# CAPSTONE PROJECT

## LOAN ELIGIBILITY PREDICTION USING MACHINE LEARNING

**Presented By:**
**1. Vedant Roy-University School of Automation and Robotics-Artificial Intelligence and Machine Learning**

# OUTLINE

- **Problem Statement**

- **Proposed System/Solution**

- **System Development Approach**

- **Data Visualization (EDA)**

- **Algorithm & Deployment**

- **Result**

- **Conclusion**

- **Future Scope**

- **References**

# PROBLEM STATEMENT

The process of loan approval in financial institutions involves significant risk assessment and decision-making. Manual evaluation of loan applications is time-consuming and prone to human error. The challenge is to develop an automated system that can predict the likelihood of loan approval based on historical data.

# PROPOSED SOLUTION

- The proposed solution leverages machine learning techniques to predict the likelihood of loan approval. The solution will consist of:

- Data Collection:
    - Historical loan application data including applicant details, loan amount, loan term, credit history, etc.
    - External data like economic indicators, regional data, etc. (if available).

- Data Preprocessing:
    - Cleaning the data to handle missing values and outliers.
    - Feature engineering to create relevant features that impact loan approval.

- Machine Learning Algorithm:
    - Implementation of algorithms such as Logistic Regression, Decision Trees, Random Forest, Gradient Boosting, and XGBoost.
    - Evaluation and comparison of these models.

- Deployment:
    - Development of an interface or application for loan approval prediction.
    - Deployment on a scalable platform.

- Evaluation:
    - Using metrics such as accuracy, precision, recall, F1-score, and AUC-ROC score to evaluate model performance.
    - Continuous monitoring and fine-tuning of the model.

edunet
foundation

# SYSTEM DEVELOPMENT APPROACH

**System Requirements**

•Python environment set up

•Jupyter Notebook or any other Python IDE

•Required Python libraries (Pandas, NumPy, Scikit-learn, Seaborn, Matplotlib, XGBoost, Imbalanced-learn)
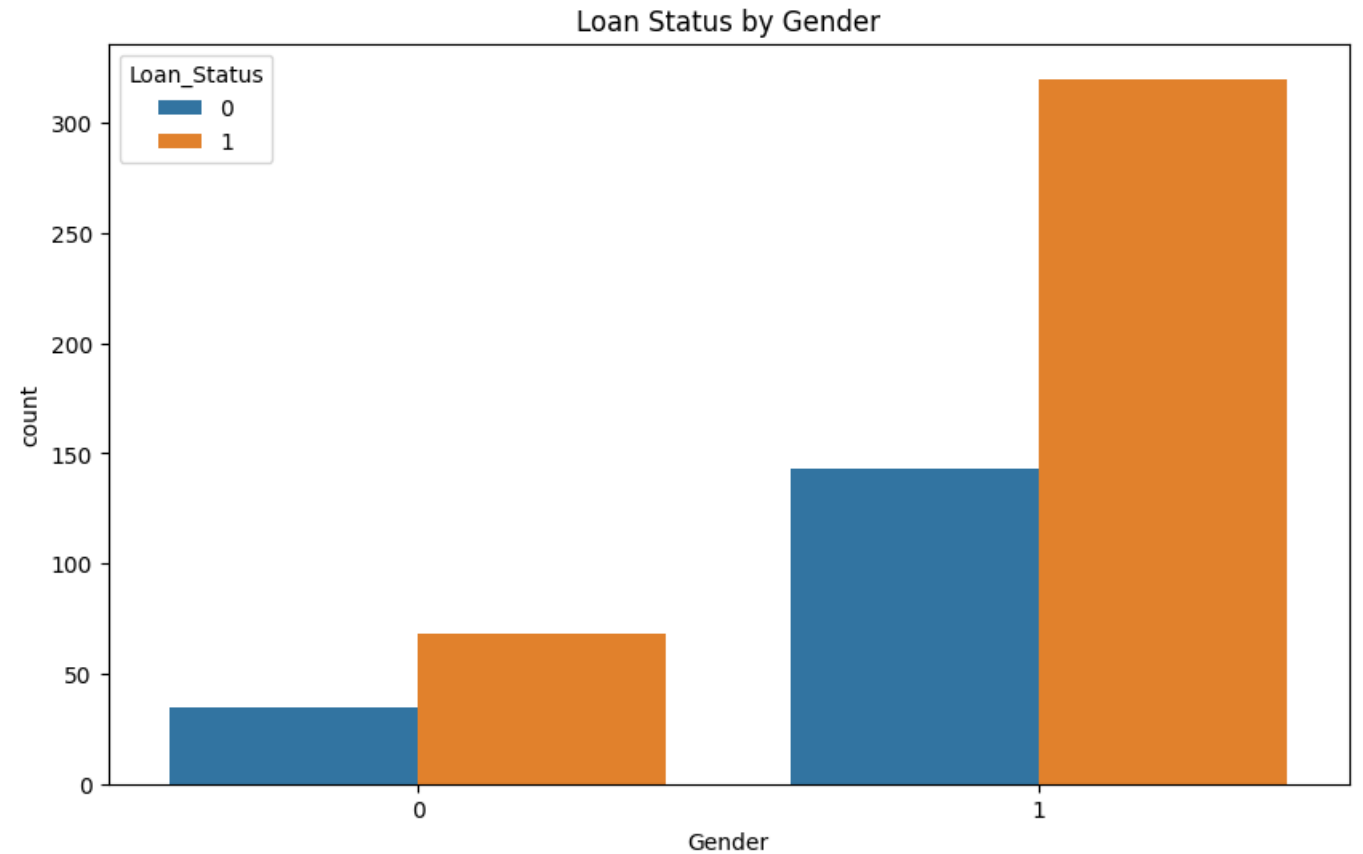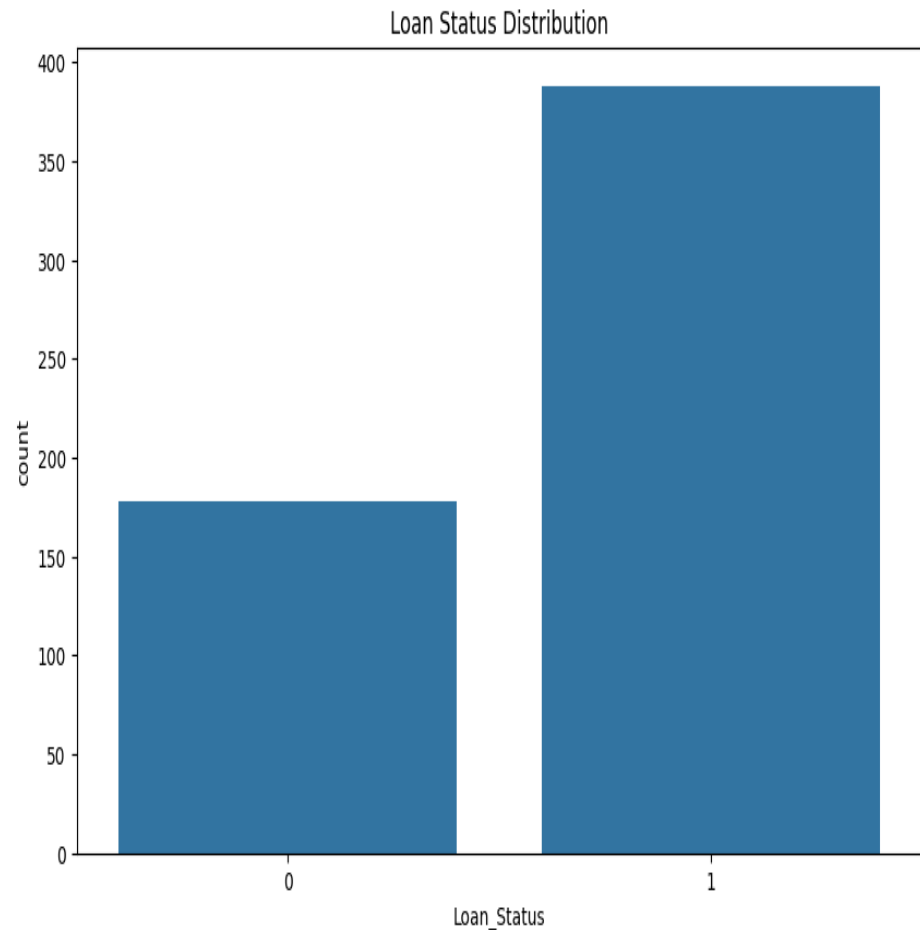
**Library Requirements**

```bash
pip install pandas numpy scikit-learn seaborn matplotlib xgboost imbalanced-learn
```
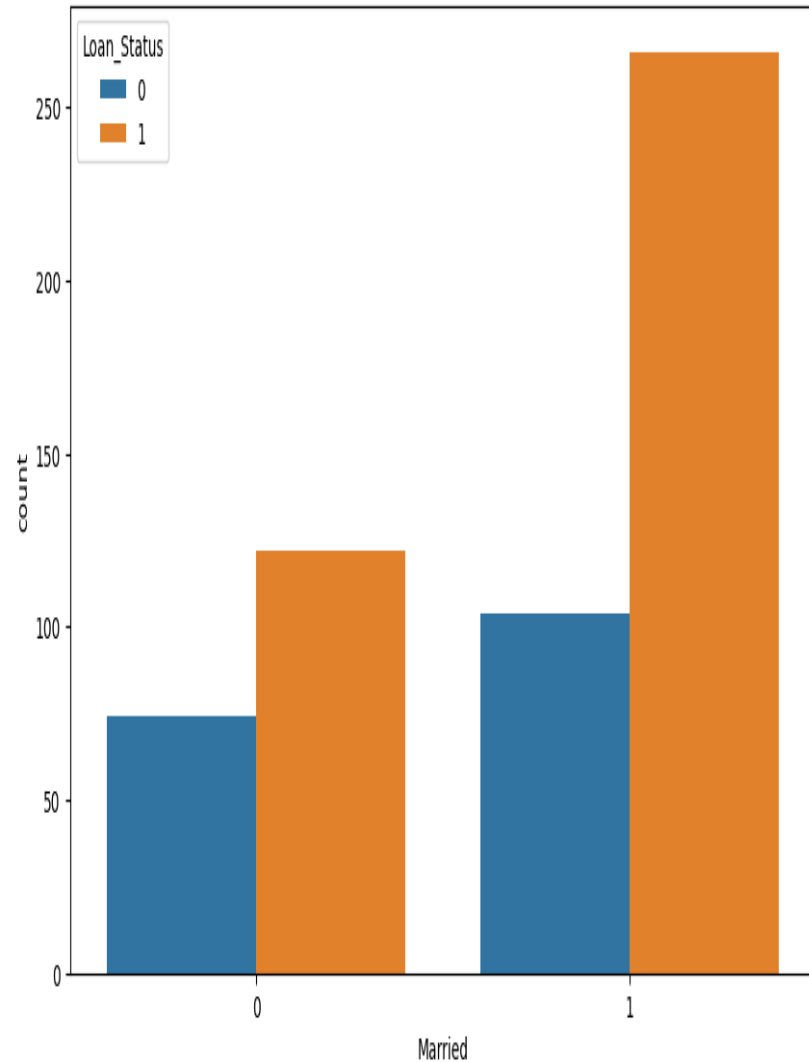
# DATA VISUALIZATION (EDA)

EDA was performed too understand the distribution and relationships of the data. The following visualizations were created to analyze the dataset:



Loan Status Distribution



Loan Status by Gender
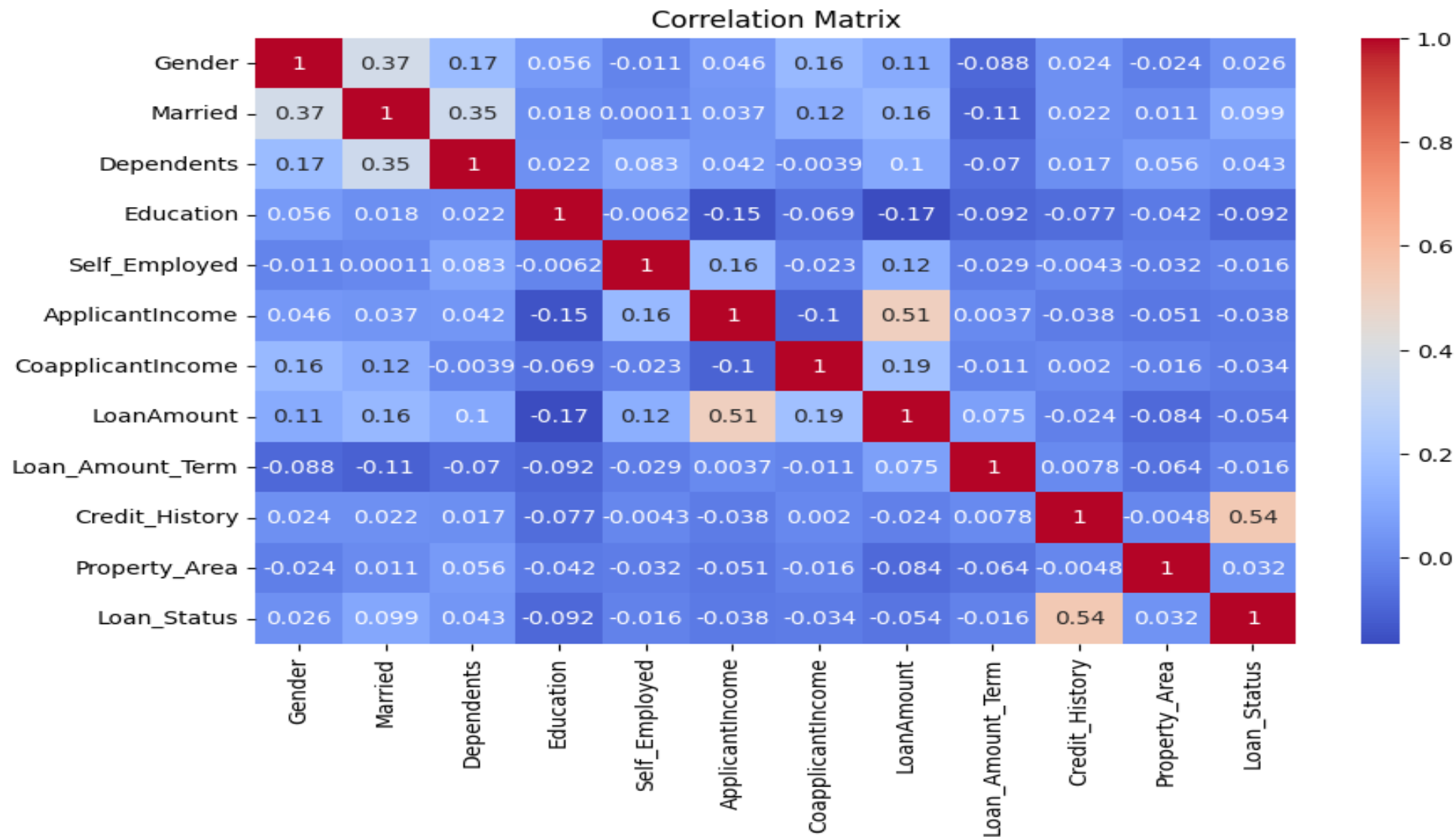
# DATA VISUALIZATION (EDA)

# DATA VISUALIZATION (EDA)

# DATA VISUALIZATION (EDA)



Correlation Matrix

# ALGORITHM & DEPLOYMENT

## Algorithm and Deployment

### Algorithm Selection

- **Logistic Regression:** Chosen for its simplicity and interpretability.

- **Decision Trees:** For capturing non-linear relationships.

- **Random Forest:** For better accuracy by reducing overfitting.

- **Gradient Boosting:** For improving model performance through boosting.

- **XGBoost:** For efficient and scalable implementation of gradient boosting.

## Data Input

- Historical loan data with features like applicant income, co-applicant income, loan amount, loan term, credit history, etc.

## Training Process

- Splitting the data into training and testing sets.

- Training each algorithm on the training set.

- Performing cross-validation and hyperparameter tuning.

## Prediction Process

- Making predictions on the test set.

- Evaluating model performance using the selected metrics.

# RESULT

The machine learning models were trained and evaluated to predict the likelihood of loan approval. Below are the detailed results for each model:
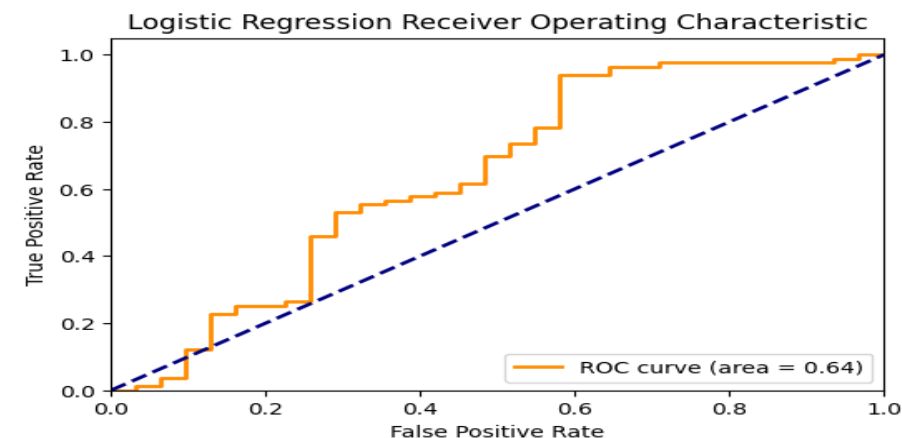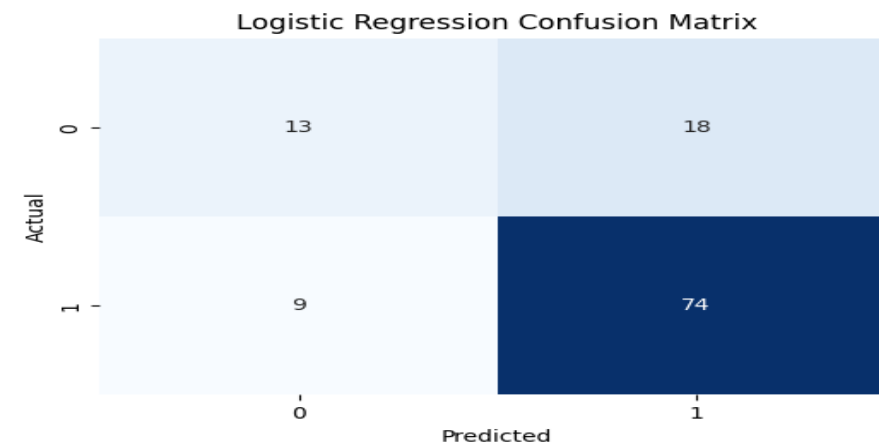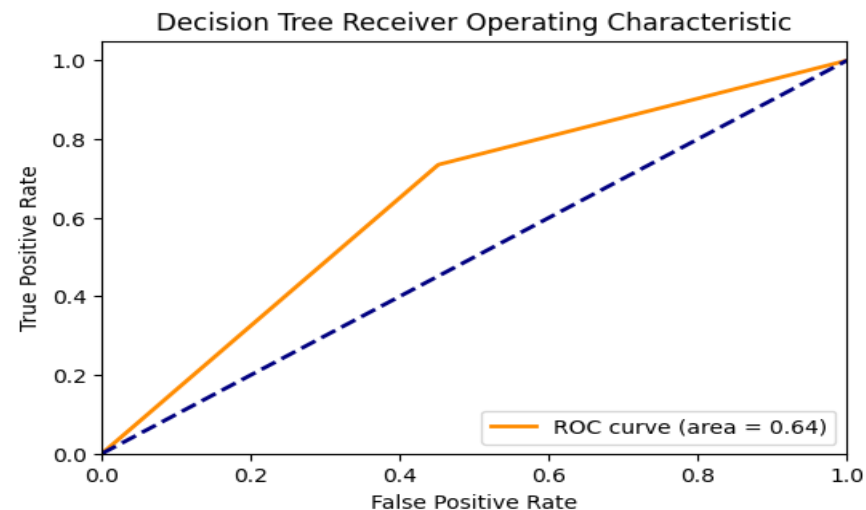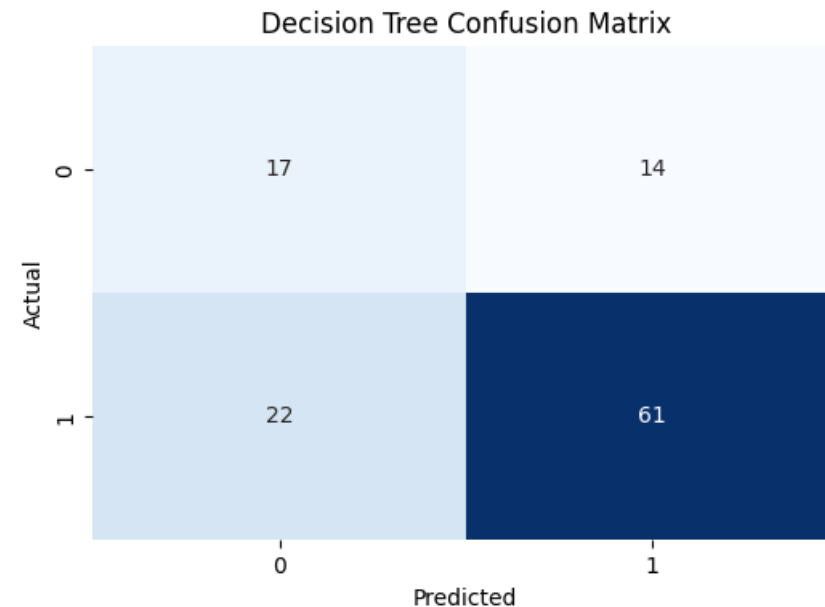
## 1. Logistic Regression

- Classification Report:

  - Precision: 0.59 for class 0, 0.80 for class 1

  - Recall: 0.42 for class 0, 0.89 for class 1

  - F1-score: 0.49 for class 0, 0.85 for class 1

  - Accuracy: 76.32%

  - AUC-ROC Score: 0.6555

- Confusion Matrix:

  - True Positives: 74

  - True Negatives: 13

  - False Positives: 18

  - False Negatives: 9



Logistic Regression Confusion Matrix



Logistic Regression Receiver Operating Characteristic

ROC curve (area = 0.64)

# RESULT

## 2. Decision Tree

- **Classification Report:**
  - Precision: 0.44 for class 0, 0.81 for class 1
  - Recall: 0.55 for class 0, 0.73 for class 1
  - F1-score: 0.49 for class 0, 0.77 for class 1
  - Accuracy: 68.42%
  - AUC-ROC Score: 0.6417

- **Confusion Matrix:**
  - True Positives: 61
  - True Negatives: 17
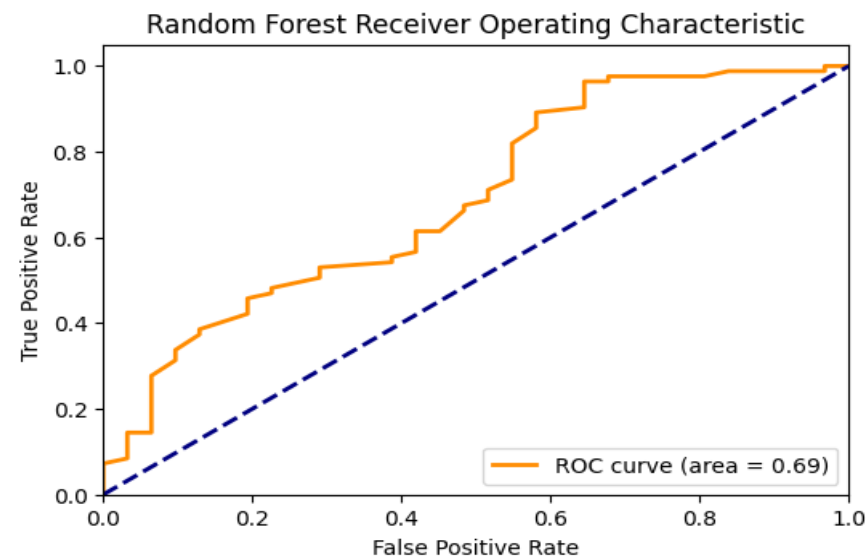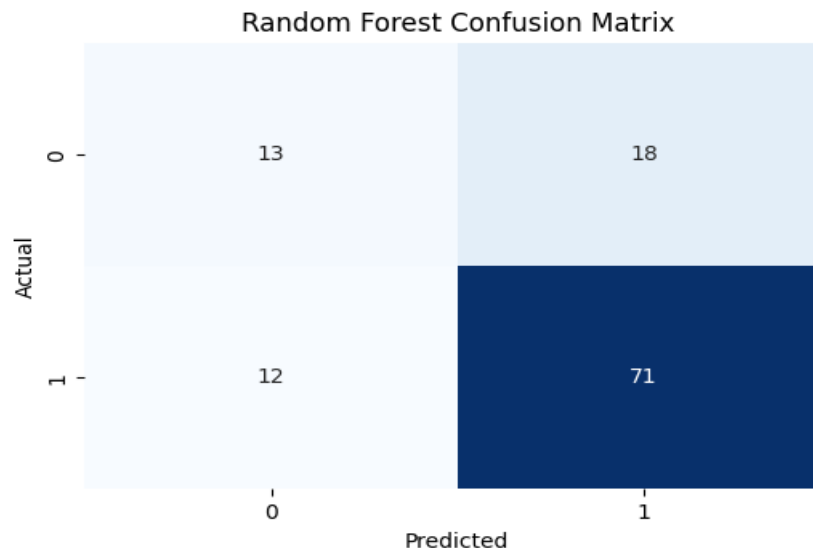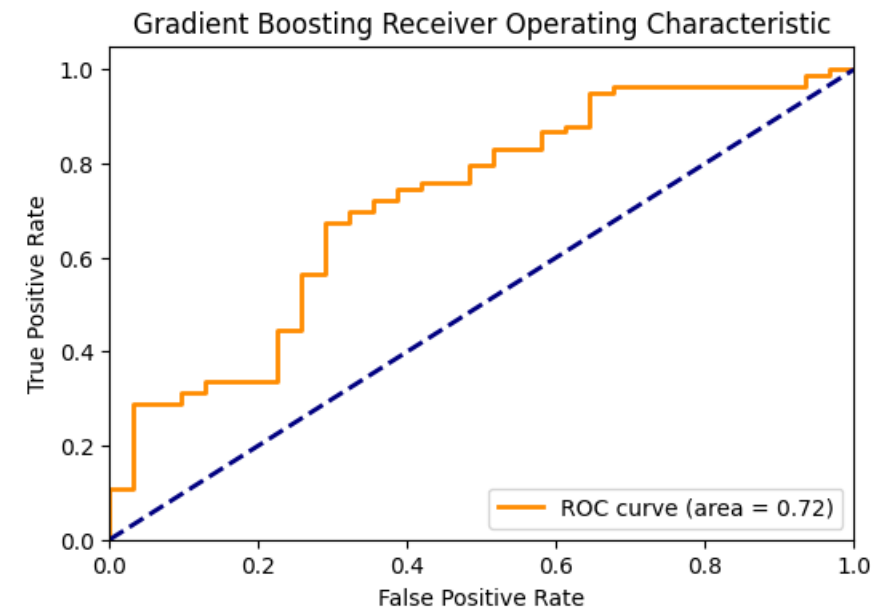  - False Positives: 14
  - False Negatives: 22



Decision Tree Confusion Matrix



Decision Tree Receiver Operating Characteristic

# RESULT

## 3. Random Forest

- **Classification Report:**

    - Precision: 0.52 for class 0, 0.80 for class 1

    - Recall: 0.42 for class 0, 0.86 for class 1

    - F1-score: 0.46 for class 0, 0.83 for class 1

    - Accuracy: 73.68%

    - AUC-ROC Score: 0.6374

- **Confusion Matrix:**

    - True Positives: 71

    - True Negatives: 13
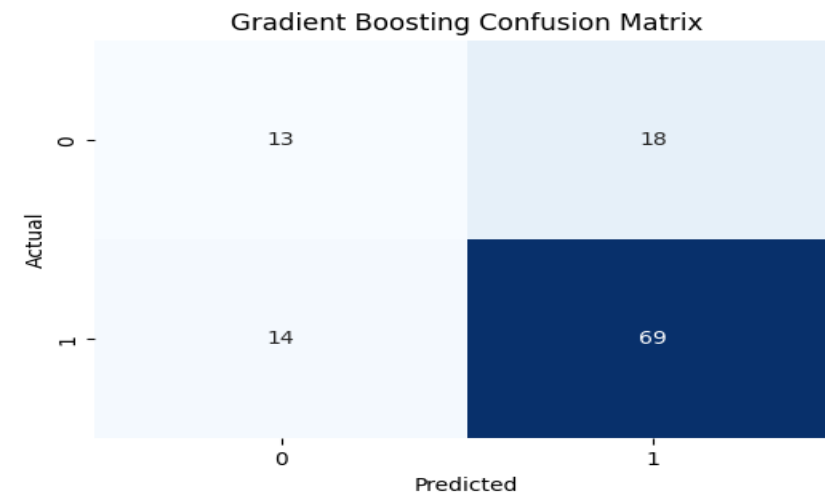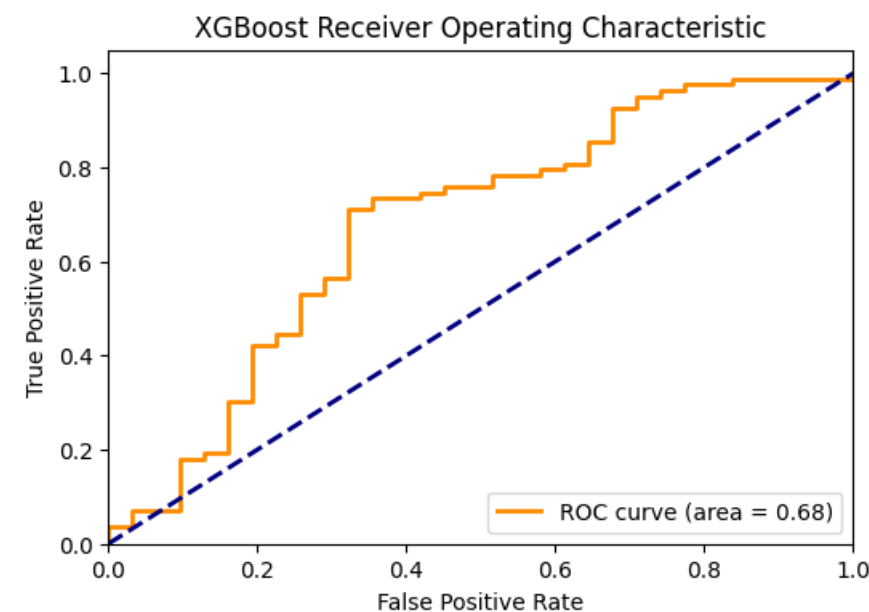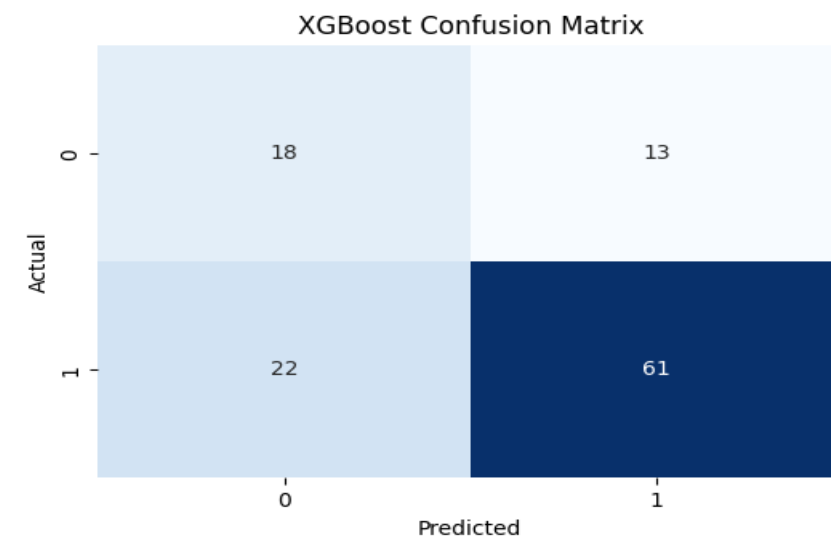
    - False Positives: 18

    - False Negatives: 12



Random Forest Confusion Matrix



Random Forest Receiver Operating Characteristic

ROC curve (area = 0.69)

# RESULT

## 4. Gradient Boosting

- **Classification Report:**

  - Precision: 0.48 for class 0, 0.79 for class 1

  - Recall: 0.42 for class 0, 0.83 for class 1

  - F1-score: 0.45 for class 0, 0.81 for class 1

  - Accuracy: 71.93%

  - AUC-ROC Score: 0.6253

- **Confusion Matrix:**

  - True Positives: 69

  - True Negatives: 13

  - False Positives: 18

  - False Negatives: 14



Gradient Boosting Confusion Matrix



Gradient Boosting Receiver Operating Characteristic

# RESULT

## 5. XGBoost

- **Classification Report:**

  - Precision: 0.45 for class 0, 0.82 for class 1

  - Recall: 0.58 for class 0, 0.73 for class 1

  - F1-score: 0.51 for class 0, 0.78 for class 1

  - Accuracy: 69.30%

  - AUC-ROC Score: 0.6578

- **Confusion Matrix:**

  - True Positives: 61

  - True Negatives: 18

  - False Positives: 13

  - False Negatives: 22



XGBoost Confusion Matrix



XGBoost Receiver Operating Characteristic

# RESULT

## Model Comparison

The results indicate that the Logistic Regression model achieved the highest accuracy of 76.32%, followed by Random Forest at 73.68%. However, in terms of the AUC-ROC score, XGBoost slightly outperformed Logistic Regression, indicating better performance in distinguishing between the classes.

# CONCLUSION

- The project successfully developed and evaluated several machine learning models to predict loan approval. Logistic Regression demonstrated the highest accuracy (0.76), indicating its reliability for this task. The Random Forest model also performed well, with an accuracy of 0.74, suggesting it as a strong candidate for deployment. Although the XGBoost model had a slightly lower accuracy (0.69), it provided the best AUC-ROC score (0.658), suggesting a robust classification capability.

- During the implementation, some challenges were encountered, such as handling missing values and class imbalances. These were addressed using data preprocessing techniques like filling missing values and employing the SMOTE algorithm for oversampling.

- Potential improvements for this project include further hyperparameter tuning of the models and exploring additional features that could improve prediction accuracy. Additionally, combining multiple models through ensemble methods might yield better results.

- Accurate loan approval predictions are crucial for financial institutions to minimize risks and streamline the approval process. The models developed in this project can significantly contribute to achieving these goals, ensuring efficient and reliable loan processing.

# FUTURE SCOPE

There are several potential enhancements and expansions for this loan approval prediction system, which can further improve its accuracy, scalability, and usability. These future improvements could include:

1. Incorporating Additional Features:

- **Credit Score:** Integrating credit scores can significantly enhance the model's predictive power, providing a more comprehensive view of an applicant's creditworthiness.

- **Employment History:** Including details about an applicant's employment history, such as job stability and length of employment, can add valuable context to the prediction.

- **Financial Indicators:** Other financial indicators, such as debt-to-income ratio, savings, and investment details, can provide a more holistic assessment of the applicant's financial health.

# REFERENCES

2. Implementing Real-Time Data Processing:

• **Real-Time Data Integration:** Enhancing the model to process real-time data can improve prediction accuracy and applicability. This would involve streaming data pipelines that continuously update the model with the latest information.

• **Dynamic Model Updating**: Implementing mechanisms for dynamically updating the model as new data comes in, ensuring the model remains relevant and accurate over time.


3. Deploying the Model as a Web Service:

• **API Integration:** Deploying the model as a web service with API endpoints allows for seamless integration with financial institution systems. This can enable automated, real-time loan approval decisions.

• **Scalable Infrastructure:** Ensuring the deployment infrastructure is scalable to handle large volumes of requests efficiently, possibly leveraging cloud services for flexibility and reliability.

• **User-Friendly Interface:** Developing a user-friendly interface for stakeholders to interact with the model, visualize results, and gain insights without needing in-depth technical knowledge.

edunet
foundation

# REFERENCES

- 1. <u>Kaggle Dataset:</u> Loan data was obtained from Kaggle. The dataset can be accessed at [Kaggle: Loan Data Set] (https://www.kaggle.com/datasets/burak3ergun/loan-data-set).

- 2. <u>Source Dataset Information:</u> The original dataset was sourced from Analytics Vidhya's practice problem, "Loan Prediction III". More details can be found at [Analytics Vidhya: Loan Prediction III] (https://datahack.analyticsvidhya.com/contest/practice-problem-loan-prediction-iii/).

- 3. <u>Logistic Regression:</u> Hosmer, D. W., & Lemeshow, S. (2000). "Applied Logistic Regression". John Wiley & Sons.

- 4. <u>Decision Trees:</u> Breiman, L., Friedman, J. H., Olshen, R. A., & Stone, C. J. (1984). "Classification and Regression Trees". Wadsworth & Brooks/Cole Advanced Books & Software.

- 5. <u>Random Forests:</u> Breiman, L. (2001). "Random Forests". Machine Learning, 45(1), 5-32. DOI: [10.1023/A:1010933404324] (https://doi.org/10.1023/A:1010933404324).

edunet
foundation

# REFERENCES

- 6. Gradient Boosting: Friedman, J. H. (2001). "Greedy Function Approximation: A Gradient Boosting Machine". Annals of Statistics, 29(5), 1189-1232. DOI: [10.1214/aos/1013203451](https://doi.org/10.1214/aos/1013203451).

- 7. XGBoost: Chen, T., & Guestrin, C. (2016). "XGBoost: A Scalable Tree Boosting System". Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 785-794. DOI: [10.1145/2939672.2939785](https://doi.org/10.1145/2939672.2939785).

- 8. SMOTE: Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). "SMOTE: Synthetic Minority Over-sampling Technique". Journal of Artificial Intelligence Research, 16, 321-357. DOI: [10.1613/jair.953](https://doi.org/10.1613/jair.953).

- 9. Model Evaluation Metrics: Fawcett, T. (2006). "An Introduction to ROC Analysis". Pattern Recognition Letters, 27(8), 861-874. DOI: [10.1016/j.patrec.2005.10.010](https://doi.org/10.1016/j.patrec.2005.10.010).

- 10. Data Preprocessing Techniques: Han, J., Kamber, M., & Pei, J. (2011). "Data Mining: Concepts and Techniques". Elsevier.

- 11. Project Code and Outputs: The complete code and outputs for this project can be accessed in the Google Colab notebook available at: https://colab.research.google.com/drive/18VV0DSwpSBA2SrK7peF-Z9Sdt-PrOGDX?usp=sharing

edunet
foundation

# COURSE CERTIFICATE (GETTING STARTED WITH ENTERPRISE- GRADE AI)

In recognition of the commitment to achieve professional excellence

Getting Started with
Enterprise-grade AI
IBM SkillBuild

IBM

## Vedant Roy

Has successfully satisfied the requirements for:

Getting Started with Enterprise-grade AI

Issued on: 11 JUL 2024
Issued by IBM

Verify: https://www.credly.com/go/hzuHkn5o

IBM.

edunet
foundation

# COURSE CERTIFICATE 2 (CLOUD COMPUTING FUNDAMENTALS)

Note- Didn't got certificate from credly website despite completing the course 100% so sharing the screenshots as proof of completion

Part- 1

# COURSE CERTIFICATE 2 (CLOUD COMPUTING FUNDAMENTALS)

Part- 2

# COURSE CERTIFICATE 2 (CLOUD COMPUTING FUNDAMENTALS)

Part- 3

# THANK YOU

edunet
foundation