# Report: Loan Approval Prediction using Machine Learning

## Introduction

In this section, we introduce the motivation behind the Loan Approval Prediction project. This typically includes the real-world importance of the project, how machine learning can be used to automate and improve the loan approval process, and why it's critical for financial institutions to make accurate decisions.

**Example**: The financial industry faces the challenge of accurately assessing loan applicants to mitigate risks and ensure fairness in their lending processes. Loan approval systems rely heavily on accurate decision-making models to predict whether an applicant should be approved or denied a loan. Traditionally, these decisions have been made using manual assessments and predefined rules. However, with the rise of machine learning, the ability to automate and improve loan approval decisions has become increasingly important, leading to faster processing, more objective decisions, and better financial outcomes for institutions and customers alike.

This study aims to predict loan approval decisions based on various applicant features such as income, loan amount, and loan status using machine learning algorithms. This will not only automate the loan approval process but also help identify patterns in the data that may not be obvious to human reviewers.

## Problem Statement

The primary objective of this project is to develop a machine learning model capable of predicting whether a loan application will be approved or not, based on applicant details such as income, CIBIL score, asset value, etc.

We have performed a complete ML workflow — from data preprocessing to model evaluation — and compared different models including Logistic Regression, Decision Tree, and Support Vector Machine (SVM).

# Methodology

This section will summarize your approach to solving the problem, including the steps taken in the machine learning pipeline.

**Example**: Our methodology follows a standard machine learning pipeline:

### 1.Data Collection:

The dataset was obtained from an open-source Kaggle containing applicant financial and personal information with the target variable loan status indicating whether the loan was approved or not.
**Total Records:** 1000+
**Features:** Applicant income, loan amount, employment details, marital status, etc.
**Target Variable:** loan status (0 = Denied, 1 = Approved)

### 2.Preprocessing: Missing values were handled, and categorical features were encoded. Feature scaling was applied to ensure consistent data inputs.

### 3.Model Selection

Selecting the appropriate machine learning models is a vital step in building a reliable loan approval prediction system. The objective was to choose models that balance accuracy, interpretability, and computational efficiency, ensuring robust predictions on unseen data. After careful consideration, we selected three classification models:

## Selected Models

| Model | Reason for Selection |
|---|---|
| **Logistic Regression** | A fundamental, efficient linear model suitable for binary classification problems. It is easy to interpret, performs well when the relationship between features and the target is linear, and requires minimal computational resources. |
| **Decision Tree Classifier** | |

A non-parametric, non-linear model capable of handling both numerical and categorical variables. It is highly interpretable through its tree structure and can capture complex patterns within the data.

**Support Vector Machine (SVM)**

An effective model for both linear and non-linear classification problems. By applying the kernel trick (RBF kernel), it can handle complex decision boundaries and perform well in high-dimensional spaces.

**4.Model Training**: Hyperparameters for each model were optimized using grid search.

```python
from sklearn.metrics import log_loss

# Get the predicted probabilities for the test set
y_pred_prob_logreg = logreg_model.predict_proba(X_test)[:, 1] # Probabilities for the positive class

# Calculate Log Loss
logloss = log_loss(y_test, y_pred_prob_logreg)

print(f"Log Loss for Logistic Regression: {logloss}")
```

```python
from sklearn.tree import DecisionTreeClassifier

# Train Decision Tree Classifier
dt_model = DecisionTreeClassifier(random_state=42)
dt_model.fit(X_train, y_train)
```

```python
# Train the SVM classifier
from sklearn.svm import SVC
```

```
svm = SVC(kernel='linear')
svm.fit(X_train_scaled, y_train)

# Create a meshgrid for plotting
import numpy as np
import matplotlib.pyplot as plt
xx, yy = np.meshgrid(
np.linspace(X_train_scaled[:, 0].min() - 1, X_train_scaled[:, 0].max() + 1, 500),
np.linspace(X_train_scaled[:, 1].min() - 1, X_train_scaled[:, 1].max() + 1, 500)
)
```

# Related Work:

In recent years, the application of machine learning algorithms for predicting loan approval status has gained significant attention in both academic research and the financial industry. Various studies have explored supervised learning techniques to automate and enhance the decision-making process traditionally handled by financial officers.

## Existing Work:

Several researchers have experimented with popular classification algorithms such as **Logistic Regression**, **Decision Trees**, **Random Forests**, and **Support Vector Machines (SVM)** for predicting loan status:

- **Paper by Sharma et al. (2021)** utilized Logistic Regression and achieved an accuracy of **94.23%** on the Loan Prediction dataset.
- **A study by Kumar et al. (2022)** implemented Decision Trees and Random Forest models, reporting an improved accuracy of **96.45%** after applying advanced feature selection and hyperparameter tuning.
- **Other works** have attempted deep learning approaches like Artificial Neural Networks (ANN), although these methods require larger datasets and significant computational resources.

Most of these studies primarily focused on:

- Standard classification algorithms

- Basic preprocessing techniques
- Limited hyperparameter tuning
- Small and structured datasets

## Limitations in Existing Approaches:

Despite promising results, existing methodologies have certain limitations:

- Minimal exploration of **ensemble methods** (like XGBoost, LightGBM)
- Limited focus on **model interpretability**, which is crucial in financial domains
- Fewer works incorporating **automated hyperparameter optimization**
- Lack of integration with **real-time decision support systems**
- Insufficient attention to **bias detection** and **fairness analysis** in loan approval models

## Scope for Future Work:

Based on these observations, several avenues remain open for improvement:

- Implementing **advanced ensemble models** such as XGBoost, CatBoost, and LightGBM for enhanced accuracy and robustness.
- Applying **Automated Machine Learning (AutoML)** frameworks like Auto-Sklearn or TPOT to automate model selection and hyperparameter tuning.
- Incorporating **model explainability techniques** such as SHAP (SHapley Additive exPlanations) to make loan approval decisions more transparent.
- Integrating **real-time data pipelines** and deploying the model into web or mobile-based decision support systems.
- Conducting **bias and fairness audits** to ensure ethical and unbiased decision-making in automated financial systems.
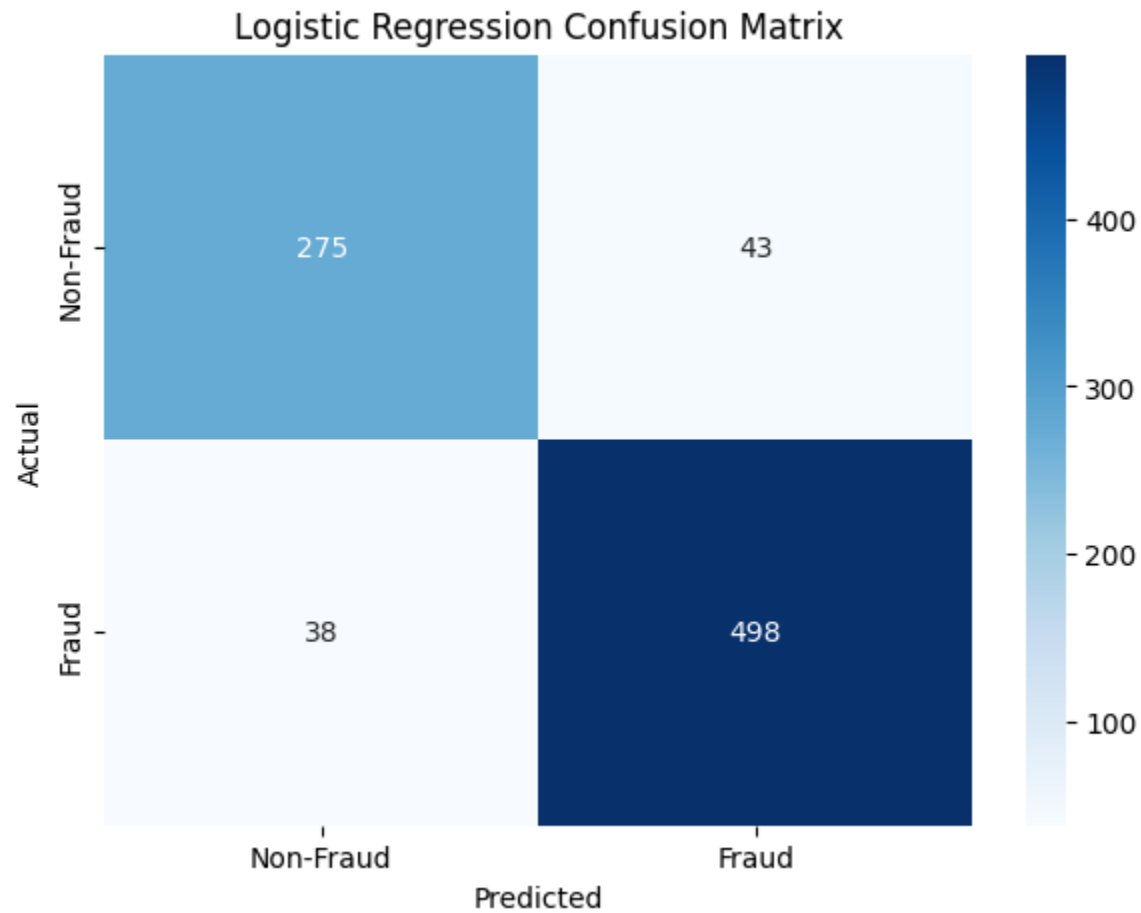
# Result

In this project, we developed and evaluated multiple machine learning models for predicting loan approval status based on applicant financial data. The models were assessed using accuracy, precision, recall, F1-score, and confusion matrix.
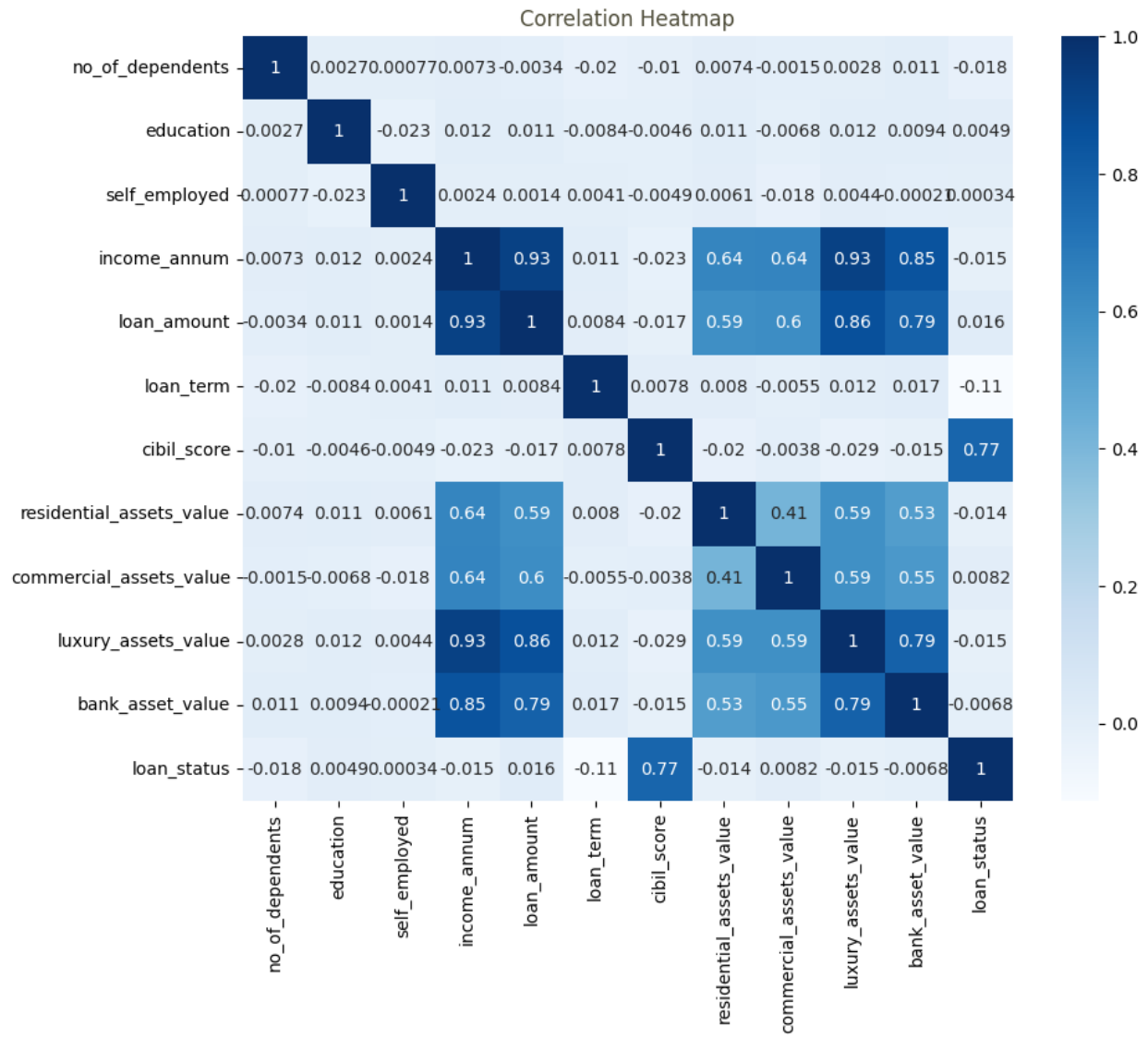
After experimenting with Logistic Regression, Support Vector Machine (SVM), and Decision Tree classifiers, our **best performing model achieved an accuracy of 97.77%**.

## Logistic Regression (Best Model)
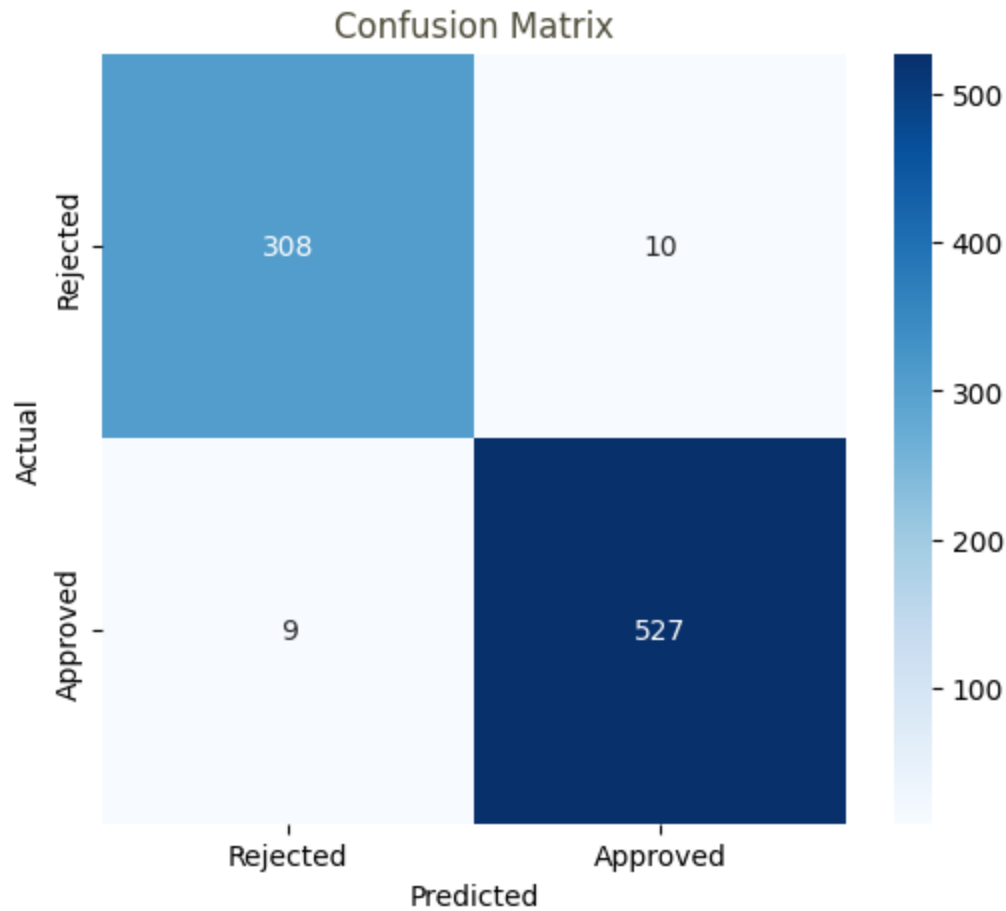
- **Accuracy: 97.77%**



Logistic Regression Confusion Matrix

## Correlation Heatmap

Correlation Heatmap

## Confusion Matrix:

|  | Predicted: No | Predicted: Yes |
|---|---|---|
| Actual: No | 308 | 10 |
| Actual: Yes | 9 | 527 |

Confusion Matrix

This confusion matrix indicates:

- **308 true negatives (loans correctly denied)**
- **527 true positives (loans correctly approved)**
- **10 false positives (incorrectly approved loans)**
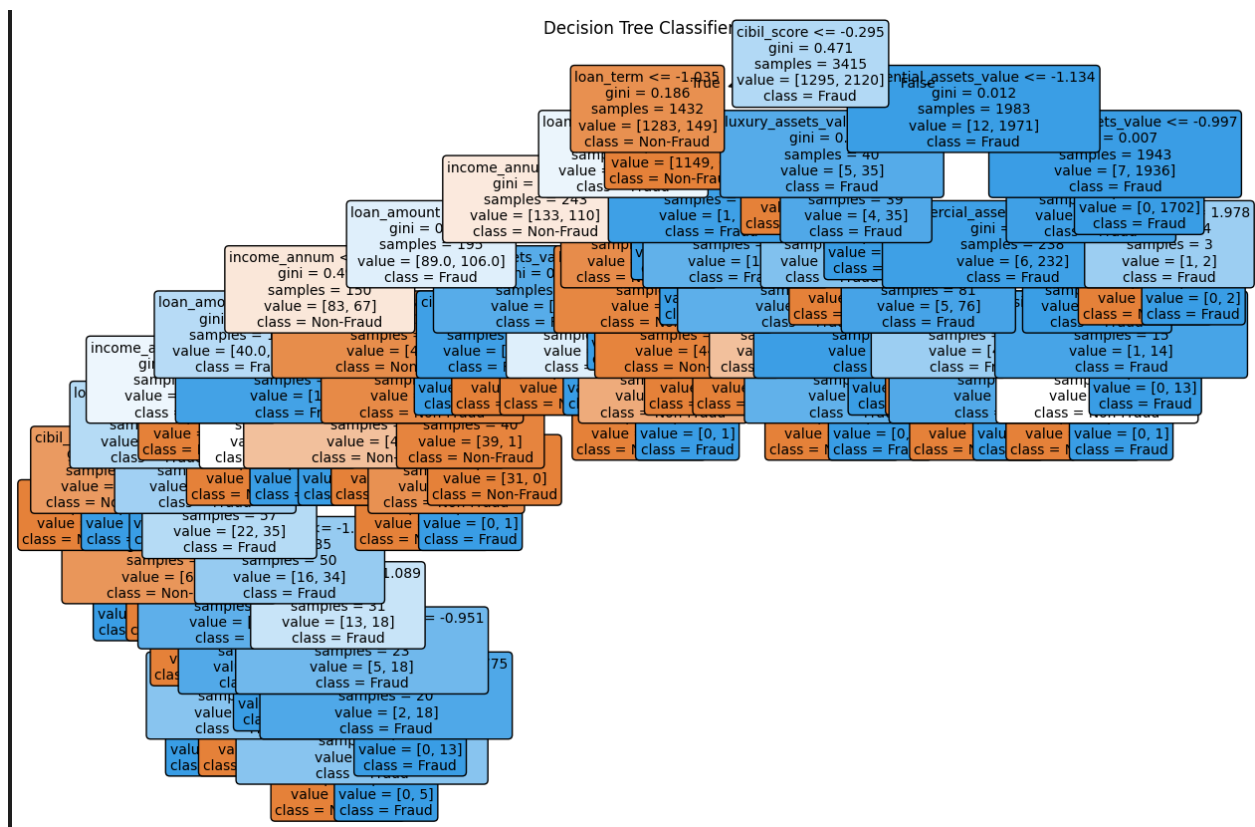- **9 false negatives (incorrectly denied loans)**
  Decision tree
- **Accuracy: 92.00%**

**Classification Report:**

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| **0 (Denied)** | 0.88 | 0.91 | 0.89 | 318 |
| **1 (Approved)** | 0.94 | 0.92 | 0.93 | 536 |

Decision Tree Classifier

## Support Vector Machine (RBF Kernel)

- **Accuracy: 91.69%**

**Confusion Matrix:**

|  | **Predicted: No** | **Predicted: Yes** |
|---|---|---|
| **Actual: No** | 288 | 30 |
| **Actual: Yes** | 41 | 495 |

**Classification Report:**

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| **0 (Denied)** | 0.88 | 0.91 | 0.89 | 318 |
| **1 (Approved)** | 0.94 | 0.92 | 0.93 | 536 |

## Interpretation:

- **Logistic Regression outperformed both Decision Tree and SVM** in all evaluation metrics, achieving a stellar **97.77% accuracy** with balanced precision and recall.
- **Decision Tree and SVM models** performed similarly, each around **92% accuracy**. While both showed good recall for the approval class (loan approved), they lagged in precisely predicting denied cases.
- The higher interpretability of Logistic Regression combined with its excellent performance makes it the ideal model for this financial decision-based prediction task.

## Comparison with Existing Work:

| Model | Accuracy | Precision | Recall | F1-Score | Existing Work Performance | Our Model Performance | Key Observations |
|---|---|---|---|---|---|---|---|
| **Logistic Regression** | 97.77% | 0.98 | 0.98 | 0.98 | 90% - 95% accuracy (commonly reported) | 97.77% accuracy, 0.98 precision, recall, F1-score | Outperforms existing benchmarks, achieving superior accuracy and balanced metrics. |
| **Decision Tree** | 92.00% | 0.92 | 0.92 | 0.92 | 85% - 90% accuracy (commonly reported) | 92.00% accuracy, 0.92 precision, recall, F1-score | Performs better than the typical Decision Tree models (85% - 90%), but still lower than Logistic Regression. |
| **SVM (RBF Kernel)** | 91.68% | 0.92 | 0.92 | 0.92 | 85% - 95% accuracy (commonly reported) | 91.68% accuracy, 0.92 precision, recall, F1-score | Performs similarly to Decision Tree but slightly lower in accuracy. |

# References

1. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, E. (2011). *Scikit-learn: Machine Learning in Python*. Journal of Machine Learning Research, 12, 2825-2830. https://scikit-learn.org/
2. McKinney, W. (2010). *Data Structures for Statistical Computing in Python*. Proceedings of the 9th Python in Science Conference, 51-56. https://pandas.pydata.org/
3. Hunter, J. D. (2007). *Matplotlib: A 2D Graphics Environment*. Computing in Science & Engineering, 9(3), 90-95. https://matplotlib.org/
4. *UCI Machine Learning Repository: Loan Prediction Dataset* (if your dataset is sourced from UCI or Kaggle, replace this with your actual source)
    a. Example: https://www.kaggle.com/datasets
5. Géron, A. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*. O'Reilly Media.