

CreditCard Fraud Prediction

Overview:

Building a machine learning model to predict the credit card fraud data. Here we are providing a dataset having the previous creditcard fraud data.using that data we are training a machine by writing the code for machine in java programing.Data is processed to the model using Linear Regression, Logistic Regression , Clustering.... After Training the model we are going to test the model using related data.

purpose:

The credit card companies are able to recognize fraudulent credit card transactions so that customers are not charged for items that they did not purchase

Literature Survey:

here is the some of the data from the dataset regarding fraud prediction

Gender,Married,Dependents,Education,Self_Employed,ApplicantIncome,CoapplicantIncome,LoanAmount,Loan_Term,Credit_History_Available,Housing,Locality,Fraud_Risk

1,0,0,1,0,5849,0,146,360,1,1,1,0
1,1,1,1,1,4583,1508,128,360,1,1,3,1
1,1,0,1,1,3000,0,66,360,1,1,1,1
1,1,0,0,1,2583,2358,120,360,1,1,1,1
1,0,0,1,0,6000,0,141,360,1,1,1,0
1,1,2,1,1,5417,4196,267,360,1,0,1,1
1,1,0,0,1,2333,1516,95,360,1,1,1,1
1,1,3,1,1,3036,2504,158,360,0,1,2,1
1,1,2,1,1,4006,1526,168,360,1,1,1,1
1,1,1,1,1,12841,10968,349,360,1,0,2,1
1,1,2,1,1,3200,700,70,360,1,0,1,1
1,1,2,1,1,2500,1840,109,360,1,0,1,1
1,1,2,1,1,3073,8106,200,360,1,0,1,1
1,0,0,1,0,1853,2840,114,360,1,1,3,1
1,1,2,1,1,1299,1086,17,120,1,1,1,1
1,0,0,1,0,4950,0,125,360,1,1,1,0
1,0,1,0,0,3596,0,100,240,1,1,1,0

0,0,0,1,0,3510,0,76,360,0,1,1,1
1,1,0,0,1,4887,0,133,360,1,1,3,1
1,1,0,1,1,2600,3500,115,12,1,1,1,1
1,1,0,0,1,7660,0,104,360,0,1,1,1
1,1,1,1,1,5955,5625,315,360,1,0,1,1
1,1,0,0,1,2600,1911,116,360,0,1,2,1
0,1,2,0,1,3365,1917,112,360,0,0,3,1
1,1,1,1,1,3717,2925,151,360,1,0,2,1
1,1,0,1,1,9560,0,191,360,1,1,2,1
1,1,0,1,1,2799,2253,122,360,1,0,2,1
1,1,2,0,1,4226,1040,110,360,1,0,1,1
1,0,0,0,0,1442,0,35,360,1,1,1,1
0,0,2,1,1,3750,2083,120,360,1,1,2,0
1,1,1,1,1,4166,3369,201,360,1,1,1,1
1,0,0,1,0,3167,0,74,360,1,1,1,1
1,0,1,1,1,4692,0,106,360,1,0,3,1
1,1,0,1,1,3500,1667,114,360,1,1,2,1
1,0,3,1,0,12500,3000,320,360,1,1,3,1
1,1,0,1,1,2275,2067,146,360,1,0,1,1
1,1,0,1,1,1828,1330,100,12,0,1,1,1
0,1,0,1,1,3667,1459,144,360,1,1,2,1
1,0,0,1,0,4166,7210,184,360,1,1,1,0
1,0,0,0,0,3748,1668,110,360,1,0,2,0
1,0,0,1,0,3600,0,80,360,1,1,1,1
1,0,0,1,0,1800,1213,47,360,1,1,1,0
1,1,0,1,1,2400,0,75,360,1,0,1,1
1,1,0,1,1,3941,2336,134,360,1,0,2,1
1,1,0,0,1,4695,0,96,12,1,1,1,1
0,0,0,1,0,3410,0,88,12,1,0,1,0
1,1,1,1,1,5649,0,44,360,1,0,1,1
1,1,0,1,1,5821,0,144,360,1,1,1,1
0,1,0,1,1,2645,3440,120,360,0,1,1,1
0,0,0,1,0,4000,2275,144,360,1,1,2,0
0,1,0,0,1,1928,1644,100,360,1,0,2,1
0,0,0,1,0,3086,0,120,360,1,0,2,0
0,0,0,1,0,4230,0,112,360,1,1,2,1
1,1,2,1,1,4616,0,134,360,1,0,1,1
0,1,1,1,1,11500,0,286,360,0,0,1,1
1,1,2,1,1,2708,1167,97,360,1,1,2,1
1,1,0,1,1,2132,1591,96,360,1,1,2,1
1,1,0,1,1,3366,2200,135,360,1,1,3,1

1,1,1,1,1,8080,2250,180,360,1,1,1,1
1,1,2,0,1,3357,2859,144,360,1,1,1,1
1,1,0,1,1,2500,3796,120,360,1,0,1,1
1,1,3,1,1,3029,0,99,360,1,1,1,1
1,1,0,0,1,2609,3449,165,180,0,1,3,1
1,1,1,1,1,4945,0,146,360,0,0,3,1
0,0,0,1,0,4166,0,116,360,0,1,2,1
1,1,0,1,1,5726,4595,258,360,1,1,2,1

Existing Problem:

It is important that credit card companies are able to recognize fraudulent credit card transactions so that customers are not charged for items that they did not purchase.

proposed solution:

Using Artificial Intelligence , Training a machine with previous data regarding fraud data there might be a very useful chances to companies to recognize fraudulent credit transactions.

Theoretical Analysis:

I have used Eclipse IDE and Weka software to train the machine and visualize the data.

Experimental Analysis:

package org1.ml;

```
import java.io.IOException;

import tech.tablesaw.api.Table;
import tech.tablesaw.plotly.Plot;
import tech.tablesaw.plotly.components.Figure;
import tech.tablesaw.plotly.components.Layout;
import tech.tablesaw.plotly.traces.BoxTrace;
import tech.tablesaw.plotly.traces.HistogramTrace;

public class Analysis {
IN: public static void main(String args[]) {
```

```

        System.out.println("Creditcard data Analysis");
    } }

```

OUT: Creditcard data Analysis

```

IN:      Table Creditcard_data =
Table.read().csv("C:\\Users\\Saketh\\eclipse-workspace\\org1.ml\\src\\main\\java\\org1.ml\\
\\fraud_dataset.csv");

        System.out.println(Creditcard_data.structure());

        System.out.println(Creditcard_data.summary());

```

OUT:

Index	Column Name	Column Type
0	Gender	INTEGER
1	Married	INTEGER
2	Dependents	INTEGER
3	Education	INTEGER
4	Self_Employed	INTEGER
5	ApplicantIncome	INTEGER
6	CoapplicantIncome	INTEGER
7	LoanAmount	INTEGER
8	Loan_Term	INTEGER
9	Credit_History_Available	INTEGER
10	Housing	INTEGER
11	Locality	INTEGER
12	Fraud_Risk	INTEGER

fraud_dataset.csv									
Summary	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_T
Count	827	827	827	827	827	827	827	827	827
sum	607	398	540	654	475	4311127	1228964	116518	
Mean	0.73397823458283	0.48125755743651766	0.6529625151148738	0.7908101571946797	0.5743651753325266	5212.970979443772	1486.0507859733966	140.89238210399026	338.128174
Min	0	0	0	0	0	150	0	9	
Max	1	1	3	1	1	81000	41667	700	
Range	1	1	3	1	1	80850	41667	691	
Variance	0.1954905709542645	0.24995095900758577	0.8757872177215122	0.16562973025990263	0.24476578900369195	31289628.52700187	7855956.815819599	6371.304384996675	5678.0975
Std. Dev	0.44214315662946146	0.4999509566023309	0.9358350376650322	0.4069763264121178	0.49473810142710045	5593.713303969186	2802.8479830022175	79.82045091952736	75.353156

```

IN:      //// Histogram

        Layout layout1=Layout.builder().title("Distribution of

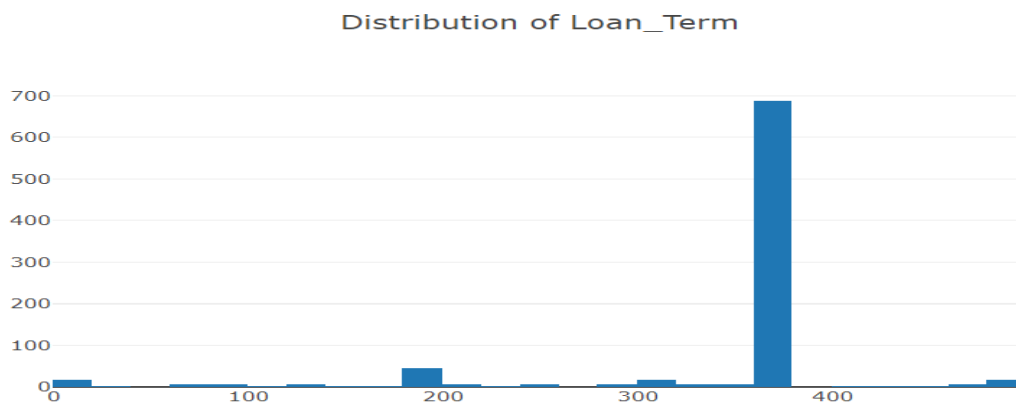
```

```

Loan_Term").build();
        HistogramTrace trace1=
HistogramTrace.builder(Creditcard_data.nCol("Loan_Term")).build();
        Plot.show(new Figure(layout1, trace1));

```

OUT:



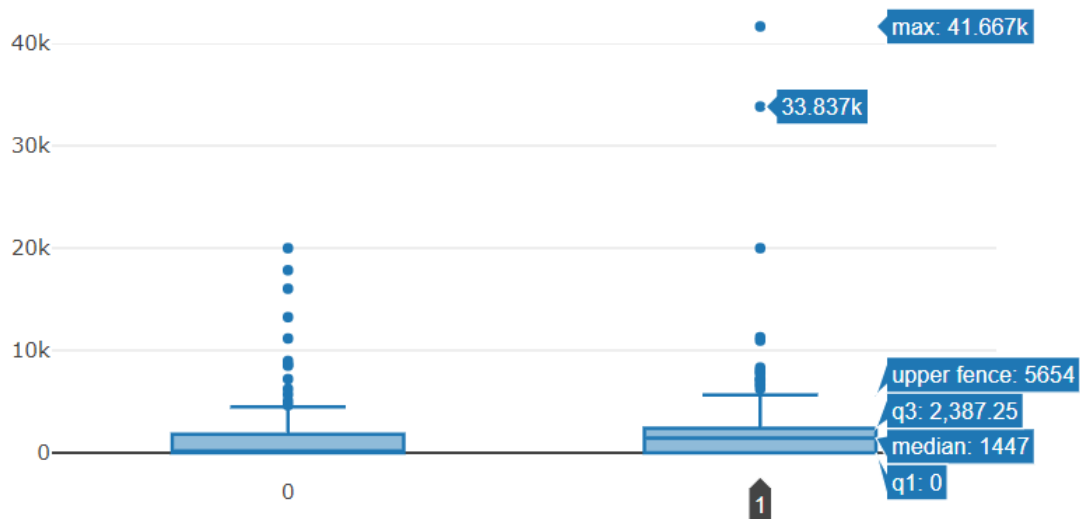
```

IN:      Layout layout3 = Layout.builder().title("fraud risk by ApplicantIncome").build();
        BoxTrace trace3
=BoxTrace.builder(Creditcard_data.categoricalColumn("Fraud_Risk"),
Creditcard_data.nCol("ApplicantIncome")).build();
        Plot.show(new Figure(layout3, trace3));

```

OUT:

fraud risk by CoapplicantIncome



IN:

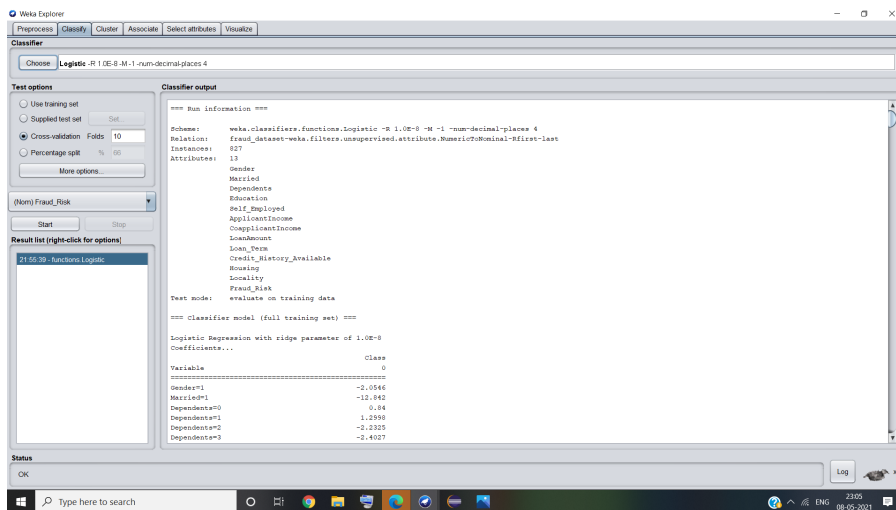
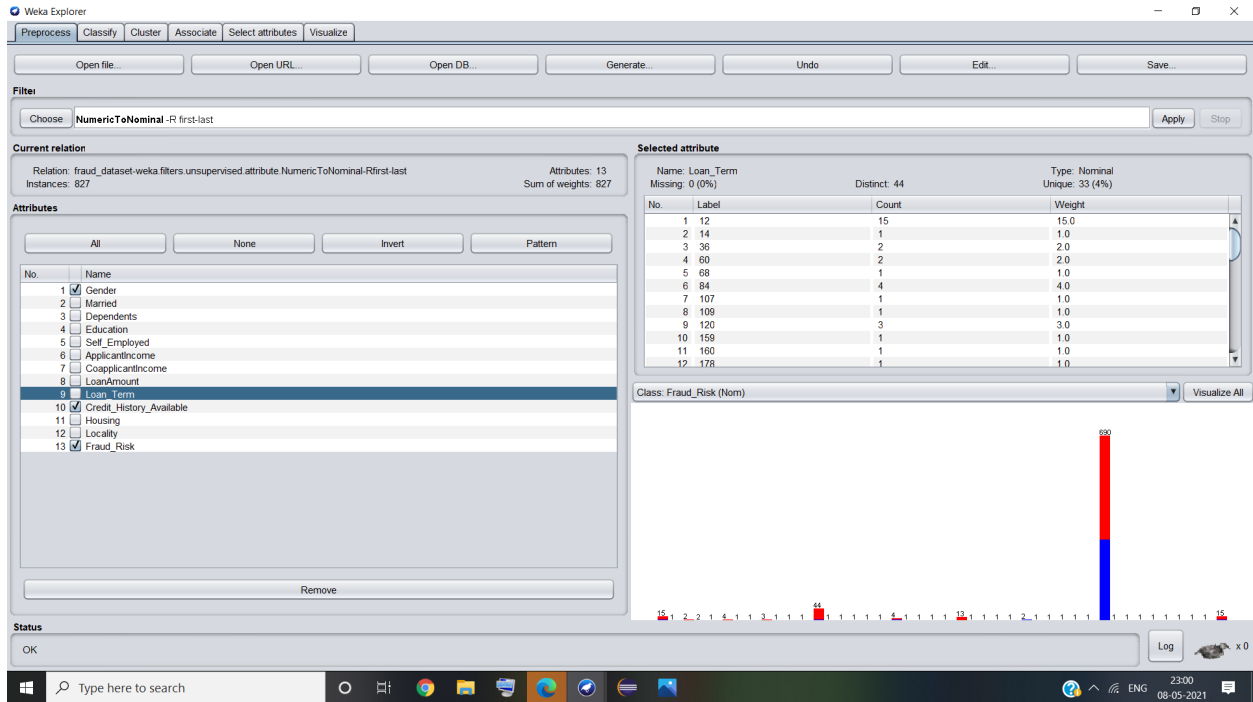
```
//linear Regression
LinearRegression lr=new LinearRegression();
lr.buildClassifier(dataset);

Evaluation lreval =new Evaluation(dataset);
lreval.evaluateModel(lr,dataset);
System.out.println(lreval.toSummaryString());
```

OUT:

Correlation coefficient	0.8428
Mean absolute error	0.1727
Root mean squared error	0.2659
Relative absolute error	35.3768 %
Root relative squared error	53.827 %
Total Number of Instances	827

// Log Regression



//Logistic using weka explorer

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **Logistic -R 1.0E-9-M-1-num-decimal-places 4**

Test options

☐ Use training set
☐ Supplied test set
☒ Cross-validation Folds **10**
☐ Percentage split % **66**
 More options...

(Nom) Fraud_Risk

Start Stop

Result list (right-click for options)

21:55:39 - Functions Logistic

Classifier output

```

Locality=3                                0.5563

Time taken to build model: 3.67 seconds

=== Evaluation on training set ===

Time taken to test model on training data: 0.11 seconds

=== Summary ===
Correctly Classified Instances      827          100 %
Incorrectly Classified Instances      0           0 %
Kappa statistic                    1.000
Mean absolute error                  0.000
Root mean squared error              0.000
Relative absolute error              0.000 %
Root relative squared error          0.000 %
Total Number of Instances          827

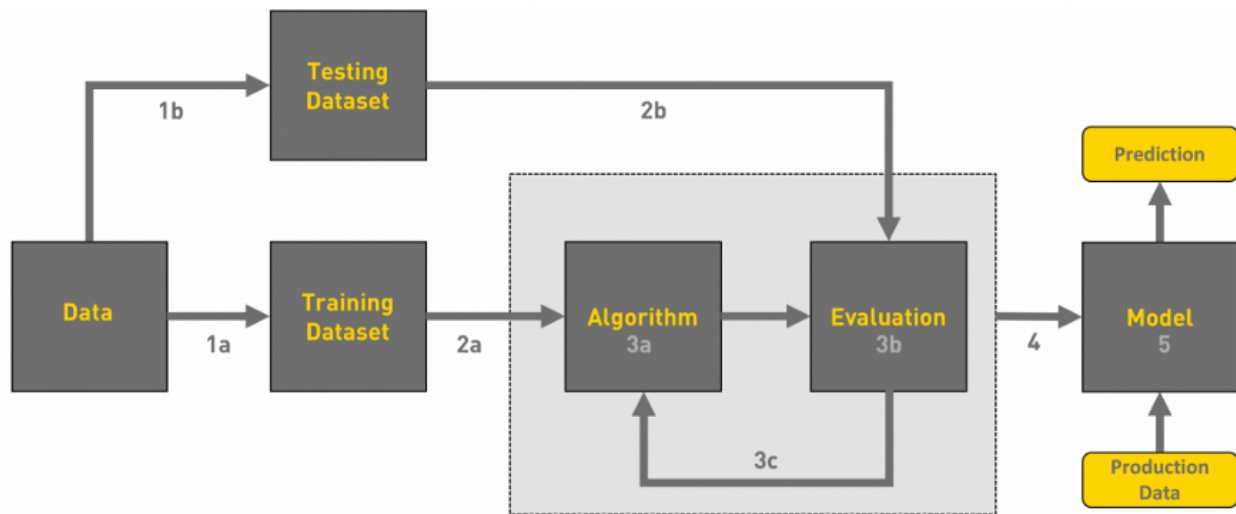
=== Detailed Accuracy By Class ===
               TP Rate  FP Rate  Precision  Recall   F-Measure  MDC     ROC Area  PRC Area  Class
               1.000    0.000    1.000    1.000    1.000    1.000    1.000    1.000    0
               1.000    0.000    1.000    1.000    1.000    1.000    1.000    1.000    1
Weighted Avg.  1.000    0.000    1.000    1.000    1.000    1.000    1.000    1.000

=== Confusion Matrix ===
      a   b   <-- classified as
350   0 |   a = 0
 0 477 |   b = 1
  
```

Status

OK Log x 0

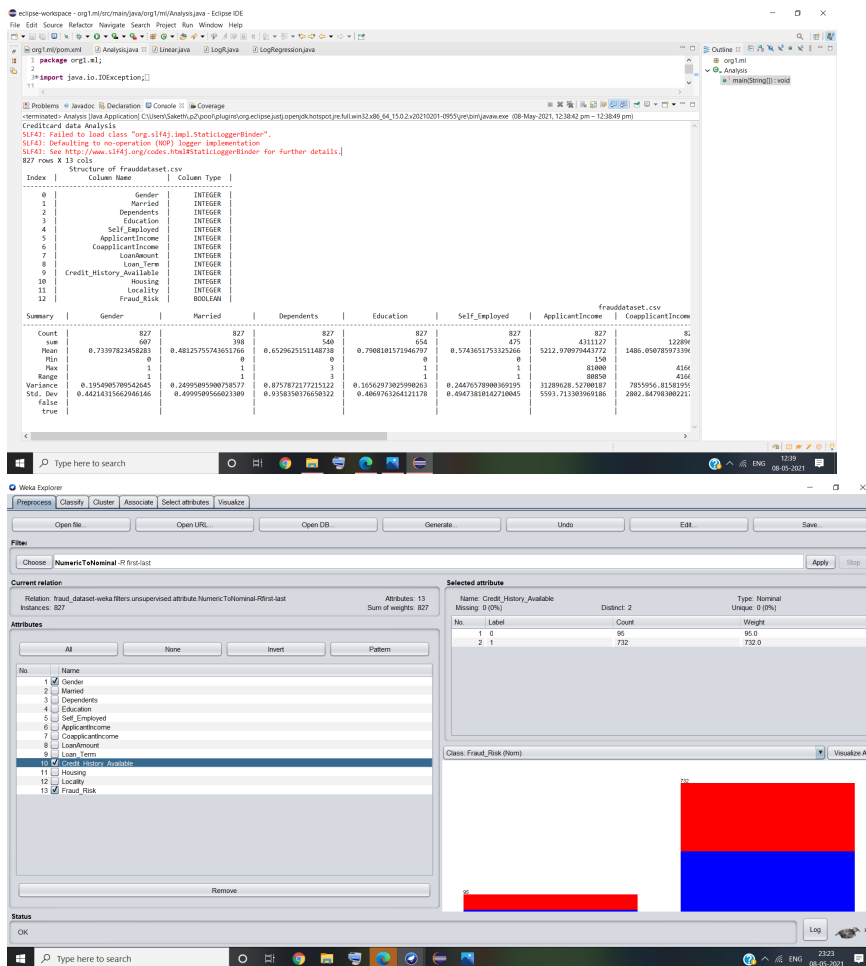
FLOW CHART:



Result:

A histogram showing the frequency of word counts in the titles of the 100 most popular books. The x-axis represents the number of words (0 to 400), and the y-axis represents the frequency (0 to 700). The distribution is highly right-skewed, with a peak frequency of approximately 680 for titles containing 350-400 words.

Number of Words (Bin)	Frequency
0-25	20
25-50	5
50-75	10
75-100	10
100-125	10
125-150	10
150-175	10
175-200	50
200-225	10
225-250	10
250-275	10
275-300	10
300-325	20
325-350	10
350-375	680
375-400	10
400-425	10
425-450	20



Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose Logistic -R 1.0E-8 -M -1 -num-decimal-places 4

Test options

☐ Use training set

☐ Supplied test set

☒ Cross-validation Folds 10

☐ Percentage split % 66

More options...

(Nom) Fraud_Risk

Start Stop

Result list (right-click for options)

21:55:39 - functions.Logistic

Classifier output

```
=== Run information ===

Scheme:      weka.classifiers.functions.Logistic -R 1.0E-8 -M -1 -num-decimal-places 4
Relation:    fraud_dataset-weka.filters.unsupervised.attribute.NumericToNominal-Rfirst-last
Instances:   827
Attributes:  13
Gender
Married
Dependents
Education
Self_Employed
ApplicantIncome
CoapplicantIncome
LoanAmount
Loan_Term
Credit_History_Available
Housing
Locality
Fraud_Risk

Test mode:   evaluate on training data

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

Logistic Regression with ridge parameter of 1.0E-8
Coefficients...

Variable          Class
-----
Gender=1          -2.0546
Married=1         -12.642
Dependents=0       0.84
Dependents=1      1.2998
Dependents=2      -2.2325
Dependents=3      -2.4027
```

Status

OK

Log

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose Logistic -R 1.0E-8 -M -1 -num-decimal-places 4

Test options

☐ Use training set

☐ Supplied test set

☒ Cross-validation Folds 10

☐ Percentage split % 66

More options...

(Nom) Fraud_Risk

Start Stop

Result list (right-click for options)

21:55:39 - functions.Logistic

Classifier output

```
Locality=3          0.5563

Time taken to build model: 3.67 seconds

=== Evaluation on training set ===

Time taken to test model on training data: 0.11 seconds

=== Summary ===

Correctly Classified Instances      827          100 %
Incorrectly Classified Instances    0           0 %
Kappa statistic                     1
Mean absolute error                 0
Root mean squared error             0
Relative absolute error             0 %
Root relative squared error         0 %
Total Number of Instances          827

=== 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     0
      1.000   0.000   1.000     1.000   1.000     1.000   1.000     1.000     1
Weighted Avg.   1.000   0.000   1.000     1.000   1.000     1.000   1.000     1.000

=== Confusion Matrix ===

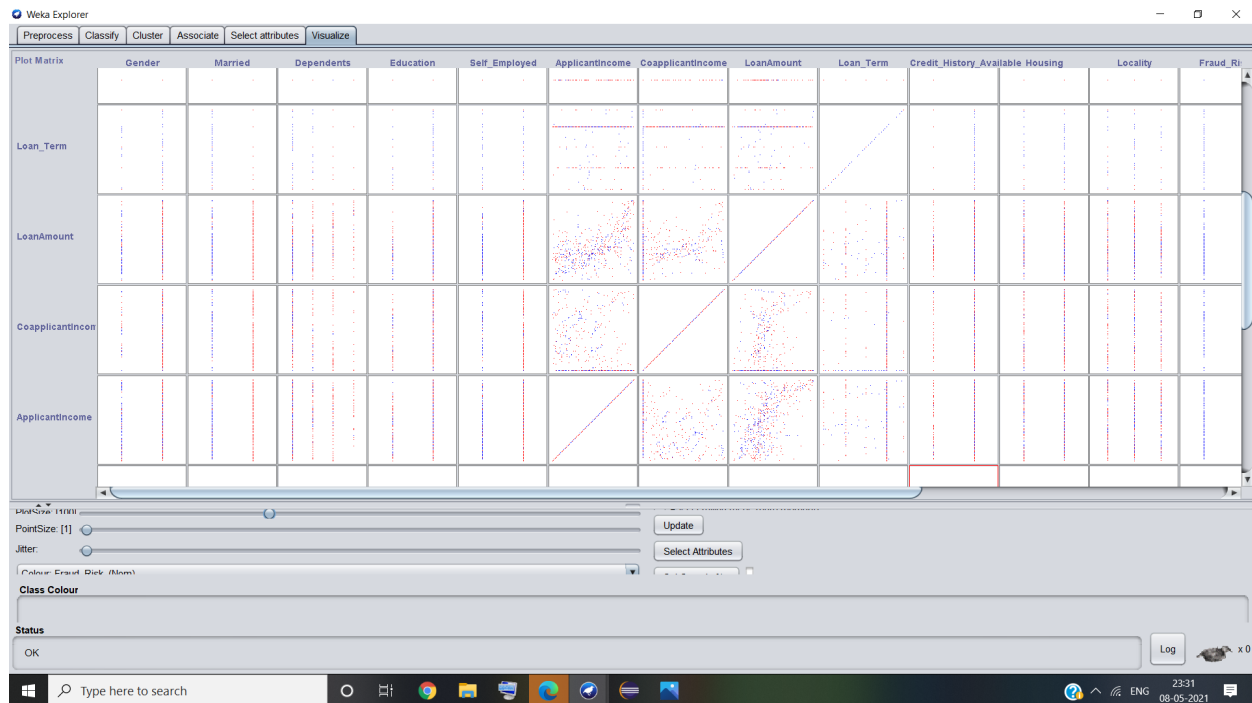
  a  b  <-- classified as
350  0  |  a = 0
  0 477 |  b = 1
```

Status

OK

Log

Visualization using weka..



Advantages and Disadvantages:

Using this model not always give the accurate data. There is no thinking power for a machine. If the trained data is different from present applicable data there would be lot of issues in dealing with the Amount.

But by using this models large data can be calculated in a short period of time.

Applications:

This data is applicable for banking sectors in the case of predicting fraud transactions

Conclusion:

here we have chosen a dataset and trained the machine with the data. by the the model can understand what is the fraud transaction and which is approved transaction. from the data provided I concluded That data is not linearly regressed.

Future scope:

In future machines are going to lead the world.Man power is going to replace with machine intelligence....