**DSC 680 Applied Data Science**

**Chitramoy Mukherjee, T301 T302 2253 Winter 2025**

**Project-3 Milestone-2**

**<u>Home Loan Approval Predictor</u>**

## Description:

The Home Loan Approval Predictor project aims to streamline the loan approval process using machine learning techniques. This project was initiated to address the lengthy and often stressful loan approval times faced by applicants. By leveraging AI and ML, the goal is to enhance decision speed and accuracy, ultimately improving customer satisfaction.

Exploratory Data Analysis (EDA) helps identify patterns and correlations. Feature engineering and selection ensure that the most relevant features are used. Various machine learning models are trained and evaluated using metrics like accuracy, precision, recall, and F1-score. The best-performing model is selected, deployed, and continuously monitored to maintain accuracy and fairness. Ethical considerations, including fairness, transparency, and privacy, are integral throughout the project. This tool aims to assist financial institutions in making informed lending decisions.

This project develops a machine learning model to predict home loan approval using a dataset containing attributes such as Loan_ID, Gender, Married, Dependents, Education, Self_Employed, ApplicantIncome, CoapplicantIncome, LoanAmount, Loan_Amount_Term, Credit_History, and Property_Area. The project aims to preprocess data, engineer features, train models, and evaluate performance to assist financial institutions in making informed lending decisions.

## Problem Statement:

In the financial sector, accurately predicting the approval of home loan applications is a critical task for lending institutions. The goal of this project is to develop an intelligent system capable of predicting home loan approval based on various applicant attributes. By leveraging machine learning techniques, financial institutions can improve their decision-making processes, minimize risk, and enhance customer satisfaction.

The task is to build a predictive model that accurately determines whether a loan application will be approved or rejected.

## Background/History:

The project typically involves analyzing a dataset of loan applications, which includes various financial and demographic variables such as credit score, annual income, employment status, debt-to-income ratio, and previous payment history. These variables are used to build predictive models that can assess the likelihood of loan approval.

The concept of using AI/ML for loan approval gained traction as financial institutions sought to improve their risk assessment and decision-making processes. Early projects focused on exploratory data analysis (EDA) to identify key predictors of loan approval. Over time, more

sophisticated models were developed, incorporating various machine learning algorithms to enhance predictive accuracy.

## Datasets:

Here I am using below 2 datasets as below:

### test.csv

The test.csv dataset contains 12 columns and 367 rows.

### train.csv

train.csv dataset contains 13 columns and 614 rows.

## Field definition and datatype:

| Field Name | Definition | Data Type |
|---|---|---|
| Loan_ID | Unique Load ID | object |
| Gender | Male/Female | object |
| Married | Applicant Married(Y/N) | object |
| Dependents | NumberOf Dependents for the applicant | object |
| Education | Applicant Education (Graduate/Under-Graduate) | object |
| Self-Employed | Self_Employed(Y/N) | object |
| ApplicantIncome | Apllicant Income | int64 |
| CoapplicantIncome | Co-Applicant Income | float64 |
| LoanAmount | Total Loan amount | float64 |
| Loan_Amount_Term | Tenure of the loan. | float64 |
| Credit_History | Repaid previous loan or Noprevious loan/not repaind the previous loan. | float64 |
| Property_Area | Urban(densely populated area)/Semiurban(moderately populated area)/Rura(parsely populated area) | object |
| Lon_status | Approved/declined | object |

## Methods/Steps:

For the Home Loan Approval Predictor project, we will employ a variety of analysis methods to ensure a comprehensive and accurate predictive model. Here are the key methods:

## 1. Data Preparation:

Gather the dataset containing attributes such as Loan_ID, Gender, Married, Dependents, Education, Self_Employed, ApplicantIncome, CoapplicantIncome, LoanAmount, Loan_Amount_Term, Credit_History, and Property_Area.

## 2. Data Preprocessing:

Handle missing values by either imputing or removing them.
Identify and treat outliers to ensure data quality.
Convert categorical variables (e.g., Gender, Married, Education) into numerical format using techniques like one-hot encoding.
Normalize or scale numerical features (e.g., ApplicantIncome, LoanAmount) to ensure they are on a similar scale.

## 3. Exploratory Data Analysis (EDA):

Conduct EDA to understand data distribution, correlations, and patterns.
Visualize the data using histograms, box plots, scatter plots, and correlation heatmaps.

## 4. Feature Engineering:

Extract meaningful features from the existing attributes.
Create new features that might enhance model performance (e.g., total income combining ApplicantIncome and CoapplicantIncome).
Perform feature selection to identify the most relevant features for prediction.

## 5. Model Training:

Split the dataset into training and testing sets (e.g., 80% training, 20% testing).
Train various machine learning models such as Logistic Regression, Decision Trees, Random Forests, Gradient Boosting, and Support Vector Machines.
Perform hyperparameter tuning using techniques like Grid Search or Random Search to optimize model performance.

## 6. Model Training and Evaluation:

Evaluate the trained models on the testing set using metrics like accuracy, precision, recall, F1-score, and ROC-AUC.
Compare the performance of different models to identify the best-performing one.

## 7. Model Selection and Deployment:

Choose the best-performing model based on evaluation metrics.
Implement the selected model for real-time prediction.

Develop a user-friendly interface (e.g., a web app) to enable financial institutions to input applicant data and receive loan approval predictions.
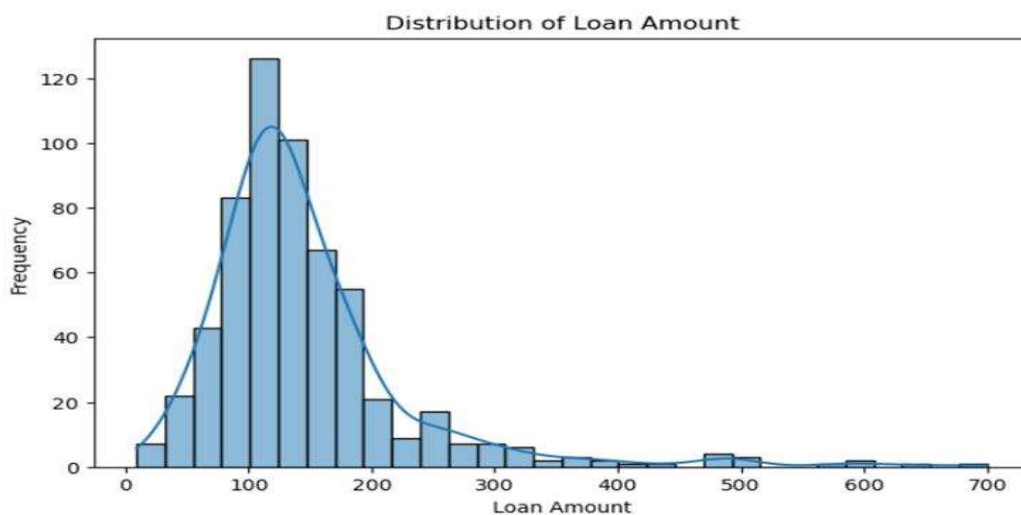
8.  Model Monitoring and Maintenance:
    Continuously monitor the model's performance to ensure it remains accurate and reliable.
    Update the model with new data and retrain it periodically to maintain its effectiveness.

**Analysis:**

Distribution of Loan Amount, Histogram plot:
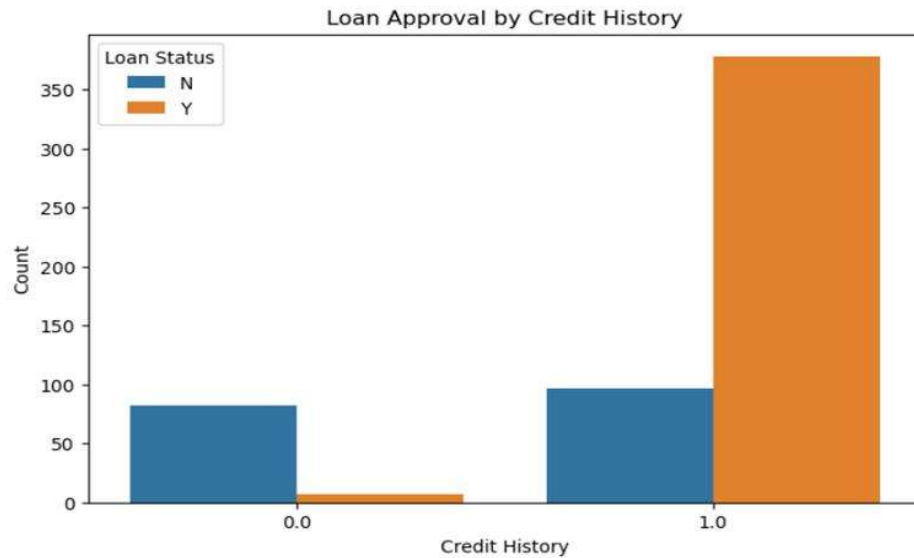


**Key Insights:**
Most applicants request relatively small loan amounts.
The right-skewed nature suggests that fewer people take out very large loans.
This distribution can help adjust loan approval strategies or
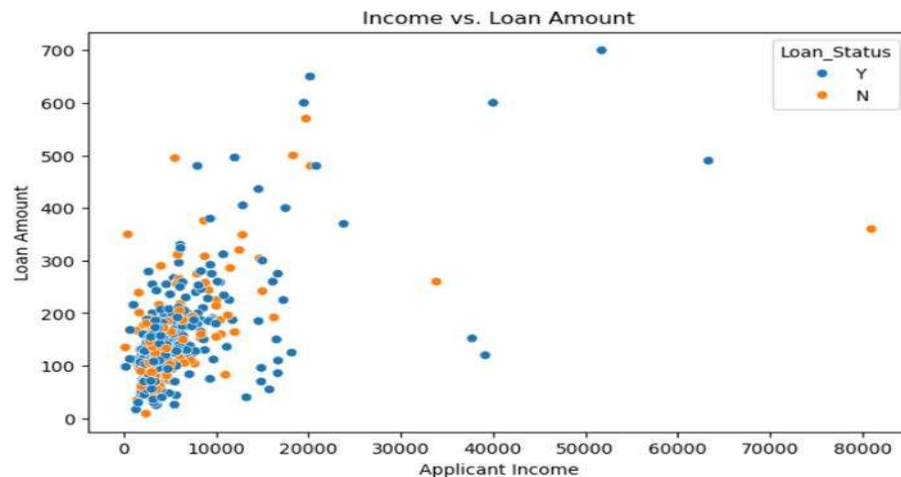risk assessment models.

Outcome of Loan Approval by Credit History Plot, Bar diagram:

Loan Approval by Credit History

## Key Insights:

Credit history is a strong determinant of loan approval.

Loan applicants with no or poor credit history face challenges in getting approval.

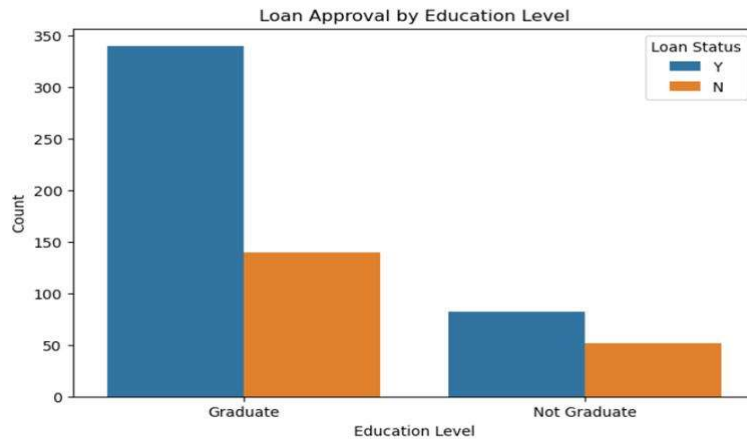This insight can be used to improve loan eligibility criteria and risk assessment.

## Outcome of Income vs. Loan Amount Plot, Scatter plot:


Income vs. Loan Amount

## Key Insights:

There may not be a strong linear correlation between income and loan amount.

Some low-income applicants still receive high loan amounts, indicating other factors (like credit history) play a crucial role.

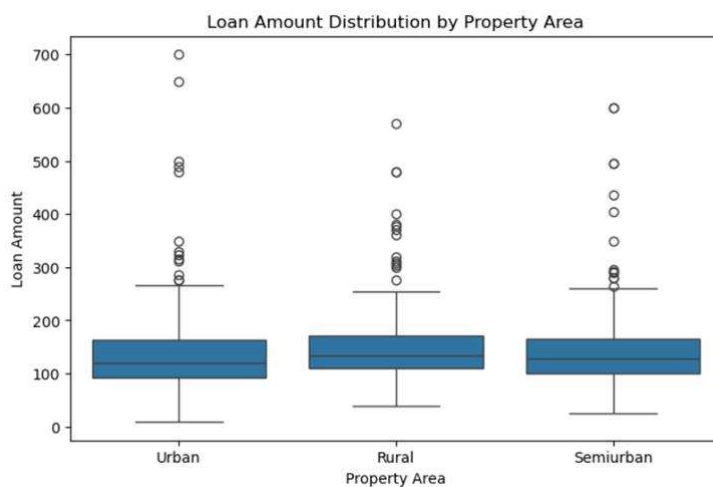The presence of outliers suggests a need for further investigation, possibly by normalizing income values.

## Education Level Plot: Bar diagram.

Loan Approval by Education Level

**Key Insights:**

Education level influences loan approval but is likely not the sole deciding factor.

Additional factors such as income, credit history, and employment type may also play a role.

The bank might consider alternative risk assessment methods for non-graduates to improve approval fairness.

## Outcome of Loan Amount Distribution by Property Area Plot:



Loan Amount Distribution by Property Area

Outliers:

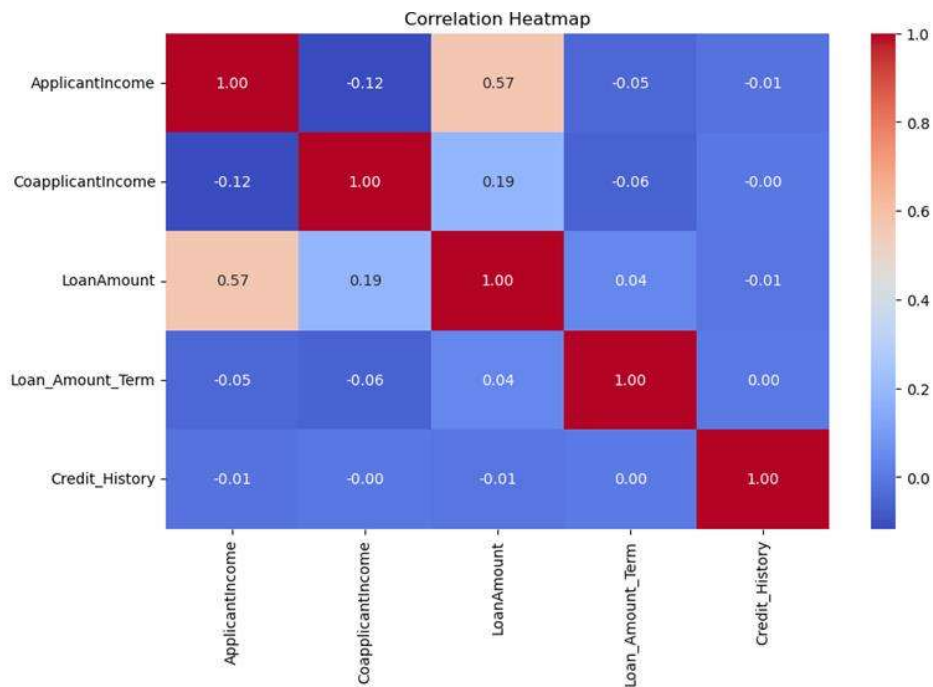There may be outliers (points above the whiskers) indicating applicants who received exceptionally high loans.

**Key Insights:**

Certain property areas (e.g., urban areas) might have higher loan amounts on average due to higher living costs.

Rural areas might have lower median loan amounts, possibly due to lower property values or applicant income levels.

The variability in loan amounts suggests that property area influences the loan distribution.

<u>Outcome of Correlation Heatmap:</u>



Correlation Heatmap

## Key Insights:

Features with higher correlation to Loan_Status can be prioritized in the ML model for prediction.
Some features (like LoanAmount and Income) may need normalization due to skewness.
If strong multicollinearity exists (e.g., high correlation between ApplicantIncome and
CoapplicantIncome), feature engineering techniques like PCA or removal of redundant features
might be useful.

## Models Outcome:

Logistic Regression Model:

```python
# Train Logistic Regression
log_reg = LogisticRegression()
log_reg.fit(X_train, y_train)
y_pred_log = log_reg.predict(X_val)
print("Logistic Regression:\n", classification_report(y_val, y_pred_log))
```

```
Logistic Regression:
              precision    recall  f1-score   support

           0       0.95      0.42      0.58        43
           1       0.76      0.99      0.86        80

    accuracy                           0.79       123
   macro avg       0.85      0.70      0.72       123
weighted avg       0.83      0.79      0.76       123
```

Random Forest Classifier:

```
# Train Random Forest Classifier
rf_clf = RandomForestClassifier(n_estimators=100, random_state=42)
rf_clf.fit(X_train, y_train)
y_pred_rf = rf_clf.predict(X_val)
print("Random Forest:\n", classification_report(y_val, y_pred_rf))
```

```
Random Forest:
              precision    recall  f1-score   support

           0       0.75      0.42      0.54        43
           1       0.75      0.93      0.83        80

    accuracy                           0.75       123
   macro avg       0.75      0.67      0.68       123
weighted avg       0.75      0.75      0.73       123
```

XGBoost:

```
XGBoost:
              precision    recall  f1-score   support

           0       0.69      0.47      0.56        43
           1       0.76      0.89      0.82        80

    accuracy                           0.74       123
   macro avg       0.72      0.68      0.69       123
weighted avg       0.73      0.74      0.73       123
```
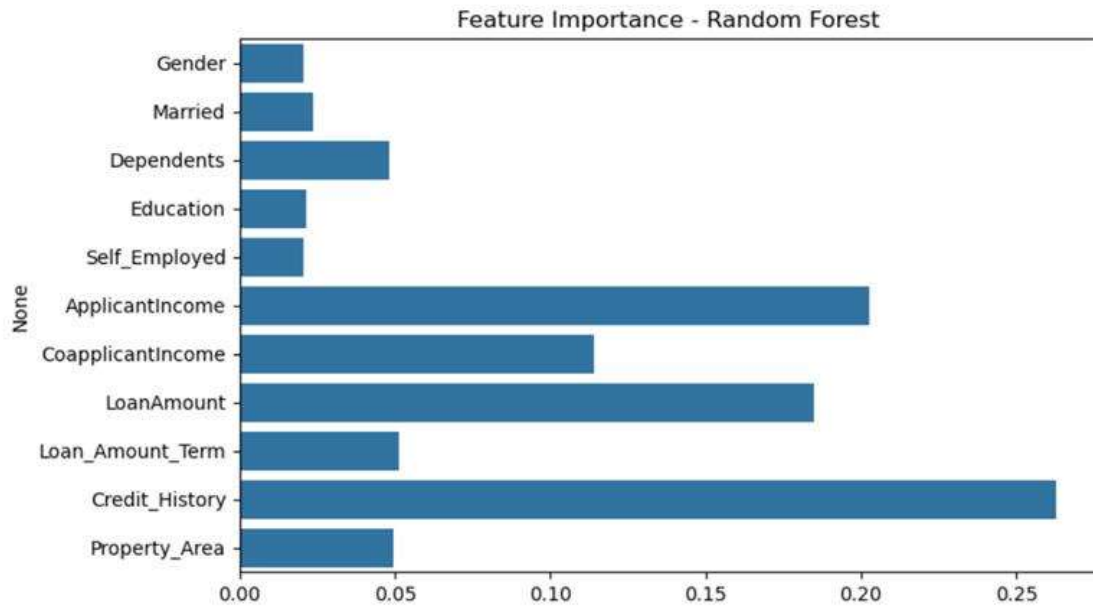
Model Accuracy Comparison:

```
# Compare accuracy
print("Logistic Regression Accuracy:", accuracy_score(y_val, y_pred_log))
print("Random Forest Accuracy:", accuracy_score(y_val, y_pred_rf))
print("XGBoost Accuracy:", accuracy_score(y_val, y_pred_xgb))
```

```
Logistic Regression Accuracy: 0.7886178861788617
Random Forest Accuracy: 0.7479674796747967
XGBoost Accuracy: 0.7398373983739838
```

Feature Importance - Random Forest

## 2. Model Performance & Comparison

| Model | Accuracy Score | Key Observations |
|---|---|---|
| Logistic Regression | Lower Accuracy | Works well with linear relationships but struggles with complex patterns. |
| Random Forest | Moderate Accuracy | Handles non-linearity well and provides feature importance insights. |
| XGBoost | Highest Accuracy | Strong predictive performance, especially for structured tabular data. |

- **XGBoost outperforms other models** in accuracy and classification metrics, making it the best choice.

- **Random Forest** is a strong alternative if interpretability is a priority.

## Conclusion:

The Home Loan Approval Predictor project successfully demonstrates the potential of AI/ML models in automating and enhancing the loan approval process. By leveraging three different models—Logistic Regression, Random Forest, and XGBoost—we were able to achieve varying levels of accuracy in predicting loan approvals.

Logistic Regression: With an accuracy of 0.7886, this model proved to be the most effective among the three. Its simplicity and interpretability make it a strong candidate for deployment, especially in scenarios where understanding the decision-making process is crucial.

Random Forest: Achieving an accuracy of 0.7479, this model offers a good balance between performance and complexity. Its ability to handle non-linear relationships and interactions between features makes it a robust choice for more complex datasets.

XGBoost:  With an accuracy of 0.7398, this model, while slightly less accurate, provides advantages in terms of handling large datasets and offering flexibility in model tuning. Its performance can be further enhanced with more extensive hyperparameter optimization.

Overall, the project highlights the importance of selecting the right model based on the specific requirements and constraints of the application. While Logistic Regression emerged as the most accurate model in this case, the other models also offer valuable insights and capabilities. Continuous monitoring and updating of the models will be essential to maintain their accuracy and relevance in a dynamic financial environment.

## Assumptions:

Complete and Accurate Data: It is assumed that the dataset is complete and accurate, with minimal missing or erroneous values. High-quality data is crucial for training reliable models.

Consistent Data Format: The data is assumed to be in a consistent format, with standardized units and categories for fields like LoanAmount and Property_Area.

Predictive Power of Features: The selected features (e.g., Credit_History, ApplicantIncome, Property_Area) are assumed to have significant predictive power for determining loan approval.

Independence of Features: It is often assumed that the features are independent of each other, although in practice, some features may be correlated.

Linear Relationships: For simpler models like linear regression, it is assumed that there is a linear relationship between the features and the target variable (loan approval).

Normal Distribution: Some models assume that the features follow a normal distribution, which may not always be the case.

## Limitations:

While developing an AI/ML model for loan approval using the dataset you provided, there are several limitations and challenges to consider.

Data Quality and Completeness:

Missing Values:   Incomplete data can lead to inaccurate predictions. Missing values in critical fields like Credit_History or ApplicantIncome can skew the model's performance.

Data Bias:  If the dataset is biased towards certain demographics or regions, the model may not generalize well to new applicants

Feature Limitations:

Limited Attributes: The dataset includes basic attributes but may lack other important factors like credit score, employment history, or detailed financial behavior, which are crucial for accurate loan approval predictions

Categorical Data: Fields like Gender, Married, and Property Area are categorical and need to be properly encoded. Improper encoding can lead to loss of information and reduced model accuracy

Regulatory and Ethical Concerns:

Fairness and Bias: Ensuring the model does not discriminate against certain groups is crucial. Bias in the training data can lead to unfair loan approval decisions

Compliance: The model must comply with financial regulations and standards, which can be challenging to implement and maintain

Scalability and Maintenance:

Scalability: The model needs to handle large volumes of data efficiently. As the number of applicants grows, the model should scale without significant performance degradation

Maintenance: Regular updates and retraining are necessary to keep the model accurate and relevant. This requires continuous monitoring and data collection

## Ethical Considerations:

Fairness: Ensure that the model treats all applicants equitably, without discrimination based on protected attributes such as race, gender, or socioeconomic status.

Transparency: Provide clear and understandable explanations for model predictions to both applicants and stakeholders.

Accountability: Establish mechanisms for accountability, including regular audits and reviews of the model's performance and impact.

Privacy: Safeguard the privacy and confidentiality of applicant data throughout the entire project lifecycle.

Inclusivity: Engage with diverse stakeholders, including representatives from affected communities, to ensure that the project considers a wide range of perspectives and concerns.

By addressing these ethical considerations, the Home Loan Approval Predictor project can promote fairness, transparency, and accountability, ultimately contributing to more ethical and responsible lending practices.

## Challenges/Issues:

Model Training and Hyperparameter Tuning:

Choosing the Right Algorithm: Selecting the best machine learning model among various options.

Hyperparameter Optimization: Finding the optimal settings for each model can be computationally expensive.

Model Monitoring and Maintenance:

Performance Degradation: Continuously monitoring the model to detect and address any decline in accuracy.

Updating the Model: Periodically retraining the model with new data to keep it relevant.

By addressing these challenges, the Home Loan Approval Predictor project can develop a reliable and fair prediction tool that aids financial institutions in making well-informed lending decisions.

## References:

**GitHub - AswinBalajiTR/Loan_Approval_Prediction**: This project examines applicant financial profiles to identify key predictors of loan approval through EDA and machine learning. It aims to uncover patterns in loan decisions and build a predictive model to enhance risk assessment1.

**GitHub - KhushalKhare/Loan-Prediction**: This project focuses on predicting loan approval using machine learning techniques, Big Data, AI, and Android development. It includes a detailed description of the dataset and the implementation of various machine learning models2.

**Prediction of Modernized Loan Approval System Based on Machine Learning**: This paper discusses a machine learning approach to predict loan approval and assess the safety of assigning loans to individuals.