \text{salary} = 46k..50k = 40 \text{ samples}.$$

$$\therefore \frac{23}{63} \quad P(C_1) = 23/63 \text{ and } P(C_2) = 40/63$$

① Entropy before split

$$H(D_1) = \sum_{i=1}^2 p_i \log_2 \left(\frac{1}{p_i} \right) = \frac{23}{63} \log_2 \frac{63}{23} + \frac{40}{63} \log_2 \frac{63}{40} = 0.9468$$

② Splitting attribute

> department

$$H(\text{department}) = \frac{30}{63} \times 0 + \frac{23}{63} \times 0 + \frac{10}{63} \times 0 + \frac{0}{63} \times 0 = 0$$

$$\text{Gain}(\text{department}) = H(D_1) - H(\text{department}) = 0.9468 - 0 = 0.9468$$

> age

$$H(\text{age}) = \frac{20}{63} \times 0 + \frac{3}{63} \times 0 + \frac{30}{63} \times 0 + \frac{10}{63} \times 0 + \frac{0}{63} \times 0 + \frac{0}{63} \times 0 = 0$$

$$\text{Gain}(\text{age}) = H(D_1) - H(\text{age}) = 0.9468 - 0 = 0.9468$$

Summary

$$\text{Gain}(\text{Department} | D_1) = 0.9468$$

$$\text{Gain}(\text{age} | D_1) = 0.9468$$

Both attributes have same gain

\therefore we choose one of them, & will place it below "salary = 46k..50k".

we choose department as attribute below "salary = 46k..50k".

From the dataset 1,



If $\text{salary} = 46k..50k$ AND $\text{department} = \text{sales}$ Then status = junior.

also = systems Then status = junior.

also = marketing Then status = senior.

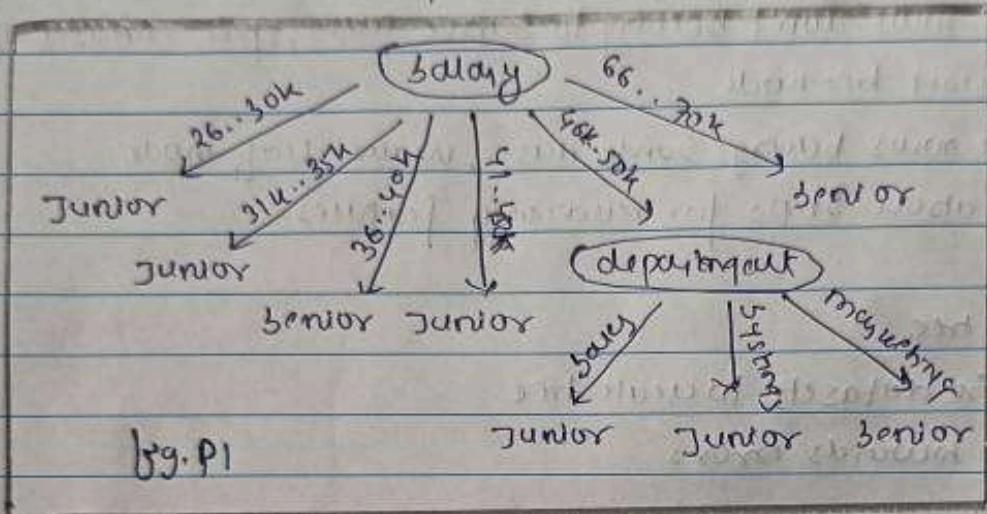


Fig. P1

- 7 ID3: ID3 stands for Iterative Dichotomies 3 and is named such because the algorithm iteratively (repeatedly) dichotomizes (divides) features into two or more groups, at each step. uses a top-down greedy approach to build decision tree. used to ID3 is for classification problems with nominal features.

principles in ID3

- uses gain and info gain to find best feature.
 - information gain
 - Entropy
 - data set as D
 - i.e. Entropy (D) = $\sum_{i=1}^n p_i \log_2 \frac{1}{p_i}$
 - $n = \text{total no. of classes}$
 - $p_i = \text{probability of class } i$
 - ratio of no. of ones to no. of zeros.
- worst info gain selected.

- randomly drawn object

- entropy

- binary classification



M. for profit gain

$$GWN(A) = \text{Entropy}(D) - \text{Entropy}(A)$$

> ID3 steps:

- (1) calculate M. gain of each feature
- (2) consider all rows don't belong to same class, split dataset D
- (3) make decision tree node
- (4) if all nodes belongs same class, make leaf node
- (5) Repeat the above steps for remaining features.

> Advantages:

- builds short tree
- searches entire datasets to weak tree
- less sensitive towards errors

> disadvantages:

- doesn't perform backtracking while searching
- may converge in locally optimal soln.

10. write a short note on the following:

a) DBSCAN:

Density Based Spatial clustering of applications with Noise
↓ population.

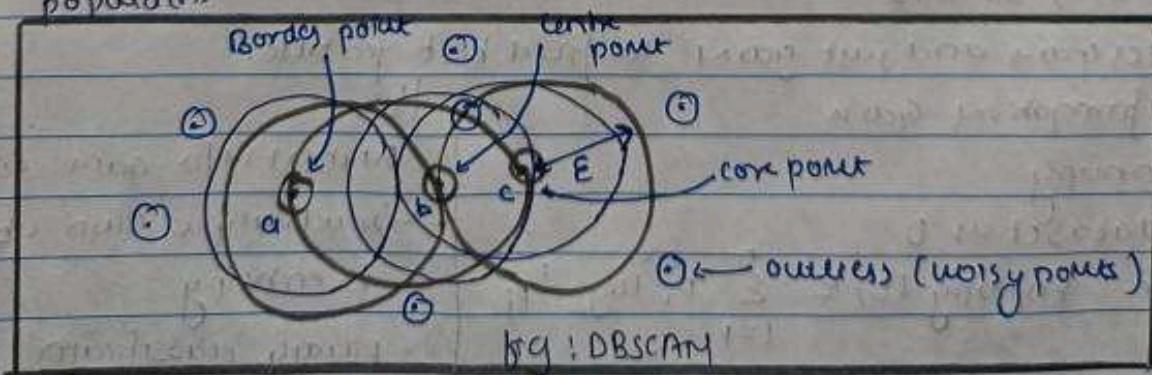


fig : DBSCAN!

(1) ϵ - neighbourhood

(2) min-points \geq must = 3

(3) core points \Rightarrow strong constituent of the clusters



- ④ Border points
- ⑤ Noise points
- directly density reachable
- ↳ any being ⑥ neighbour or core point
- possible through core point
- for example $c \rightarrow b$, $b \rightarrow a \therefore c \rightarrow a$
 $a \rightarrow b$, $b \rightarrow c$ but $a \rightarrow c$ because a isn't core point
- this method is based on the notion of density.
 the basic idea is to continue growing the given cluster, as long as the density in the neighbourhood exceeds some threshold, i.e. for each data point within a given cluster, the radius of a given cluster has to contain at least a minimum no. of points.

b) Web structure mining:

- web structure mining is used for creating a model of web organization. It is the process of analyzing the nodes and connection structure of a website using graph theory.
- There are two things that can be obtained from this:
 - ① the structure of website in terms of how it connected
 - ② document structure of website itself
- The web structure mining can be used to discover link structure of hyperlinks. It is used to identify that the web pages are either linked by informal or direct link connection.
- Example: Page Rank algorithm, CLEVER, HITS).

● Page Rank:

- The Page Rank (PR) algorithm is applicable in web pages.
- Page Rank is an algorithm used by Google search to rank websites in their search engine results.
- Page rank algorithm is designed to measure the effectiveness of

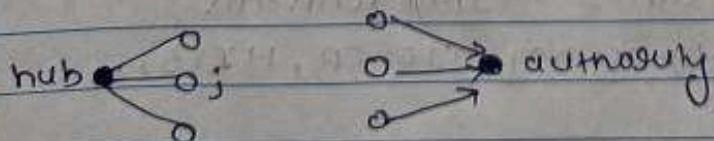


search engines and improve their efficiency.

- It's a way of measuring - the importance of website pages
- page Rank is used to prioritize the pages returned from a traditional search engine using keyword searchingly.
- page rank is calculated based on the no. of pages that point to it.

• Hyperlink-Reduced Topic Search (HITS)

- HITS, also known as Hubs and authorities, developed by John Kleinberg is a link analysis algorithm that rates web pages. It was precursor to Page Rank.
- The idea behind Hubs and Authority stemmed from a pathway insight into the heart of web pages when the Internet was originally forming.
- The HITS algorithm treats www as directed graph $G(V, E)$, where V is the set of vertices representing pages and E is a set of edges corresponds to links.
- Based on mutually recursive facts: Hubs point to lots of authorities. Authorities are pointed to by lots of hubs.
 - Authority: valuable & informative web page usually pointed to by many hyperlinks.
 - Hub: A webpage that points to many authority page is itself a resource and is called a hub.
- Authorities and hubs reinforce one another.



big! Hub and authority

A good authority pointed by many good hubs.

A good hub points to many good authorities.