## ⌄ Project Report: Loan Approval Prediction System

## 1. Introduction/Project Overview

This project aims to develop a system to predict loan approval status based on various applicant and loan characteristics. The goal is to assist financial institutions in automating and streamlining the loan approval process, reducing manual effort and potential biases.

## ⌄ 2. Data Preparation

### Subtask:

Load the dataset, handle missing values, encode categorical variables, separate features and target, and split the data into training and testing sets.

```python
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.preprocessing import LabelEncoder
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix

# Assuming the dataset is named 'loan_approval_dataset.csv' and is in the /content/ directory
try:
    df = pd.read_csv('/content/loan_approval_dataset.csv')
except FileNotFoundError:
    print("Error: 'loan_approval_dataset.csv' not found. Using placeholder data.")
    # Using a larger placeholder dataframe for demonstration purposes if the file is not found
    data = {
        'Loan_ID': [f'LP00{i:04d}' for i in range(1, 101)],
        'Gender': np.random.choice(['Male', 'Female'], 100),
        'Married': np.random.choice(['Yes', 'No'], 100),
        'Dependents': np.random.choice([0, 1, 2, 3], 100),
        'Education': np.random.choice(['Graduate', 'Not Graduate'], 100),
        'Self_Employed': np.random.choice(['Yes', 'No'], 100),
        'ApplicantIncome': np.random.randint(1500, 10000, 100),
        'CoapplicantIncome': np.random.randint(0, 5000, 100),
        'LoanAmount': np.random.randint(50, 500, 100),
        'Loan_Amount_Term': np.random.choice([120, 180, 360, 480], 100),
        'Credit_History': np.random.choice([0, 1], 100),
        'Property_Area': np.random.choice(['Urban', 'Rural', 'Semiurban'], 100),
        'Loan_Status': np.random.choice(['Y', 'N'], 100, p=[0.7, 0.3]) # More balanced
    }
    df = pd.DataFrame(data)
    # Introduce some NaNs for testing
    for col in ['LoanAmount', 'Credit_History', 'Self_Employed', 'Dependents']:
        df.loc[df.sample(frac=0.05).index, col] = np.nan


# Handle missing values - Using mode for categorical and median for numerical
for column in ['Gender', 'Married', 'Dependents', 'Self_Employed', 'Loan_Amount_Term', 'Credit_History']:
    if df[column].isnull().any():
        df[column] = df[column].fillna(df[column].mode()[0])

if df['LoanAmount'].isnull().any():
    df['LoanAmount'] = df['LoanAmount'].fillna(df['LoanAmount'].median())


# Encode categorical variables
categorical_cols = ['Gender', 'Married', 'Dependents', 'Education', 'Self_Employed', 'Property_Area']
df = pd.get_dummies(df, columns=categorical_cols, drop_first=True)

# Encode target variable
if 'Loan_Status' in df.columns:
    le = LabelEncoder()
    df['Loan_Status'] = le.fit_transform(df['Loan_Status'])


# Separate features and target variable
if 'Loan_Status' in df.columns:
    X = df.drop(['Loan_ID', 'Loan_Status'], axis=1) # Dropping Loan_ID as it's an identifier
    y = df['Loan_Status']
else:
    # Assuming a placeholder for X if Loan_Status is not in columns
```

```
    X = df.drop(['Loan_ID'], axis=1)
    y = pd.Series([0] * len(X), name='Loan_Status') # Placeholder target


# Split data into training and testing sets
# Using stratify to handle potential imbalance in the target variable
# Check if stratification is possible
if len(y.unique()) > 1 and y.value_counts().min() >= 2:
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42, stratify=y)
    print("\nData split with stratification.")
else:
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
    print("\nData split without stratification (target variable has only one class or insufficient samples).")


display(X_train.head())
display(y_train.head())
```

⊋ Error: 'loan_approval_dataset.csv' not found. Using placeholder data.

Data split with stratification.

| | ApplicantIncome | CoapplicantIncome | LoanAmount | Loan_Amount_Term | Credit_History | Gender_Male | Married_Yes | Dependents |
|---|---|---|---|---|---|---|---|---|
| 70 | 2706 | 733 | 136.0 | 480 | 0.0 | True | True | |
| 64 | 8307 | 3081 | 499.0 | 360 | 1.0 | True | True | |
| 15 | 1518 | 2672 | 196.0 | 360 | 1.0 | True | True | |
| 48 | 6338 | 921 | 251.0 | 120 | 0.0 | True | False | |
| 17 | 1997 | 4798 | 237.0 | 480 | 1.0 | False | True | |

| | Loan_Status |
|---|---|
| 70 | 0 |
| 64 | 1 |
| 15 | 1 |
| 48 | 1 |
| 17 | 1 |

**dtype:** int64

## ⌄ 3. Model Selection and Training

### Subtask:

Choose appropriate machine learning models (e.g., Logistic Regression, Decision Tree, Random Forest) for the prediction task and train them on the training data.

```
# Instantiate the classifiers
logistic_regression_model = LogisticRegression()
decision_tree_model = DecisionTreeClassifier()
random_forest_model = RandomForestClassifier()

# Train the classifiers
logistic_regression_model.fit(X_train, y_train)
decision_tree_model.fit(X_train, y_train)
random_forest_model.fit(X_train, y_train)

print("Models trained successfully.")
```

⊋ /usr/local/lib/python3.12/dist-packages/sklearn/linear_model/_logistic.py:465: ConvergenceWarning: lbfgs failed to conve
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
Models trained successfully.

## ⌄ 4. Model Evaluation

## Subtask:

Evaluate the performance of the trained models using appropriate metrics such as accuracy, precision, recall, and F1-score. Analyze the results and compare the performance of different models.

```
# Make predictions on the test set
lr_predictions = logistic_regression_model.predict(X_test)
dt_predictions = decision_tree_model.predict(X_test)
rf_predictions = random_forest_model.predict(X_test)

# Evaluate Logistic Regression
print("Logistic Regression Performance:")
print("Accuracy:", accuracy_score(y_test, lr_predictions))
print("Classification Report:\n", classification_report(y_test, lr_predictions))
print("Confusion Matrix:\n", confusion_matrix(y_test, lr_predictions))

# Evaluate Decision Tree
print("\nDecision Tree Performance:")
print("Accuracy:", accuracy_score(y_test, dt_predictions))
print("Classification Report:\n", classification_report(y_test, dt_predictions))
print("Confusion Matrix:\n", confusion_matrix(y_test, dt_predictions))

# Evaluate Random Forest
print("\nRandom Forest Performance:")
print("Accuracy:", accuracy_score(y_test, rf_predictions))
print("Classification Report:\n", classification_report(y_test, rf_predictions))
print("Confusion Matrix:\n", confusion_matrix(y_test, rf_predictions))
```

```
Logistic Regression Performance:
Accuracy: 0.75
Classification Report:
              precision    recall  f1-score   support

           0       0.50      0.40      0.44         5
           1       0.81      0.87      0.84        15

    accuracy                           0.75        20
   macro avg       0.66      0.63      0.64        20
weighted avg       0.73      0.75      0.74        20

Confusion Matrix:
 [[ 2  3]
 [ 2 13]]

Decision Tree Performance:
Accuracy: 0.75
Classification Report:
              precision    recall  f1-score   support

           0       0.50      0.40      0.44         5
           1       0.81      0.87      0.84        15

    accuracy                           0.75        20
   macro avg       0.66      0.63      0.64        20
weighted avg       0.73      0.75      0.74        20

Confusion Matrix:
 [[ 2  3]
 [ 2 13]]

Random Forest Performance:
Accuracy: 0.8
Classification Report:
              precision    recall  f1-score   support

           0       1.00      0.20      0.33         5
           1       0.79      1.00      0.88        15

    accuracy                           0.80        20
   macro avg       0.89      0.60      0.61        20
weighted avg       0.84      0.80      0.75        20

Confusion Matrix:
 [[ 1  4]
 [ 0 15]]
```

## 5. Hyperparameter tuning

## Subtask:

Optimize the hyperparameters of the chosen models to improve their performance.

```
# Define parameter grids for each model
lr_param_grid = {'C': [0.001, 0.01, 0.1, 1, 10, 100], 'penalty': ['l2']} # Using l2 penalty as it's default and common

dt_param_grid = {'max_depth': [None, 5, 10, 15, 20],
                 'min_samples_split': [2, 5, 10],
                 'min_samples_leaf': [1, 2, 4]}

rf_param_grid = {'n_estimators': [50, 100, 200],
                 'max_depth': [None, 5, 10, 15],
                 'min_samples_split': [2, 5, 10],
                 'min_samples_leaf': [1, 2, 4]}

# Instantiate GridSearchCV for each model
lr_grid_search = GridSearchCV(LogisticRegression(max_iter=1000), lr_param_grid, cv=5, scoring='accuracy')
dt_grid_search = GridSearchCV(DecisionTreeClassifier(), dt_param_grid, cv=5, scoring='accuracy')
rf_grid_search = GridSearchCV(RandomForestClassifier(), rf_param_grid, cv=5, scoring='accuracy')

# Fit GridSearchCV to the training data
print("Performing Grid Search for Logistic Regression...")
lr_grid_search.fit(X_train, y_train)
print("Performing Grid Search for Decision Tree...")
dt_grid_search.fit(X_train, y_train)
print("Performing Grid Search for Random Forest...")
rf_grid_search.fit(X_train, y_train)

# Print the best hyperparameters
print("\nBest hyperparameters for Logistic Regression:", lr_grid_search.best_params_)
print("Best hyperparameters for Decision Tree:", dt_grid_search.best_params_)
print("Best hyperparameters for Random Forest:", rf_grid_search.best_params_)
```

```
Performing Grid Search for Logistic Regression...
/usr/local/lib/python3.12/dist-packages/sklearn/linear_model/_logistic.py:465: ConvergenceWarning: lbfgs failed to conv
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
/usr/local/lib/python3.12/dist-packages/sklearn/linear_model/_logistic.py:465: ConvergenceWarning: lbfgs failed to conv
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
/usr/local/lib/python3.12/dist-packages/sklearn/linear_model/_logistic.py:465: ConvergenceWarning: lbfgs failed to conv
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
/usr/local/lib/python3.12/dist-packages/sklearn/linear_model/_logistic.py:465: ConvergenceWarning: lbfgs failed to conv
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
/usr/local/lib/python3.12/dist-packages/sklearn/linear_model/_logistic.py:465: ConvergenceWarning: lbfgs failed to conv
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
/usr/local/lib/python3.12/dist-packages/sklearn/linear_model/_logistic.py:465: ConvergenceWarning: lbfgs failed to conv
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
/usr/local/lib/python3.12/dist-packages/sklearn/linear_model/_logistic.py:465: ConvergenceWarning: lbfgs failed to conv
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
/usr/local/lib/python3.12/dist-packages/sklearn/linear_model/_logistic.py:465: ConvergenceWarning: lbfgs failed to conv
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
  n_iter_i = _check_optimize_result(
```

## ˅ 6. Model Evaluation with Tuned Hyperparameters

### Subtask:

Evaluate the models with the best hyperparameters found by GridSearchCV on the test set and compare their performance.

```python
# Get the best models from GridSearchCV
best_lr_model = lr_grid_search.best_estimator_
best_dt_model = dt_grid_search.best_estimator_
best_rf_model = rf_grid_search.best_estimator_

# Make predictions with the best models on the test set
best_lr_predictions = best_lr_model.predict(X_test)
best_dt_predictions = best_dt_model.predict(X_test)
best_rf_predictions = best_rf_model.predict(X_test)

# Evaluate the best Logistic Regression model
print("Best Logistic Regression Performance:")
print("Accuracy:", accuracy_score(y_test, best_lr_predictions))
print("Classification Report:\n", classification_report(y_test, best_lr_predictions))
print("Confusion Matrix:\n", confusion_matrix(y_test, best_lr_predictions))

# Evaluate the best Decision Tree model
print("\nBest Decision Tree Performance:")
print("Accuracy:", accuracy_score(y_test, best_dt_predictions))
print("Classification Report:\n", classification_report(y_test, best_dt_predictions))
print("Confusion Matrix:\n", confusion_matrix(y_test, best_dt_predictions))

# Evaluate the best Random Forest model
print("\nBest Random Forest Performance:")
print("Accuracy:", accuracy_score(y_test, best_rf_predictions))
print("Classification Report:\n", classification_report(y_test, best_rf_predictions))
print("Confusion Matrix:\n", confusion_matrix(y_test, best_rf_predictions))
```

```
Best Logistic Regression Performance:
Accuracy: 0.8
Classification Report:
               precision    recall  f1-score   support

           0       1.00      0.20      0.33         5
           1       0.79      1.00      0.88        15

    accuracy                           0.80        20
   macro avg       0.89      0.60      0.61        20
weighted avg       0.84      0.80      0.75        20


Confusion Matrix:
 [[ 1  4]
 [ 0 15]]

Best Decision Tree Performance:
Accuracy: 0.6
Classification Report:
               precision    recall  f1-score   support

           0       0.33      0.60      0.43         5
           1       0.82      0.60      0.69        15

    accuracy                           0.60        20
   macro avg       0.58      0.60      0.56        20
weighted avg       0.70      0.60      0.63        20


Confusion Matrix:
 [[3 2]
 [6 9]]

Best Random Forest Performance:
Accuracy: 0.8
Classification Report:
               precision    recall  f1-score   support

           0       1.00      0.20      0.33         5
           1       0.79      1.00      0.88        15

    accuracy                           0.80        20
   macro avg       0.89      0.60      0.61        20
weighted avg       0.84      0.80      0.75        20


Confusion Matrix:
 [[ 1  4]
 [ 0 15]]
```

## ⌄ 7. Project documentation

### Subtask:

Document the entire project process, including data preparation steps, model selection and training details, evaluation results, and conclusions. Align this documentation with the weekly activity logs and sections provided in your document.

**Reasoning**: I need to document the project process based on the previous steps and the assumed structure from the PDF. I will create a markdown report summarizing the project, data preparation, model training, tuning, and evaluation.

```
# Insert the actual evaluation results into the project report string
project_report_with_results = """
# Project Report: Loan Approval Prediction System

## 1. Introduction/Project Overview

This project aims to develop a system to predict loan approval status based on various applicant and loan characteristics. T

## 2. Data Preparation

The dataset used for this project contains information about loan applicants, including demographic details, income, loan am

**Steps Performed:**
*   **Loading Data:** The dataset was loaded into a pandas DataFrame.
*   **Handling Missing Values:** Missing values in categorical columns (Gender, Married, Dependents, Self_Employed, Loan_Amo
*   **Encoding Categorical Variables:** Categorical features were converted into numerical format using one-hot encoding. Th
*   **Encoding Target Variable:** The target variable 'Loan_Status' ('Y'/'N') was encoded into numerical labels (1/0) using
*   **Splitting Data:** The dataset was split into training (80%) and testing (20%) sets using `train_test_split`, with stra

## 3. Model Selection and Training

Based on the nature of the binary classification problem (predicting loan approval), the following machine learning models w

*   **Logistic Regression:** A linear model suitable for binary classification.
*   **Decision Tree:** A non-linear model that partitions the data based on features.
*   **Random Forest:** An ensemble method that combines multiple decision trees to improve robustness and accuracy.

These models were instantiated and trained on the preprocessed training data (`X_train`, `y_train`).

## 4. Hyperparameter Tuning

To optimize the performance of the selected models, hyperparameter tuning was performed using `GridSearchCV`. The following

*   **Logistic Regression:**
    *   `C`: [0.001, 0.01, 0.1, 1, 10, 100]
    *   `penalty`: ['l2']
*   **Decision Tree:**
    *   `max_depth`: [None, 5, 10, 15, 20]
    *   `min_samples_split`: [2, 5, 10]
    *   `min_samples_leaf`: [1, 2, 4]
*   **Random Forest:**
    *   `n_estimators`: [50, 100, 200]
    *   `max_depth`: [None, 5, 10, 15]
    *   `min_samples_split`: [2, 5, 10]
    *   `min_samples_leaf`: [1, 2, 4]

The tuning process aimed to find the combination of hyperparameters that yielded the best accuracy on the training data thro

**Best Hyperparameters Found:**
*   **Logistic Regression:** C=0.001, penalty='l2'
*   **Decision Tree:** max_depth=10, min_samples_leaf=4, min_samples_split=10
*   **Random Forest:** n_estimators=50, max_depth=None, min_samples_leaf=1, min_samples_split=10

## 5. Model Evaluation

The performance of the trained models, both before and after hyperparameter tuning, was evaluated on the test set (`X_test`,

*   **Accuracy:** The proportion of correctly classified instances.
*   **Precision:** The ability of the model to correctly identify positive instances (approved loans).
*   **Recall:** The ability of the model to find all positive instances.
*   **F1-score:** The harmonic mean of precision and recall.
*   **Confusion Matrix:** A table summarizing the performance of the classification model.

**Evaluation Results (After Hyperparameter Tuning):**

**Best Logistic Regression Performance:**
Accuracy: 0.8
Classification Report:
              precision    recall  f1-score   support
```

```
            0       1.00      0.20      0.33         5
            1       0.79      1.00      0.88        15

     accuracy                           0.80        20
    macro avg       0.89      0.60      0.61        20
 weighted avg       0.84      0.80      0.75        20
```

```
Confusion Matrix:
 [[ 1  4]
 [ 0 15]]
```

**Best Decision Tree Performance:**
Accuracy: 0.6
Classification Report:
```
              precision    recall  f1-score   support

            0       0.33      0.60      0.43         5
            1       0.82      0.60      0.69        15

     accuracy                           0.60        20
    macro avg       0.58      0.60      0.56        20
 weighted avg       0.70      0.60      0.63        20
```

```
Confusion Matrix:
 [[3 2]
 [6 9]]
```

**Best Random Forest Performance:**
Accuracy: 0.8
Classification Report:
```
              precision    recall  f1-score   support

            0       1.00      0.20      0.33         5
            1       0.79      1.00      0.88        15

     accuracy                           0.80        20
    macro avg       0.89      0.60      0.61        20
 weighted avg       0.84      0.80      0.75        20
```

```
Confusion Matrix:
 [[ 1  4]
 [ 0 15]]
```

**Comparison:**

Both the optimized Logistic Regression and Random Forest models achieved the highest accuracy of 0.80 on the test set. The D

## 6. Conclusion

This project successfully built and evaluated machine learning models for loan approval prediction. The data was preprocesse

**Next Steps:**
*   Address the class imbalance issue using techniques such as oversampling the minority class, undersampling the majority c
*   Investigate the convergence warning for Logistic Regression and consider data scaling or increasing `max_iter`.
*   Explore other models or techniques that might be more robust to imbalanced data (e.g., Gradient Boosting, SVM with class
*   Consider feature engineering to potentially improve model performance.
*   Deploy the best performing model, taking into account the need to balance predicting both approved and non-approved loar
"""

print(project_report_with_results)

```
    # Project Report: Loan Approval Prediction System

    ## 1. Introduction/Project Overview

    This project aims to develop a system to predict loan approval status based on various applicant and loan characteristi

    ## 2. Data Preparation

    The dataset used for this project contains information about loan applicants, including demographic details, income, lo

    **Steps Performed:**
    *   **Loading Data:** The dataset was loaded into a pandas DataFrame.
    *   **Handling Missing Values:** Missing values in categorical columns (Gender, Married, Dependents, Self_Employed, Loa
    *   **Encoding Categorical Variables:** Categorical features were converted into numerical format using one-hot encodir
    *   **Encoding Target Variable:** The target variable 'Loan_Status' ('Y'/'N') was encoded into numerical labels (1/0) u
    *   **Splitting Data:** The dataset was split into training (80%) and testing (20%) sets using `train_test_split`, witr

    ## 3. Model Selection and Training

    Based on the nature of the binary classification problem (predicting loan approval), the following machine learning moc
```

```
*    **Logistic Regression:** A linear model suitable for binary classification.
*    **Decision Tree:** A non-linear model that partitions the data based on features.
*    **Random Forest:** An ensemble method that combines multiple decision trees to improve robustness and accuracy.

These models were instantiated and trained on the preprocessed training data (`X_train`, `y_train`).

## 4. Hyperparameter Tuning

To optimize the performance of the selected models, hyperparameter tuning was performed using `GridSearchCV`. The follo

*    **Logistic Regression:**
   *    `C`: [0.001, 0.01, 0.1, 1, 10, 100]
   *    `penalty`: ['l2']
*    **Decision Tree:**
   *    `max_depth`: [None, 5, 10, 15, 20]
   *    `min_samples_split`: [2, 5, 10]
   *    `min_samples_leaf`: [1, 2, 4]
*    **Random Forest:**
   *    `n_estimators`: [50, 100, 200]
   *    `max_depth`: [None, 5, 10, 15]
   *    `min_samples_split`: [2, 5, 10]
   *    `min_samples_leaf`: [1, 2, 4]

The tuning process aimed to find the combination of hyperparameters that yielded the best accuracy on the training data

**Best Hyperparameters Found:**
*    **Logistic Regression:** C=0.001, penalty='l2'
*    **Decision Tree:** max_depth=10, min_samples_leaf=4, min_samples_split=10
*    **Random Forest:** n_estimators=50, max_depth=None, min_samples_leaf=1, min_samples_split=10

## 5. Model Evaluation

The performance of the trained models, both before and after hyperparameter tuning, was evaluated on the test set (`X_t

*    **Accuracy:** The proportion of correctly classified instances.
```

# 8. Presentation and reporting

## Subtask:

Prepare a presentation or report summarizing your project, findings, and conclusions.

## Summary:

### Data Analysis Key Findings

- The initial attempt to load the dataset failed, and a placeholder dataset was used for the analysis.
- Missing values were handled by imputing the mode for categorical features and the median for numerical features.
- Categorical features were one-hot encoded, and the target variable was label encoded.
- The data was split into training (80%) and testing (20%) sets with stratification.
- Logistic Regression, Decision Tree, and Random Forest models were trained on the preprocessed data.
- Hyperparameter tuning using GridSearchCV was performed for all three models, identifying the best parameters.
- After tuning, both Logistic Regression and Random Forest achieved an accuracy of 0.80 on the test set, while the Decision Tree had an accuracy of 0.60.
- A significant challenge observed was the difficulty of Logistic Regression and Random Forest in predicting the minority class (non-approved loans), resulting in lower precision and f1-score for this class. The Decision Tree performed slightly better in predicting the minority class.

### Insights or Next Steps

- Address the class imbalance issue using techniques like oversampling, undersampling, or using evaluation metrics more suitable for imbalanced datasets (e.g., ROC AUC) to improve the prediction of the minority class.
- Investigate the convergence warning for Logistic Regression by considering data scaling or increasing the maximum number of iterations.
- Explore other models or techniques that might be more robust to imbalanced data (e.g., Gradient Boosting, SVM with class weights).
- Consider feature engineering to potentially improve model performance.
- Deploy the best performing model, taking into account the need to balance predicting both approved and non-approved loans effectively.