# CREDIT_RISK_CUSTOMERS PROJECT

## *DOCUMENTATION*

- **Abstract/Executive Summary :-**

The Project is about Credit Risk of the Customer according to various attributes of the customer. We need to predict about the target Class depending on the features that whether he comes under good class or bad class. There are 1000 data points and 21 features in the dataset of the project. The EDA is done to find the relation of the various features with the target and predictions are done using various Machine Learning Algorithms. After performing all the EDA process and Machine Learning Algorithms, the following were the key findings:-

➢ The Best Model to predict about the credit risk in this project is Random Forest Random Search CV.
➢ Generally there are very less number of customers who are risky for granting the credit. The duration for most of the loans are between 10 and 25. The maximum customers have taken credit between 1000 to 5000. Customers between age of 20 and 30 are more in a need of credit.
➢ Customers who have high savings and more employment duration are very less risky for granting the credit and most of the customers who take credit for luxurious items.
➢ Higher Loan amount is generally cleared on time. Customer who are in middle age and those have less years of employment are generally prone towards the credit.

## ● Introduction :-

The project is of *Fintech industry*. It is a <u>classification</u> problem.

The following is the problem statement. We need to predict whether the customer is risky or not for providing the credit

***Problem Statement:-Predict if customers are risky or not for credit. There are 21 attributes of the customers a prediction needs to be done whether the customer is risky or not. The Target in the dataset is class, whether the customer comes under good or bad class.***

## ● Data Collection and Preprocessing :-

The data is collected from *Kaggle*. There are <u>1000 data points and 21 features</u> in the dataset.

The following are the features of the dataset:-***This dataset classifies people described by a set of attributes as good or bad credit risks.***

- ➢ checking_status - Status of existing checking account
- ➢ duration - Duration in months
- ➢ credit_history - credits taken, paid back duly, delays, critical accounts
- ➢ purpose - Purpose of the credit
- ➢ credit_amount - Amount of credit
- ➢ savings_status - Status of savings account/bond
- ➢ employment - Present employment, in number of years
- ➢ installment_commitment - Installment rate in percentage of disposable income
- ➢ personal_status - sex and marital data
- ➢ other_parties - Other debtors / guarantors
- ➢ residence_since -   How many years customer is residing.
- ➢ property_magnitude  - Type of Property
- ➢ age – Age of the customer
- ➢ other_payment_plans – Any other Payment Plan
- ➢ housing – Customer has own or rented house or has free accomodation
- ➢ existing_credits – Number of existing credit
- ➢ job – skilled or unskilled worker
- ➢ num_dependents  - number of dependants
- ➢ own_telephone – possession of telephone or not
- ➢ foreign_worker – Foreign worker or not
- ➢ class – Customer comer under which class for credit – good or bad

There were some inconsistencies in the data type of the features which was dealt and the data was made consistent. There were no null values in the data. Data was Standardised as well as Normalised using Scaling Techniques like Standard Scaler and MinMax Scaler. Encoding was also done so that the Machine Learning Models can give predictions with huge accuracy.

## ● Methodology:-

The problem statement was the classification problem statement so the first ML Algorithm which was used was Logistic Regression. Later the other Algorithms which where used were KNN, Decision Tree Classifier, Support Vector Classifier Random Forest Classifier and Boosting Algorithms. The Algorithms were Hyper Parameter Tuned as well so that the Accuracy of the Model be high. The Performance Metrics used is ACCURACY.
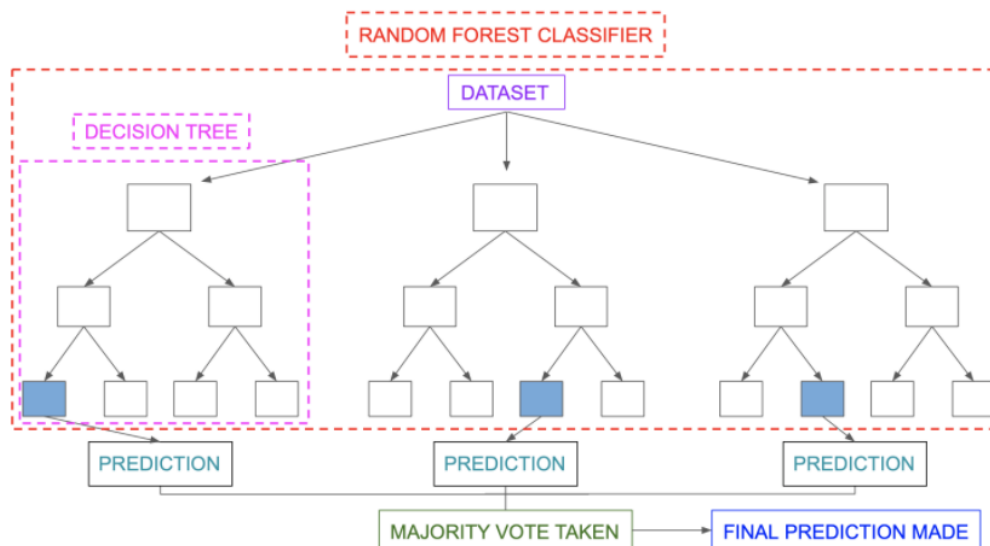
The Result of all the Machine Learning Algorithms are as follows:-

| Model | Training Accuracy | Testing Accuracy |
|---|---|---|
| Logistic Regression | 71.60 | 75.2 |
| KNN | 72.00 | 71.2 |
| SVC | 70.53 | 72.8 |
| Decision Tree | 100.00 | 70.4 |
| Random Forest | 84.90 | 82.4 |
| Cat Boost | 96.80 | 82.4 |
| XG Boost | 100.00 | 78.8 |
| Gardient Boost | 83.60 | 87.4 |

The above results are *Hyper Parameter Tuned*.

From the table we can see that the Best Fit Model is Random Forest Random Search CV.

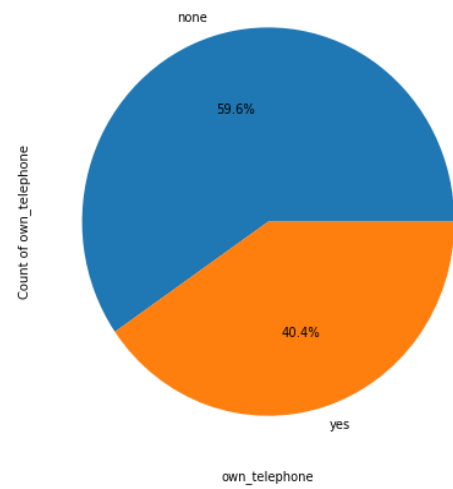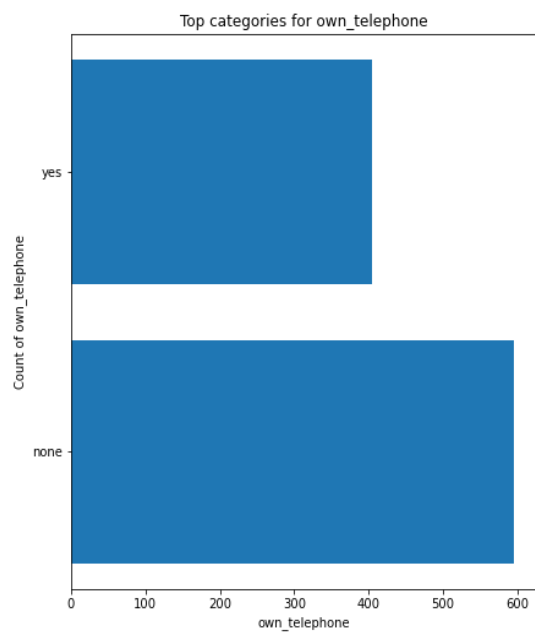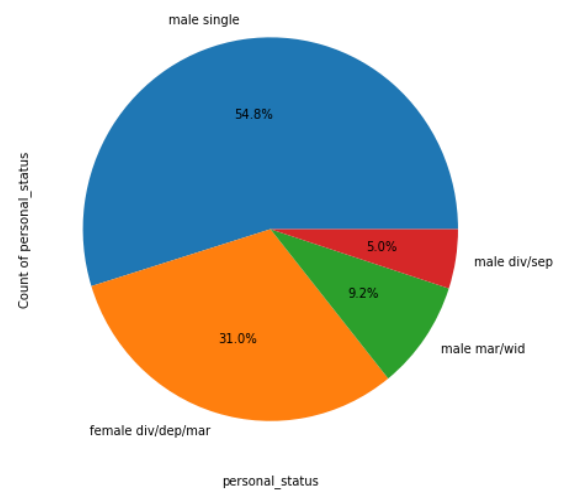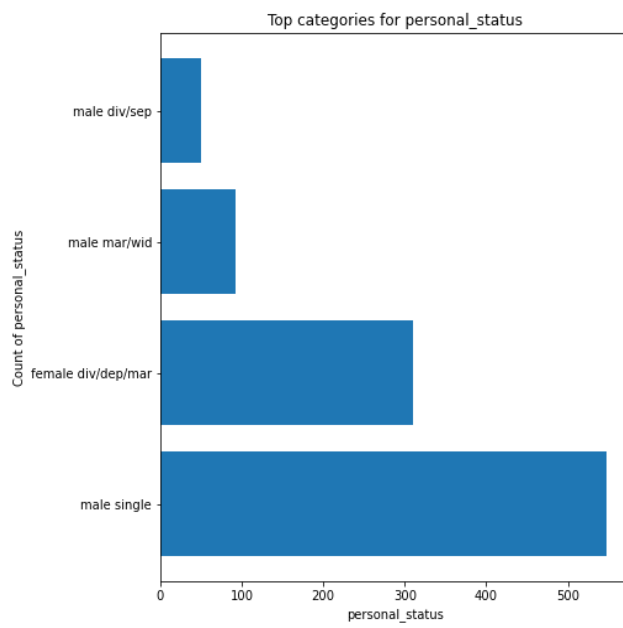## FLOW CHART:-

● **Results and Conclusion:-**

## EDA RESULTS

### Top categories for categorical features in the data :-

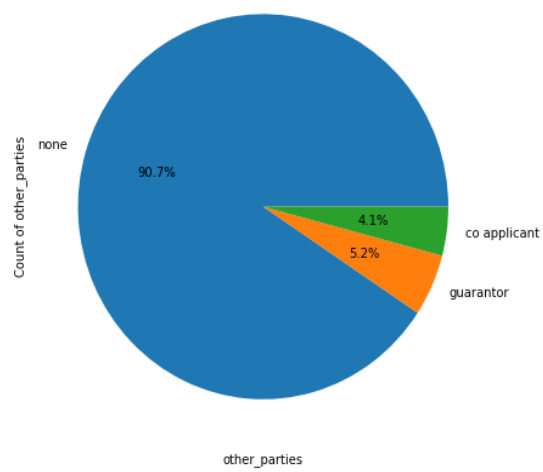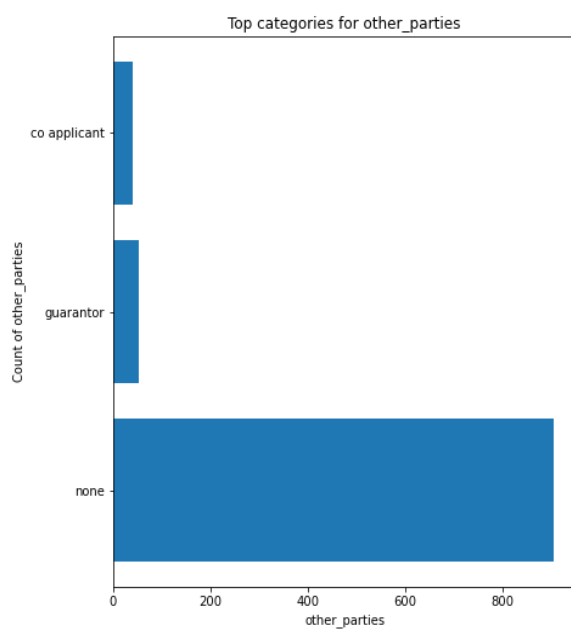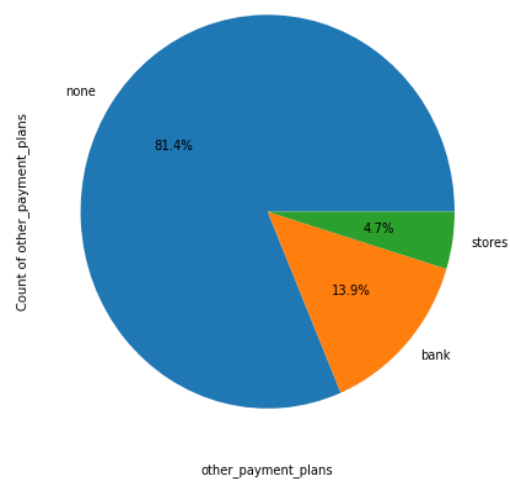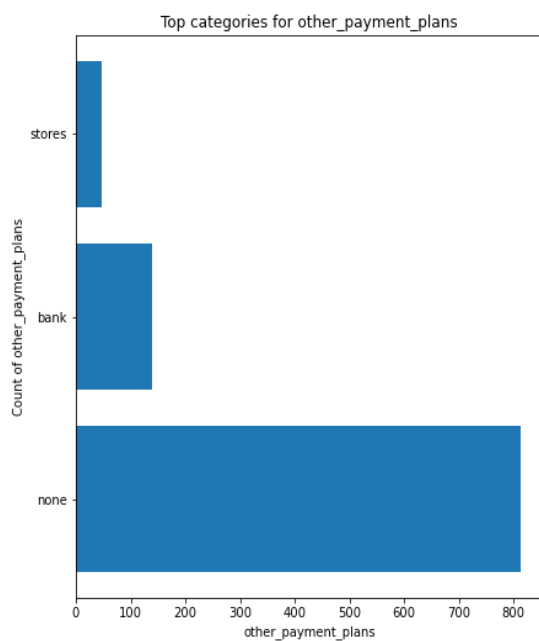Top categories for employment



Top categories for credit_history

## Top categories for existing_credits



## Top categories for residence_since

Top categories for purpose

Top categories for property_magnitude

## Top categories for personal_status



## Top categories for own_telephone

## Top categories for other_payment_plans



## Top categories for other_parties

Top categories for num_dependents



Top categories for job

Top categories for installment_commitment



Top categories for housing

Top categories for foreign_worker



Top categories for savings_status

✓ **Inference from above charts**

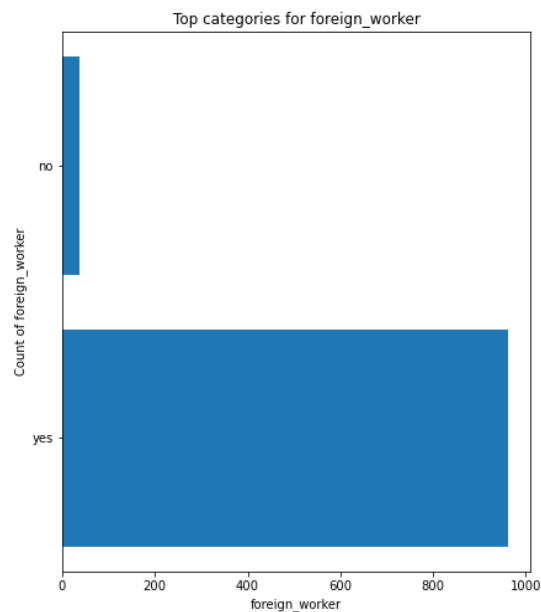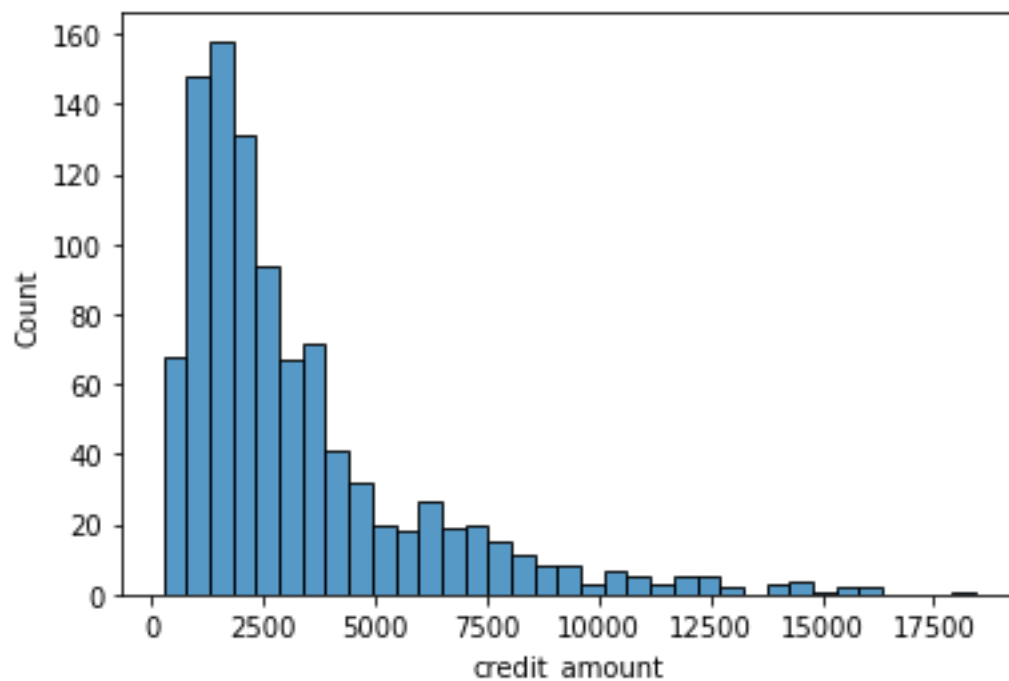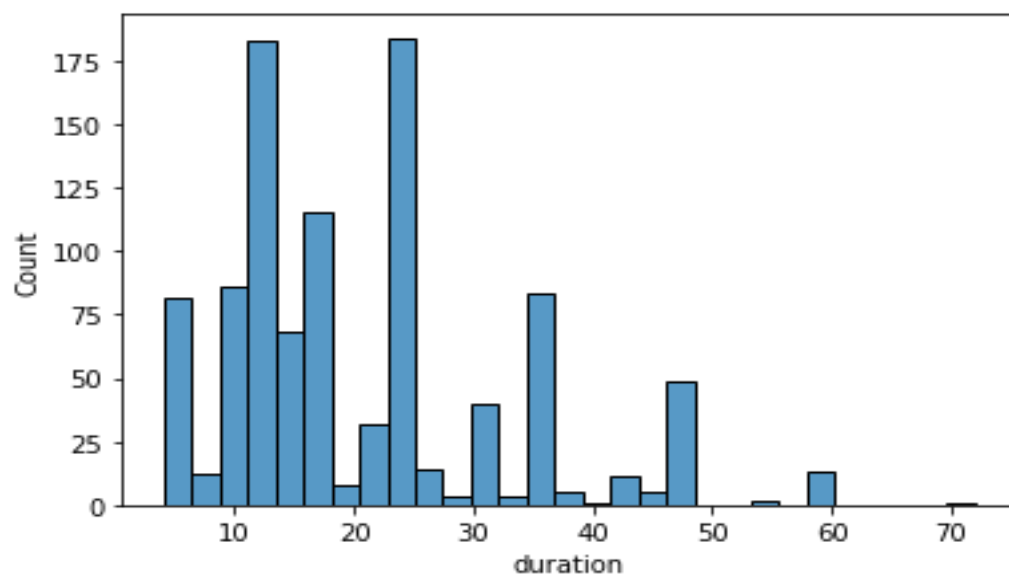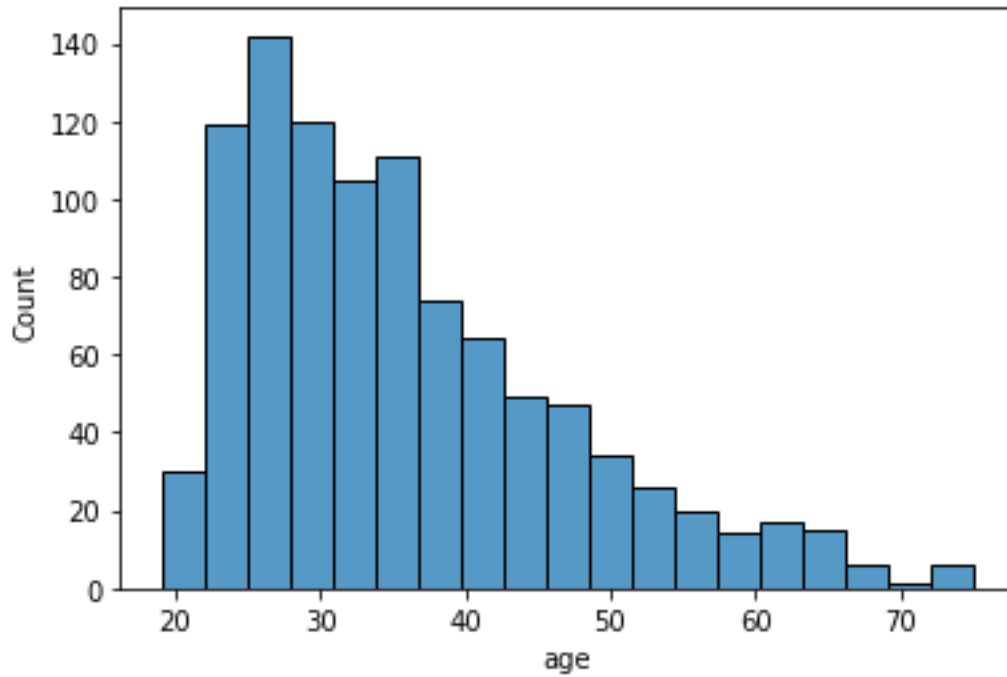➢ Almost 40% people have checking_status as no checking, which means there will be no account checking for them if they are availing for credit. 27.5% people have less than 0 in their account checking. 27% have between 0 to 200 and only 6% people have more than 200 in their account checking.

➢ People who are existing paid or critical/another existing credit are more prone towards credit.

➢ People are taking more credit for purchasing Radio, TV, New Car, Furnitures and equipments.

➢ Customers who have less savings are more prone towards credit as they don't have liquidity.

- ➤ Customers who have more tenure of employment are more prone towards credit which shows they need more money.
- ➤ Almost 50% customers have installment commitment of 4
- ➤ Male single are in need of more credit which is surprising as bachelor male have less expenses.
- ➤ In more than 90% cases there is no involvement of other parties in loan cases.
- ➤ Customers who have more tenure of residence are more prone towards credit.
- ➤ Customers who are in possession of Car and Real Estate are more prone towards credit
- ➤ Customers who are owning house are more inclined towards credit
- ➤ Maximum customers are having only one existing loan credit
- ➤ 70% customers have good credit and 30% have bad credit
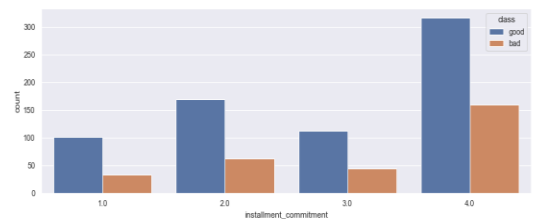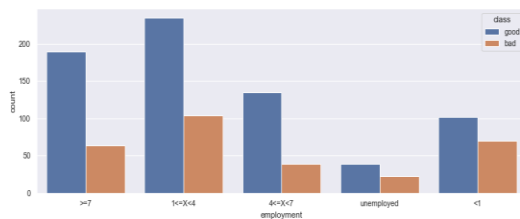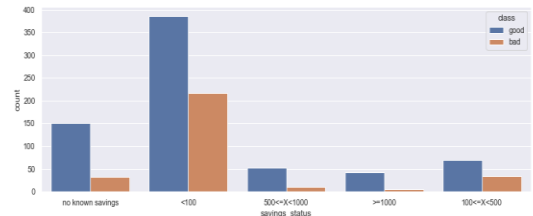
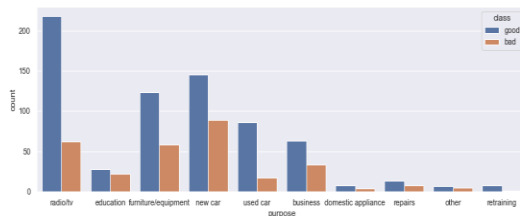*Top categories for continuous features in the data :-*

✓ **Inference from above charts**
➢ The duration for most of the loans are between 10 and 25
➢ The maximum customers have taken credit between 1000 to 5000
➢ Customers between age of 20 and 30 are more in a need of credit

*Relation of categorical and continuous columns with the Target:-*

✓ **Inference from above charts**

➢ We can see that number of people in class good is 700 and in class bad is 300.

➢ Generally for all the checking status customers who are not risky for credit are more than number of risky customers but most of those who have no checking as checking status are very less risky customer and have class good and credit can be granted to them without any second thought.

➢ Most of the customers having credit history as existing paid are not risky for credit. Same is the case for critical/other existing credit. In customers who have credit history as no credits/all paid & all paid there are more number of cu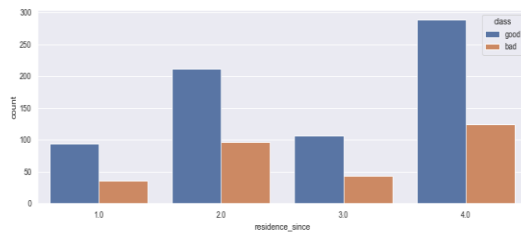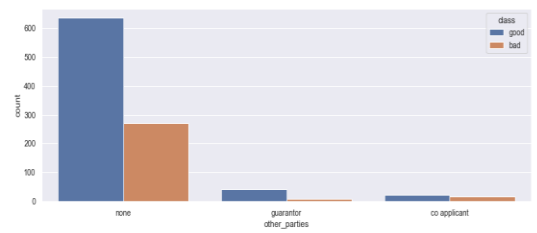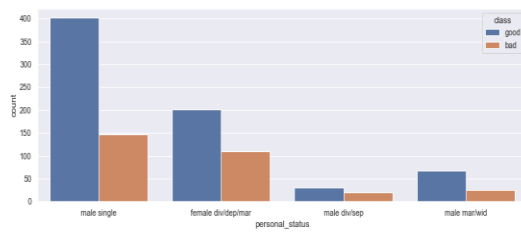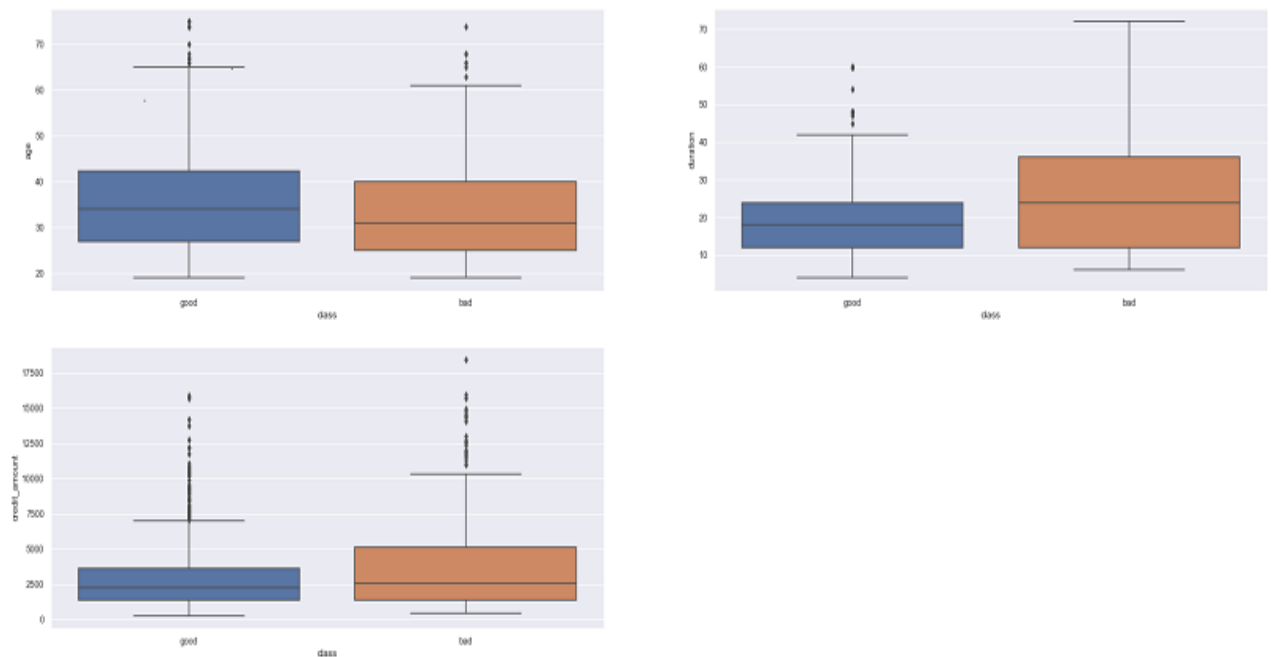stomers who are in bad class, hence they are risky but a small credit can be given to them because they will pay back the money later and that credit will not go in NPA because they are in class of no credits or all paid which clearly shows they return the money but not in time.

➢ Most of the customers who take credit for radio, TV, furniture, new car and retraining generally pay of the loan and they are not risky customers.

➢ Customers who are having savings_status between 500 to 1000 are generally very less risky customers.

➢ Customers who are employed between 1 and 4 are generally less risky. Same is the case for employment greater than 7.

➢ Most of the customers who are having installment commitment of 4 are less risky.

➢ Customers who are residing for a longer period of time are less risky.

➢ Customers who own the house are generally less risky as compared to those who does not posses a property.

➢ Customer who posses car and also investment in real estate are less risky.

➢ Individual loan account customers are less risky than loan accounts where co-applicant and guarantor are involved.

➢ Customers who have existing credit of less than 3 are more risky to customers who have more existing credits.

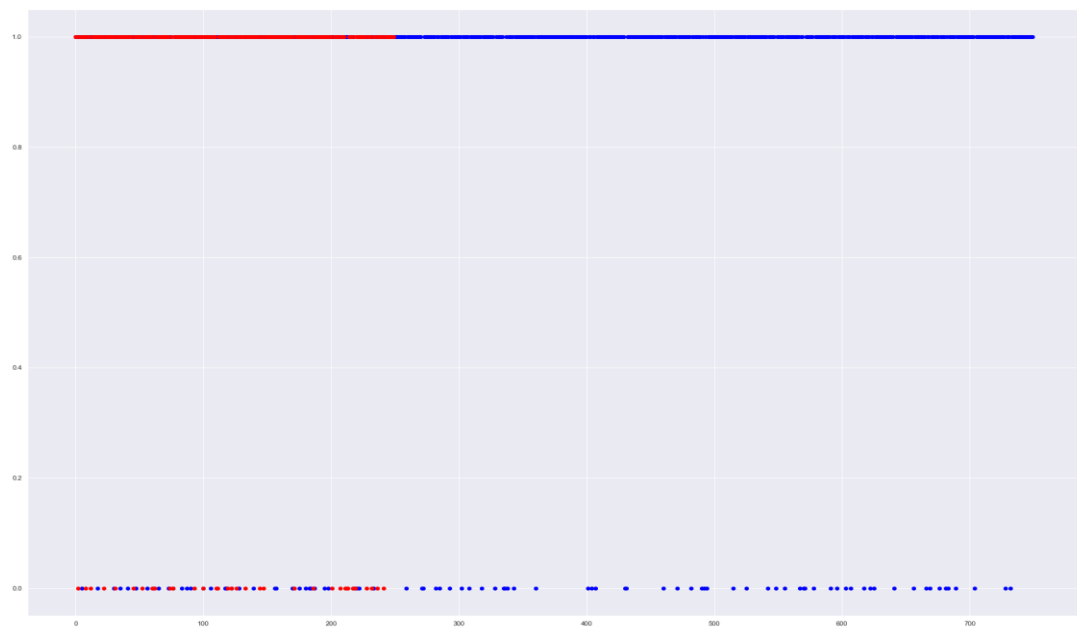➢ Customers who are skilled are less risky compared to others.

- Customers whose age is between 25 to 45 are who are need of credit. Customers who are above 40 years of age are very less risky compared to other age groups.
- Customers where duration of credit is less are less risky. And as the duration of credit increases the risk of the customer also increases.
- Customers who take credit more than 3500 are very risky compared to customers who take less credit amount.
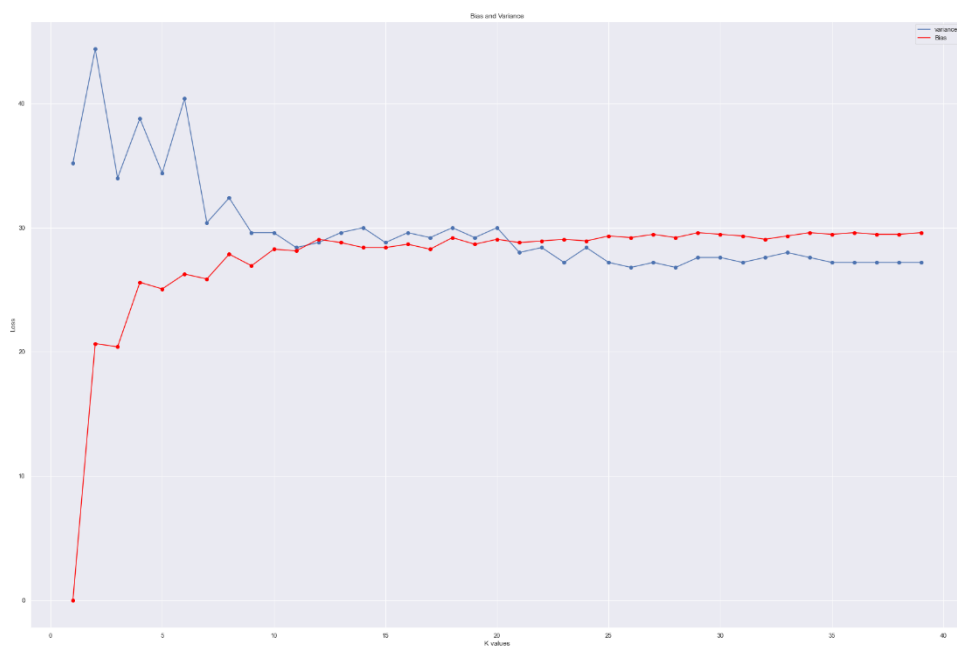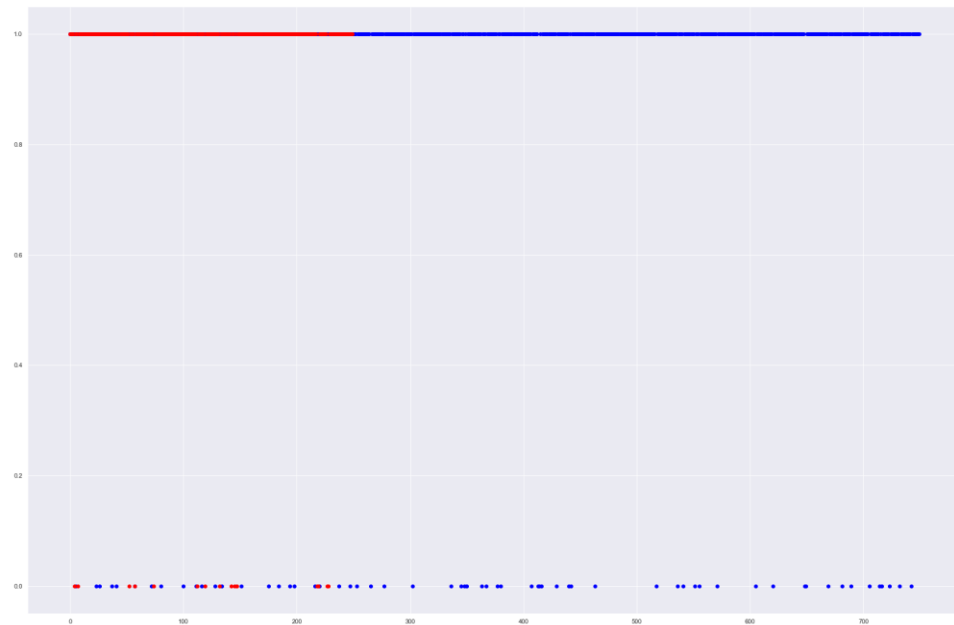
## MACHINE LEARNING RESULTS:-

Below are the scatter plots of all the models involved in the prediction of the customer's risk towards the credit.

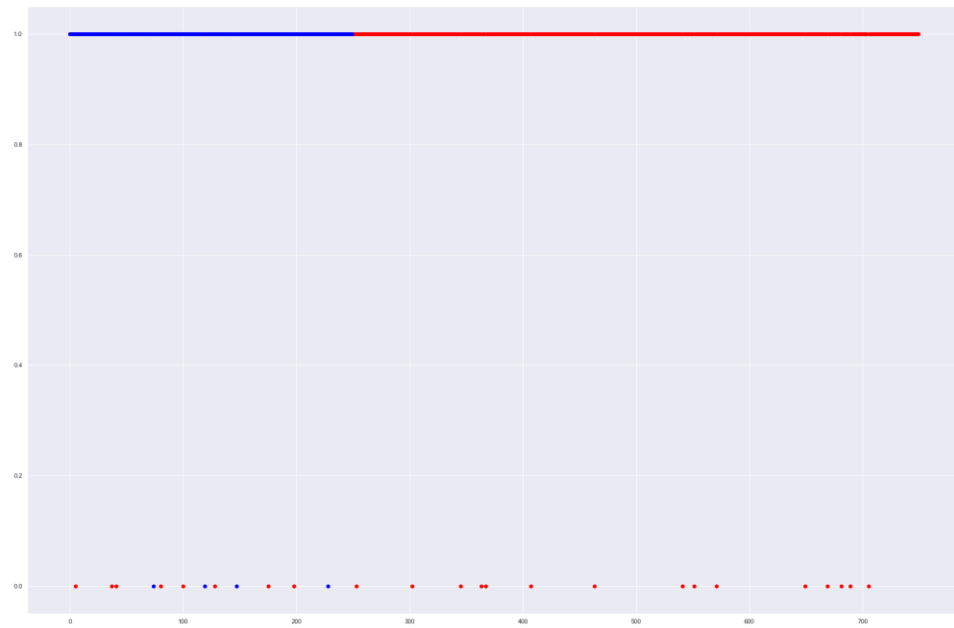The Blue colour is for Training Data and Red colour is for Test Data.
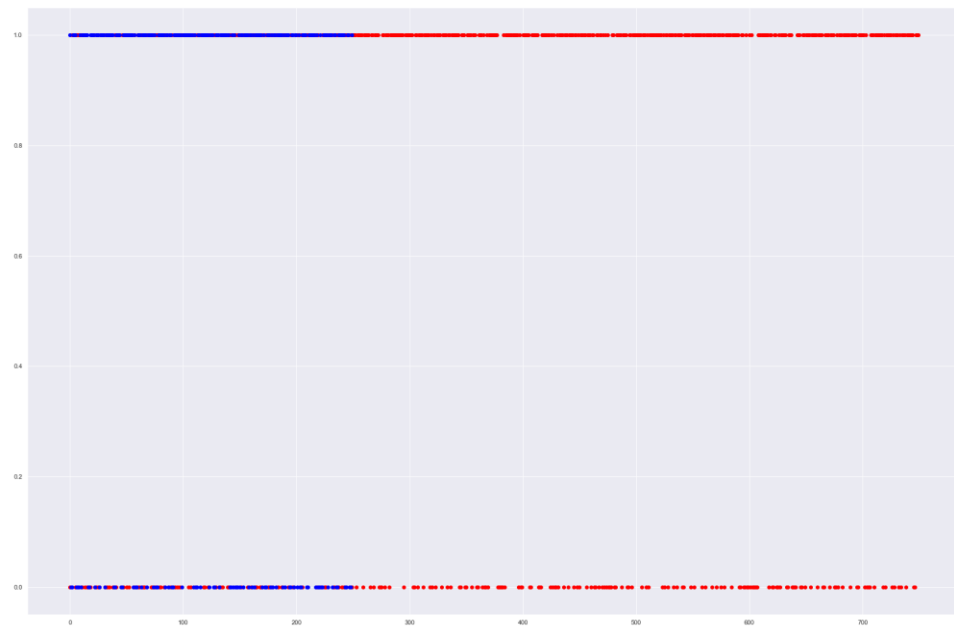
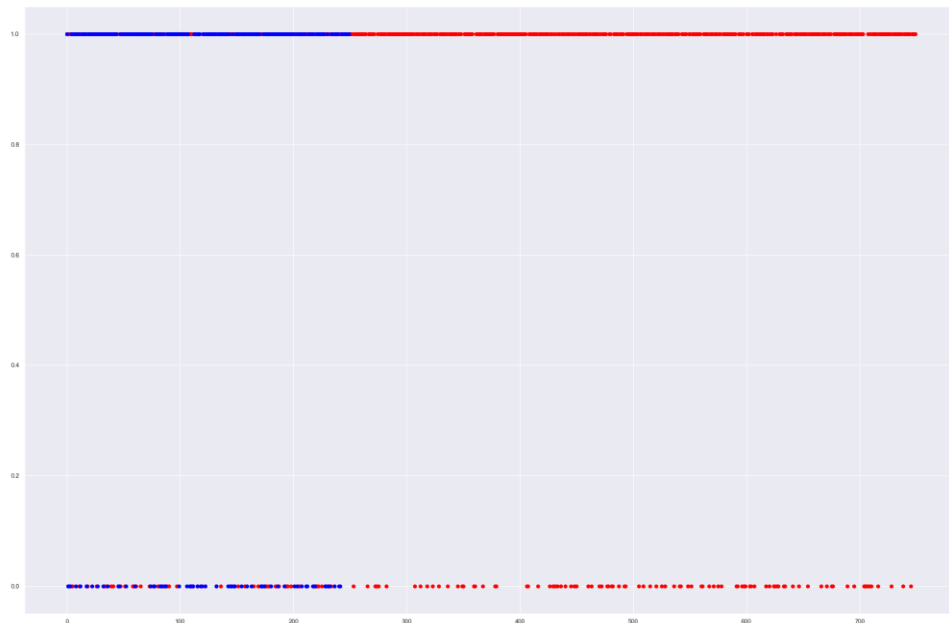## LOGISTIC REGRESSION

# KNN & KNN BV Trade Off

# SVM



# Decision Tree

**Random Forest**



*The Best Prediction for the Problem Statement is Random Forest Random Search CV.*

● **References:-**

➢ Random Forest Flow Chart:-
https://www.google.com/search?q=Random+forest+random+search+cv+flow+chart&rlz=1C1CHBD_enIN1056IN1056&sxsrf=APwXEdeHtBsPdqdn3BMvLplpTTtGrzoYJA:1686565153893&source=lnms&tbm=isch&sa=X&ved=2ahUKEwjhro3nwL3_AhUIyDgGHXBODfEQ_AUoAXoECAIQAw&biw=1280&bih=512&dpr=1.5#imgrc=hPgBAWOXUzu0VM

➢ Dataset Source:-
https://www.kaggle.com/datasets/ppb00x/credit-risk-customers

➢ https://www.google.com/search?q=credit+risk&tbm=isch&ved=2ahUKEwjH9PjXmL7_AhVSLLcAHfV7CPgQ2-cCegQIABAA&oq=credit+risk&gs_lcp=CgNpbWcQAzIECCMQJzIFCAAQgAQyBQgAEIAEMgUIABCABDIFCAAQgAQyBQgAEIAEMgUIABCABDIFCAAQgAQyBQgAEIAEMgUIABCABDoHCAAQGBCABFDYBVi3FGCbGWgAcAB4AIABvAKIAfkQkgEHMC42LjQuMZgBAKABAaoBC2d3cy13aXotaW1nwAEB&sclient=img&ei=SE2HZMewGdLY3LUP9fehwA8&bih=569&biw=1280&rlz=1C1CHBD_enIN1056IN1056#imgrc=PqVdVup3nHUU8M