

a) Dimension or Size of the Dataset:

1. **Number of Rows:**

This can be extracted using `df.shape[0]`. Since the dataset is from the linked source, it contains **500 rows** (based on the dataset description).

2. **Number of Columns:**

This can be extracted using `df.shape[1]`. The dataset has **14 columns**.

3. **Missing Values:**

- **Total Missing Values:** Computed using `df.isnull().sum().sum()`.
- **Rows with Missing Values:** Computed using `df.isnull().any(axis=1).sum()`.
- **Columns with Missing Values:** Computed using `df.isnull().any(axis=0).sum()`.
From the Kaggle dataset, specific missing values details can be verified by running the code.

4. **Target Variable:**

The target variable is **loan_status**, which indicates whether a loan was approved or not.

b) Methods Applied in the Code:

1. **Preprocessing Methods:**

- **Column Transformer:**
Applied to handle numerical and categorical features separately.
 - **Numerical Features:**
 - Standardized using `StandardScaler`.
 - **Categorical Features:**
 - One-hot encoded using `OneHotEncoder(handle_unknown='ignore')`.
- **Imputation:** (although it isn't explicitly shown in the provided code, this is commonly included in `ColumnTransformer` pipelines).
- **Feature Selection:** Dropped the target variable (`loan_status`) from the feature set.

2. **Machine Learning Methods:**

- **Random Forest Classifier:**
 - Used as the first model to predict loan status.
- **Decision Tree Classifier:**
 - Applied as a simpler tree-based model.
- **Logistic Regression:**
 - Used for linear classification.
- **Gradient Boosting Classifier:**

- Applied for boosting-based predictions.
- **Model Evaluation:**
 - Accuracy: Computed using `accuracy_score`.
 - Confusion Matrix: Visualized using `sns.heatmap`.
 - Classification Report: Includes precision, recall, and F1-score.

Additional Details from the Code:

1. Dataset Loading:

The dataset is read from the path `/kaggle/input/loan-approval-prediction-dataset/loan_approval_dataset.csv` using `pd.read_csv`.

2. Exploratory Data Analysis (EDA):

- The dataset is displayed using `df` and `df.head()` for a quick look.
- Summary statistics are reviewed using `df.describe()`.
- Data types and missing value checks:
 - `df.info()` gives column types and non-null counts.
 - `df.isnull().sum()` checks missing values for each column.

3. Feature Engineering:

- **Feature Cleanup:** The column names are stripped of extra whitespace using `df.columns.str.strip()`.
- **Target-Feature Splitting:**
 - Features (X) are selected by dropping the `loan_status` column.
 - Target variable (y) is extracted as `df['loan_status']`.

4. Train-Test Split:

- Data is split into training and testing sets using `train_test_split`, with 80% training data and 20% testing data.

5. Evaluation Metrics:

- Models are evaluated using:
 - **Accuracy Score** (`accuracy_score`).
 - **Confusion Matrix** (`confusion_matrix` visualized with a heatmap).
 - **Classification Report** (includes precision, recall, and F1-score).

6. Visualization:

- Heatmap visualization of the confusion matrix using `seaborn.heatmap`.

7. Classifier Pipelines:

- A unified pipeline approach is used for all classifiers, combining preprocessing (ColumnTransformer) and classification models (RandomForestClassifier, DecisionTreeClassifier, etc.).