

CAR RESALE VALUE PREDICTION USING MACHINE LEARNING TECHNIQUES

CONTENTS

1. INTRODUCTION

a. OVERVIEW b. PURPOSE

2. LITERATURE SURVEY

a. EXISTING PROBLEM b. PROPOSED SOLUTION

3. THEORITICAL ANALYSIS

a. BLOCK DIAGRAM b. SOFTWARE DESIGNING

4. EXPERIMENTAL INVESTIGATIONS

5. RESULT

6. ADVANTAGES AND DISADVANTAGES

7. APPLICATIONS

8. CONCLUSION

9. APPENDIX a. SOURCE CODE

1.INTRODUCTION

OVERVIEW

The prices of the new manufactured cars in the automobile industry by the manufacturer with some extra cost which are incurred by the Government in the form of taxes. So, customers buying a new car can be assured of the money which they are investing to be worthy. But due to the increased price of new cars and the problems of customers to buy new cars due to improper financial status, used car sales are on increase in demand. Due to the high raise in demand for the used cars, there is a need for a car resale price prediction system to effectively determine the quality and worth of the car using different features, which are very much required while assessing the price of resale car.

PURPOSE

Client will be able to efficiently determine the price that the chosen resale vehicle must cost. So, that people can invest money according to their preferences. In this project, we suggest a solution using Machine Learning algorithms implemented in JAVA.

2. LITERATURE SURVEY

EXISTING PROBLEM

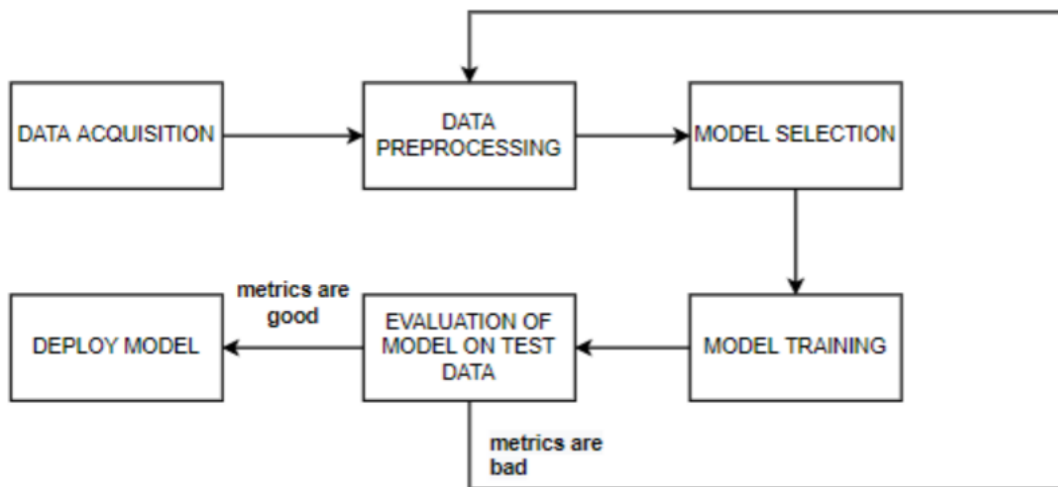
Existing models use Linear Regression algorithm, Decision Tree algorithms as the basic models. Problem with these models is that, it takes a huge amount of time for training the model, and if any attributes are categorical, proper encoding techniques must be used, and yet due to high correlated data, this technique produces mid-level percentage of accuracy (70 - 80%), and if overfitting took place if we tried to remove the duplicate values.

PROPOSED SOLUTION

Proposed model is built on Random Forest regressor algorithm, instead of using conventional Linear Regression, and Decision Trees models. This model takes less amount of training and the correlation coefficient is high, and low error metrics, thus making suitable for modelling.

3.THEROTICAL ANALYSIS

BLOCK DIAGRAM



SOFTWARE DESIGNING

1. JAVA 16.0
2. WEKA 3.8
3. MAVEN
4. ECLIPSE IDE

4. EXPERIMENTAL INVESTIGATIONS

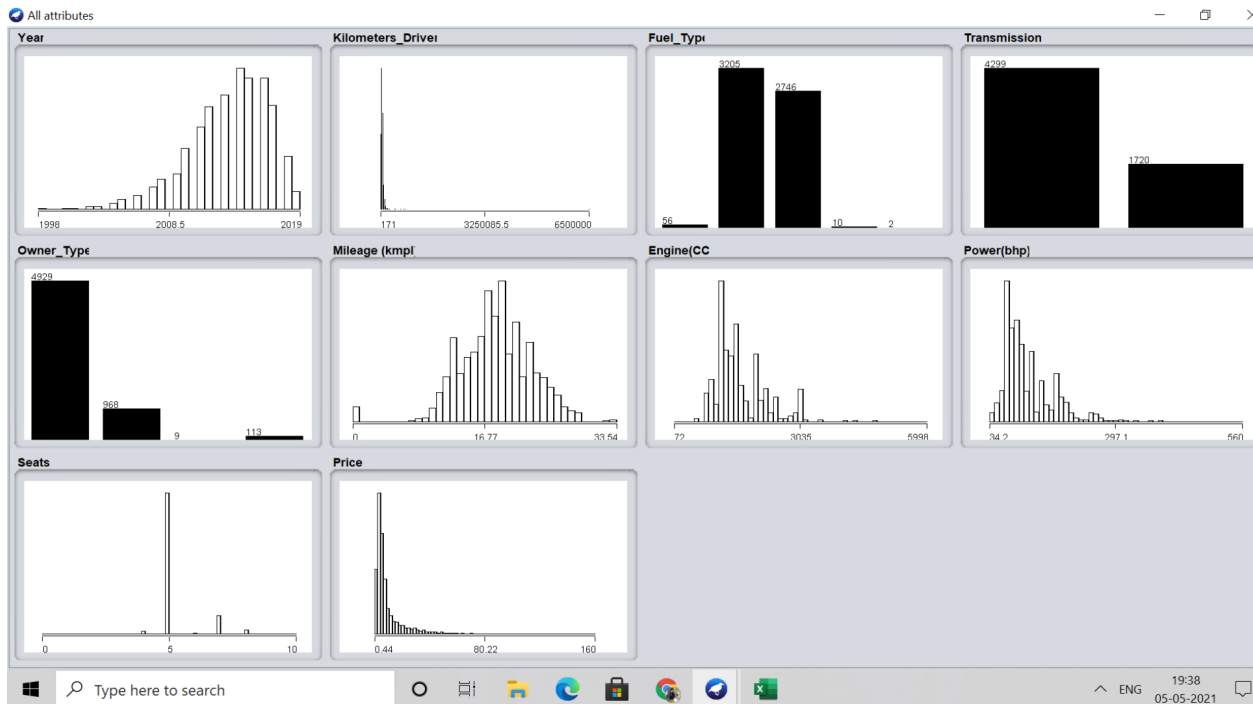
There are separate datasets, one for the training, and other for testing the built model.

TRAINING DATASET: Training data consist of 12 independent variables and 1 dependent variable. Many attributes consisted of much missing data. The following table depicts the details of train data and test data.

Train Data		
S.no	Attribute	Missing data
1	Name	0 (0%)
2	Location	0 (0%)
3	Year	0 (0%)
4	Kilometers_Driven	0 (0%)
5	Fuel_Type	0 (0%)
6	Transmission	0 (0%)
7	Owner_Type	0 (0%)
8	Mileage (kmpl)	0 (0%)
9	Engine (CC)	36 (1%)
10	Power (bhp)	143 (2%)
11	Seats	42 (1%)
12	New_Price	5192 (86%)
13	Price	0 (0%)

Test Data		
S.no	Attribute	Missing data
1	Name	0 (0%)
2	Location	0 (0%)
3	Year	0 (0%)
4	Kilometers_Driven	0 (0%)
5	Fuel_Type	0 (0%)
6	Transmission	0 (0%)
7	Owner_Type	0 (0%)
8	Mileage (kmpl)	0 (0%)
9	Engine (CC)	10 (1%)
10	Power (bhp)	32 (3%)
11	Seats	11 (1%)
12	New_Price	1052 (85%)

There are total of 6020 rows and 13 columns in the training dataset.



There are total of 1235 rows and 12 columns in the testing dataset.

Meaning of each attribute is shown below

Name	The brand and model of the car
Location	The location in which the car is being sold or is available for purchase
Year	The year or edition of the model
Kilometers_Driven	The total kilometers driven in the car by the previous owner(s) in KM
Fuel_Type	The type of fuel used by the car
Transmission	The type of transmission used by the car
Owner_Type	Whether the ownership is Firsthand, Second hand or other
Mileage	The standard mileage offered by the car company in kmpl or km/kg
Engine	The displacement volume of the engine in cc
Power	The maximum power of the engine in bhp
Seats	The number of seats in the car
New_Price	Price of new model
Price	The price of the used car in INR Lakhs

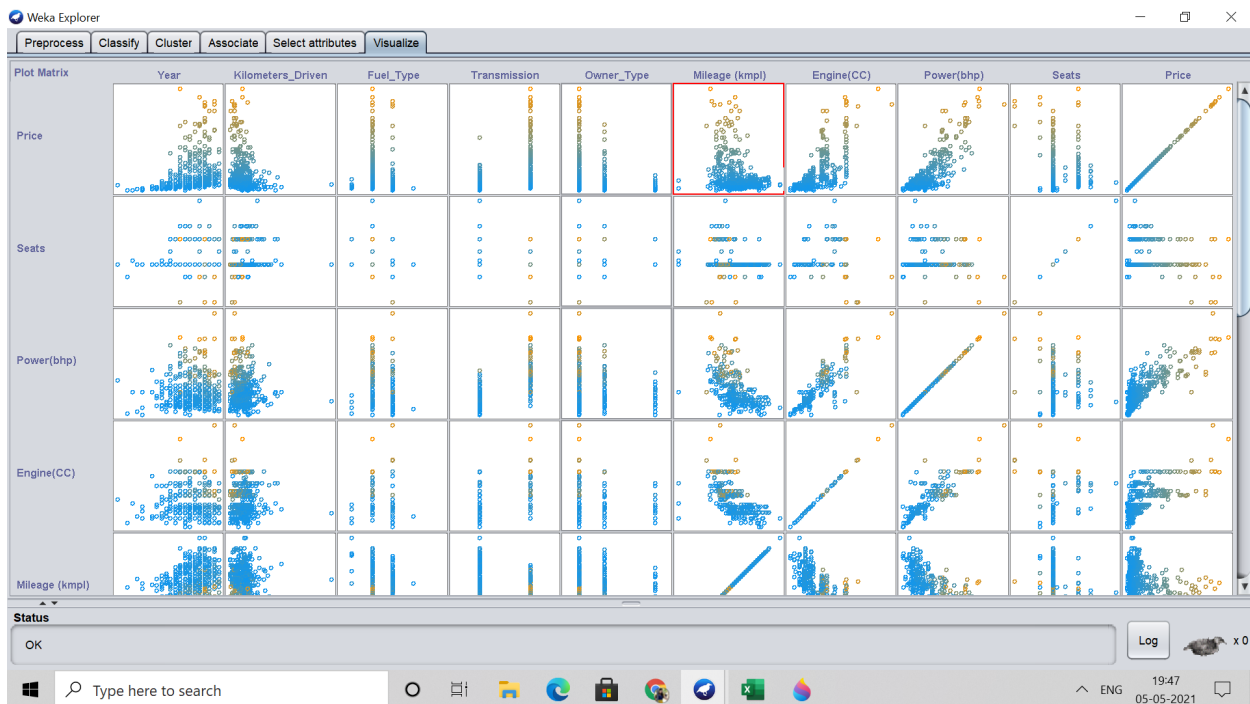
DATA CLEANING:

As depicted in the train dataset properes table, there are many missing values in the dataset. These missing will affect the accuracy in building the model. So, there is a need to remove the missing data, or to fill the missing data. "RemoveMissingValue" filter has been used in-order to replace missing value with median. And in for categorical variables such as Seats, etc. Missing values were replaced by mode.

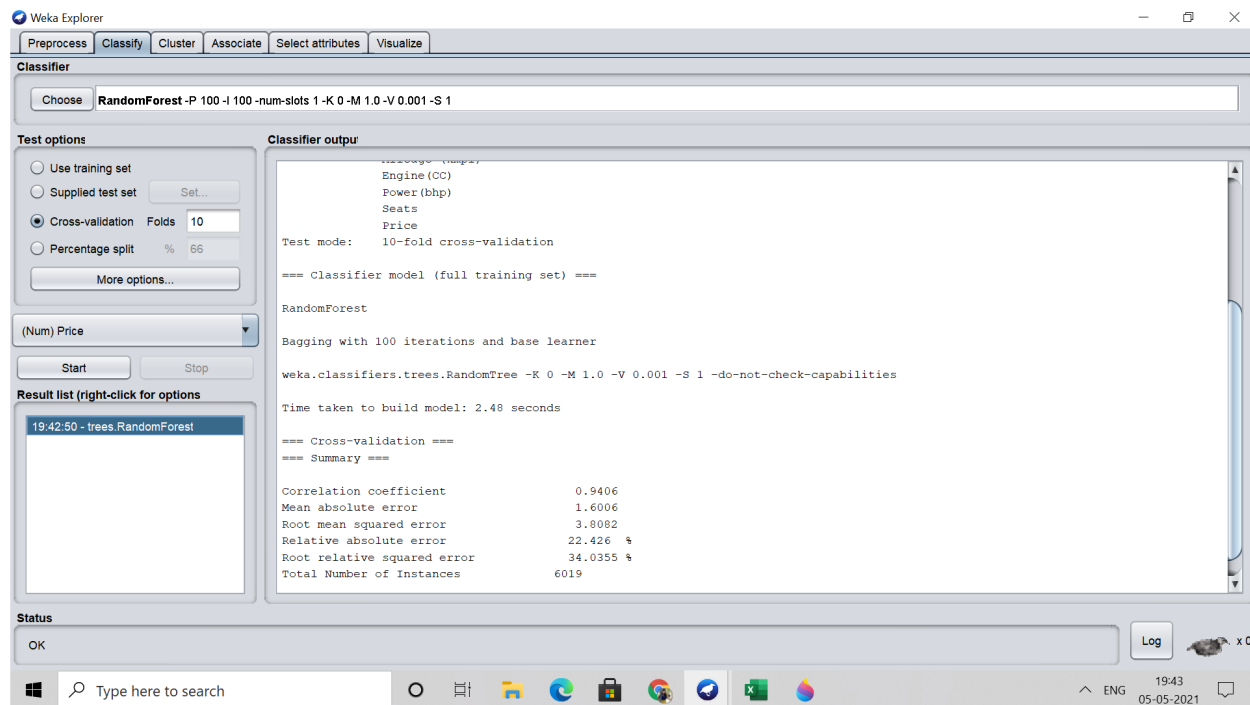
DATA VISUALIZATION:

Visualizing the data is the most crucial part in analyzing the model. We use WEKA explorer to visualize the data in the form of Histogram plots. Below figures shows the modified data after using the "RemoveMissingValue" filter

and removing unwanted attributes such "name","Location","New_Price".

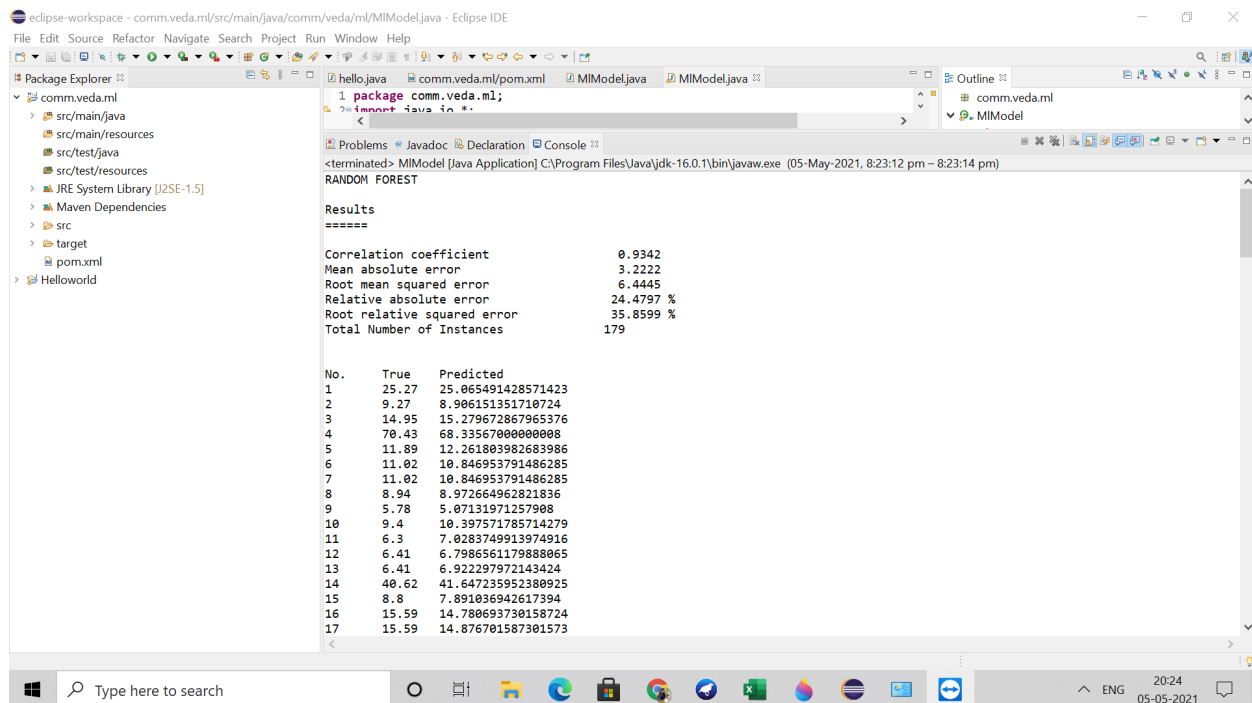


Random Forest regressor is used ,the model showed an accuracy measure of 94%, with the following metric evaluations show below



5. RESULT

From the above analysis, the most suitable algorithm for building the model is to use Random Forest regressor. The below figure shows the output of the random forest regressor, which was run on the test dataset. And actual values to that of the predicted values is also being shown in the figure.



6. ADVANTAGES AND DISADVANTAGES

Advantages:

1. Accuracy.
2. Generic model i.e. can be used for any type of raw data.
3. Robust.
4. Feature Engineering is not needed.

Disadvantages:

1. Sometimes may undergo overfitting, if dataset is small.

7. APPLICATIONS

This model can be used in following sectors:

1. Petrol price prediction
2. Land resale value prediction

8. CONCLUSION

The best accuracy was found out in Random Forest Regressor (94.4%) while linear regression algorithm shows 84.4% accuracy. While new features can be created via feature engineering which may help in predicting the target variable. overall, Random Forest Regressor classifier provides the best result in terms of accuracy for the given dataset, without any feature engineering needed. Because of its simplicity and the fact that it can be implemented relatively easy and quick.

11. APPENDIX

SOURCE CODE

comm.veda.ml/pom.xml

```
<project xmlns="http://maven.apache.org/POM/4.0.0"
xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
xsi:schemaLocation="http://maven.apache.org/POM/4.0.0
https://maven.apache.org/xsd/maven-4.0.0.xsd">
  <modelVersion>4.0.0</modelVersion>
  <groupId>org</groupId>
  <artifactId>comm.veda.ml</artifactId>
  <version>0.0.1-SNAPSHOT</version><dependencies>
  <dependency>
    <groupId>nz.ac.waikato.cms.weka</groupId>
    <artifactId>weka-stable</artifactId>
    <version>3.8.0</version>
  </dependency>
  <dependency>
    <groupId>tech.tablesaw</groupId>
    <artifactId>tablesaw-core</artifactId>
    <version>0.38.1</version>
  </dependency><dependency>
    <groupId>tech.tablesaw</groupId>
    <artifactId>tablesaw-jsplot</artifactId>
    <version>0.38.1</version>
  </dependency><!-- Thanks for using https://jar-download.com --></dependencies>
  <properties>
    <maven.compiler.source>1.8</maven.compiler.source>
    <maven.compiler.target>1.8</maven.compiler.target>
  </properties>
</project>
```

MIModel.java

```
package comm.veda.ml;
import java.io.*;
import java.io.BufferedReader;
import java.io.FileReader;
import java.util.Random;

import weka.classifiers.Classifier;
import weka.classifiers.Evaluation;
import weka.classifiers.bayes.NaiveBayes;
import weka.classifiers.functions.LinearRegression;
import weka.classifiers.trees.RandomForest;
import weka.core.Instances;
public class MIModel {
    private static String getCurrPath() {
        String
path=System.getProperty("user.dir")+"\\src\\main\\java\\comm\\veda\\ml";
        //System.out.println(path);
        return path.replace("\\", "\\");
    }
    public static String trainPath="";
    public static String testPath="";
    private static void p6(Instances testDf, Instances trainDf) throws Exception {
        // TODO Auto-generated method stub
        System.out.println("-----RANDOM
FOREST-----");
        RandomForest rf=new RandomForest();
        rf.setNumIterations(200);
        Evaluation evaluation = new Evaluation(trainDf);
        evaluation.crossValidateModel(rf, testDf, 10, new Random(1));
        rf.buildClassifier(testDf);
        System.out.println(evaluation.toSummaryString("\nEvaluation
Results\n*****\n", true));
        System.out.println();
        System.out.println("S.No.\tTrueValue\tPredictedValue");
```

```

        for(int i=0;i<testDf.numInstances();i++) {
            String trueValue;
            trueValue=testDf.instance(i).toString(testDf.classIndex());
            double predValue=rf.classifyInstance(testDf.instance(i));
            System.out.println((i+1)+"\t"+trueValue+"\t"+predValue);
        }
    }

    public static void main(String args[]) throws Exception{
        String trainFileName="\\\\\\train.arff";
        String testFileName="\\\\\\test1.arff";
        trainPath=getCurrPath()+trainFileName;
        //System.out.println(trainPath);
        testPath=getCurrPath()+testFileName;
        BufferedReader brTrain=new BufferedReader(new FileReader(trainPath));
        Instances trainDf=new Instances(brTrain);
        BufferedReader brTest=new BufferedReader(new FileReader(testPath));
        Instances testDf=new Instances(brTest);
        testDf.setClassIndex(testDf.numAttributes()-1);
        trainDf.setClassIndex(trainDf.numAttributes()-1);
        brTrain.close();
        brTest.close();
        p6(testDf,trainDf);
    }

}

```