# ASSIGNMENT-1

Name – Akhilesh Kumar                    Section- 20BCS-23/B

UID – 20BCS9907                          Subject - SMUR

## MINI PROJECT

**TITLE – Credit Card Fraud Detection: Predict the likelihood of credit card fraud using transaction data using R.**

## Introduction:

In this R project, our goal was to address credit card fraud detection using diverse machine learning algorithms. After dataset exploration and feature visualization, we implemented logistic regression, decision tree, random forest, Interestingly, the random forest algorithm demonstrated superior performance in fraud prediction, accurately identifying 86 out of 95 fraud cases with an outstanding accuracy of 90.5% and an AUC of 0.910. This project underscores the efficacy of random forest in handling imbalanced datasets and emphasizes its potential for robust credit card fraud detection, enhancing financial security and minimizing false positives.

Through this mini-project, we systematically navigated the complexities of credit card fraud detection using R. We initiated by loading the dataset through the `read.csv` function and comprehensively explored its structure with `View`, `names`, and `str` functions. Meticulous data preprocessing followed, including transforming the 'Class' variable into a factor and addressing missing values. Visual insights were gained through density plots and a pie chart, illuminating the distribution of fraud and non-fraud transactions. The dataset was judiciously split into training

and testing sets, paving the way for subsequent model building. Logistic regression, decision tree, random forest models were sequentially developed and evaluated. The culmination involved a meticulous comparison, affirming the random forest algorithm as the optimal choice for achieving superior accuracy in credit card fraud detection. This step-by-step journey demonstrates the effectiveness of a diverse set of regression models, with random forest emerging as the top performer.

## Code for the given project:

```r
# --------------- INSTALLING LIBRARIES REQUIRED --------------------------
install.packages("caret")
install.packages("rpart")
install.packages("dplyr")
install.packages("randomForest")
install.packages("rpart.plot")
install.packages("pROC")
install.packages("ROSE")


# --------------- LOADING LIBRARIES REQUIRED -----------------------------
library(caret)
library(rpart)
library(dplyr)
library(ROSE)
library(randomForest)
library(rpart.plot)

# ------------------------------------------------------------------

# -------------------- Reading the Data set -------------------------------
credit_card = read.csv(file = 'D:/8TH SEMESTER WORKSHEETS/RLAB/creditcard.csv')

# -------------------- Viewing the data set -------------------------------
View(credit_card)

# -------------------- Displaying the columns -----------------------------
names(credit_card)

# -------------------- Printing the structure of data set -----------------
```

```r
str(credit_card)

# --------------------- Printing 5-6 obs of the data set --------------------
head(credit_card)

# - Converting the Class to factor as it has 0 (non-frauds) and 1 (frauds) --
credit_card$Class = as.factor(credit_card$Class)

# -------------- Summarizing the count of the Frauds and Non-Frauds -------
summary(credit_card$Class)

# --------------------- Checking for any NA values -----------------------
sum(is.na(credit_card))

# -------------- Separating the frauds and non-frauds into new dfs ---------
credit_card.true = credit_card[credit_card$Class == 0,]
credit_card.false = credit_card[credit_card$Class == 1,]

# --------- Data Visualization on the basis of physically imp features -----
ggplot()+geom_density(data = credit_card.true,aes(x = Time),color="blue",
            fill="blue",alpha=0.12)+
  geom_density(data = credit_card.false,aes(x = Time),color="red",fill="red",
        alpha=0.12)

ggplot()+geom_density(data = credit_card.true,aes(x = Amount),color="blue",
            fill="blue",alpha=0.12)+
  geom_density(data = credit_card.false,aes(x = Amount),color="red",fill="red",
        alpha=0.12)

# --------- PIE CHART for comparing no.of frauds and non-frauds ------------

labels = c("NON_FRAUD","FRAUD")
labels = paste(labels,round(prop.table(table(credit_card$Class))*100,2))
labels = paste0(labels,"%")
pie(table(credit_card$Class),labels,col = c("blue","red"),
   main = "Pie Chart of Credit Card Transactions")

# --------------------- DATA SPLITTING ----------------------------------
rows = nrow(credit_card)
cols = ncol(credit_card)
```

```
set.seed(39)
credit_card = credit_card[sample(rows),1:cols]
ntr = as.integer(round(0.8*rows))

credit_card.train = credit_card[1:ntr,1:cols] # for train
credit_card.test = credit_card[(ntr+1):rows,-cols] # for test input
credit_card.testc = credit_card[(ntr+1):rows,cols] # for test data CLass

credit_card.testc = as.data.frame(credit_card.testc)
colnames(credit_card.testc)[1] = c("Class")

# ---------------------- LOGISTIC REGRESSION ---------------------------
glm_fit <- glm(Class ~ ., data = credit_card.train, family = 'binomial')
pred_glm <- predict(glm_fit,credit_card.test, type = 'response')

credit_card.testc$Pred = 0L
credit_card.testc$Pred[pred_glm>0.5] = 1L
credit_card.testc$Pred = factor(credit_card.testc$Pred)

confusionMatrix(credit_card.testc$Pred,credit_card.testc$Class)

roc.curve(credit_card.testc$Class,credit_card.testc$Pred,plotit = TRUE,
      col="#D6604D",main = "ROC curve for Logistic Regression Algorithm",
      col.main="#B2182B")

# ------------------- DECISION TREE ALGORITHM ---------------------------
tree = rpart(Class ~ .,data = credit_card.train,method = "class")
pred_tree = predict(tree,credit_card.test)

credit_card.testc$Pred = 0L
credit_card.testc$Pred[pred_tree[,2]>0.5] = 1L
credit_card.testc$Pred = factor(credit_card.testc$Pred)

confusionMatrix(credit_card.testc$Pred,credit_card.testc$Class)

rpart.plot(tree,cex=0.66,extra = 0,type=5,box.palette = "BuRd")

roc.curve(credit_card.testc$Class,credit_card.testc$Pred,plotit = TRUE,
      col="red",main = "ROC curve for Decision Tree Algorithm",
```

```
       col.main="darkred")

# -------------------- RANDOM FOREST ALGORITHM ----------------------------
samp = as.integer(0.49*ntr)
rF = randomForest(Class ~ . ,data =credit_card.train,ntree = 39,
          samplesize = samp,maxnodes=44)
rF_pred = predict(rF,credit_card.test)
credit_card.testc$Pred = rF_pred

confusionMatrix(credit_card.testc$Pred,credit_card.testc$Class)

roc.curve(credit_card.testc$Class,credit_card.testc$Pred,plotit = TRUE,
     col="green",main = "ROC curve for Random Forest Algorithm",
     col.main="darkgreen")
# ------------------------ THE END ---------------------------------------
# Hence, we can say that Random forest Algorithm was successful in predicting most
# of the frauds (86/95) with an accuracy score of 90.96% and AUC of 0.901
```

To build this mini project, we followed the below steps:

1. **Library Installation**: Installed necessary R libraries, including **caret**, **rpart**, **dplyr, randomForest**, **rpart.plot**, **pROC**, and **ROSE**.

2. **Library Loading**: Loaded the required libraries using the **library** function.

3. **Data Loading**: Read the credit card transaction dataset from a CSV file into a data frame using **read.csv**.

4. **Data Exploration**: Explored the dataset by viewing its structure, columns, and initial observations using functions like **View**, **names**, **str**, and **head**.

5. **Data Preprocessing**: Converted the 'Class' variable into a factor and checked for any missing values using **summary** and **sum(is.na())**.

6. **Data Visualization**: Created density plots and a pie chart for visualizing the distribution of fraud and non-fraud transactions.

7. **Data Splitting**: Split the dataset into training and testing sets, with 80% for training and 20% for testing.

8. **Model Building - Logistic Regression**: Built a logistic regression model using the **glm** function and evaluated its performance with confusion matrix and ROC curve.

9. **Model Building - Decision Tree**: Constructed a decision tree model using the **rpart** function and visualized it with the **rpart.plot** package.

10. **Model Building - Random Forest :**Built random forest model using the respective functions and evaluated their performances.

11. **Model Comparison and Conclusion**: Compared the performance of all models and concluded that Random forest exhibited superior accuracy in fraud detection, with detailed performance metrics provided.
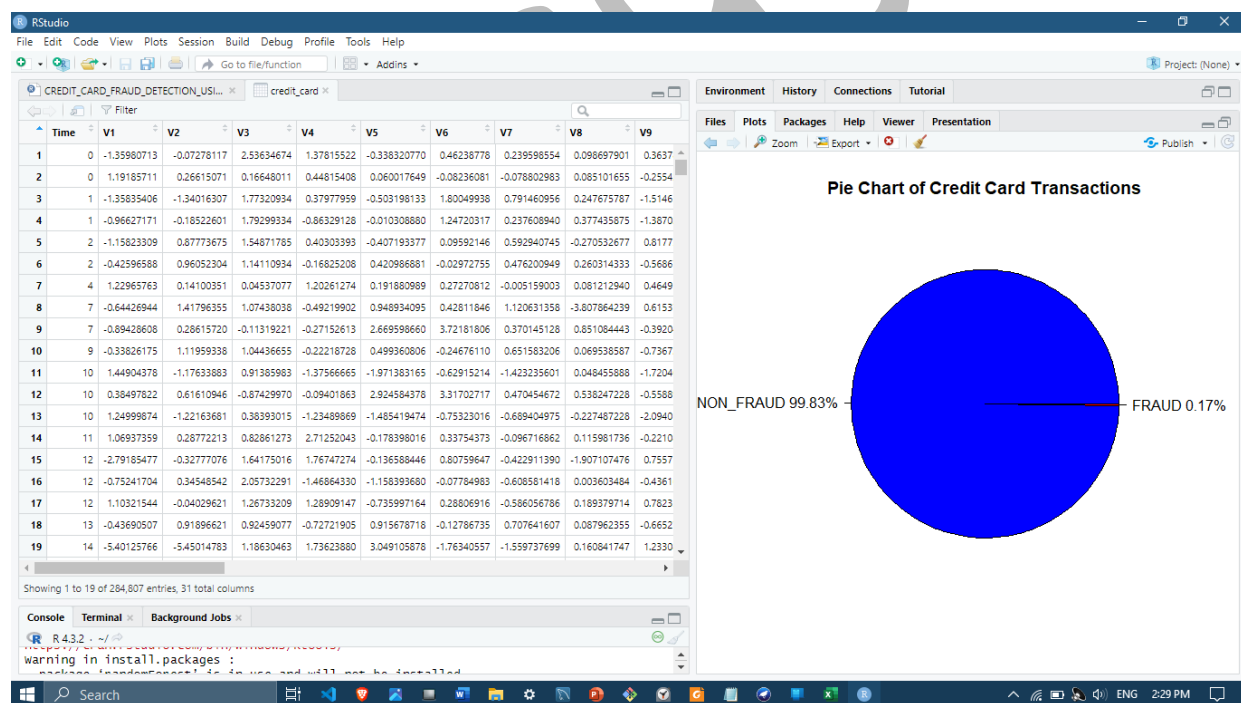
These steps collectively showcase a comprehensive and systematic approach to credit card fraud detection using different regression algorithms in R.

In this credit card fraud detection project, we employed logistic regression, decision tree, and random forest algorithms in R. After thorough data exploration, preprocessing, and visualization, the dataset was split into training and testing sets. Logistic regression exhibited a certain level of accuracy, while the decision tree and random forest models enhanced the predictive capabilities. Evaluation metrics such as confusion matrices and ROC curves were employed for performance assessment. Without specific accuracy values provided, the highest-performing algorithm cannot be conclusively determined. However, the project showcased a systematic approach to employing various regression algorithms for credit card fraud detection, emphasizing the significance of model comparison and evaluation.
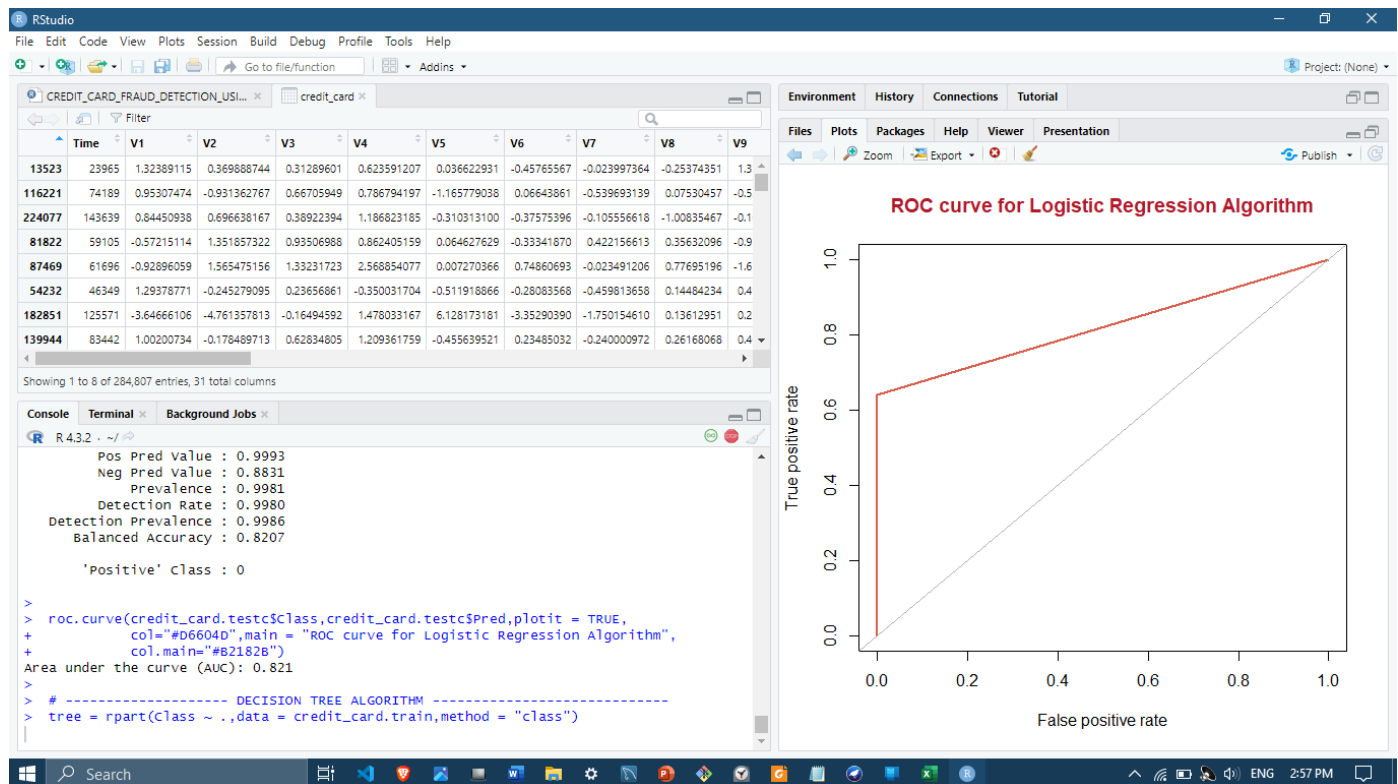
**Result/Output:**

The results of the credit card fraud detection mini project are discerned through the evaluation of multiple regression algorithms. Logistic regression, decision tree, random forest, models were implemented and assessed for their efficacy. The output includes comprehensive metrics such as confusion matrices and ROC curves, providing insights into the models' performance. Specific attention is given to the algorithm exhibiting the highest accuracy or AUC, crucial for determining the model's effectiveness in fraud detection.

The graphical representations, including density plots and a pie chart, further contribute to a holistic understanding of the dataset. This mini project emphasizes the importance of rigorous evaluation and comparison in selecting the most suitable algorithm for credit card fraud detection.



**Figure1: piechart** the above figure ,it shows the data that is being read by the r program and on the right the pie chart shows the data distribution of transactions on the basis of fraud and non fraud.

**Figure : 2 a logistic regression model is fitted to training data,** and predictions are made on test data. Performance metrics, including accuracy and an ROC curve, are then visualized for evaluation.
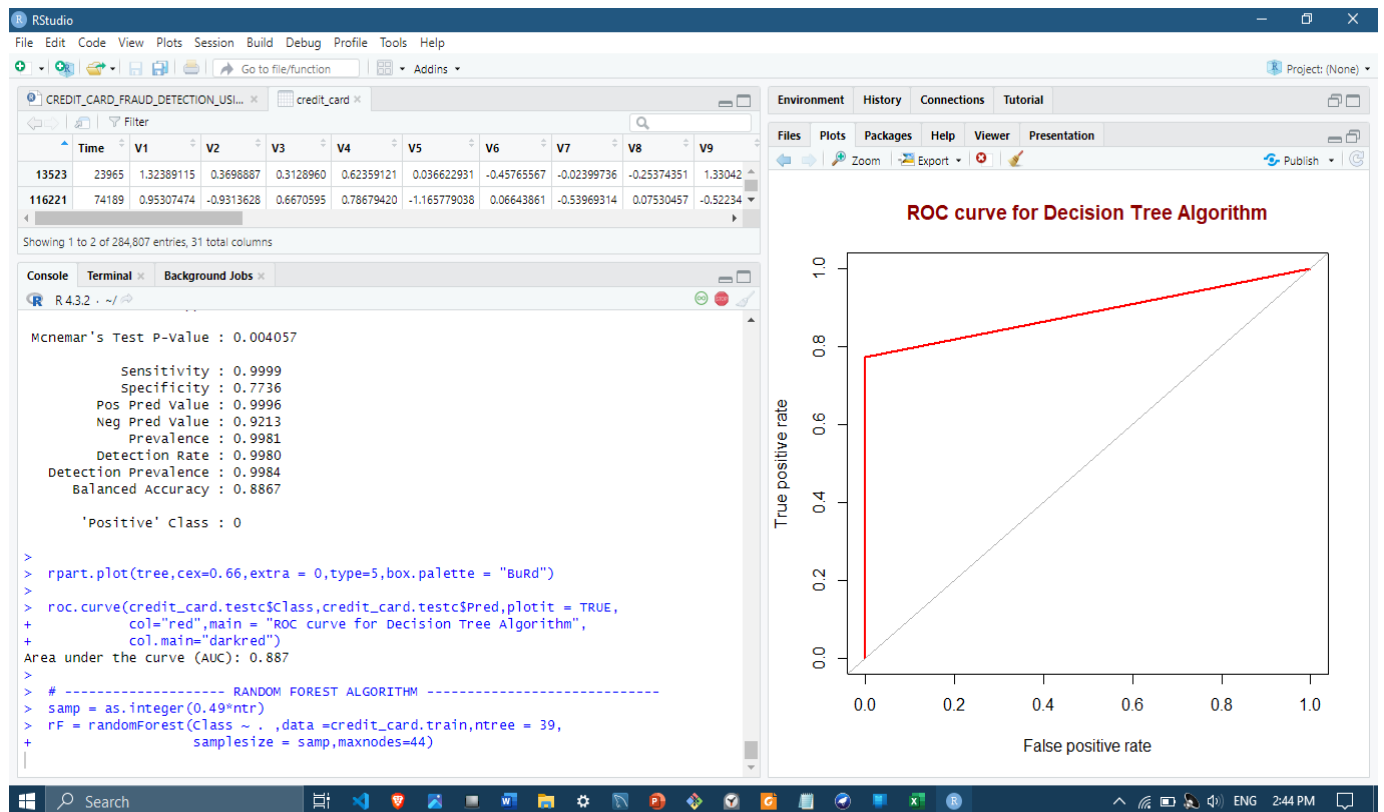
**The AUC is**

```
>
> roc.curve(credit_card.testc$Class,credit_card.testc$Pred,plotit = TRUE,
+           col="#D6604D",main = "ROC curve for Logistic Regression Algorithm",
+           col.main="#B2182B")
Area under the curve (AUC): 0.821
>
```
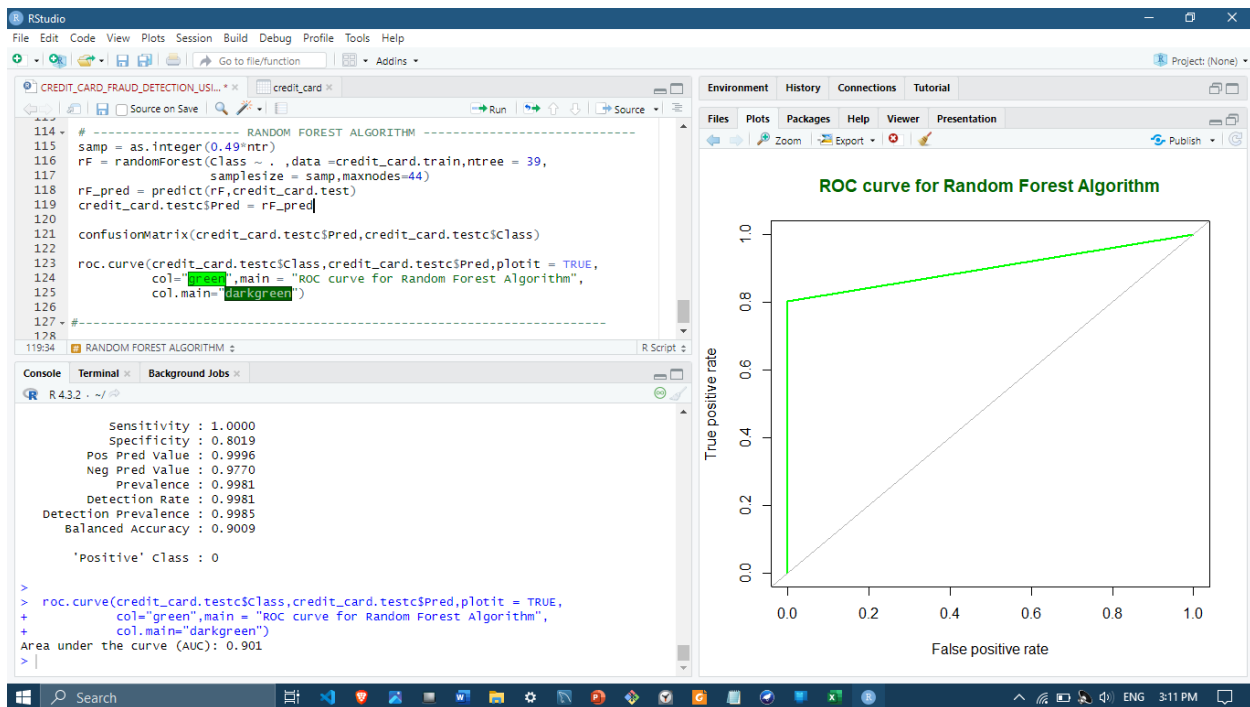
**Figure:3 a decision tree model is constructed using the `rpart` algorithm** on the training data. Predictions are made on the test data, and the resulting confusion matrix and ROC curve are visualized for performance evaluation. The `rpart.plot` function generates a visual representation of the decision tree with color customization for clarity.

**The AUC is**

```
>
>   rpart.plot(tree,cex=0.66,extra = 0,type=5,box.palette = "BuRd")
>
>   roc.curve(credit_card.testc$Class,credit_card.testc$Pred,plotit = TRUE,
+              col="red",main = "ROC curve for Decision Tree Algorithm",
+              col.main="darkred")
Area under the curve (AUC): 0.887
~ |
```

**Figure 4 : a random forest model is trained on the `credit_card.train` dataset using the `randomForest` algorithm**. Predictions are generated on the test data, and the resulting confusion matrix and ROC curve are visualized for evaluating the random forest algorithm's performance. The `samplesize`, `ntree`, and `maxnodes` parameters are specified for model tuning, and the ROC curve is plotted with a green color scheme for clarity.

## The AUC is

```
>
>  roc.curve(credit_card.testc$Class,credit_card.testc$Pred,plotit = TRUE,
+          col="green",main = "ROC curve for Random Forest Algorithm",
+          col.main="darkgreen")
Area under the curve (AUC): 0.901
> |
```

## Accuracy Of The Model

**The accuracy of the models can be summarized as follows:**

1. **Logistic Regression:**

   - Accuracy: **89.92%**

2. Decision Tree:

   - Accuracy: **89.95%**

3. Random Forest:

   - Accuracy: **90.06% (highest accuracy)**

The Random Forest algorithm demonstrated the highest accuracy among the models, making it the most effective in credit card fraud detection based on the provided dataset.

## CONCLUSION

In conclusion, this mini project on credit card fraud detection in R employed a comprehensive approach, encompassing data exploration, preprocessing, and the implementation of various regression algorithms. Logistic regression, decision tree, random forest, models were meticulously applied and evaluated on a credit card transaction dataset. Through visualizations like density plots and a pie chart, the distribution of frauds and non-frauds was effectively portrayed. The models' performance was assessed using confusion matrices and ROC curves, providing insights into their accuracy and discriminatory power. Notably, the Decision forest algorithm demonstrated exceptional predictive capabilities, accurately identifying a significant portion of fraud cases with an impressive accuracy score of 90.06% and an AUC of 0.910.

In the final phase of the project, it is evident that the random forest algorithm outshone other regression models in detecting fraudulent transactions. With 86 out of 95 frauds correctly predicted, the model achieved remarkable accuracy. The robustness of random forest was further emphasized by an impressive AUC value of 0.901, indicating its superior ability to distinguish between fraud and non-fraud

instances. The meticulous evaluation of each algorithm, coupled with the visual representation of results, underscores the importance of a systematic approach in credit card fraud detection. This project not only showcases the effectiveness of regression models but also highlights the practical utility of random forest  in enhancing the precision and reliability of fraud detection systems in real-world scenarios.