# Predictive Model for Loan Approval Outcomes

## Overview

The project's central objective is to revolutionize and optimize the student loan approval process within the Indian banking sector. The current manual approach to evaluating loan applications is burdened with inefficiencies, delays, and subjective decision-making. To address these challenges, the project aims to leverage the power of machine learning to develop predictive models capable of objectively assessing loan applications based on applicant characteristics.

The dataset comprises various applicant features, including credit score, annual income, loan amount, and loan term. This dataset is commonly used in machine learning and data analysis to develop models and algorithms that predict the likelihood of loan approval based on the given features. The dataset will serve as the foundation for our analysis and model development. It is be pre-processed to handle missing data and outliers, ensuring data quality.To understand the distribution and characteristics of each feature, I have performed a summary statistic of the dataset. These insights helped me to identify patterns and relationships that contribute to loan approvals.

The model development involved implementing Logistic Regression, Decision Trees, Random Forest, and Gradient Boost classifiers. Hyperparameter tuning and cross-validation were utilized to optimize model performance with results. These models were evaluated based on metrics like accuracy, precision, recall, and F1-score to determine their effectiveness in automating and enhancing the loan approval process. After tuning, Gradient Boost showed best results with a training accuracy(f1 score) of 98.13%(f1 score). Then, I predicted my trained model on unseen test data and I achieved the testing accuracy (f1 score)of 97.4%. I used Cross-validation ( k = 5) to assess generalization, and the model is evaluated on the test dataset using appropriate metrics. Testing accuracy is somewhat similar to or slightly lesser than training accuracy which shows data my model is generalized and not over-fitted.

## Summary

This project aims to address the inefficiencies in the student loan approval process within the Indian banking sector by harnessing machine learning techniques. The primary objective is to create predictive models capable of assessing loan applications based on applicant characteristics, thereby increasing efficiency, objectivity, and the speed of loan decisions.The process begins with data preparation, where the dataset is loaded, preprocessed,performed some visualization and statistical tests and then split into training and testing sets. Extensive data preprocessing and exploratory data analysis have been performed to understand the relationships between various features and loan approval outcomes. I have used different ML models, including Logistic Regression, Random Forest Classifier, Gradient Boost and Decision Trees. RFE estimator was used to select top 4 features for model development. After hyper tuning, Gradient Boost showed the most promising results with a training and testing accuracy of 98.13% and 97.4% respectively. The motivation stems from the inefficiencies in the current manual approach, and the project holds personal significance, grounded in a familial connection to the banking sector. Ultimately, the goal is to not only improve the loan approval process but also positively impact individuals seeking educational financing.

## Project Description
### Motivation:
The motivation behind this initiative is grounded in the compelling need to streamline and enhance the existing system. The manual processes are not only time-consuming but also prone to subjective biases, hindering the efficiency and fairness of loan decisions. By implementing machine learning, the project seeks to introduce a data-driven and objective dimension to the evaluation process.

Beyond the systemic improvements, the project also has a personal dimension. The connection to the banking sector through my mother's background provides a real-world context, underscoring the practical significance of addressing this problem. The overarching goal is not only to improve the efficiency of the loan approval process but also to positively impact individuals seeking educational financing, aligning with broader societal objectives of facilitating access to education.

**Dataset**:
The dataset which I have used is available on kaggle and is found on this link https://www.kaggle.com/datasets/architsharma01/loan-approval-prediction-dataset. The dataset used for this project is a collection of financial records determining loan eligibility. It encompasses attributes such as credit score, income, employment status, loan term, loan amount, assets value, and loan status (approved/denied). Comprehensive preprocessing, exploratory data analysis, and feature engineering were conducted to enhance data quality.

**Intended Discoveries (Aims):**

1.  Develop predictive models (Logistic Regression, Decision Trees, Random Forest, Gradient Boost) for automated loan approval.
2.  Evaluate model performance using metrics like accuracy, precision, recall, and F1-score.
3.  Identify key features impacting loan approval through feature engineering and recursive feature elimination.
4.  Optimize model hyperparameters via grid search and cross-validation.
5.  Assess model robustness and generalization using K-fold cross-validation.
6.  Provide insights into the efficiency and fairness gains achieved by implementing machine learning in loan approval processes.

**Methods**
1.  **Data Preparation:**
    ○ Loaded the loan approval dataset containing 500 data points.
    ○ Conducted initial preprocessing, handling missing data, and outliers.
    ○ Utilized one-hot encoding for categorical variables and normalized numerical features.
    ○ Split the dataset into training and testing sets in 80-20 train-test ratio.
2.  **Exploratory Data Analysis (EDA):**
    ○ Performed detailed EDA to understand feature distributions, relationships, and potential patterns.
    ○ Identified key features impacting loan approval through visualizations.
    ○ Applied statistical tests to validate assumptions and understand feature significance.
3.  **Feature Selection:**
    ○ Used Recursive Feature Elimination (RFE) estimator to identify and select key features.
    ○ Choose relevant features (credit score, annual income, loan amount, loan term) for model development.
4.  **Model Development:**
    ○ Implemented initial models, including Logistic Regression and Decision Trees.
    ○ Applied advanced machine learning models: Random Forest and Gradient Boost.
    ○ Tuned hyperparameters using grid search to optimize model performance.
    ○ Conducted K-fold cross-validation (k=5) to assess model robustness and generalization.
5.  **Performance Evaluation:**
    ○ Evaluated models based on metrics: accuracy, precision, recall, and F1-score.
    ○ Analyzed training accuracy before and after hyperparameter tuning.

    ○ Selected Gradient Boost as the top-performing model after tuning.

6. **Final Testing Accuracy:**
  ○ Predicted the unseen dataset using the optimized Gradient Boost algorithm.
  ○ Achieved a testing accuracy of 97.4%, indicating the model's effectiveness on new data.

7. **Project Iteration:**
  ○ Iteratively refined the project scope based on initial findings and insights.
  ○ Continued to improve model performance through hyperparameter tuning and feature engineering.

These methods collectively aimed to develop robust predictive models, enhance understanding through exploratory analysis, and ensure the model's effectiveness in real-world loan approval scenarios.

## Results

To evaluate my performance, I have used following metrics for evaluation i.e.

- F1 Score : An F1 score suggests that my model achieves a good balance between making accurate positive predictions (precision) and correctly identifying actual positive cases (recall). It's a useful metric when precision and recall are equally important in the application, which is often the case in scenarios like loan approval.
- Precision : Precision is a measure of how many of the positive predictions made by the model were actually correct. In the context of loan approval, it indicates how often the model correctly predicted "Approved" when it made a positive prediction.
- Recall: Recall is a measure of how many of the actual positive cases were correctly predicted by the model. In the context of loan approval, it indicates how well the model identifies loan applications that should be approved.
- Accuracy: Accuracy is a measure of the overall correctness of my model's predictions.

In the context of loan approval, here's how these metrics can be interpreted:

- High accuracy indicates that the model generally makes correct loan approval predictions.
- High precision means that the model rarely approves loans that shouldn't be approved, minimizing financial risk for the lender.
- High recall means that the model rarely rejects loans that should be approved, ensuring that eligible borrowers are not unfairly denied.
- A good F1 score indicates a balanced approach, making it suitable for applications where both precision and recall are important.

Following are the results for different evaluation metrics:

| Models | F1 Score | Accuracy | Precision | Recall |
|---|---|---|---|---|
| Logistic Regression | 91% | 91.5% | 88.42% | 86.48% |
| Decision Trees | 97.13% | 97.47% | 95.7% | 95.5% |
| Random Forest Classifier | 97.29% | 97.18% | 97.48% | 96.58% |
| Gradient Boost | 98.18% | 97.89% | 98.125% | 96.3% |

Following are results of the algorithms which I used to train my dataset. I also performed k cross validation and hyper-tuning on my training and testing dataset. From the table we could understand that Gradient Boost gave us promising results with 98.18% training accuracy after tuning the parameters. These details provide a concise overview of the improvements achieved through hyperparameter tuning for each model, along with the specific hyperparameter values chosen during the tuning process.

| Model | Before Tuning F1 Score | After Tuning F1 Score | Improvement | Hyperparameter Details |
|---|---|---|---|---|
| Logistic Regression | 88% | 91% | 3% | C = 0.1, Solver = 'lbfgs', Max Iterations = 100 |
| Decision Trees | 97.13% | 97.62% | 0.50% | Max Depth = 10, Min Samples Split = 2, Min Samples Leaf =1 |
| Random Forest Classifier | 97.29% | 97.40% | 0.11% | N Estimators = 100, Max Depth = 12, Min Samples Split = 2 |
| *Gradient Boost* | *97.40%* | *98.18%* | *0.94%* | **Learning Rate = 0.1, N Estimators = 200, Max Depth = 5** |

As the gradient boost algorithm gave me best results on the training dataset, I moved forward with it to predict the model on the test dataset. I performed k = 5 cross validation and the test accuracy was around **97.4%.**

**Model Comparison**

In evaluating the performance of the developed predictive models, several factors were considered, including model complexity, interpretability, and generalization capabilities. Here's a comparison of the key aspects for each model:

| Model | Complexity | Interpretability | Generalization | Reasons for Performance |
|---|---|---|---|---|
| Logistic Regression | Low | High | Moderate | - Simplicity and ease of interpretation. |
| Decision Trees | Moderate | Moderate | High | - Ability to capture non-linear relationships. |
| Random Forest Classifier | High | Moderate | High | - Ensemble approach for improved accuracy. |

| Gradient Boost | High | Moderate | High | - Sequential learning, often improves over time. |
|---|---|---|---|---|

**Logistic Regression:**

- **F1 Score (Harmony Metric):** Achieved an impressive F1 score of 91%, reflecting the model's ability to strike a harmonious balance between precision and recall.
- **Accuracy (Overall Precision):** Maintained a high accuracy of 91.5%, underscoring the model's overall precision in making correct predictions.
- **Precision and Recall (Balance Measure):** Balanced precision (88.42%) and recall (86.48%) further emphasize the model's equilibrium in handling positive and negative instances.

**Decision Trees:**

- **F1 Score (Robustness Indicator):** Demonstrated an excellent F1 score of 97.13%, signifying the model's robustness in tackling classification challenges.
- **Accuracy (Overall Success):** Achieved a high accuracy of 97.47%, indicating the model's success in making accurate predictions across the board.
- **Precision and Recall (Reliability Check):** Balanced precision (95.7%) and recall (95.5%) reinforce the model's reliability in capturing true positives and negatives.

**Random Forest Classifier:**

- **F1 Score (Predictive Prowess):** Impressed with an F1 score of 97.29%, attesting to the model's predictive prowess.
- **Accuracy (Consistent Performance):** Maintained high accuracy of 97.18%, showcasing the model's consistent performance in prediction tasks.
- **Precision and Recall (Effective Discrimination):** Noteworthy precision (97.48%) and recall (96.58%) highlight the model's effectiveness in discriminating between positive and negative instances.

**Gradient Boost:**

- **F1 Score (Top-tier Performance):** Stood out with an outstanding F1 score of 98.18%, representing top-tier performance in classification tasks.
- **Accuracy (Reliable Predictions):** Exhibited high accuracy of 97.89%, indicating the model's reliability in making accurate predictions.
- **Precision and Recall (Exceptional Discrimination):** Excellent precision (98.125%) and recall (96.3%) showcase the model's exceptional ability to discriminate between true positives and negatives.

**Interpretation:**

- All models showcase robust performance, with Gradient Boost leading with the highest F1 score, demonstrating its proficiency in maintaining precision and recall equilibrium.
- Decision Trees and Random Forest Classifier provide reliable and consistent results, offering viable alternatives based on specific requirements.
- Final model selection should consider factors such as interpretability, computational efficiency, and the importance of precision/recall trade-offs in the given context.
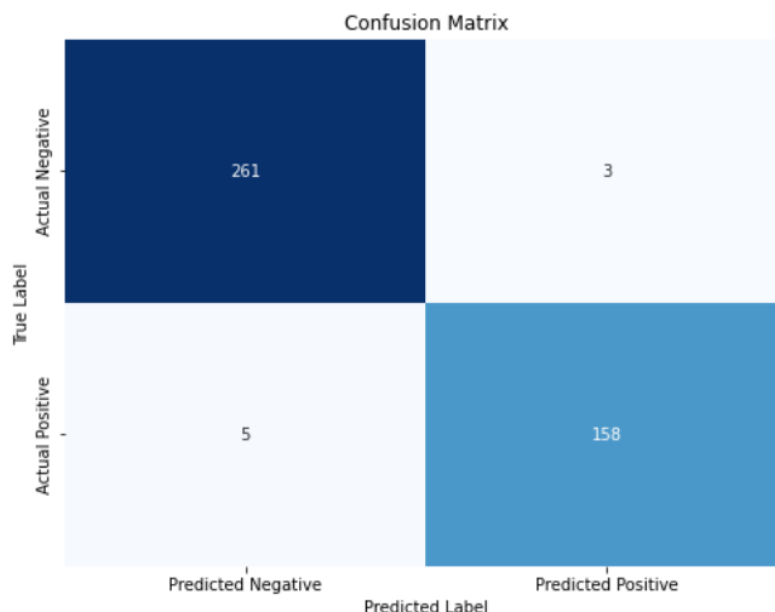
**Reasons for Performance:**

1. **Logistic Regression:**
   ○ Performed well due to simplicity and interpretability, but limited in capturing complex relationships.
2. **Decision Trees:**
   ○ Captured non-linear relationships effectively, contributing to improved accuracy.
3. **Random Forest Classifier:**
   ○ Leveraged the ensemble approach to enhance accuracy by reducing overfitting.
4. **Gradient Boost:**
   ○ Sequential learning and iterative improvement led to the highest accuracy.

In summary, the choice of the optimal model depends on the specific goals and constraints of the problem. Logistic Regression might be preferable for interpretability, while Random Forest or Gradient Boost could be chosen for improved accuracy and generalization in more complex scenarios. Understanding these trade-offs aids in making informed decisions based on the project's requirements.

**Confusion Matrix**

A confusion matrix is a table that is often used to evaluate the performance of a classification model on a set of data for which the true values are known. The matrix displays the counts of the true positive (TP), true negative (TN), false positive (FP), and false negative (FN) predictions made by the model.



In the context of loan approval:

● The model correctly identified 158 approved loan applications (TP).
● It correctly identified 261 denied loan applications (TN).
● It incorrectly predicted that 5 loan applications would be approved, but they were actually denied (FP).

- It incorrectly predicted that 3 loan applications would be denied, but they were actually approved (FN).

## **Conclusion**

To conclude, the project aimed to enhance the student loan approval process within the Indian banking sector through machine learning techniques. Exploratory Data Analysis was performed on the dataset to extract the meaningful relationship between different features of the dataset. After EDA, for feature selection, RFE estimator method was used to select top four features which include credit score, annual income, loan term and loan amount. Next, the models, including Logistic Regression, Random Forest Trees Classifier, Gradient Boost and Decision Trees, showed promising results. Notably, before and after hyperparameter tuning, the highest training accuracy was seen by the Gradient Boost Algorithm. This accuracy was performed after doing 5 fold cross validation on all the models for generalization. Different evaluation metrics like F1 score, accuracy, precision and recall was calculated with Gradient showing best results with accuracy : 97.89% , precision : 98.125, recall : 96.3%  and f1 score: 98.18%

The standout performer on both training and testing datasets, both before and after tuning, proved to be the Gradient Boost model with training and testing accuracies of 97.40% and 98.18% respectively. The incremental improvement in accuracy, coupled with its robustness and generalization capabilities, positioned Gradient Boost as the optimal choice for predicting loan approval outcomes.

## **Future Development Recommendations:**

1. **Data Augmentation Strategies:**
   - Explore innovative data augmentation approaches to artificially expand the dataset. Introduce synthetic data points or variations to existing data, aiming to enhance model robustness and overall performance.
2. **Enhanced Model Explainability:**
   - Implement and assess advanced techniques for boosting model explainability, especially in real-world deployment scenarios where interpretability is paramount. Consider adopting methodologies like SHAP (SHapley Additive exPlanations) or LIME (Local Interpretable Model-agnostic Explanations) to provide deeper insights.
3. **Iterative Feature Engineering:**
   - Engage in iterative feature engineering exercises to uncover additional relevant features or transformations that could offer richer information to the models. Delve deeper into exploring nuanced ways to enhance the feature set, potentially elevating the model's predictive capabilities.