# LOAN ELIGIBILITY PREDICTION USING SUPERVISED LEARNING

**(Logistic Regression, Linear Regression, Naïve Bayes, KNN, SVM, Gradient Descent)**

Pavithra Shankar Babu

# CONTENTS

- Introduction

- Business Question

- Problem Statement & Problem Solution

- Software Used

- Proposed Work

- Methodology

# INTRODUCTION

Loans are the core business of banks. The main profit comes directly from the loan's interest. The loan companies grant a loan after an intensive process of verification and validation. However, they still don't have assurance if the applicant can repay the loan with no difficulties. In this Project, a predictive model ha been built to predict if an applicant is eligible for the loan or not.

# BUSINESS QUESTION

"What factors most accurately predict whether a loan applicant is likely to repay their loan on time, and how can this information be leveraged to make informed decisions regarding loan approval or rejection?"

# PROBLEM STATEMENT

- Intensive time Consumption process of verification and validation.
- Human errors can be introduced during the validation process.
- No cross referencing previous loan records.

# PROBLEM SOLUTION

- Our Machine learning model calculates all the parameters given and predicts if the applicant is eligible for loan or not in very less time.
- Time required for verification, and validation reduces significantly.

# SOFTWARE USED

**Libraries Used:**

o Pandas

o Numpy

o Seaborn

o Matplotlib.pyplot

o Sklearn.preprocessing

o Sklearn.svm

o Sklearn.linear_model.LinearRegression
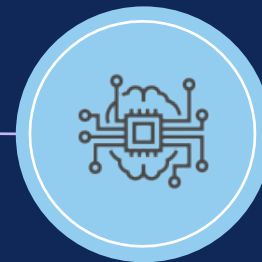
**Language Used:**

o Python

# PROPOSED WORK



**Step 1**

Data Input

**Step 2**

Data Preparation

**Step 3**

ML Model building & Training
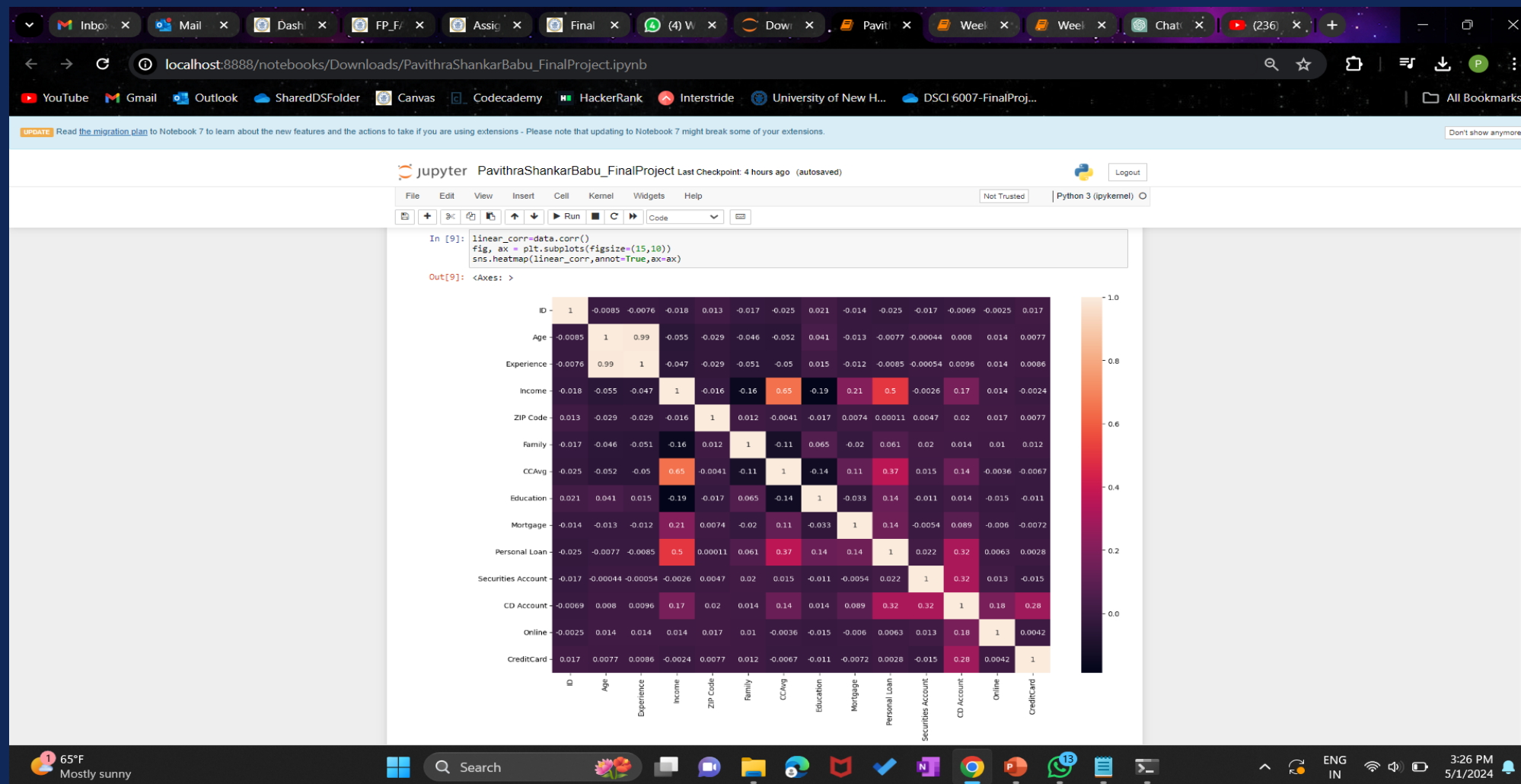
**Step 4**

Model Evaluation

**Step 5**

Prediction & Results

# METHODOLOGY

o Data Collection: A single dataset containing historical loan application data, from Kaggle is taken including features such as applicant demographics, financial information, credit history, and loan outcome (eligible or not eligible).

o Analyzing Data: Explore the data to understand its structure, distributions, and relationships between variables. Determine which features are relevant for predicting loan eligibility and identify any patterns or trends in the data.

o Data Cleaning: Clean the data by handling missing values, outliers, and inconsistencies. Ensure that all features are in the appropriate format for analysis and preprocessing, such as numerical encoding for categorical variables.

o Model Building: Implement various supervised learning algorithms for loan eligibility prediction, including K-Nearest Neighbors (KNN), Support Vector Machines (SVM), Logistic Regression, Linear Regression, Naive Bayes, and Gradient Descent. Train each model on the training data and tune hyperparameters as needed.

o Evaluating Performance Metrics of Models: Evaluate the performance of each model using appropriate metrics such as accuracy, precision, recall, F1-score, and ROC-AUC. Compare the performance of different algorithms to determine which one works best for predicting loan eligibility in this specific context. Tune the selected model further if necessary to improve its performance.

# RESULT

## Correlation Matrix:

# Outliers:



```python
In [7]:  #Checking the outliers of the columns using boxplot

         plt.figure(figsize=(20,15))
         plt.subplot(3,3,1)
         sns.boxplot(x='Personal Loan',y='Age',data=data)
         plt.subplot(3,3,2)
         sns.boxplot(x='Personal Loan',y='Experience',data=data)
         plt.subplot(3,3,3)
         sns.boxplot(x='Personal Loan',y='Income',data=data)
         plt.subplot(3,3,4)
         sns.boxplot(x='Personal Loan',y='CCAvg',data=data)
         plt.subplot(3,3,5)
         sns.boxplot(x='Personal Loan',y=data['Mortgage'].loc[data['Mortgage']!=0],data=data)
```

```
Out[7]: <Axes: xlabel='Personal Loan', ylabel='Mortgage'>
```

Thus, it is clear that CCAvg,Income and Mortgage column has outliers.

## Linear Regression:

```
In [22]: #Linear
         linear_reg(X,y)

The linear model prediction is 84.18533367404265%
The confusion matrix is
[[945   1]
 [ 70 484]]
ROC value for linear model is 93.6294563466925%
```

## SVM:

```
In [25]: #SVM
         svm_fun(X,y)

The KNN model prediction is 98.0%
The confusion matrix is
[[944   2]
 [ 28 526]]
the Classification report is
              precision    recall  f1-score   support

           0       0.97      1.00      0.98       946
           1       1.00      0.95      0.97       554

    accuracy                           0.98      1500
   macro avg       0.98      0.97      0.98      1500
weighted avg       0.98      0.98      0.98      1500

ROC value for svm model is 97.3672159424825%
```

## Logistic Regression:

```
In [26]: #Logistic Regression
         logistic_reg(X,y)

         #Since in the input data, we have more value for personal Loan as 0 than 1,
             #we must consider the 0 class level f1 score - Here it is 96% good that it would predict who would get the Personal Loan

         The Logistic model prediction is 97.26666666666667%
         The confusion matrix is
         [[941   5]
          [ 36 518]]
         the Classification report is
                       precision    recall  f1-score   support

                    0       0.96      0.99      0.98       946
                    1       0.99      0.94      0.96       554

             accuracy                           0.97      1500
            macro avg       0.98      0.96      0.97      1500
         weighted avg       0.97      0.97      0.97      1500

         ROC value for logistic model is 96.486631913968%
```

## Naïve Bayes:

```
In [27]: #Naive Bayes
         naive_bayes(X,y)

         #Since in the input data, we have more value for personal Loan as 0 than 1,
             #we must consider the 0 class level f1 score - Here it is 91% good that it would predict who would get the Personal Loan

         The Naive Bayes model prediction is 93.60000000000001%
         The confusion matrix is
         [[943   3]
          [ 93 461]]
         the Classification report is
                       precision    recall  f1-score   support

                    0       0.91      1.00      0.95       946
                    1       0.99      0.83      0.91       554

             accuracy                           0.94      1500
            macro avg       0.95      0.91      0.93      1500
         weighted avg       0.94      0.94      0.93      1500

         ROC value for linear model is 91.447935827708115%
```

# KNN:

```
In [28]:  #KNN
          knn(X,y,3)

          #Since in the input data, we have more value for personal Loan as 0 than 1,
             #we must consider the 0 class level f1 score - Here it is 96% good that it would predict who would get the Personal Loan
```

```
The KNN model prediction is 97.2%
The confusion matrix is
[[936  10]
 [ 32 522]]
the Classification report is
              precision    recall  f1-score   support

           0       0.97      0.99      0.98       946
           1       0.98      0.94      0.96       554

    accuracy                           0.97      1500
   macro avg       0.97      0.97      0.97      1500
weighted avg       0.97      0.97      0.97      1500

ROC value for linear model is 96.58337213118507%
```

Here, 10 people predicted by the model built that they will get loan, dont get loan and 32 who we predict dont get loan actually got loan.
Recall=522/(32+522)=0.94

In KNN, we see that the Precession value,ROC and f1-score of 1 is higher than compared to Logistic and Naive bayes, So I recommend to follow KNN Algorithm.

# SUMMARY

Estimating the probability of default for approved loans: This can help the bank assess the risk associated with approving a loan and make informed decisions about loan terms and interest rates.

---

Predicting the loan amount: Given a set of applicant characteristics, the model can predict the amount of loan that is likely to be approved, providing valuable information to both the bank and the applicant. By developing this accurate predictive model, financial institutions can streamline their loan approval processes, minimize risk, and improve customer satisfaction.

# THANK YOU

Pavithra Shankar Babu

pshan4@unh.newhaven.edu