



American International University-Bangladesh (AIUB)

Department of Computer Science

Faculty of Science & Technology (FST)

Summer 21-22

Section: A

Data Warehouse and Data Mining

PROJECT

A Report submitted

By

SN	Student Name	Student ID
1	SAJIDUL HASAN	18-38627-2
2	TURSHIN ARA ASHTARY	18-38593-2
3	ABDULLAH AL MAKSUD	18-38582-2

Under the supervision of

Tohedul Islam
Assistant Professor

Table of Contents

Task1	2
introduction.....	2
Applying Naïve Bayes Classifiar	4
Applying KNN classifiar	6
Result.....	6
Discussion	7
Task2	7
Introduction	7
Result.....	9
Discussion	9
Task 3	10
Introduction	11
Discussion	13
Result	13

TASK 1:

Here is Our Dataset (Supervised Learning Dataset)

Introduction: Data mining using labelled data is known as supervised learning .where labelled data is a specially designated attribute and the aim is to use the data given to predict the value of that attribute for instance that has not to get been seen.

Targeted Features are: Housing In London Monthly Variable

Attributes: 7 attributes in this dataset

1. date
2. area
3. average_price
4. code

5. houses_sold
6. no_of_crimes
7. borough_flag

There is a total of 13549 instances of these 7 attributes and all these instances were used for classification. Here are the graphical details of the attributes:

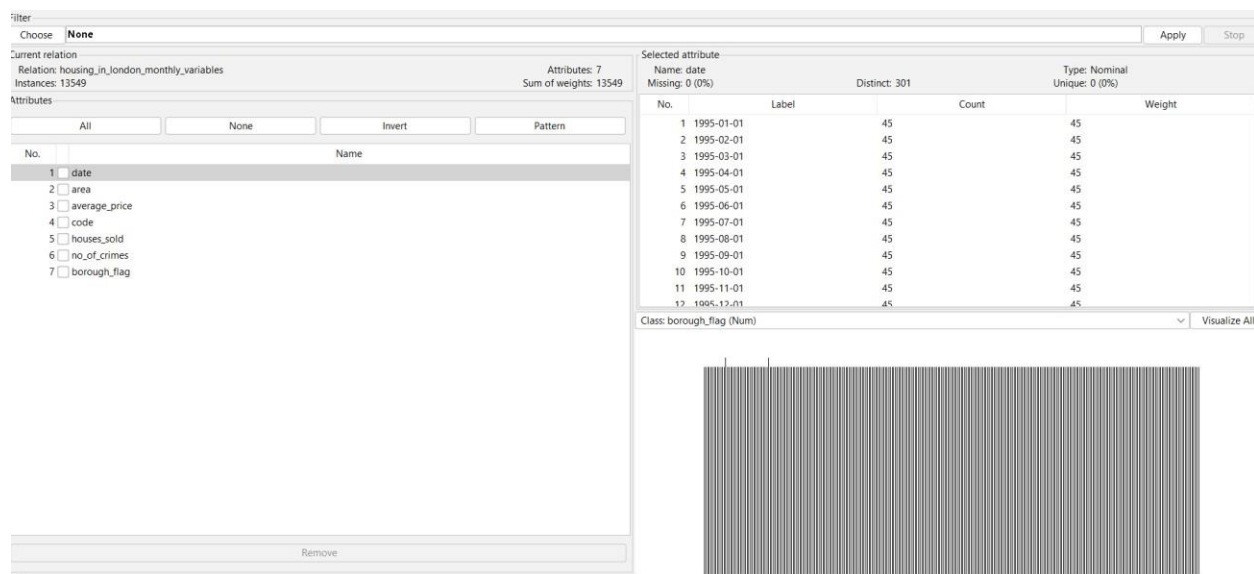


Fig 1: Selected dataset

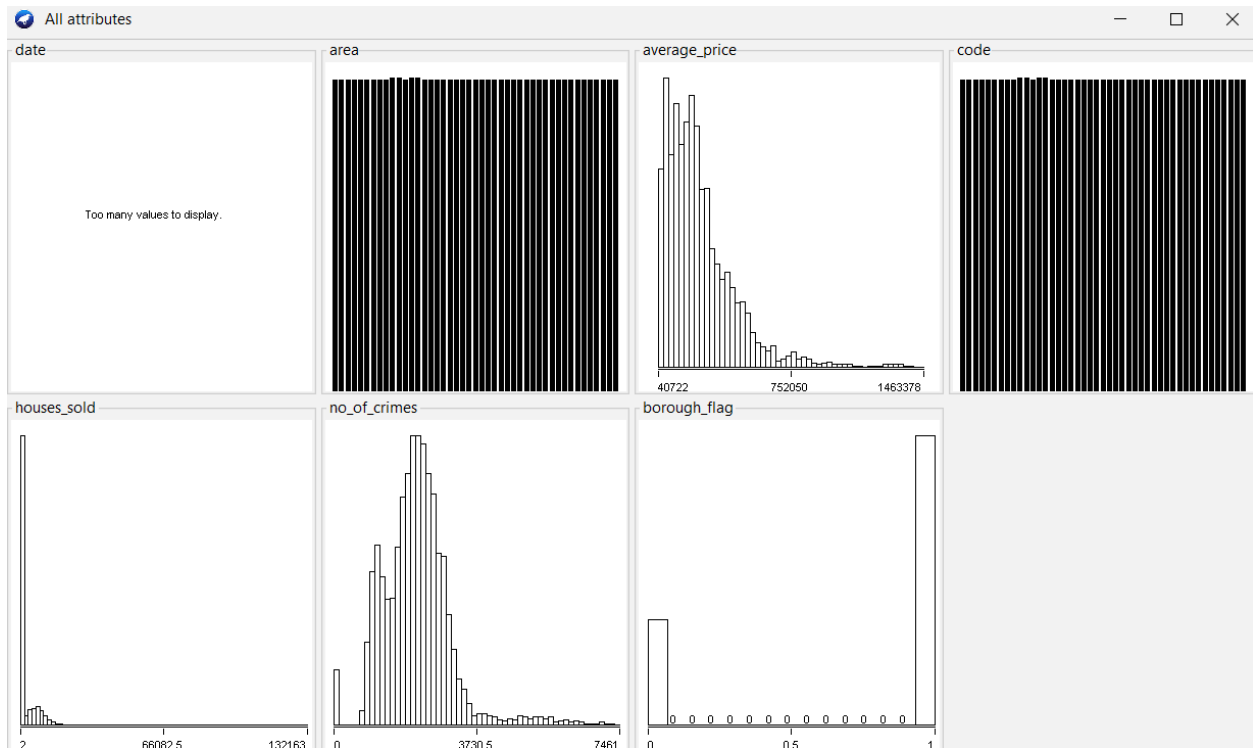


Fig 2: Details of all attribute

Applying Naïve Bayes Theorem: Naive Bayes algorithm is a supervised learning algorithm, which is based on the Bayes theorem and is used for solving classification problems. It is mainly used in text classification that includes a high-dimensional training dataset. It is a probabilistic classifier, which means it predicts on the basis of the probability of an object. While classifying the selected dataset, the Naïve Bayes format was selected from the Bayes folder.

```

=== Run information ===

Scheme:      weka.classifiers.bayes.NaiveBayes
Relation:    housing_in_london_monthly_variables
Instances:   13549
Attributes:  7
    date
    area
    average_price
    code
    houses_sold
    no_of_crimes
    borough_flag
Test mode:   evaluate on training data

=== Classifier model (full training set) ===

Naive Bayes Classifier

Attribute      Class
              city of london  barking and dagenham  barnet  bexley  brent  bromley  camden
              (0.02)          (0.02)          (0.02)  (0.02)  (0.02)  (0.02)  (0.02)
=====
date
1995-01-01      2.0          2.0          2.0      2.0      2.0      2.0      2.0
1995-02-01      2.0          2.0          2.0      2.0      2.0      2.0      2.0
1995-03-01      2.0          2.0          2.0      2.0      2.0      2.0      2.0
1995-04-01      2.0          2.0          2.0      2.0      2.0      2.0      2.0
1995-05-01      2.0          2.0          2.0      2.0      2.0      2.0      2.0
1995-06-01      2.0          2.0          2.0      2.0      2.0      2.0      2.0
1995-07-01      2.0          2.0          2.0      2.0      2.0      2.0      2.0
1995-08-01      2.0          2.0          2.0      2.0      2.0      2.0      2.0
1995-09-01      2.0          2.0          2.0      2.0      2.0      2.0      2.0

```

Fig-3

```

=== Evaluation on training set ===

Time taken to test model on training data: 8.2 seconds

=== Summary ===

Correctly Classified Instances      13537          99.9114 %
Incorrectly Classified Instances     12          0.0886 %
Kappa statistic                     0.9991
Mean absolute error                  0.002
Root mean squared error              0.013
Relative absolute error              4.4907 %
Root relative squared error          8.8265 %
Total Number of Instances           13549

=== Detailed Accuracy By Class ===

TP Rate  FP Rate  Precision  Recall  F-Measure  MCC  ROC Area  PRC Area  Class
1.000    0.000    1.000     1.000    1.000     1.000  1.000    1.000    city of london
1.000    0.000    1.000     1.000    1.000     1.000  1.000    1.000    barking and dagenham
1.000    0.000    1.000     1.000    1.000     1.000  1.000    1.000    barnet
1.000    0.000    1.000     1.000    1.000     1.000  1.000    1.000    bexley
1.000    0.000    1.000     1.000    1.000     1.000  1.000    1.000    brent
1.000    0.000    0.997     1.000    0.998     0.998  1.000    1.000    bromley
1.000    0.000    1.000     1.000    1.000     1.000  1.000    1.000    camden
1.000    0.000    0.984     1.000    0.992     0.992  1.000    1.000    croydon
1.000    0.000    1.000     1.000    1.000     1.000  1.000    1.000    ealing
0.993    0.000    0.997     0.993    0.995     0.995  1.000    0.997    enfield
0.997    0.000    0.993     0.997    0.995     0.995  1.000    0.996    tower hamlets
1.000    0.000    1.000     1.000    1.000     1.000  1.000    1.000    greenwich
1.000    0.000    0.997     1.000    0.998     0.998  1.000    1.000    barknaw

```

Fig-4: Applying Naïve Bayes Classifier


```

=== Evaluation on training set ===

Time taken to test model on training data: 22.66 seconds

=== Summary ===

Correctly Classified Instances      13545      99.9705 %
Incorrectly Classified Instances      4      0.0295 %
Kappa statistic      0.9997
Mean absolute error      0.0002
Root mean squared error      0.0036
Relative absolute error      0.3611 %
Root relative squared error      2.4755 %
Total Number of Instances      13549

=== Detailed Accuracy By Class ===

      TP Rate    FP Rate    Precision    Recall    F-Measure    MCC      ROC Area    PRC Area    Class
1.000    0.000    1.000    1.000    1.000    1.000    1.000    1.000    E090000001
1.000    0.000    1.000    1.000    1.000    1.000    1.000    1.000    E090000002
1.000    0.000    1.000    1.000    1.000    1.000    1.000    1.000    E090000003
1.000    0.000    1.000    1.000    1.000    1.000    1.000    1.000    E090000004
1.000    0.000    1.000    1.000    1.000    1.000    1.000    1.000    E090000005
1.000    0.000    1.000    1.000    1.000    1.000    1.000    1.000    E090000006
1.000    0.000    1.000    1.000    1.000    1.000    1.000    1.000    E090000007
1.000    0.000    1.000    1.000    1.000    1.000    1.000    1.000    E090000008
1.000    0.000    1.000    1.000    1.000    1.000    1.000    1.000    E090000009
0.997    0.000    0.997    0.997    0.997    0.997    0.998    0.993    E090000010
0.997    0.000    0.997    0.997    0.997    0.997    0.998    0.993    E090000030
1.000    0.000    1.000    1.000    1.000    1.000    1.000    1.000    E090000011
0.997    0.000    0.997    0.997    0.997    0.997    0.998    0.993    E090000012
0.997    0.000    0.997    0.997    0.997    0.997    0.998    0.993    E120000008

```

[illegible]

RESULT:

Classifier	Accuracy
Naïve Bayes	99.91%
KNN	99.97%

Discussion:

After applying two types of classifiers, the highest percentage of correctly classified instances is for the naïve Bayes classifier with 99.9114%. After that comes the KNN classifier with 99.9705%. The KNN classifier is considered the best classifier for the dataset.

From the above discussion, we saw that the accuracy rate of the KNN classifier is higher than the Naïve Bayes

Task 2:

Introduction:

Supervised learning data(from task-1 Data set) to Test data set:

We make the test data set from the supervised learning data set we used before. We take here 2709 instances & 7 attribute in this data set.

Attribute are:

1. date
2. area
3. average_price
4. code
5. houses_sold
6. no_of_crimes
7. borough_flag

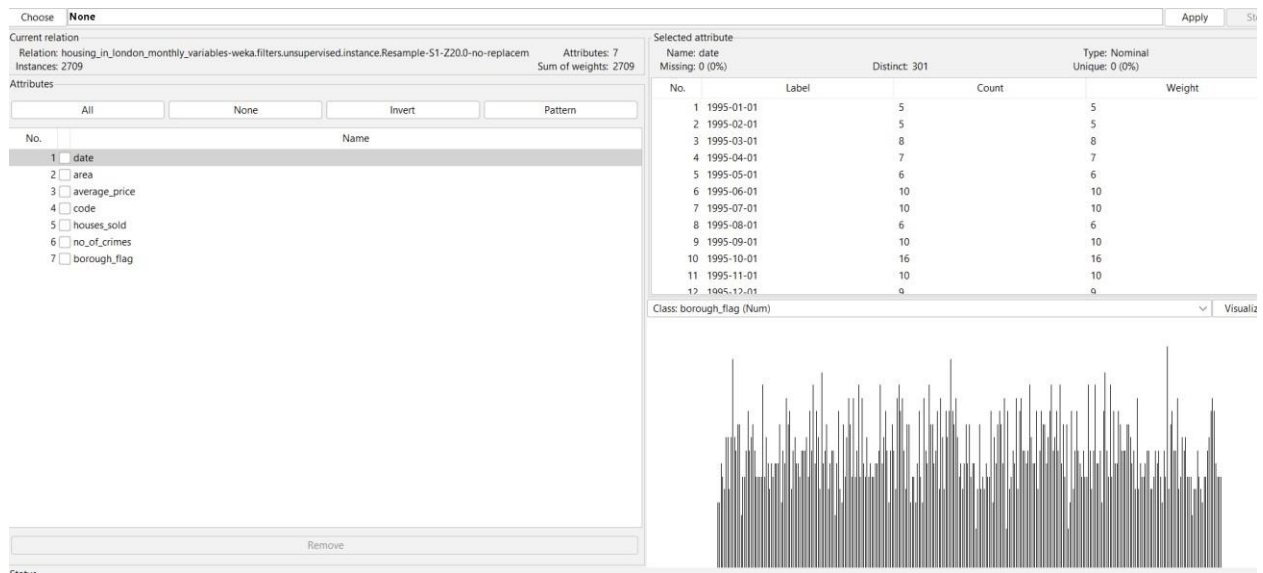


Fig 5: Selected dataset

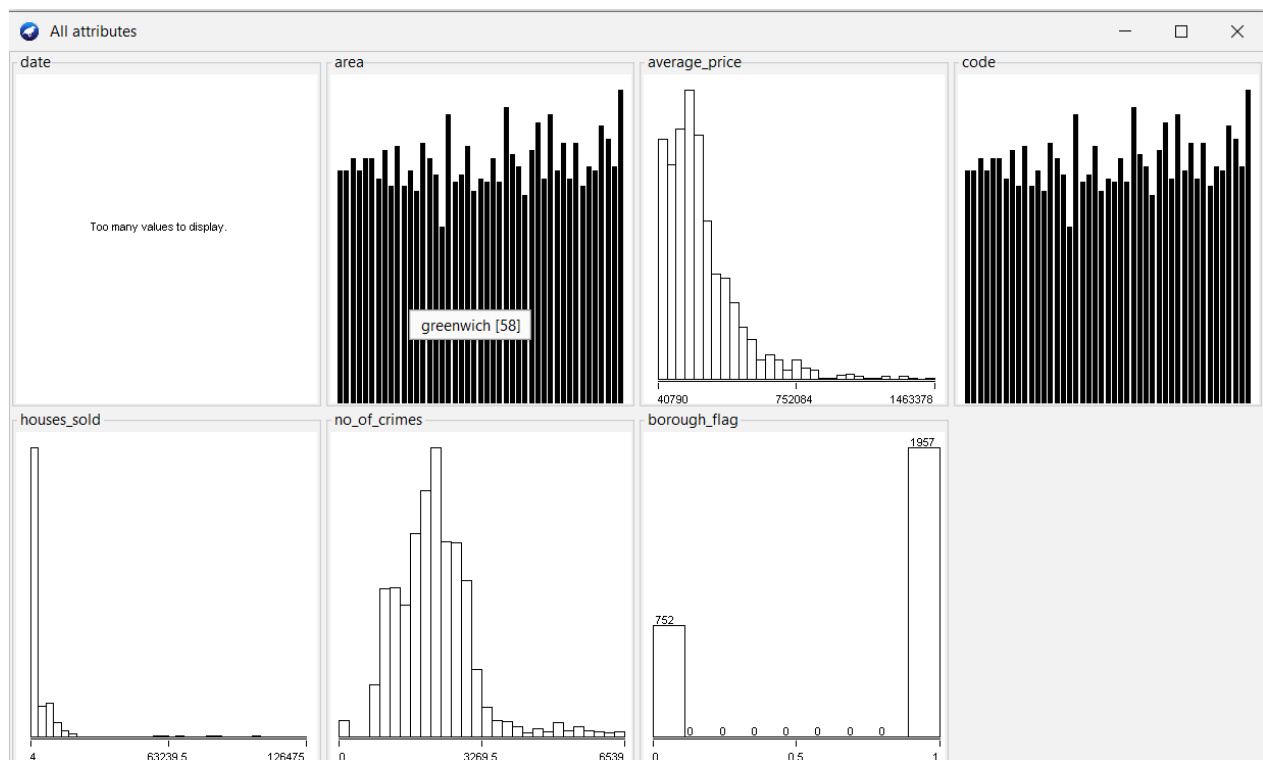


Fig 6: Details of all attribute

=== Evaluation on test set ===

Time taken to test model on supplied test set: 2.48 seconds

=== Summary ===

Correctly Classified Instances	2707	99.9262 %
Incorrectly Classified Instances	2	0.0738 %
Kappa statistic	0.9992	
Mean absolute error	0	
Root mean squared error	0.0057	
Relative absolute error	0.1045 %	
Root relative squared error	3.8757 %	
Total Number of Instances	2709	

=== Detailed Accuracy By Class ===

TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000001
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000002
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000003
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000004
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000005
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000006
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000007
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000008
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000009
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000010
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000030
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000011
0.981	0.000	0.981	0.981	0.981	0.981	1.000	0.971	E09000012
0.985	0.000	0.985	0.985	0.985	0.984	1.000	0.978	E12000008
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000013
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000014
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000015
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000016
1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	E09000017

Fig 7: Applying J48 classifier

evaluating data mining models. Since the test set data already contains known values for the property I want to predict, it is easy to determine if the model's assumptions are correct. So here I carefully separate the data into test/training and apply them in Weka after that it gives the error table and the accuracy rate is decreasing than before that's how I get my proper result and accuracy model.

I use here j48 classifier & it shows 99.62% accuracy whereas in the previous when we use it in the training dataset using KNN classifier it shows 99.97% accuracy.

TASK-3:

Introduction:

Unsupervised Learning is a machine learning technique in which the users do not need to supervise the model. Instead, it allows the model to work on its own to discover patterns and information that was previously undetected. It mainly deals with unlabelled data.

Unsupervised learning is helpful for finding useful insights from the data. Unsupervised learning is much similar to human learning to think by their own experiences, which makes it closer to the real AI.

Unsupervised learning works on unlabeled and uncategorized data which makes unsupervised learning more important.

I have chosen the “accidental-deaths-in-USA-monthly” dataset. I will also use K means clustering Algorithm.

About the dataset: In this report, the used “accidental-deaths-in-usa-

monthly”, a CSV dataset file [Later converted into .arff], collected from

Kaggle.com.

RESULT :

Applying K means clustering Algorithm.

K-Means Clustering is an Unsupervised Learning algorithm, which groups the unlabeled dataset into different clusters. Here K defines the number of pre-defined clusters that need to be created in the process, as if $K=2$, there will be two clusters, and for $K=3$, there will be three clusters, and so on. It allows us to cluster the data into different groups and is a convenient way to discover the categories of groups in the unlabeled dataset on its own without the need for any training.

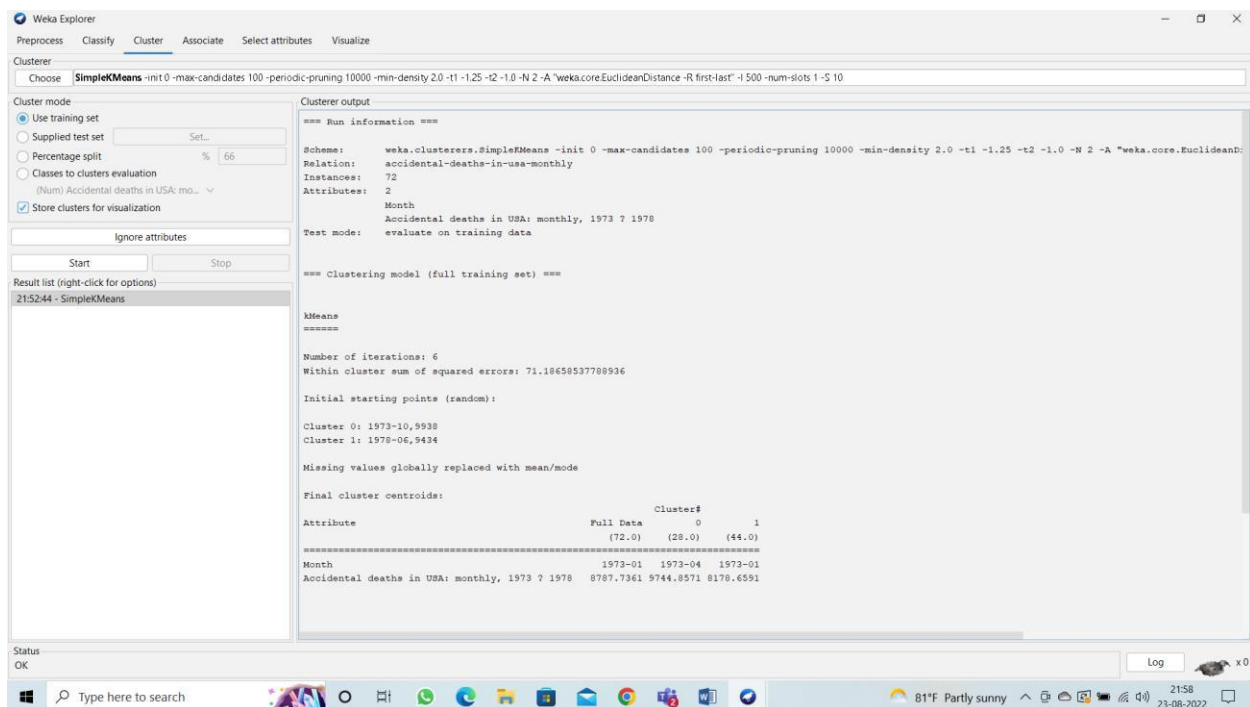


Fig-8

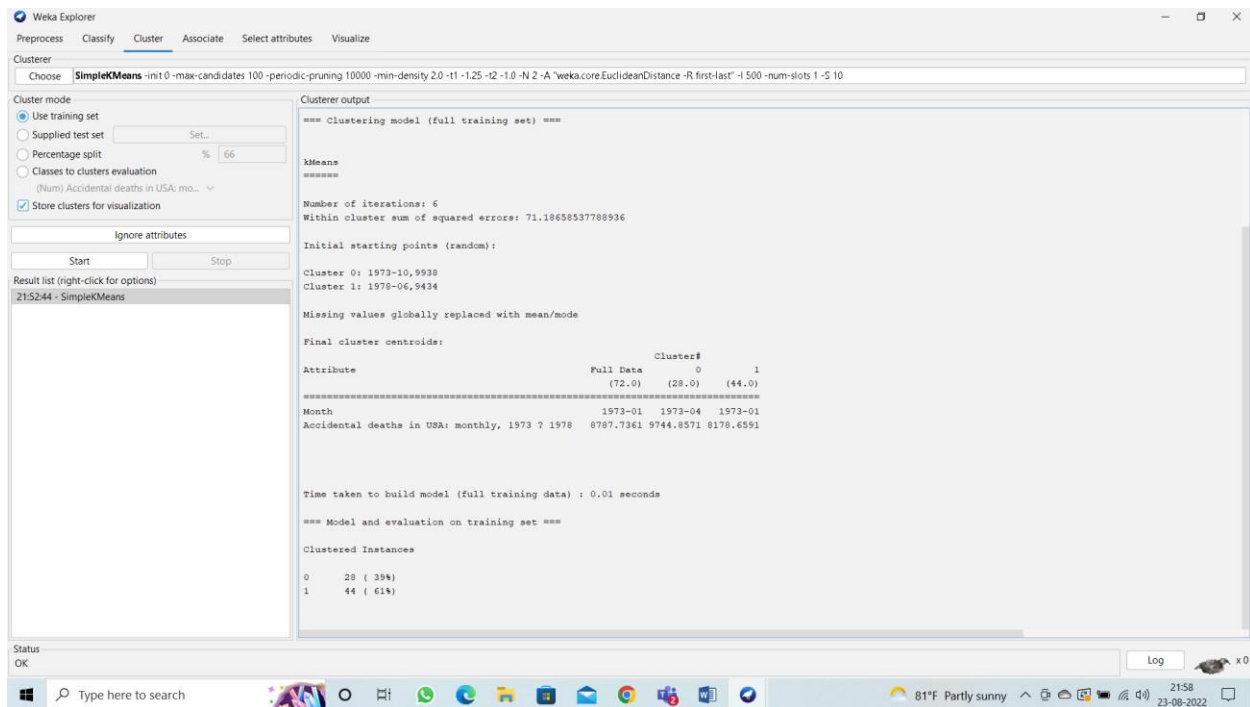


Fig-9

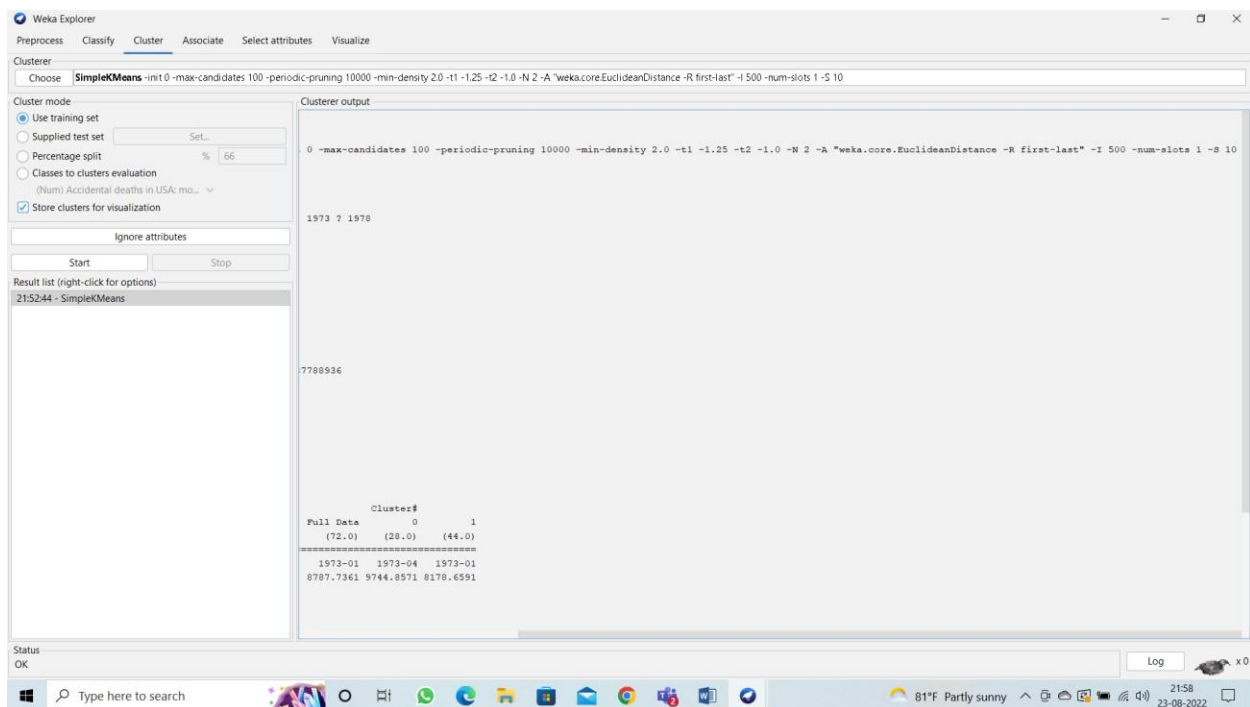


Fig-10

Figures 1,2,3,: K means clustering Algorithm

Here is the summary of the K means clustering Algorithm result:

Number of iterations: 6

Initial starting points (random):

Cluster 0: 1973-10,9938

Cluster 1: 1978-06,9434

Time taken to build model (full training data) : 0.01 seconds

Clustered Instances

0 28(39%)

1 44(61%)

Discussion:

I have chosen a proper unsupervised dataset. I convert csv file to arff file and then apply K means Clustering Algorithm. By default, the value of K is 2 so there are two results. For 0 instances it is 39% whereas for 1 is 61%. I also find the final cluster of centroids also. By that I have completed K means clustering for unsupervised data