**Aim:** Creation of a Data warehouse.

→ Build Data Ware house / Dart Mart (using open source tools like pentaho Data. Integration. Tool, pentaho

. Business Analytics, or other data ware house tools like Microsoft - SSIS, Information, Business objects etc...

→ Design multi - dimensional data models namely star, Snowflake and fact constellation schemas for any one enterprise.

(Eg :- Banking, Insurance, Finance, Healthcare, manufacturing, Automobile, Sales etc).

→ Write ETL scripts and implement using data warehouse tools.

→ perform various OLAP operations such slice, dice, rollup, drillup and pivot.

**Objectivies:**

Data warehousing is a technique of gathering and analyzing data from many sources to get valuable, business, insights. Typically a data ware house integrates and analyzes business data from many sources. Data ware housing is a vital component of business intelligence.

## Preprocess :-

The data that is collected from the field contains many unwanted things that leads to wrong analysis. Those the data must be preprocessed to meet the requirements of the type of analysis you are setting. This is the done in the preprocessing module.

## Classifiers :-

Classifiers in WEKA are the models for Predicting nominal or numeric quantities. The learning schemes available in WEKA include decision trees and lists, instance - based classifiers, classificers include bagging, boosting, Stacking, error - correcting output codes and locally weighted learning.

## WEKA :-

WEKA (Waikota Environment for Knowledge Analysis) is a popular suite of machine learning software written in Java, developed at the university of waikota, Newzealand. weka is free Software available under the GNO General Public license.

Weka - an open source software provides tools for data preprocessing implementation of several machine learning algorithms and visualization tools so that you can develop machine learning techniques and apply them to real - world data mining problems.

## Explorer: It is an environment for exploring data

Explorer consists of several tabs. They are:

→ Preprocess: It is the first step in machine learning is to preprocess the data. It is used to select the data file, process it and make it fit for applying the various machine learning algorithms.

→ Classify: The classify tab provides you several machine learning algorithms for the classification of your data Such as Linear Regression, Logistic Regression.

→ Cluster: under the cluster tab there are several clustering algorithm provided - such as simple k Means, Filtered Clusterer, Hierarchical cluster.

→ Associate: under the Associate tab you would find Apriori filtered Associator and FP Growth.

• Select Attributes Tab:

Select Attributes allows you feature selection based on several algorithms such as classifier, subset Eval, Principal components.

• Visualize Tab:

The visualize option allows you to visualize your processed data for analysis.

→ Simple CLI: It provides a simple command - line interface and allows direct execution of weka commands.

→ Experimenter: It is an environment for performing experiment and conducting statistical tests between learning scheme. knowledge +

→ Knowledge Flow: It is a Java-Beans based interface for setting up running machine learning experiments.

Trees J48 Classifier:

It is an algorithm to generate a decision tree that is generated by (4.5. It is also known as a statistical classifier. For decision tree classification, we need a database. wet
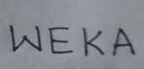
weather nominal:

In weka, attributes can be nominal or numeric. The value of a nominal attribute is represented by a word: sunny, overcast and rainy for the outlook attribute; yes and no for the play attribute.

Steps Required :-

1. open WEKA you can see 5 tabs on the right side of the application. these are: Explorer, Experimenter, knowledge flow, work bench, Simple CLI.

2. Click on "Explorer".

3. On preprocess. Click on "Open file".

4. Go to "C:\ program files\ weka - 3-8-6\ data", select "weather. nominal. arff" and click on open.

5. Click on "Classify" and then click on choose.

6. You will see the following options. select J48 & Click on "start".

7. Click on the resulted list to see the visual.

8. click on the resulted list and click on visualize tree option. outcome of the experiment :-

⇒ Out come of the Experiment :-

| Program | visualization | Tools | Help | — ☐ X Application |
|---|---|---|---|---|
| | | | | Explorer |
| | | | | Experimenter |
| | WEKA | | | Knowledge Flow |
| | | | | Work bench |
| | | | | Simple CLI |

weka Explorer                                    – □ X

Preprocess  classify  cluster  Associate  selectio  Visualiz

openfile   openURL   openDB

---

Look in  data

→ airlens
→ breakcase
→ CPU            Weathernominal

---

Preprocess  classify  cluster  Associate   Selection  visualize

openfile   openURL   openDB   Generate

File

Choose   None

Select different

| No | label | count | weight |
|----|-------|-------|--------|
| 1  | sunny | 5     | 5      |
| 2  | over  | 4     | 4      |
| 3  | rainy | 5     | 5      |

5                                5



13/10/22

**Aim:** Demonstrate performing classification on data.

→ Load each dataset into weka and run id3, J48 classification algorithm. Study the classifier output. Compute entropy values, kappa statistic.

→ Extract if then rules from the decision tree generated by the classifier, observe the confusion matrix.

→ Load each dataset into weka and perform Naive-bays classification and K-Nearest Neighbour Classification. Interpret the results obtained.

→ Plot ROC curves

→ Compare classification result of ID3, J48, Naive-Bayes and K-NN classifiers for each dataset and reduce which classifier is Performing best and poor for each dataset and justify.

**Objectives:**

The ultimate objective of classification is to relate a variable of interest with observed variables. The actual variable of interest is meant to be of "Qualitative" type. The algorithm required for performing the classification is known as the classifier.

**Zero R :-**

→ Zero R is the simplest classification method which relies on the target and ignores all predictors.

→ Zero R classifier simply predicts the majority category.

→ Although there is no predictability power in Zero R it is useful for determining a baseline performance as a benchmark for other classification methods.

**One R :-**

→ This method is used in the sequential learning algorithm for leaving the rules.

→ It returns a single rule that covers at least some examples.

→ However, what makes it really powerful is its ability to create relations among the attributes given. Hence covering a larger hypothesis space.

**Explorer:** It is an environment for exploring data.

**Simple CLI:** It provides a simple command-line-interface and allows direct execution of weka commands.

**Experimenter: Knowledge Flow:**

It is a Java-Beans based interface to setting up and running machine learning experiments.

## Experimenter:

It is an environment for performing experiment and conducting statistical tests between learning scheme.

## Preprocess:

It is the first step in machine learning to preprocess the data. It is used to select the data file preprocessing and make it fit for applying the various machine learning algorithms.

## Classify:

The classify tab provides you several machine learning algorithms for, the classification of your data such as linear Regression, logistic Regression.

## Test options:

Before you run the classification algorithm, you need to set test options. Set test options in the 'Test Options' box. The test options that available are.

### 1. Use training set:-

Evaluates the classifier on how well it Predicts the class of the instances it was trained on.

### 2. Supplied test set:

Evaluates the classifier on how well it predicts the class of a set of instances loaded from a file. clicking on the 'set.' button brings up a dialog allowing you to choose the file to test on.

### 3. Cross validation:

Evaluates the classifier by cross-validation using the number of folds that are entered in the 'Folds' text field.

### 4. percentage Spilt:

Evaluates the classifier on how well it predicts a certain percentage of the data, which is held out for testing. The amount of data held out depends on the value entered in the '%' field.

### Steps Required:

1. Open weka you can see 5 tabs on the right side of the application. These are explorer, experimentor, knowledge flow, workbench, simple CLI.

2. Click on 'explorer'.

3. You can see classify tab click on the classify button.

4. You can observe choose, test options etc..

5. In test option you can see cross-validation folds. set it as 10.

6. Right click on choose option, then select the zero R algorithm or one R algorithm.

7. Then click Start button.

8. Zero R algorithm or one R algorithm will execute and it gives the output.

Output :-

Zero R

| Pre process Classify Cluster Associate Select attribute Visualiz | -□X |
|---|---|
| choose zero R - B 6 | Classifier output: |
| Test options | Correct classifier 9 64.285 |
| • use training set | incorrect classifier 5 35.740 |
| • Supplied test | |
| • Cross-Validation Fold 10 | |
| • percentage split % 56 | |
| 21·36·38 — rules· zero | |

One R :-

| Pre process Classify cluster Associate Select Visualize | -□X |
|---|---|
| Choose one-R-D6 | classify output |
| Test option | Correctly classifier instance |
| • use training set | 6    42·85 71 %. |
| • supplied test | Incorrectly classifier instance |
| • cross validation 10 | 8    57·1429 %. |
| • Percentage split% 56 | |
| 21-36·38 - rules-zero R | |
| 21-36·38 -rules-one R | |

¢
13/10/22

Write a program of cluster analysis using simple k-means algorithm python programming language.

## Cluster Analysis:

Cluster Analysis is a statistical method for processing data. It works by organizing items into groups, or clusters on the basis of how closely associated they are.

## K-means algorithm:

k-means algorithm is a simple two steps clustering process. The first step is cluster assignment and the second one is the move centroid step. However, this unsupervised algorithm can easily create, implement and handle massive datasets.

## Steps involved in k-means Algorithm:

Step 1: Select the number k to decide the number of clusters.

Step 2: Select random k points or centroids.

Step 3: Assign each data point to their closest centroid, which will form the predefined k clusters.

Step 4: Calculate the variance and place a new centroid of each cluster.

Step 5: Repeat the third steps, which means reassign each datapoint to the new closest centroid of each cluster.

Step 6: If any reassignment occurs, then go to step-4, else go to FINISH.

Step 7: The model is ready.

# K-means Algorithm using Python programming.

```python
import numpy as nm
import matplotlib.pyplot as mtp
import pandas as pd
dataset = pd.read_csv('/content/Sample_data/Mall_
                        customers.csv')
X = dataset.iloc[:, [3,4]].values
from sklearn.cluster import KMeans
wcss_list = []
for i in range (1,11):
    Kmeans = kmeans(n_clusters = i, init =
                    'k-means++', random_state=42)
    Kmeans.fit(x)
    wcss_list.append(kmeans.inertia_)
mtp.plot(range(1,11), wcss_list)
mtp.title('The Elbow Method Graph')
mtp.xlabel('Number of Clusters(k)')
mtp.ylabel('wcss_list')
mtp.show()

kmeans = kmeans mode(n_clusters=5, init = 'k-means++'
                    random_state=42)

y_predict = kmeans.fit_predict(x)
mtp.scatter(x[y_predict = 0,0], x[y_predict=0,1],
            s=100, c='blue', label = 'cluster 1')
mtp.scatter(x[y_predict = 1,0], x[y_predict = 1,1],
            s=100, c='green', label = 'cluster 2')
mtp.scatter(x[y_predict = 2,0], x[y_predict=2,1],
            s=100, c='red', label = 'cluster 3')
mtp.scatter(kmeans.cluster_centers_[:, 0],
            kmeans.cluster_centers_[:,1], s=300,
            c='yellow', label = 'centroid')
mtp.title('clusters of customers')
mtp.xlabel('Annual Income (k$')
mtp.ylabel('Spending score (1-100)')
mtp.legend()
mtp.show()
```
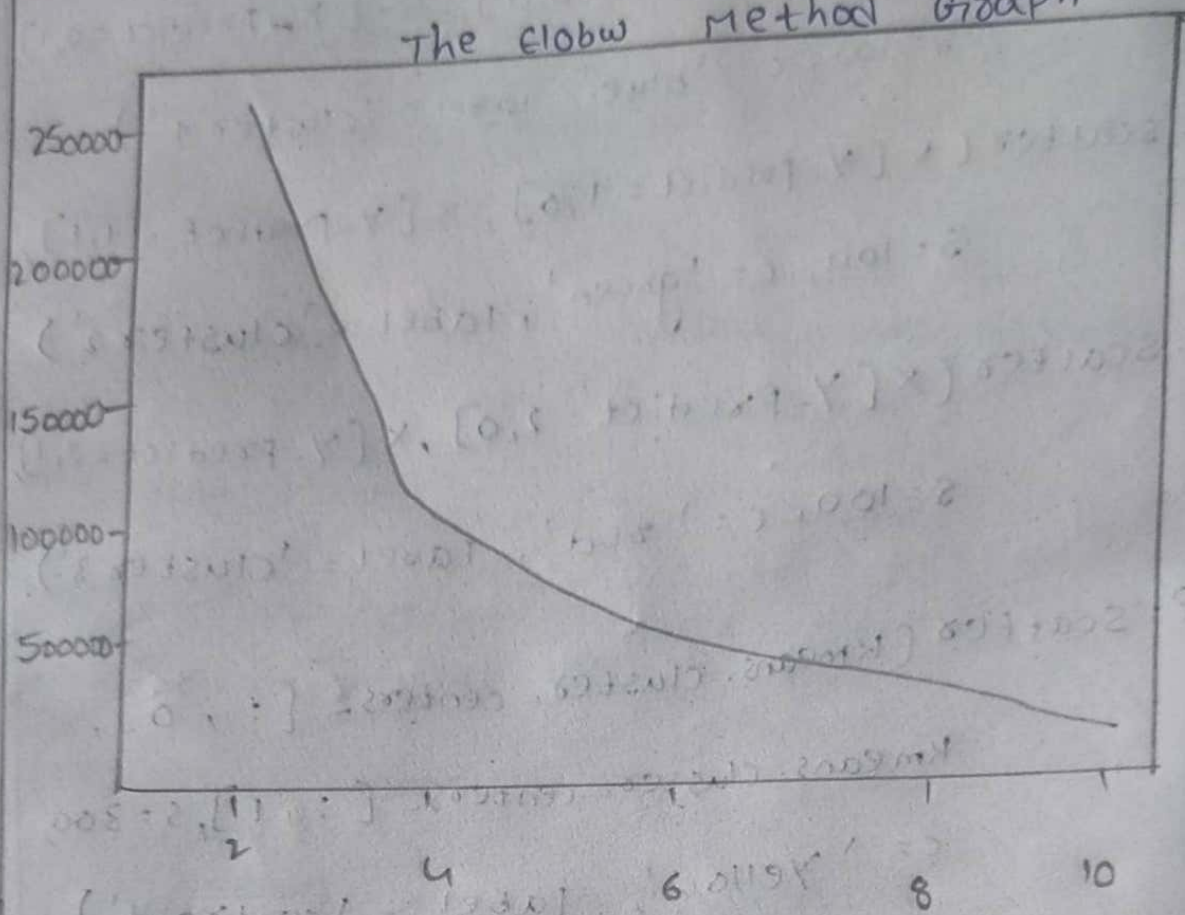
Out put

# The Elobw Method Graph



250000

200000

150000

100000

50000

2    4    6    8    10

Number of clusters(k)

Clusters of customers



Spending score (1-100)

100

80

60

40

20

0

13/10/22

20   40   60   80   100   120   140

Annual Income (k $)

• Cluster 1
• Cluster 2
• Cluster 3
  Centroid