

# Project

2023-04-05

## LOAN RISK PREDICTION MODEL USING DECISION TREE

### DATA LOADING

```
data=read.csv("credit_risk_dataset.csv")
```

```
head(data)
```

```
##   person_age person_income person_home_ownership person_emp_length loan_intent
## 1         22         59000                RENT             123    PERSONAL
## 2         21          9600                 OWN              5    EDUCATION
## 3         25          9600            MORTGAGE              1    MEDICAL
## 4         23        65500                RENT              4    MEDICAL
## 5         24        54400                RENT              8    MEDICAL
## 6         21          9900                 OWN              2    VENTURE
##   loan_grade loan_amnt loan_int_rate loan_status loan_percent_income
## 1          D    35000         16.02           1           0.59
## 2          B     1000         11.14           0           0.10
## 3          C     5500         12.87           1           0.57
## 4          C    35000         15.23           1           0.53
## 5          C    35000         14.27           1           0.55
## 6          A     2500          7.14           1           0.25
##   cb_person_default_on_file cb_person_cred_hist_length
## 1                          Y                        3
## 2                          N                        2
## 3                          N                        3
## 4                          N                        2
## 5                          Y                        4
## 6                          N                        2
```

### DATA PREPROCESSING

```
#Checking for any missing values in the dataset.
sum(is.na(data))
```

```
## [1] 4011
```

There are missing values in the dataset. We need to remove the missing values from the dataset.

```
#Removing rows with missing values
data=na.omit(data)
sum(is.na(data))
```

```
## [1] 0
```

All the missing values have been removed from the dataset.

```
#Checking the datatype of each column in R
str(data)
```

```
## 'data.frame': 28638 obs. of 12 variables:
## $ person_age : int 22 21 25 23 24 21 26 24 24 21 ...
## $ person_income : int 59000 9600 9600 65500 54400 9900 77100 78956 83000 10000 ...
## $ person_home_ownership : chr "RENT" "OWN" "MORTGAGE" "RENT" ...
## $ person_emp_length : num 123 5 1 4 8 2 8 5 8 6 ...
## $ loan_intent : chr "PERSONAL" "EDUCATION" "MEDICAL" "MEDICAL" ...
## $ loan_grade : chr "D" "B" "C" "C" ...
## $ loan_amnt : int 35000 1000 5500 35000 35000 2500 35000 35000 35000 1600 ...
## $ loan_int_rate : num 16 11.1 12.9 15.2 14.3 ...
## $ loan_status : int 1 0 1 1 1 1 1 1 1 1 ...
## $ loan_percent_income : num 0.59 0.1 0.57 0.53 0.55 0.25 0.45 0.44 0.42 0.16 ...
## $ cb_person_default_on_file : chr "Y" "N" "N" "N" ...
## $ cb_person_cred_hist_length: int 3 2 3 2 4 2 3 4 2 3 ...
## - attr(*, "na.action")= 'omit' Named int [1:3943] 40 51 58 60 63 71 72 85 86 88 ...
## ..- attr(*, "names")= chr [1:3943] "40" "51" "58" "60" ...
```

```
#convert the categorical variables to factors
data$cb_person_default_on_file=as.factor(data$cb_person_default_on_file)
data$person_home_ownership=as.factor(data$person_home_ownership)
data$loan_intent=as.factor(data$loan_intent)
data$loan_grade=as.factor(data$loan_grade)
```

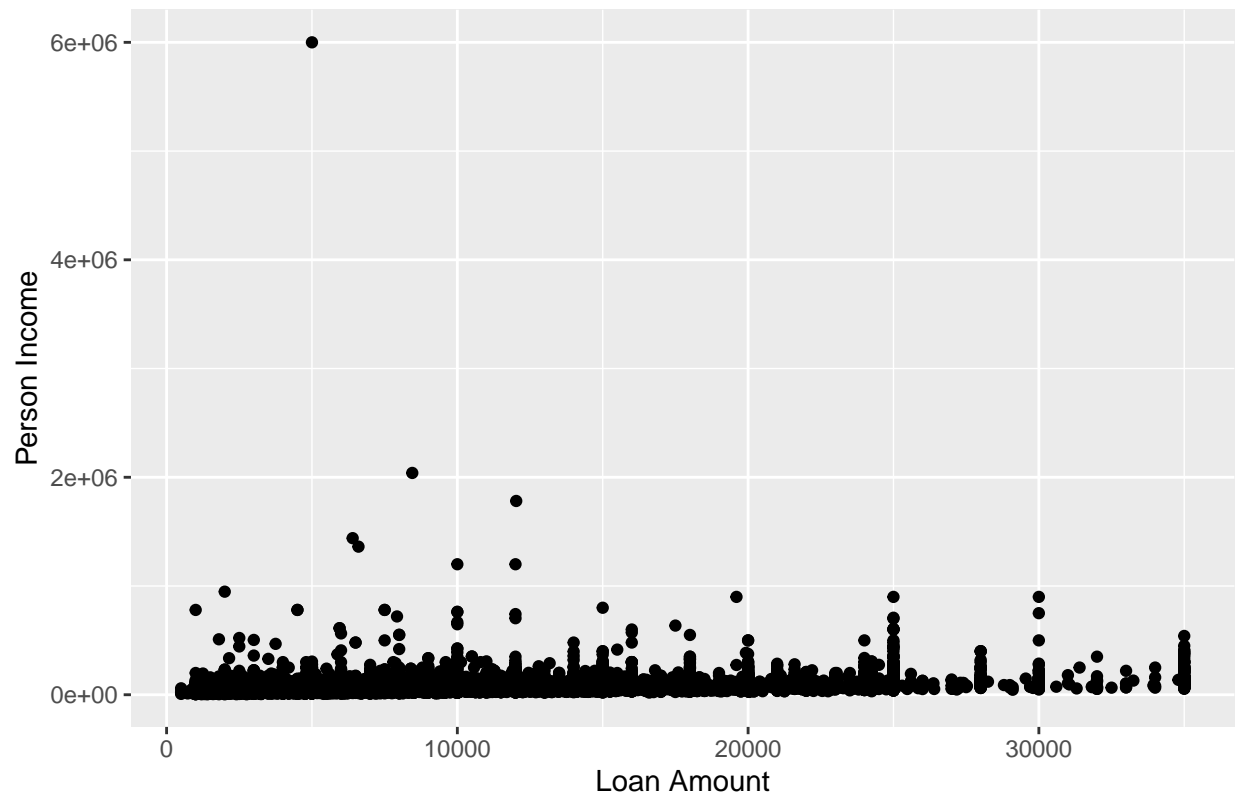
All the datatypes are correct, So no need of changing. The data pre processing is done.

## DATA VISUALISATION

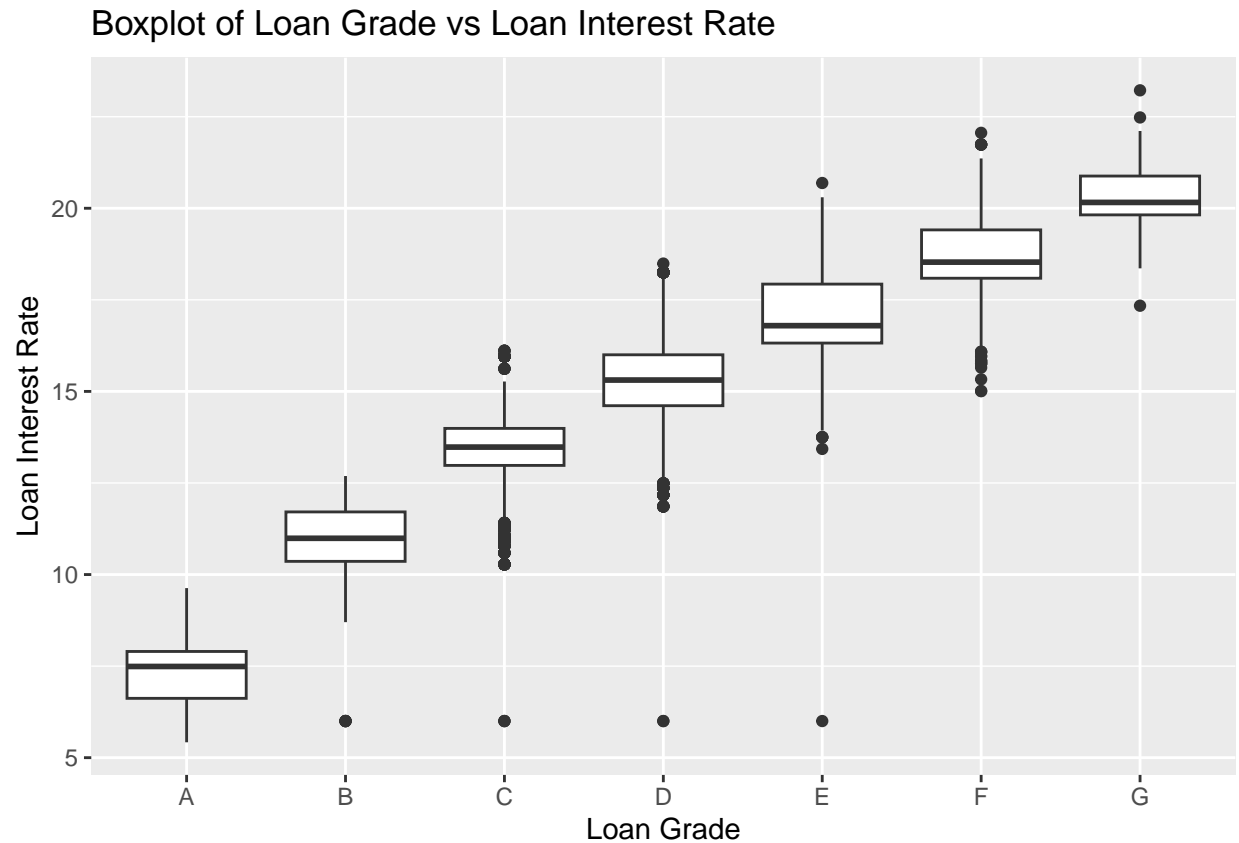
### PLOTTING A SCATTERPLOT BETWEEN THE INCOME AND LOAN AMOUNT

```
library(ggplot2)
pt=ggplot(data=data)+geom_point(aes(y = data$person_income, x = data$loan_amnt))+labs(y = "Person Income", x = "Loan Amount")
print(pt)
```

Scatterplot of Person Income vs Loan Amount



```
pt2=ggplot(data=data)+geom_boxplot(aes(x = data$loan_grade, y = data$loan_int_rate))+labs(x = "Loan Grade", y = "Loan Interest Rate")
plot(pt2)
```



We can see that the Interest Rates are considerably increasing with the level of Loan Grades.

## DATA SPLITTING

We have a large dataset of nearly 30,000 observations. So we can use more data to train. The preferable ratio will be 90% to training data and 10% to testing data.

```
library(caret)
```

```
## Loading required package: lattice
```

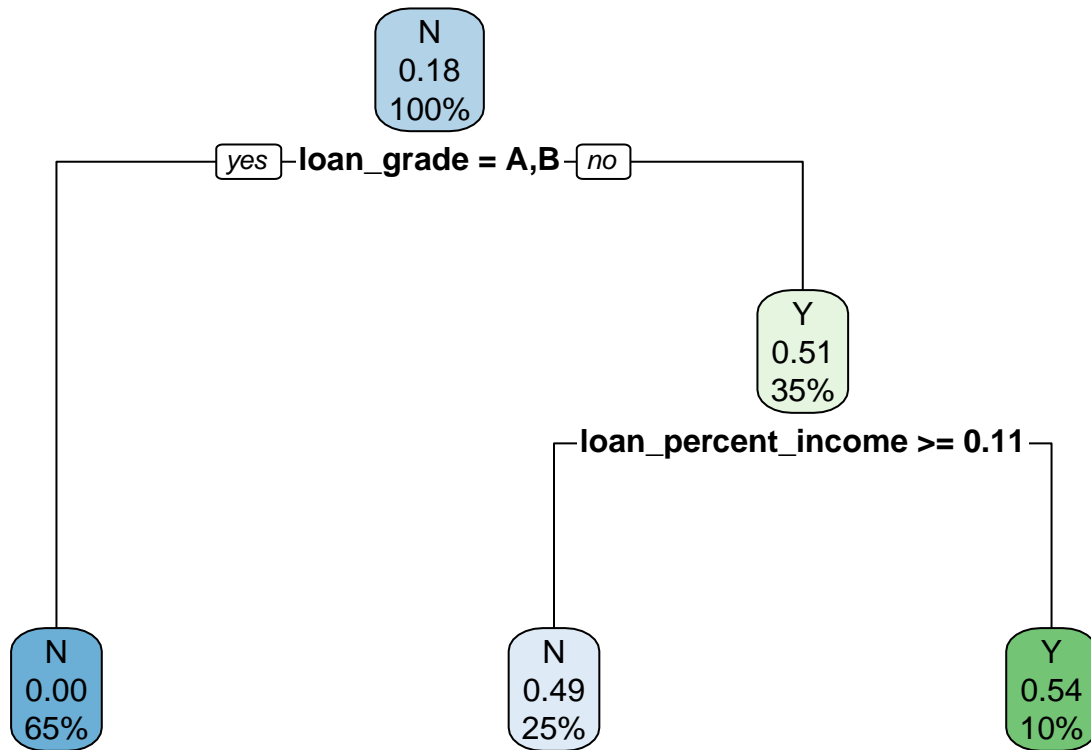
```
train_index <- createDataPartition(data$loan_status, p = 0.9, list = FALSE)
train_data <- data[train_index, ]
test_data <- data[-train_index, ]
```

## BUILDING THE FIRST MODEL

```
library(rpart) #required library
library(rpart.plot)
variable=data$cb_person_default_on_file
```

```
model1= rpart(cb_person_default_on_file ~ ., data = train_data)
```

```
rpart.plot(model1)
```



## MODEL-1 PREDICTION AND EVALUATION

```
predict1 =predict(model1, test_data, type = "class")
```

```

# Compute confusion matrix and accuracy for the Decision Tree model
confmat1 = table(test_data$cb_person_default_on_file, predict1)
acc1 = sum(diag(confmat1)) / sum(confmat1)
print(paste("Accuracy:", round(acc1, 3)*100,"%"))

```

```
## [1] "Accuracy: 82.5 %"
```

## MODEL-2 BUILDING - NAIVE BAYES

```

library(e1071)
model2 = naiveBayes(cb_person_default_on_file ~ ., data = train_data)
pred2 = predict(model2, test_data)

```

## MODEL-2 EVALUATION

```
# Compute confusion matrix for Naive Bayes Model
confmat2 <- table(test_data$cb_person_default_on_file, pred2)
print(confmat2)
```

```
##      pred2
##          N      Y
## N 1886  453
## Y   64  460
```

```
# Compute accuracy for the model
acc2 <- sum(diag(confmat2)) / sum(confmat2)
print(paste("Accuracy:", round(acc2, 3)*100, "%"))
```

```
## [1] "Accuracy: 81.9 %"
```