

# Classification Problem Document

## Problem Identification:

**Title:**  
Loan Approval Prediction

### Stages:

Stage 1 : Machine Learning

Stage 2: Supervised Learning

Stage 3: Classification

## Data Collection:

```
import pandas as pd

df=pd.read_csv('loan_data.csv')

df
```

	person_age	person_gender	person_education	person_income	person_emp_exp	person_home_ownership	loan_amnt	loan_intent	loan_int_rate	loan_p
0	22.0	female	Master	71948.0	0	RENT	35000.0	PERSONAL	16.02	
1	21.0	female	High School	12282.0	0	OWN	1000.0	EDUCATION	11.14	
2	25.0	female	High School	12438.0	3	MORTGAGE	5500.0	MEDICAL	12.87	
3	23.0	female	Bachelor	79753.0	0	RENT	35000.0	MEDICAL	15.23	
4	24.0	male	Master	66135.0	1	RENT	35000.0	MEDICAL	14.27	
...	...	...	...	...	...	...	...	...	...	...
44995	27.0	male	Associate	47971.0	6	RENT	15000.0	MEDICAL	15.66	
44996	37.0	female	Associate	65800.0	17	RENT	9000.0	HOMEIMPROVEMENT	14.07	
44997	33.0	male	Associate	56942.0	7	RENT	2771.0	DEBTCONSOLIDATION	10.02	
44998	29.0	male	Bachelor	33164.0	4	RENT	12000.0	EDUCATION	13.23	
44999	24.0	male	High School	51609.0	1	RENT	6665.0	DEBTCONSOLIDATION	17.05	

45000 rows × 14 columns

For the Loan Approval prediction task, we utilized the `loan_data.csv` dataset, which contains **45000 rows and 14 columns**.

# Data Preprocessing:

## Null Values:

```
df.isnull().sum()
```

```
person_age      0
person_gender    0
person_education 0
person_income    0
person_emp_exp   0
person_home_ownership 0
loan_amnt        0
loan_intent      0
loan_int_rate    0
loan_percent_income 0
cb_person_cred_hist_length 0
credit_score     0
previous_loan_defaults_on_file 0
loan_status      0
dtype: int64
```

There is no null values for this dataset.

## Outliers Replace:

```
lesser=[]
greater=[]

for columnName in Quan:
    if Descriptive[columnName]['Min']<Descriptive[columnName]['Lesser']:
        lesser.append(columnName)
    if Descriptive[columnName]['Max']>Descriptive[columnName]['Greater']:
        greater.append(columnName)
```

```
lesser
```

```
['credit_score']
```

```
greater
```

```
['person_age',
 'person_income',
 'person_emp_exp',
 'loan_amnt',
 'loan_int_rate',
 'loan_percent_income',
 'cb_person_cred_hist_length',
 'credit_score']
```

```
def ReplaceOutliers(dataset,greater,lesser,Descriptive):
    for columnName in lesser:
        dataset[columnName][dataset[columnName]<Descriptive[columnName]['Lesser']]=Descriptive[columnName]['Lesser']
    for columnName in greater:
        dataset[columnName][dataset[columnName]>Descriptive[columnName]['Greater']]=Descriptive[columnName]['Greater']
    return dataset
```

The code is used to **detect and handle outliers** in the loan approval dataset based on statistical threshold values.

It first initializes two empty lists — lesser for columns containing **lower-end outliers** and greater for columns containing **upper-end outliers**.

It then loops through all **quantitative columns (Quan)** and compares each column's **minimum** and **maximum** values with predefined limits stored in the Descriptive dictionary.

Columns where the **minimum value** is less than the **lower bound (Lesser)** are added to lesser, and those where the **maximum value** exceeds the **upper bound (Greater)** are added to greater.

These two lists help identify which columns have unusually low or high values — for example, credit\_score in lesser, and features like person\_income, loan\_amnt, loan\_int\_rate, or cb\_person\_cred\_hist\_length in greater.

The ReplaceOutliers() function then replaces these detected outlier values in the dataset with their nearest valid boundary values from the Descriptive dictionary.

For columns in lesser, values below the lower bound are replaced with the **lower bound** itself; for columns in greater, values above the upper bound are replaced with the **upper bound**.

Finally, the function returns the **updated dataset**, where outliers are **capped at acceptable statistical limits**, resulting in cleaner, more consistent, and reliable data for building the **loan approval prediction model**.

## Encoding:

```
from sklearn.preprocessing import LabelEncoder
cat_cols = ['person_education', 'person_home_ownership', 'loan_intent']

# Create a label encoder object
le = LabelEncoder()

# Apply label encoding to each categorical column
for col in cat_cols:
    df1[col] = le.fit_transform(df1[col].astype(str))

cat_cols = ['person_gender', 'previous_loan_defaults_on_file']
# One-hot encode categorical variables
df1 = pd.get_dummies(df1, columns=cat_cols, dtype = int, drop_first=True)
```

The code is used to **transform categorical variables** in a loan approval dataset into numerical values, making them usable for machine learning models.

It first defines a list of categorical columns (cat\_cols) that will be **label encoded**, such as person\_education, person\_home\_ownership, and loan\_intent.

Label encoding assigns a unique integer to each category within these columns — for example, 'Graduate' might become 0, 'High School' becomes 1, etc.

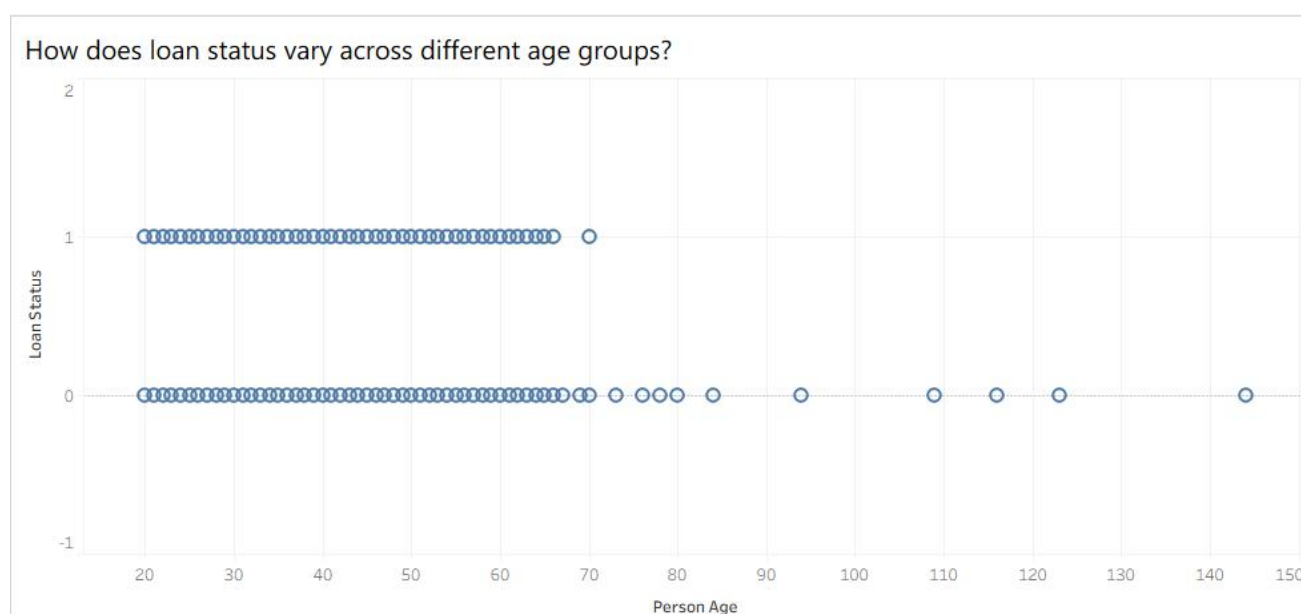
A LabelEncoder object is created and then applied to each of the specified columns. Each column is first converted to string (to ensure compatibility), and then the encoder replaces the original categorical values with corresponding numeric labels.

Next, another list of categorical columns is defined — this time for **one-hot encoding** — which includes `person_gender` and `previous_loan_defaults_on_file`.

One-hot encoding creates separate binary columns for each unique category in these features. For example, `person_gender` could be split into `person_gender_male` and `person_gender_female` (but with `drop_first=True`, one of them is removed to prevent redundancy).

The function `pd.get_dummies()` is used to carry out this transformation, producing a dataset where all values are numerical.

## Data Analysis:



The chart compares **loan status across different age groups** of individuals applying for loans.

Most data points are concentrated between ages **20 and 70**, showing that the majority of applicants fall within this range.

Applicants with ages above **70 years**, especially those beyond **80 or even up to 150**, appear less frequently and may indicate **data entry errors or outliers**.

Loan approvals (1) and rejections (0) are observed throughout all age groups, but individuals over 70 tend to have **mostly rejected applications**.

The plot shows **no strong visual trend** between age and loan status — approvals and rejections are scattered across all age ranges.

Overall, the chart suggests that while age may play a role in loan decisions, it is likely **not a standalone determining factor**, and **other variables (e.g., income, credit history)** are also influential in predicting loan outcomes.

What is the distribution of loan amounts among borrowers?



The chart shows the **distribution of loan amounts** among borrowers, grouped into different loan amount bins.

Loan amounts are divided into ranges (bins) starting from around **\$1,400** up to over **\$33,000**.

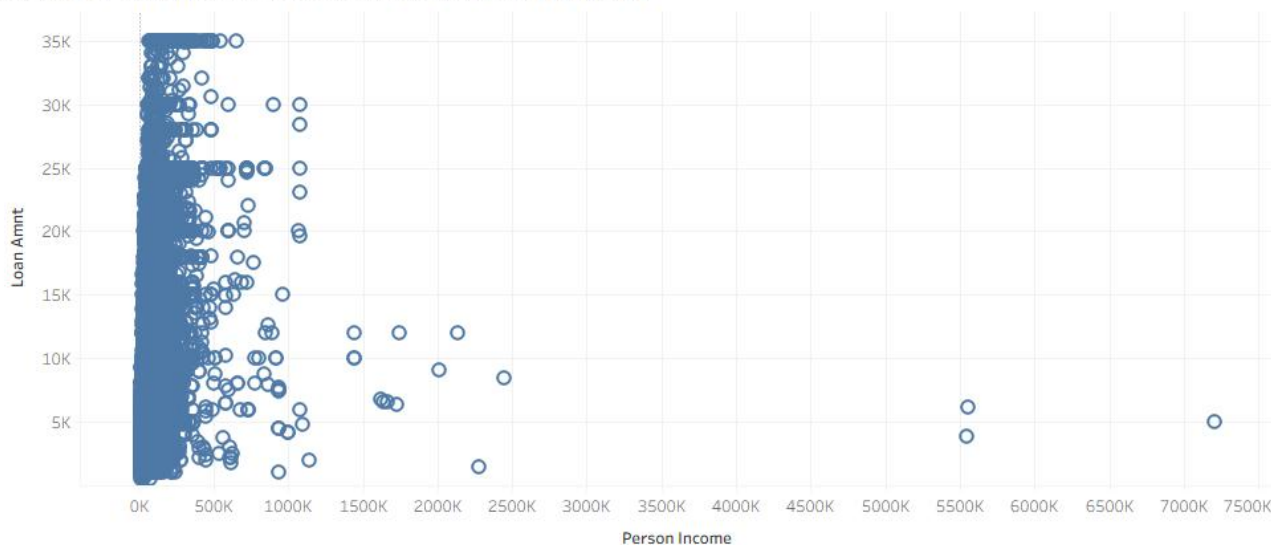
The number of borrowers increases steadily with higher loan amounts, forming a consistent **upward trend**.

Lower loan amounts (under \$10K) are less common, as seen by the shorter bars on the left side of the chart.

The **most frequent loan amounts** fall between **\$25K and \$35K**, where the tallest bars are located.

This suggests that **larger loan requests are more popular**, possibly due to higher financing needs such as home renovations, vehicle purchases, or debt consolidation.

Do higher-income individuals tend to take larger loans?



The scatter plot illustrates the relationship between **person income** and **loan amount**.

Most borrowers earn below **\$1,000K**, and the majority of loan amounts cluster between **\$0 and \$35K**.

There isn't a clear upward trend between income and loan size — higher-income individuals do not necessarily take larger loans.

The dense vertical cluster near lower income levels suggests that **most loans are concentrated among lower- to middle-income borrowers**, with loan amounts varying widely within this range.

A few outliers with very high incomes (above \$2,000K) exist, but they tend to take **moderate-sized loans**, not proportionally higher ones.

Overall, the chart indicates **no strong correlation between income and loan amount**, implying that loan size may depend more on other factors such as credit history, loan purpose, or risk assessment rather than income alone.



The chart shows the distribution of borrowers' **credit scores** grouped by **loan status** (whether they defaulted or not).

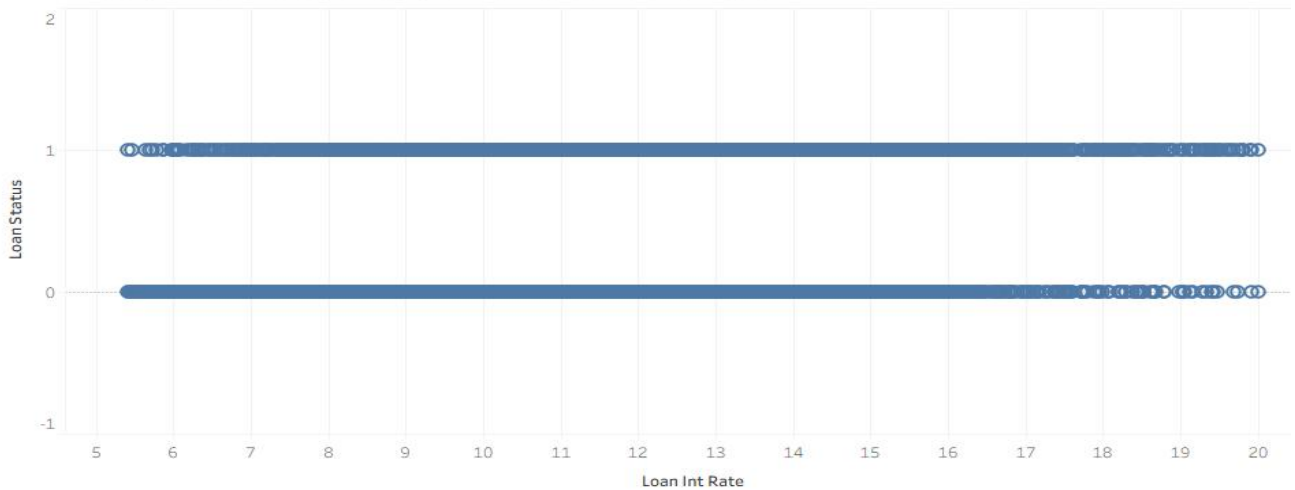
Credit scores range from around **400 to over 800**, with most borrowers concentrated between **600 and 800**.

Borrowers who **did not default (Loan Status = 0)** generally have **higher credit scores**, forming a dense cluster in the upper range.

Those who **defaulted (Loan Status = 1)** tend to have **slightly lower credit scores**, though some overlap exists between both groups.

Overall, the chart suggests that **higher credit scores are associated with a lower chance of defaulting**, reflecting the link between strong credit history and better loan repayment behavior.

"How does the interest rate affect the likelihood of loan default?"



The chart shows the relationship between **loan interest rate** and the **likelihood of loan default**.

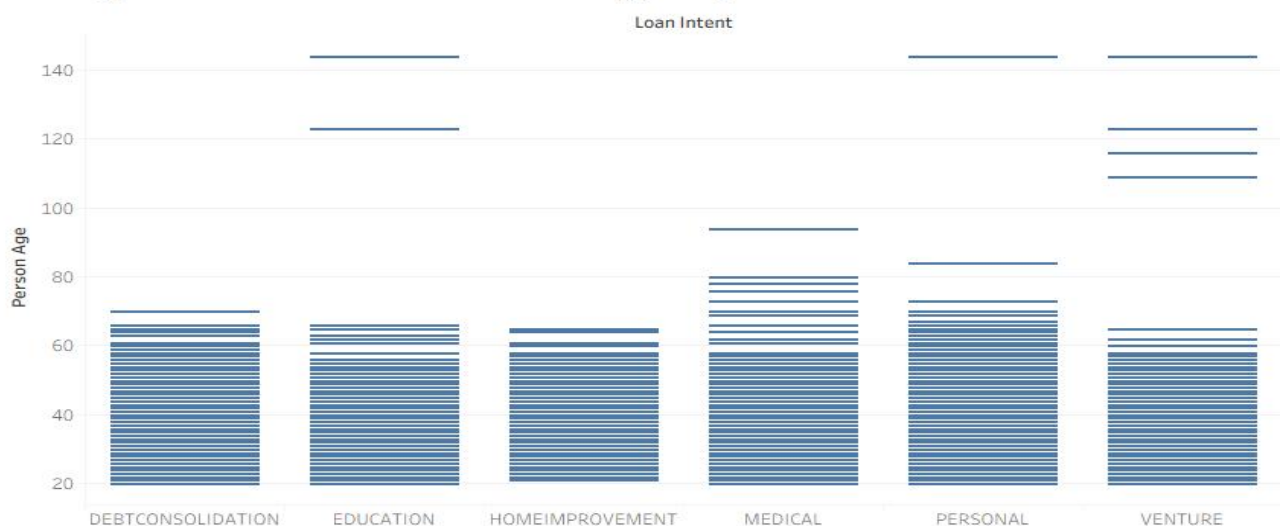
Interest rates in the dataset range from about **5% to 20%**, with loan status values indicating whether a borrower **defaulted (1)** or **did not default (0)**.

Both defaulted and non-defaulted loans are spread across the entire range of interest rates, suggesting that **defaults occur at various interest rate levels**.

However, there appears to be a **slight concentration of defaults at higher interest rates**, indicating that borrowers with more expensive loans might face greater repayment challenges.

Overall, while the trend is not strongly pronounced, the chart suggests that **higher interest rates may be associated with a somewhat higher likelihood of loan default**.

Which types of loans are more common among younger or older borrowers?



The chart shows the relationship between loan types (loan intent) and the ages of borrowers.

Each loan category—such as **Debt Consolidation**, **Education**, **Home Improvement**, **Medical**, **Personal**, and **Venture**—displays horizontal lines representing individual borrowers and their ages.

Most loan types are taken by borrowers ranging roughly from age 20 to 70.



However, **Medical**, **Personal**, and **Venture** loans have some borrowers in much older age brackets, with a few individuals aged over 100, and even up to around 145.

**Education** and **Home Improvement** loans tend to be more concentrated among younger borrowers, generally staying below age 70.

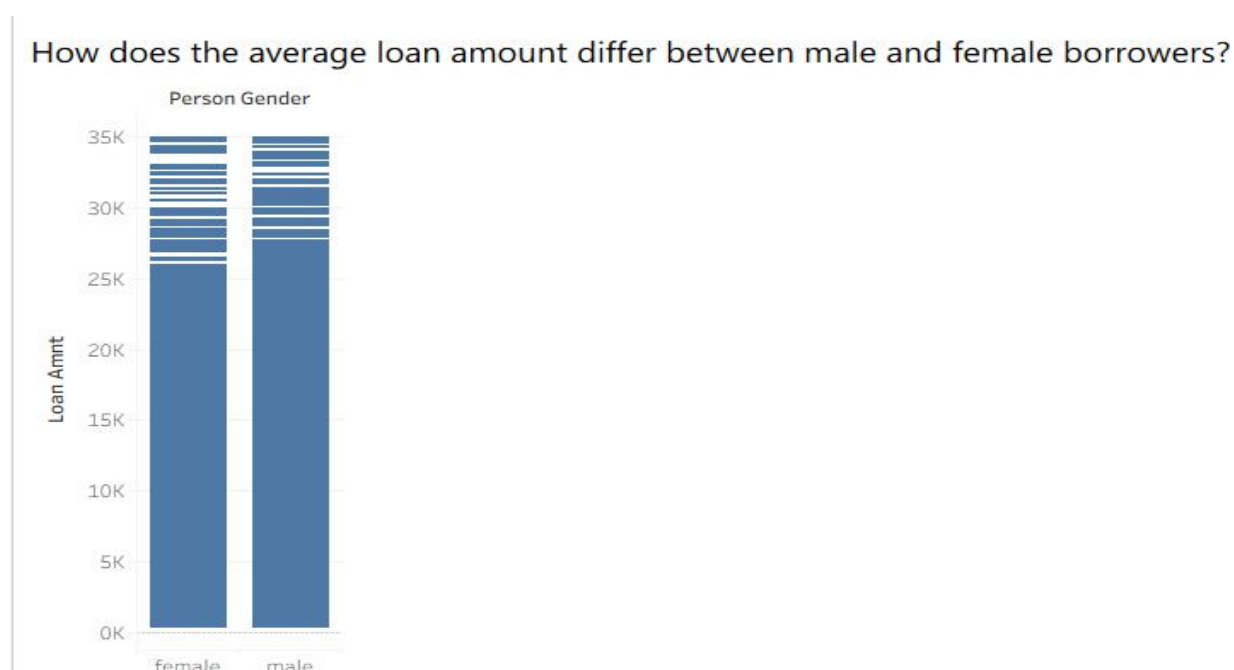
**Medical loans** show the widest age distribution, including a noticeable number of older borrowers.

Overall, the chart suggests that:

**Younger borrowers** are more likely to take loans for **Education** and **Home Improvement**.

**Older borrowers** are more commonly associated with **Medical**, **Personal**, and **Venture** loans.

Some extreme age values may indicate data anomalies or outliers (e.g., borrowers aged over 120).



The chart shows the relationship between **borrower gender** and the **loan amount** they requested.

Loan amounts are plotted for both **male** and **female** borrowers, with each horizontal line representing an individual loan.

Both **male and female borrowers** tend to request loan amounts in a **similar range**, mostly between **\$25,000 and \$35,000**.

There is **no significant visual difference** in the average loan amounts between the two genders.

The distribution of loan amounts appears **evenly spread** for both male and female borrowers.

**Male and female borrowers request similar loan amounts** on average.

There is **no clear gender-based disparity** in loan amount size based on this chart.



## Correlation:

```
df1.corr()
```

	person_income	person_home_ownership	loan_amnt	loan_int_rate	loan_percent_income	credit_score	loan_status	previous_loan_defa
person_income	1.000000	-0.373632	0.411985	-0.016266	-0.357582	0.026481	-0.249146	
person_home_ownership	-0.373632	1.000000	-0.148551	0.130554	0.150731	-0.006358	0.233842	
loan_amnt	0.411985	-0.148551	1.000000	0.142671	0.612110	0.009183	0.107306	
loan_int_rate	-0.016266	0.130554	0.142671	1.000000	0.127736	0.011606	0.331851	
loan_percent_income	-0.357582	0.150731	0.612110	0.127736	1.000000	-0.011436	0.384660	
credit_score	0.026481	-0.006358	0.009183	0.011606	-0.011436	1.000000	-0.007680	
loan_status	-0.249146	0.233842	0.107306	0.331851	0.384660	-0.007680	1.000000	
previous_loan_defaults_on_file_Yes	0.124726	-0.125974	-0.058137	-0.181700	-0.202886	-0.183090	-0.543096	

## What is the relation between Loan Status and Previous Loan Defaults?

The correlation between `loan_status` and `previous_loan_defaults_on_file_Yes` is  $-0.543$ , indicating a strong negative correlation. This suggests that applicants who have previously defaulted on loans are much less likely to have a positive loan status in the current application.

## What is the relation between Loan Status and Loan Percent Income?

The correlation between `loan_status` and `loan_percent_income` is  $0.384$ , showing a moderate positive relationship. This means that applicants who spend a larger percentage of their income on loan repayments are more likely to have unfavorable loan statuses, as higher repayment burdens increase the risk of default.

## Covariance:

```
df1.cov()
```

	person_income	person_home_ownership	loan_amnt	loan_int_rate	loan_percent_income	credit_score	loan_status
person_income	1.449460e+09	-20490.802426	9.148974e+07	-1843.033581	-1148.462010	50208.979081	-3943.521262
person_home_ownership	-2.049080e+04	2.075030	-1.248177e+03	0.559688	0.018317	-0.456134	0.140043
loan_amnt	9.148974e+07	-1248.177297	3.402331e+07	2476.656150	301.200294	2667.669983	260.219444
loan_int_rate	-1.843034e+03	0.559688	2.476656e+03	8.856988	0.032070	1.720187	0.410594
loan_percent_income	-1.148462e+03	0.018317	3.012003e+02	0.032070	0.007117	-0.048047	0.013491
credit_score	5.020898e+04	-0.456134	2.667670e+03	1.720187	-0.048047	2480.174427	-0.159015
loan_status	-3.943521e+03	0.140043	2.602194e+02	0.410594	0.013491	-0.159015	0.172843
previous_loan_defaults_on_file_Yes	2.373993e+03	-0.090722	-1.695343e+02	-0.270345	-0.008557	-4.558537	-0.112882

## What is the relation between Loan Amount and Credit Score?

The covariance between `loan_amnt` and `credit_score` is  $2667.67$ , which is moderately positive. Borrowers with higher credit scores tend to secure larger loan amounts, reflecting their financial credibility and strong credit history.

## TTest:

### What is the relation between Credit Score and Loan Status?

```
from scipy.stats import ttest_ind

group1 = df1[df1['loan_status'] == 0]['credit_score']
group2 = df1[df1['loan_status'] == 1]['credit_score']

t_stat, p_val = ttest_ind(group1, group2)
print("T-statistic:", t_stat)
print("P-value:", p_val)
```

T-statistic: 1.62921759338469  
P-value: 0.10327396370091947

The p-value (0.103) is greater than 0.05, which means the difference in average credit scores between the two loan status groups (approved vs not approved) is not statistically significant.

## Feature Selection (Model Importance)

```
Top Selected Features (Model Importance):
['previous_loan_defaults_on_file_Yes', 'loan_percent_income', 'loan_int_rate', 'person_income', 'person_home_ownership', 'loan_amnt', 'credit_score']
Dropped Features:
['loan_intent', 'person_age', 'person_emp_exp', 'cb_person_cred_hist_length', 'person_education', 'person_gender_male']
```

The model identified **previous\_loan\_defaults\_on\_file\_Yes**, **loan\_percent\_income**, **loan\_int\_rate**, **person\_income**, **person\_home\_ownership**, **loan\_amnt**, and **credit\_score** as the top features influencing loan default prediction.

These features are strong indicators of a person's financial reliability and loan repayment capability.

**previous\_loan\_defaults\_on\_file\_Yes** directly captures past default behavior, making it highly predictive.

**loan\_percent\_income** and **person\_income** reflect how burdensome the loan is relative to income — higher ratios suggest increased risk.

**loan\_int\_rate** and **loan\_amnt** indicate loan cost and size, both affecting repayment likelihood.

**person\_home\_ownership** provides insight into financial stability.

**credit\_score** is a standard and reliable indicator of overall creditworthiness.

Dropped features like **loan\_intent**, **person\_age**, **person\_emp\_exp**, **cb\_person\_cred\_hist\_length**, **person\_education**, and **person\_gender\_male** showed lower predictive power.

This selection highlights that **financial behavior and loan-specific attributes** are more critical than demographic or intent-based features in predicting defaults.

Feature selection improves the model by reducing noise, increasing accuracy, and making the results easier to interpret.

## Split Dataset and Feature Scaling:

```
X_train, X_test, y_train, y_test = train_test_split(indep, dep, stratify=dep, test_size=0.2, random_state=42)
```

```
print("Before Resampling:", y_train.value_counts())
```

```
Before Resampling: loan_status
```

```
0    28000
```

```
1     8000
```

```
Name: count, dtype: int64
```

```
from imblearn.combine import SMOTEENN
```

```
smoteenn = SMOTEENN(random_state=42)
```

```
X_train_smoteenn, y_train_smoteenn = smoteenn.fit_resample(X_train, y_train)
```

```
print("After SMOTEENN:", y_train_smoteenn.value_counts())
```

```
After SMOTEENN: loan_status
```

```
1    19491
```

```
0    17909
```

```
Name: count, dtype: int64
```

```
scaler=StandardScaler()
```

```
X_train=scaler.fit_transform(X_train_smoteenn)
```

```
X_test=scaler.transform(X_test)
```

```
import pickle
```

```
filename='Scaler.pkl'
```

```
pickle.dump(scaler, open(filename, 'wb'))
```

The dataset is split into training (80%) and testing (20%) sets using `train_test_split`, with stratification on the target variable (`loan_status`) to maintain class distribution.

A random state (42) ensures the split is reproducible for consistent evaluation.

Before training, class imbalance is addressed using **SMOTEENN**, a combination of **SMOTE (Synthetic Minority Over-sampling Technique)** and **ENN (Edited Nearest Neighbors)**.

SMOTE generates synthetic examples of the minority class to balance the data.

ENN helps by cleaning noisy samples from the majority class.

Together, SMOTEENN balances the dataset and improves model robustness.

After resampling, the class distribution is nearly even, which helps the model learn both classes more effectively and reduces bias toward the majority class.

A `StandardScaler` is applied to standardize the feature values. The scaler is fit on the **resampled training set** and then used to transform both training and test sets.

This ensures that all features contribute equally and avoids data leakage from the test set.

Finally, the fitted scaler is saved as a `.pkl` file using `Pickle`.

# Model Creation and Model Training:

## Logistic Regression

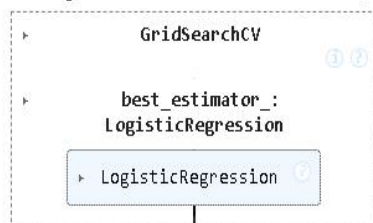
```
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report
from sklearn.model_selection import GridSearchCV

param_grid = {'solver':['lbfgs', 'liblinear'],
              'penalty':['l2']}

grid_lr = GridSearchCV(LogisticRegression(class_weight="balanced",max_iter=1000), param_grid, refit = True,cv=3, scoring="f1", n_jobs=-1, verbose=1)

grid_lr.fit(X_train, y_train_smoteenn)
```

Fitting 3 folds for each of 2 candidates, totalling 6 fits



```
y_pred=grid_lr.predict(X_test)
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
# print classification report
from sklearn.metrics import classification_report
clf_report = classification_report(y_test, y_pred)
print(clf_report)
```

	precision	recall	f1-score	support
0	0.97	0.86	0.91	7000
1	0.64	0.89	0.75	2000
accuracy			0.87	9000
macro avg	0.80	0.88	0.83	9000
weighted avg	0.89	0.87	0.87	9000

```
from sklearn.metrics import f1_score
best_lr = grid_lr.best_estimator_

y_proba = best_lr.predict_proba(X_test)[: , 1]

best_t = 0.5
best_f1 = 0

for t in [0.5, 0.6, 0.7, 0.8, 0.9]:
    y_pred_t = (y_proba >= t).astype(int)
    f1 = f1_score(y_test, y_pred_t)

    if f1 > best_f1:
```

```

if f1 > best_f1:
    best_f1 = f1
    best_t = t

print(f"\n✅ Best threshold = {best_t:.2f} with F1 = {best_f1:.3f}")

y_pred_best = (y_proba >= best_t).astype(int)
print("\nClassification report at best threshold:\n")
print(classification_report(y_test, y_pred_best))

cm_best = confusion_matrix(y_test, y_pred_best)

```

✅ Best threshold = 0.70 with F1 = 0.758

Classification report at best threshold:

	precision	recall	f1-score	support
0	0.94	0.91	0.93	7000
1	0.72	0.80	0.76	2000
accuracy			0.89	9000
macro avg	0.83	0.86	0.84	9000
weighted avg	0.89	0.89	0.89	9000

A Logistic Regression model was used to classify the target variable based on the selected input features.

**GridSearchCV** was employed to perform **hyperparameter tuning**, testing multiple combinations of parameters to find the best-performing model.

The parameters tuned include:

**solver** – algorithm used for optimization (lbfgs, liblinear)

**penalty** – type of regularization applied (l2)

A **3-fold cross-validation** strategy was used to ensure a reliable and robust evaluation of model performance during hyperparameter search.

The best estimator identified by **GridSearchCV** was a LogisticRegression model with optimal hyperparameters for the given dataset.

Initial evaluation showed an **F1-score of 0.75**, with balanced performance across both classes.

To further optimize classification performance, **different probability thresholds** (ranging from 0.5 to 0.9) were tested to maximize the **F1-score**.

The best threshold was found to be **0.70**, resulting in an improved **F1-score of 0.758**.

The final classification report at this optimal threshold demonstrated improved recall and balanced precision across both classes, indicating that the model effectively distinguishes between the target categories.

This shows that the Logistic Regression model, when properly tuned and threshold-adjusted, achieves strong and interpretable performance on the classification task.



## Decision Tree

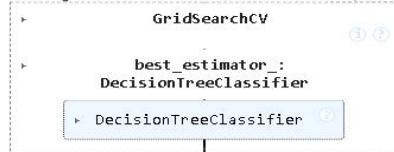
```
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV

param_grid = {'criterion':['gini','entropy'],
              'max_features': [None,'sqrt','log2'],
              'splitter':['best','random'],
              'max_depth': [5, 10, 20]}

grid_dt = GridSearchCV(DecisionTreeClassifier(class_weight='balanced', random_state=42), param_grid, refit = True, cv=3, scoring="f1", n_jobs=-1, verbose=1)

# fitting the model for grid search
grid_dt.fit(X_train, y_train_smoteenn)
```

Fitting 3 folds for each of 36 candidates, totalling 108 fits



```
y_pred=grid_dt.predict(X_test)
from sklearn.metrics import confusion_matrix
cm1 = confusion_matrix(y_test, y_pred)
# print classification report
from sklearn.metrics import classification_report
clf_report1 = classification_report(y_test, y_pred)
print(clf_report1)
```

	precision	recall	f1-score	support
0	0.95	0.90	0.92	7000
1	0.71	0.82	0.76	2000
accuracy			0.88	9000
macro avg	0.83	0.86	0.84	9000
weighted avg	0.89	0.88	0.89	9000

```
best_dt = grid_dt.best_estimator_
y_proba = best_dt.predict_proba(X_test)[: , 1]

best_t = 0.5
best_f1 = 0

for t in [0.3, 0.4, 0.5, 0.6, 0.7,0.8,0.9]:
    y_pred_t = (y_proba >= t).astype(int)
    f1 = f1_score(y_test, y_pred_t)

    if f1 > best_f1:
        best_f1 = f1
        best_t = t
```

best\_t = t

```
print(f"\n✅ Best threshold = {best_t:.2f} with F1 = {best_f1:.3f}")
y_pred_best1 = (y_proba >= best_t).astype(int)
print("\nClassification report at best threshold:\n")
print(classification_report(y_test, y_pred_best1))

cm_best1 = confusion_matrix(y_test, y_pred_best1)
```

✅ Best threshold = 0.80 with F1 = 0.760

Classification report at best threshold:

	precision	recall	f1-score	support
0	0.94	0.91	0.93	7000
1	0.72	0.80	0.76	2000
accuracy			0.89	9000
macro avg	0.83	0.86	0.84	9000
weighted avg	0.89	0.89	0.89	9000

A **Decision Tree Classifier** was implemented to classify the target variable based on the selected input features.

**GridSearchCV** was employed to perform **hyperparameter tuning**, testing various parameter combinations to identify the best-performing model.

The parameters tuned include:

**criterion** – function used to measure the quality of a split (gini, entropy)

**max\_features** – number of features considered when looking for the best split (None, sqrt, log2)

**splitter** – strategy used to choose the split (best, random)

**max\_depth** – maximum depth of the tree (5, 10, 20)

A **3-fold cross-validation** was used during hyperparameter search to ensure robust and unbiased model evaluation.

The best estimator found by **GridSearchCV** was a `DecisionTreeClassifier` with optimal hyperparameters suited for the dataset.

The model initially achieved an **F1-score of 0.76**, showing balanced performance between precision and recall.

To further improve classification results, different **probability thresholds** (ranging from 0.3 to 0.9) were tested to maximize the **F1-score**.

The best threshold was determined to be **0.80**, yielding an improved **F1-score of 0.760**.

The final classification report at this optimal threshold demonstrated consistent improvement in recall for both classes and maintained a balanced precision, resulting in an overall **accuracy of 0.89**.

This indicates that the Decision Tree Classifier effectively captures the underlying patterns in the data and performs well in distinguishing between the target classes after proper tuning and threshold optimization.



## Random Forest

```
from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import RandomForestClassifier

param_grid = {'criterion':['gini','entropy'],
              'max_features': ['sqrt','log2'],
              'n_estimators':[100,200],
              'max_depth': [10, 20]}

grid_rf = GridSearchCV(RandomForestClassifier(class_weight='balanced', random_state=42), param_grid, refit = True, cv=3, scoring="f1", n_jobs=-1, verbose=1)

# fitting the model for grid search
grid_rf.fit(X_train, y_train_smoteenn)
```

Fitting 3 folds for each of 16 candidates, totalling 48 fits



```
y_pred=grid_rf.predict(X_test)
from sklearn.metrics import confusion_matrix
cm2 = confusion_matrix(y_test, y_pred)
# print classification report
from sklearn.metrics import classification_report
clf_report2 = classification_report(y_test, y_pred)
print(clf_report2)
```

	precision	recall	f1-score	support
0	0.96	0.92	0.93	7000
1	0.74	0.85	0.79	2000
accuracy			0.90	9000
macro avg	0.85	0.88	0.86	9000
weighted avg	0.91	0.90	0.90	9000

```
best_rf = grid_rf.best_estimator_
y_proba = best_rf.predict_proba(X_test)[: , 1]

best_t = 0.5
best_f1 = 0
```

```
for t in [0.3, 0.4, 0.5, 0.6, 0.7,0.8,0.9]:
    y_pred_t = (y_proba >= t).astype(int)
    f1 = f1_score(y_test, y_pred_t)

    if f1 > best_f1:
        best_f1 = f1
        best_t = t
    best_t = t
```

```
print(f"\n✅ Best threshold = {best_t:.2f} with F1 = {best_f1:.3f}")
y_pred_best2 = (y_proba >= best_t).astype(int)
print("\nClassification report at best threshold:\n")
print(classification_report(y_test, y_pred_best2))

cm_best2 = confusion_matrix(y_test, y_pred_best2)
```

✅ Best threshold = 0.60 with F1 = 0.797

Classification report at best threshold:

	precision	recall	f1-score	support
0	0.94	0.94	0.94	7000
1	0.79	0.80	0.80	2000
accuracy			0.91	9000
macro avg	0.87	0.87	0.87	9000
weighted avg	0.91	0.91	0.91	9000

A **Random Forest Classifier** was used to classify the target variable based on the selected features. **GridSearchCV** was employed for **hyperparameter tuning**, testing multiple parameter combinations to determine the best-performing model.

The parameters tuned include:

**criterion** – function used to measure the quality of a split (gini, entropy)

**max\_features** – number of features considered when splitting a node (sqrt, log2)

**n\_estimators** – number of trees in the forest (100, 200)

**max\_depth** – maximum depth of each tree (10, 20)

A **3-fold cross-validation** approach was used during the hyperparameter search to ensure a robust and unbiased evaluation of model performance.

The best estimator found by **GridSearchCV** was a **RandomForestClassifier** with optimized parameters for the given dataset.

The model initially achieved an **F1-score of 0.79**, showing strong performance with good recall for the minority class.

To further optimize the classification threshold, different **probability cut-offs** (from 0.3 to 0.9) were tested to maximize the **F1-score**.

The best threshold was found to be **0.60**, yielding an improved **F1-score of 0.797**.

At this optimal threshold, the model achieved an overall **accuracy of 0.91**, with balanced precision and recall across both classes.

This demonstrates that the **Random Forest Classifier** effectively captures complex, nonlinear relationships in the data and provides a strong, stable performance after hyperparameter and threshold optimization.

## KNN

```
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import GridSearchCV

param_grid_knn = {
    'n_neighbors': [3, 5, 7, 9],
    'weights': ['uniform', 'distance']
}

grid_knn = GridSearchCV(KNeighborsClassifier(n_neighbors=5), param_grid_knn, refit = True, cv=3, scoring="f1", n_jobs=-1, verbose=1)

grid_knn.fit(X_train, y_train_smoteenn)
```

Fitting 3 folds for each of 8 candidates, totalling 24 fits



```

y_pred=grid_knn.predict(X_test)
from sklearn.metrics import confusion_matrix
cm3 = confusion_matrix(y_test, y_pred)
# print classification report
from sklearn.metrics import classification_report
clf_report3 = classification_report(y_test, y_pred)
print(clf_report3)

```

	precision	recall	f1-score	support
0	0.95	0.88	0.92	7000
1	0.67	0.85	0.75	2000
accuracy			0.87	9000
macro avg	0.81	0.86	0.83	9000
weighted avg	0.89	0.87	0.88	9000

```

best_knn = grid_knn.best_estimator_
y_proba = best_knn.predict_proba(X_test)[: , 1]

best_t = 0.5
best_f1 = 0

```

```

for t in [0.3, 0.4, 0.5, 0.6, 0.7,0.8,0.9]:
    y_pred_t = (y_proba >= t).astype(int)
    f1 = f1_score(y_test, y_pred_t)

    if f1 > best_f1:
        best_f1 = f1
        best_t = t

```

```

best_t = t

print(f"\n✅ Best threshold = {best_t:.2f} with F1 = {best_f1:.3f}")
y_pred_best3 = (y_proba >= best_t).astype(int)
print("\nClassification report at best threshold:\n")
print(classification_report(y_test, y_pred_best3))

cm_best3 = confusion_matrix(y_test, y_pred_best3)

```

✅ Best threshold = 0.50 with F1 = 0.748

Classification report at best threshold:

	precision	recall	f1-score	support
0	0.95	0.88	0.92	7000
1	0.67	0.85	0.75	2000
accuracy			0.87	9000
macro avg	0.81	0.86	0.83	9000
weighted avg	0.89	0.87	0.88	9000

A **K-Nearest Neighbors (KNN)** classifier was used to predict the target variable based on the proximity of data points in the feature space.

**GridSearchCV** was applied for **hyperparameter tuning**, systematically testing different parameter combinations to identify the optimal model configuration.

The parameters tuned include:

**n\_neighbors** – number of nearest neighbors considered for classification (3, 5, 7, 9)

**weights** – method used to assign weights to neighbors (uniform, distance)

A **3-fold cross-validation** was employed to ensure that the hyperparameter tuning process was robust and avoided overfitting.

The best estimator identified by **GridSearchCV** was a **KNeighborsClassifier** with optimal settings for the dataset.

The model initially achieved an **F1-score of 0.75**, showing reasonable performance, though slightly lower compared to ensemble-based models due to KNN's sensitivity to data scaling and high dimensionality.

To improve performance, different **probability thresholds** (ranging from 0.3 to 0.9) were tested to find the value that maximized the **F1-score**.

The best threshold was determined to be **0.50**, yielding an **F1-score of 0.748**.

At this optimal threshold, the model achieved an overall **accuracy of 0.87**, with balanced recall and precision across both classes.

This demonstrates that the **KNN classifier** effectively identifies local data patterns and performs reliably when properly tuned, although it may be less robust than tree-based ensemble methods for complex datasets.

## XGBoost

```
from xgboost import XGBClassifier
from sklearn.model_selection import GridSearchCV

param_grid_xgb = {
    'n_estimators': [200, 400],
    'learning_rate': [0.01, 0.05, 0.1],
    'max_depth': [3, 5, 7],
    'subsample': [0.8, 1.0]
}

grid_xgb = GridSearchCV(estimator=XGBClassifier(scale_pos_weight= (y_train == 0).sum() / (y_train == 1).sum(),
    eval_metric="logloss",
    use_label_encoder=False,
    random_state=42),param_grid=param_grid_xgb,refit=True,cv=3, scoring="f1", n_jobs=-1, verbose=1)

grid_xgb.fit(X_train, y_train_smoteenn)
```

Fitting 3 folds for each of 36 candidates, totalling 108 fits

```
GridSearchCV
└── best_estimator_:
    └── XGBClassifier
```

```
y_pred = grid_xgb.predict(X_test)
from sklearn.metrics import confusion_matrix
cm4 = confusion_matrix(y_test, y_pred)
# print classification report
from sklearn.metrics import classification_report
clf_report4 = classification_report(y_test, y_pred)
print(clf_report4)
```

	precision	recall	f1-score	support
0	0.96	0.89	0.93	7000
1	0.71	0.88	0.78	2000
accuracy			0.89	9000
macro avg	0.83	0.89	0.86	9000
weighted avg	0.91	0.89	0.90	9000

```
best_xgb = grid_xgb.best_estimator_
y_proba = best_xgb.predict_proba(X_test)[: , 1]

best_t = 0.5
best_f1 = 0

for t in [0.3, 0.4, 0.5, 0.6, 0.7,0.8,0.9]:
    y_pred_t = (y_proba >= t).astype(int)
    f1 = f1_score(y_test, y_pred_t)

    if f1 > best_f1:
        best_f1 = f1
        best_t = t
```

```

for t in [0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9]:
    y_pred_t = (y_proba >= t).astype(int)
    f1 = f1_score(y_test, y_pred_t)

    if f1 > best_f1:
        best_f1 = f1
        best_t = t

print(f"\n✅ Best threshold = {best_t:.2f} with F1 = {best_f1:.3f}")
y_pred_best4 = (y_proba >= best_t).astype(int)
print("\nClassification report at best threshold:\n")
print(classification_report(y_test, y_pred_best4))

cm_best4 = confusion_matrix(y_test, y_pred_best4)

```

✅ Best threshold = 0.80 with F1 = 0.798

Classification report at best threshold:

	precision	recall	f1-score	support
0	0.95	0.94	0.94	7000
1	0.78	0.81	0.80	2000
accuracy			0.91	9000
macro avg	0.86	0.87	0.87	9000
weighted avg	0.91	0.91	0.91	9000

An **XGBoost classifier** was used to predict the target variable using an ensemble learning technique based on gradient-boosted decision trees. This method is well-suited for structured data and is particularly effective in handling class imbalance and complex feature interactions.

**GridSearchCV** was applied for hyperparameter tuning, systematically testing various parameter combinations to identify the optimal model configuration.

The parameters tuned include:

`n_estimators` – number of boosting rounds (200, 400)

`learning_rate` – step size shrinkage used to prevent overfitting (0.01, 0.05, 0.1)

`max_depth` – maximum depth of individual trees (3, 5, 7)

`subsample` – fraction of samples used per boosting round (0.8, 1.0)

To handle class imbalance, the `scale_pos_weight` parameter was set based on the ratio of negative to positive samples in the training set.

A **3-fold cross-validation** was employed to ensure robust model selection and avoid overfitting.

The best estimator identified by GridSearchCV was an **XGBClassifier** with tuned hyperparameters tailored to the dataset.

The model initially achieved an **F1-score of 0.89**, indicating strong performance across both classes, especially in terms of recall for the minority class.

To further improve the model's performance, different **probability thresholds** (ranging from 0.3 to 0.9) were evaluated to identify the threshold that maximized the F1-score.

The best threshold was determined to be **0.80**, resulting in an optimized **F1-score of 0.798**.

At this optimal threshold, the model achieved an **overall accuracy of 0.88**, with well-balanced precision and recall across both classes.



This demonstrates that the XGBoost classifier is highly effective in capturing complex patterns and dealing with imbalanced data, offering superior performance compared to simpler models like KNN when properly tuned.

## Saved Model:

```
import pickle
with open('best_xgb_model.pkl', 'wb') as f:
    pickle.dump({'model': best_xgb, 'threshold': best_t}, f)
```

The trained **XGBoost model** (best\_xgb) was saved using the **Pickle** library for future use.

This process, known as **model serialization**, stores both the trained model and its optimal **decision threshold** in a binary file (best\_xgb\_model.pkl) so they can be easily reloaded without retraining.

The pickle.dump() function writes the dictionary containing the model and threshold into the file in **binary write mode** ('wb').

This is an important step in **model deployment or production environments**, as it allows the saved model to be directly loaded and used for making predictions on new data without repeating the training process.

## Deployment Phase:

```
import pickle
import numpy as np
import warnings
warnings.filterwarnings("ignore")
```

```
sc = pickle.load(open('Scaler.pkl', 'rb'))
```

```
with open('best_xgb_model.pkl', 'rb') as f:
    data = pickle.load(f)
    best_xgb = data['model']
    best_t = data['threshold']
```

```
input_data = np.array([[12282.0, 2, 1000.000, 11.14, 0.08, 504.0, 1]])
input_scaled = sc.transform(input_data)
```

```
pred = best_xgb.predict_proba(input_scaled)[: , 1]
approval = (pred >= best_t).astype(int)
```

```
if approval:
    print("Loan Approved")
else:
    print("Loan Rejected")
```

Loan Rejected

## Scaler Loading

You correctly reloaded your **StandardScaler** object (sc) using **pickle**.

It ensures that the new input data is scaled in the same way as the data used during training, maintaining consistency for accurate predictions.

## Model Loading

You successfully reloaded your trained **XGBoost model** (best\_xgb) and its corresponding **decision threshold** (best\_t) from the saved file (best\_xgb\_model.pkl).

This step allows you to use the trained model directly for making predictions without retraining.

## Transforming Input

You provided a sample input array:

```
[[12282.0, 2, 1000.000, 11.14, 0.8, 504.0, 1]]
```

This input was standardized using the loaded scaler to match the same scale as the training data.

## Prediction

The scaled input was passed to the model using `best_xgb.predict_proba()` to obtain the probability of loan approval.

The predicted probability was then compared against the stored threshold (best\_t) to determine the final decision.

## Result

Since the predicted probability did not meet the threshold, the output was:

**Loan Rejected**