

# Recsify Technologies Assignment

## Machine Learning

**Perform all the various steps of machine learning like data exploration, feature engineering and model building.**

Step-by-Step Explanation:

### Step 1: Importing Libraries

We use libraries such as pandas for data manipulation, matplotlib and seaborn for visualization, and scikit-learn for machine learning and evaluation metrics. These libraries are essential for data preprocessing, model training, and evaluation.

### Step 2: Load and Preprocess the Data

We load the dataset using pandas. Preprocessing involves handling missing values and encoding categorical variables:

- **\*Handling Missing Values\***: We fill missing values for numerical columns with the median and for categorical columns with the mode. This ensures that the dataset is complete and can be used for model training without errors due to missing values.
- **\*Encoding Categorical Variables\***: Categorical variables are converted to numerical values using LabelEncoder. This is necessary because machine learning models typically require numerical input.

### Step 3: Define Target Variable and Features

We define the target variable, which is the column we want to predict (Loan\_Status in this case). Features are the columns used to make the prediction. We exclude the Loan\_ID column as it is an identifier and not useful for prediction.

### Step 4: Split the Data

We split the dataset into training and testing sets using an 80-20 split. The training set is used to train the model, and the testing set is used to evaluate its performance.

#### Step 5: Train the RandomForest Model

We use the RandomForest algorithm, a robust and widely used machine learning model, to train on the training data. The model learns patterns and relationships in the data to make predictions.

#### Step 6: Evaluate Model Performance

Model performance is evaluated using several metrics:

- Confusion Matrix\*\*: This matrix shows the number of true positive, true negative, false positive, and false negative predictions. It helps understand how well the model is distinguishing between the different classes.
- Classification Report\*\*: This report provides precision, recall, F1-score, and support for each class. These metrics give a detailed insight into the model's performance for each class.
  - \*Precision\*: The proportion of true positive predictions among all positive predictions.
  - \*Recall\*: The proportion of true positive predictions among all actual positives.
  - \*F1-Score\*: The harmonic mean of precision and recall, providing a single metric that balances both concerns.
  - \*Support\*: The number of actual occurrences of each class in the dataset.
- \*ROC AUC Score\*: The Area Under the Receiver Operating Characteristic curve. It measures the ability of the model to distinguish between the classes. A higher score indicates better performance.

#### Step 7: Feature Importance

We determine the importance of each feature in making predictions using the feature importance attribute of the RandomForest model. Features with higher importance scores have a more significant impact on the model's predictions.

## Step 8: Generate PDF Report

We compile the confusion matrix, classification report, and feature importance into a PDF report. This report provides a comprehensive summary of the model's evaluation, making it easy to share and review.

## Results Interpretation

### Confusion Matrix

The confusion matrix indicates the accuracy of the model's predictions. High values along the diagonal represent a high number of correct predictions, while high off-diagonal values indicate many errors.

### Classification Report

The classification report breaks down the model's performance for each class:

- **\*Precision\*** indicates how many of the predicted positive instances are actually positive.
- **\*Recall\*** measures how many actual positive instances are correctly predicted by the model.
- **\*F1-Score\*** provides a balance between precision and recall.
- **\*Support\*** shows the number of instances for each class in the test set.

### ROC AUC Score

The ROC AUC score is a single number summary of the model's ability to distinguish between classes. A score of 1.0 represents perfect discrimination, while a score of 0.5 suggests no better than random guessing.

## Feature Importance

The feature importance plot identifies which features contribute most to the model's decisions. This can provide insights into which factors are most influential in determining loan status, helping to understand the model's behavior and potentially inform feature selection in future models.

## Conclusion

By following this process, we can develop a machine learning model to predict loan risk, evaluate its performance using various metrics, and identify the key features influencing its predictions. The generated PDF report serves as a comprehensive summary of the analysis and results, aiding in decision-making and communication with stakeholders.

## **Submitted By:**

**Hanuman Chandra Shekar Reddy.**