

```
import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

import seaborn as sns


from sklearn.model_selection import train_test_split, GridSearchCV

from sklearn.preprocessing import LabelEncoder, StandardScaler

from sklearn.impute import SimpleImputer

from sklearn.ensemble import RandomForestClassifier

from sklearn.linear_model import LogisticRegression

from xgboost import XGBClassifier

from sklearn.metrics import accuracy_score, classification_report, confusion_matrix, roc_auc_score,
roc_curve

import warnings

warnings.filterwarnings("ignore")


# Assuming the dataset is in the default Colab sample data directory
try:

    # Try reading the CSV file - ensure the path is correct and the file is a valid CSV

    from google.colab import files

    uploaded = files.upload()

import io

df = pd.read_csv(/loan_approval_dataset.csv)


except FileNotFoundError:
```

```
print("Error: The file '/cuda-keyring_1.1-1_all.deb' was not found. Please ensure the file is uploaded or provide the correct path.")
```

```
exit()
```

```
except Exception as e:
```

```
print(f"An error occurred while reading the CSV file: {e}")
```

```
print("Please ensure the file is a valid CSV and the path is correct.")
```

```
exit()
```

```
df.head()
```

```
df.info()
```

```
df.describe()
```

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

```
df['Loan_Status'].value_counts(normalize=True)
```

```
# Separate numerical and categorical columns
```

```
num_cols = df.select_dtypes(include=['int64', 'float64']).columns
```

```
cat_cols = df.select_dtypes(include=['object']).columns
```

```
# Impute missing numerical values with mean
```

```
df[num_cols] = df[num_cols].apply(lambda x: x.fillna(x.mean()))
```

```
# Impute missing categorical values with mode
```

```
df[cat_cols] = df[cat_cols].apply(lambda x: x.fillna(x.mode()[0]))
```

```
le = LabelEncoder()
```

```
for col in cat_cols:
```

```

df[col] = le.fit_transform(df[col])

X = df.drop(['Loan_ID', 'Loan_Status'], axis=1) # Drop unnecessary columns
y = df['Loan_Status']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

scaler = StandardScaler()

X_train = scaler.fit_transform(X_train)

X_test = scaler.transform(X_test)

lr = LogisticRegression()

lr.fit(X_train, y_train)

y_pred_lr = lr.predict(X_test)

rf = RandomForestClassifier()

rf.fit(X_train, y_train) # This was y_train

y_pred_rf = rf.predict(X_test)

xgb = XGBClassifier(use_label_encoder=False, eval_metric='logloss')

xgb.fit(X_train, y_train)

y_pred_xgb = xgb.predict(X_test)

def evaluate_model(y_test, y_pred, model_name):

    print(f"\n--- {model_name} ---")

    print("Accuracy:", accuracy_score(y_test, y_pred))

    print("Classification Report:\n", classification_report(y_test, y_pred))

    print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))

evaluate_model(y_test, y_pred_lr, "Logistic Regression")

```

```
evaluate_model(y_test, y_pred_rf, "Random Forest")
```

```
evaluate_model(y_test, y_pred_xgb, "XGBoost")
```

```
y_prob = xgb.predict_proba(X_test)[:,-1]
```

```
fpr, tpr, thresholds = roc_curve(y_test, y_prob)
```

```
plt.figure(figsize=(8,6))
```

```
plt.plot(fpr, tpr, label='XGBoost')
```

```
plt.plot([0,1],[0,1], 'k--')
```

```
plt.xlabel('False Positive Rate')
```

```
plt.ylabel('True Positive Rate')
```

```
plt.title('ROC Curve - XGBoost')
```

```
plt.legend()
```

```
plt.show()
```

```
print("ROC-AUC Score:", roc_auc_score(y_test, y_prob))
```

```
importances = xgb.feature_importances_
```

```
features = X.columns
```

```
plt.figure(figsize=(10,6))
```

```
sns.barplot(x=importances, y=features)
```

```
plt.title("Feature Importance - XGBoost")
```

```
plt.show()
```

Here's a **summary of your Loan Approval Prediction project** based on the uploaded notebook:

□ Loan Approval Prediction Project – Summary

✓ Project Goal:

Build a machine learning model to automatically predict whether a loan should be approved or not, based on applicant data (demographics, income, credit history, etc.).

🔧 Tools & Libraries Used:

- **Data Manipulation:** pandas, numpy
 - **Visualization:** matplotlib, seaborn
 - **ML Models:** LogisticRegression, RandomForestClassifier, XGBClassifier
 - **Evaluation:** accuracy_score, classification_report, roc_auc_score, confusion_matrix, roc_curve
 - **Preprocessing:** LabelEncoder, StandardScaler, SimpleImputer
 - **Environment:** Google Colab (with file upload)
-

🔍 Key Steps Executed:

1. Data Loading

- File upload via Google Colab.
- Attempt to read loan_approval_dataset.csv.

2. Data Inspection

- Used .head(), .info(), .describe(), .isnull().sum() to understand the dataset.
- Checked target variable distribution: Loan_Status.

3. Data Cleaning

- Handled missing values:
 - Numerical: filled with mean.
 - Categorical: filled with mode.

4. Encoding

- All categorical columns encoded using LabelEncoder.

5. Feature Selection

- Dropped Loan_ID and Loan_Status from features.
- Split features (X) and target (y).

6. Train-Test Split

- 80/20 split using train_test_split.

7. Feature Scaling

- Standardized features using StandardScaler.

8. Model Training

Trained 3 classification models:

- **Logistic Regression**
- **Random Forest**
- **XGBoost**

⚠ Small error: rf.fit(X_train, y_test) should be rf.fit(X_train, y_train).

9. Model Evaluation

- Accuracy, classification report, and confusion matrix for all models using a custom evaluate_model() function.
- ROC Curve plotted for XGBoost.
- ROC-AUC Score printed.

10. Feature Importance

- XGBoost feature importances plotted using Seaborn.



Visual Output:

- **ROC Curve** for XGBoost model
- **Feature Importance Chart** (top predictive features)

✓ **Overall Result:**

You have successfully built and compared three machine learning models to predict loan approvals and visualized their performance using evaluation metrics and feature importance.

Would you like me to fix the RandomForestClassifier bug and return a clean Jupyter Notebook version or a PDF report of the entire project?