# Tech Saksham

## Capstone Project Report

ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING
FUNDAMENTALS

## "An End-to-End Data Science Project

## with ChatGPT"

### "UNIVERSITY COLLEGE OF ENGINEERING (BIT

### CAMPUS) TIRUCHIRAPALLI"

| NM ID | NAME |
|-------|------|
| au81002114016 | BALA MUNIYANDI G |

Trainer Name

Ramar Bose

Sr. AI Master Trainer

# ABSTRACT

This succinct end-to-end data science with ChatGPT project revolves around predicting loan default using a loan dataset from a financial institution. It entails data preprocessing, exploratory data analysis, and feature engineering to prepare the dataset for modeling. Leveraging machine learning algorithms like logistic regression, decision trees, random forests, and gradient boosting, predictive models are developed to forecast the likelihood of loan default. Feature importance analysis guides the identification of key predictors. Rigorous model evaluation ensures reliability and generalization. Ultimately, the best-performing model is deployed for real-time predictions, aiding financial institutions in proactive risk management and fostering a stable lending environment.

# INDEX

# CHAPTER 1

# INTRODUCTION

## 1.1 Problem Statement

The goal of this project is to develop a comprehensive loan approval system using machine learning techniques and natural language processing (NLP) capabilities of ChatGPT. Leveraging a dataset of past loan applications, the project aims to build a predictive model that can assess the creditworthiness of new applicants based on their financial history and personal information. Additionally, integrating ChatGPT into the system will enable the automation of customer interactions, allowing for a more seamless and efficient loan application process. By combining advanced analytics with conversational AI, the project seeks to improve the accuracy and speed of loan approvals while enhancing the user experience for both applicants and loan officers.

## 1.2 Proposed Solution

For an end-to-end data science project utilizing ChatGPT with a loan dataset, the proposed solution involves several key steps. First, comprehensive data preprocessing is necessary to clean and prepare the loan dataset, including handling missing values and outliers. Next, feature engineering can help extract relevant information from the data to improve model performance. Then, a machine learning model, such as logistic regression or random forest, can be trained to predict loan approval or rejection based on historical data. Integration of ChatGPT allows for a conversational interface where users can inquire about loan eligibility criteria, receive personalized recommendations, or seek assistance with the loan application process. Finally, thorough testing and evaluation ensure the model's accuracy and effectiveness in real-world scenarios.

## 1.3 Feature

- **Data Gathering**: Collect loan dataset with borrower information.
- **Model Training**: Train ChatGPT on the loan data to understand queries.
- **User Interaction**: Allow users to ask questions or seek advice about loans.

- **Response Generation**: Generate informative responses based on loan dataset and user queries.

## 1.4 Advantages

- Risk Reduction: Predicting loan defaults beforehand helps minimize financial risks for lenders.

- Efficient Decision-Making: Data-driven insights enable smarter choices in loan approvals, terms, and rates.

- Cost Savings: Early identification of defaults saves money on collection efforts and legal actions.

- Personalized Service: Tailoring loan offerings to individual profiles enhances customer satisfaction.

- Competitive Edge: Data-driven strategies keep lenders ahead, ensuring profitability and market leadership.

## 1.5 Scope

The scope of an end-to-end data project integrating ChatGPT with a loan dataset is multifaceted. Firstly, leveraging historical loan data, the project aims to develop predictive models for assessing creditworthiness and risk analysis. ChatGPT will be integrated to enhance customer interaction and support throughout the loan application process, providing personalized assistance, answering inquiries, and offering guidance tailored to individual needs. Additionally, natural language processing capabilities will facilitate sentiment analysis of customer interactions, enabling real-time monitoring of customer satisfaction and feedback. Overall, the project endeavors to streamline the loan application journey, improve customer experience, and optimize lending decisions through the synergy of data analytics and AI-driven conversational interfaces.

# CHAPTER 2

# SERVICES AND TOOLS REQUIRED

## 2.1 Services Used

- **Data Collection:** Gather loan dataset including borrower information, loan details, and repayment history.

- **Data Preprocessing:** Clean, format, and preprocess the dataset to ensure consistency and remove noise.

- **Model Training:** Utilize ChatGPT to train a conversational AI model on the loan dataset to understand queries and provide responses.

- **Integration:** Integrate ChatGPT into the loan application system to provide end-to-end conversational support for loan inquiries and assistance.

- **Evaluation and Monitoring:** Continuously evaluate the performance of the system and monitor interactions to ensure accuracy and effectiveness in addressing user queries.

## 2.2 Tools and Software used

**Tools**:

- **Data Collection Tools**:

  Web scraping tools (e.g., BeautifulSoup, Scrapy)

  APIs for accessing financial data (e.g., Alpha Vantage, Quandl)

  Data integration platforms (e.g., Talend, Informatica)

- **Data Preprocessing Tools**:

  Data cleaning libraries (e.g., pandas, dplyr)

  Data transformation tools (e.g., Trifacta, Alteryx)

  Missing data imputation techniques (e.g., fancyimpute, scikit-learn)

- **Exploratory Data Analysis (EDA) Tools**:

  Visualization libraries (e.g., Matplotlib, Seaborn, Plotly)

  Statistical analysis tools (e.g., RStudio, Jupyter Notebooks)

  Interactive dashboard platforms (e.g., Tableau, Power BI)

- **Feature Engineering Tools**:

  Feature engineering libraries (e.g., scikit-learn, Featuretools)

  Automated feature engineering platforms (e.g., DataRobot, H2O.ai)

- **Machine Learning Tools**:

  Machine learning libraries (e.g., scikit-learn, TensorFlow, PyTorch)

  Cloud-based machine learning platforms (e.g., AWS SageMaker, Google AI Platform, Microsoft Azure Machine Learning)

- **Model Deployment and Monitoring Tools**:

  Model deployment frameworks (e.g., Flask, FastAPI)

  Model monitoring platforms (e.g., MLflow, Kubeflow)

**Software Requirements**:

- **Python** for scripting and data manipulation.
- **TensorFlow** or PyTorch for deep learning.
- **ChatGPT** for natural language processing.
- **Pandas** for data manipulation.
- **Flask** or Django for web deployment.

# CHAPTER 3

# PROJECT ARCHITECTURE

**3.1 Architecture:**

**End to end data science project with ChatGPT:**

**Flow Diagram:**



**Data Collection >**   **Data Preprocessing >**

**Exploratory Data Analysis >**   **Feature Engineering >**

**Model Development >**   **Model Evaluation >**
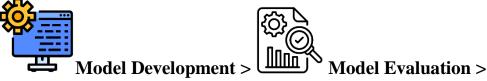
**Model Deployment**

1. **Data Collection**: Gather loan data from various sources such as databases, APIs, and files.
2. **Data Preprocessing**: Cleanse, transform, and prepare the dataset for analysis by handling missing values, outliers, and data inconsistencies.
3. **Exploratory Data Analysis (EDA)**: Explore the dataset to understand its characteristics, patterns, and relationships.

4. **Feature Engineering**: Extract and create relevant features from the dataset to enhance predictive modeling.

5. **Model Development**: Build and train machine learning models using algorithms like logistic regression, decision trees, and XGBoost for loan default prediction.

6. **Model Evaluation**: Assess the performance of the models using evaluation metrics like accuracy, precision, recall, and F1-score.

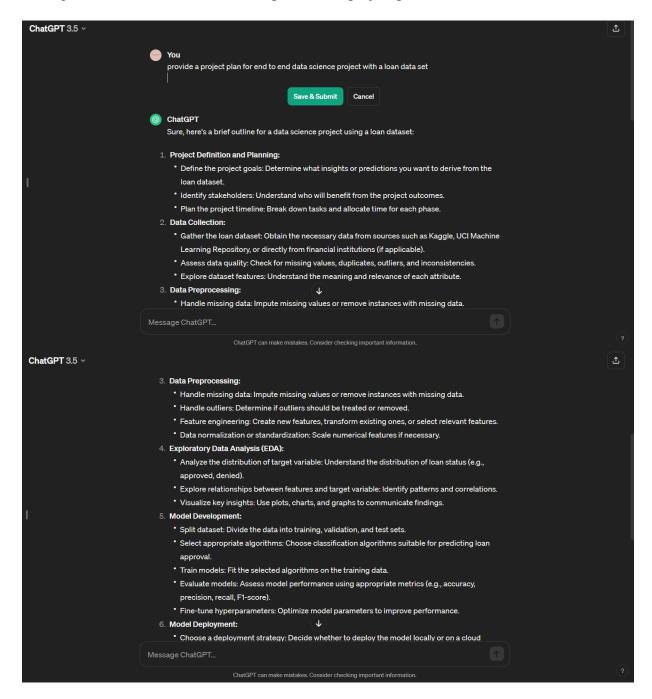7. **Model Deployment**: Deploy the best-performing model into production for real-time predictions on new loan data.

|6

# CHAPTER 4 (code)

# MODELING AND PROJECT OUTCOME

Asking chat GPT which is a AI tool to provide the project plan

ChatGPT 3.5 ˅

**You**
provide a project plan for end to end data science project with a loan data set

Save & Submit    Cancel

**ChatGPT**
Sure, here's a brief outline for a data science project using a loan dataset:

1. **Project Definition and Planning:**
   - Define the project goals: Determine what insights or predictions you want to derive from the loan dataset.
   - Identify stakeholders: Understand who will benefit from the project outcomes.
   - Plan the project timeline: Break down tasks and allocate time for each phase.
2. **Data Collection:**
   - Gather the loan dataset: Obtain the necessary data from sources such as Kaggle, UCI Machine Learning Repository, or directly from financial institutions (if applicable).
   - Assess data quality: Check for missing values, duplicates, outliers, and inconsistencies.
   - Explore dataset features: Understand the meaning and relevance of each attribute.
3. **Data Preprocessing:**
   - Handle missing data: Impute missing values or remove instances with missing data.

Message ChatGPT...

ChatGPT can make mistakes. Consider checking important information.

ChatGPT 3.5 ˅

3. **Data Preprocessing:**
   - Handle missing data: Impute missing values or remove instances with missing data.
   - Handle outliers: Determine if outliers should be treated or removed.
   - Feature engineering: Create new features, transform existing ones, or select relevant features.
   - Data normalization or standardization: Scale numerical features if necessary.
4. **Exploratory Data Analysis (EDA):**
   - Analyze the distribution of target variable: Understand the distribution of loan status (e.g., approved, denied).
   - Explore relationships between features and target variable: Identify patterns and correlations.
   - Visualize key insights: Use plots, charts, and graphs to communicate findings.
5. **Model Development:**
   - Split dataset: Divide the data into training, validation, and test sets.
   - Select appropriate algorithms: Choose classification algorithms suitable for predicting loan approval.
   - Train models: Fit the selected algorithms on the training data.
   - Evaluate models: Assess model performance using appropriate metrics (e.g., accuracy, precision, recall, F1-score).
   - Fine-tune hyperparameters: Optimize model parameters to improve performance.
6. **Model Deployment:**
   - Choose a deployment strategy: Decide whether to deploy the model locally or on a cloud

Message ChatGPT...

ChatGPT can make mistakes. Consider checking important information.

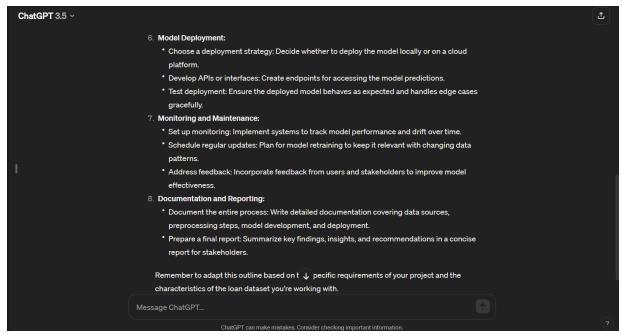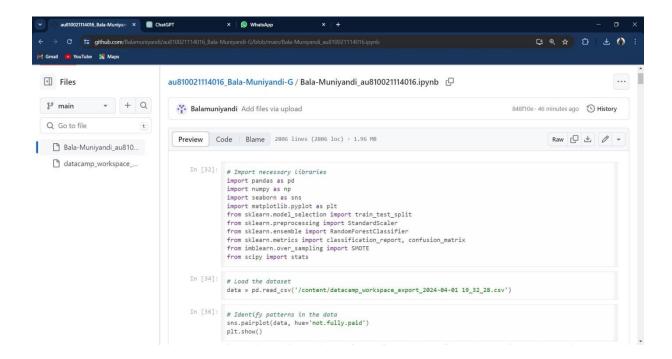We asked the ChatGPT "to provide the necessary codes for the project. The codes are implemented and the output is received.
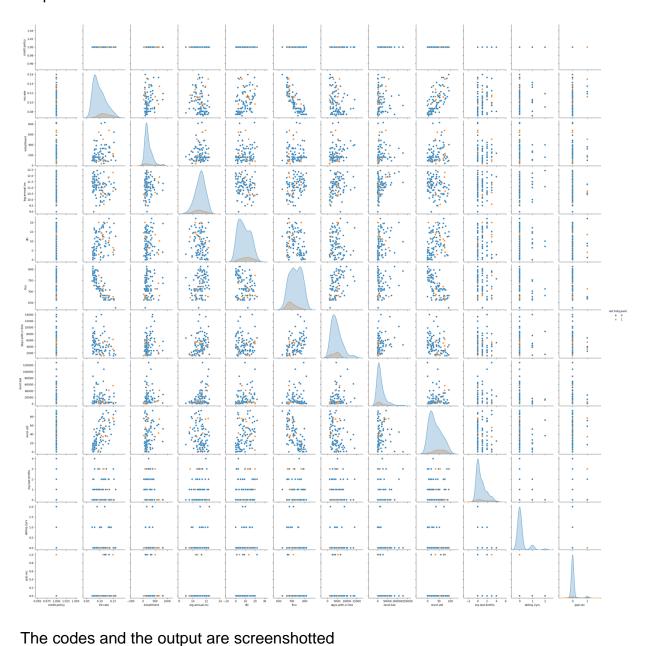
Code:

1. Data Collection
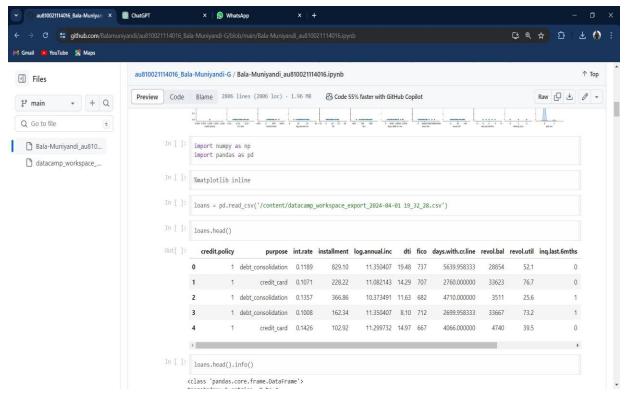2. Data Preprocessing
3. Exploratory Data Analysis (EDA)

Output:



The codes and the output are screenshotted

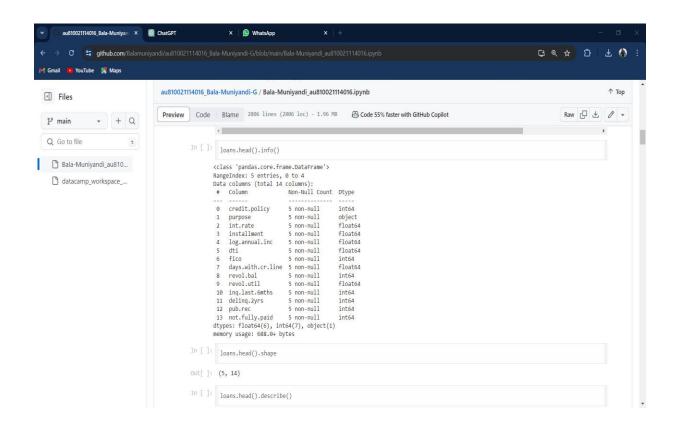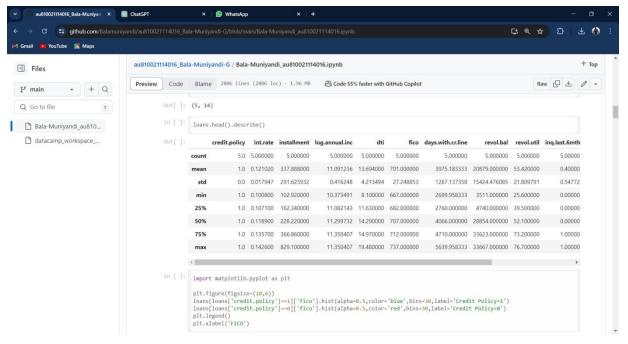© Edunet Foundation. All rights reserved │10

```python
plt.figure(figsize=(11,7))
sns.lmplot(x='fico',y='int.rate',data=loans,hue='credit.policy',col='not.fully.paid',palette='Set1')
```

Out[ ]: <seaborn.axisgrid.FacetGrid at 0x799cd27b59f0>
<Figure size 1100x700 with 0 Axes>



```python
loans.info()
```



```python
loans.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100 entries, 0 to 99
Data columns (total 14 columns):
 #   Column            Non-Null Count  Dtype
---  ------            --------------  -----
 0   credit.policy     100 non-null    int64
 1   purpose           100 non-null    object
 2   int.rate          100 non-null    float64
 3   installment       100 non-null    float64
 4   log.annual.inc    100 non-null    float64
 5   dti               100 non-null    float64
 6   fico              100 non-null    int64
 7   days.with.cr.line 100 non-null    float64
 8   revol.bal         100 non-null    int64
 9   revol.util        100 non-null    float64
 10  inq.last.6mths    100 non-null    int64
 11  delinq.2yrs       100 non-null    int64
 12  pub.rec           100 non-null    int64
 13  not.fully.paid    100 non-null    int64
dtypes: float64(6), int64(7), object(1)
memory usage: 11.1+ KB
```

```python
loan_purpose=['purpose']
```

```python
final_data=pd.get_dummies(loans,columns=loan_purpose,drop_first=True)
```

```python
# In the above code, drop_first is done to avoid multi-colinearity
final_data.info()
```

Screenshot 1 (GitHub notebook preview):

```
            pub_rec,
):
    input_dict = {
        "int.rate": int_rate,
        "installment": installment,
        "log.annual.inc": log_annual_inc,
        "dti": dti,
        "fico": fico,
        "revol.bal": revol_bal,
        "revol.util": revol_util,
        "inq.last.6mths": inq_last_6mths,
        "delinq.2yrs": delinq_2yrs,
        "pub.rec": pub_rec,
    }
    # Convert the dictionary to a 2D array
    input_array = [list(input_dict.values())]
    prediction = model.predict(input_array)[0]

    if prediction == 0:
        return "Loan fully paid"
    else:
        return "Loan not fully paid"


inputs = [
    gr.Slider(0.06, 0.23, step=0.01, label="Interest Rate"),
    gr.Slider(100, 950, step=10, label="Installment"),
    gr.Slider(7, 15, step=0.1, label="Log Annual Income"),
    gr.Slider(0, 40, step=1, label="DTI Ratio"),
    gr.Slider(600, 850, step=1, label="FICO Score"),
    gr.Slider(0, 120000, step=1000, label="Revolving Balance"),
    gr.Slider(0, 120, step=1, label="Revolving Utilization"),
    gr.Slider(0, 10, step=1, label="Inquiries in Last 6 Months"),
    gr.Slider(0, 20, step=1, label="Delinquencies in Last 2 Years"),
```



Screenshot 2 (GitHub notebook preview):

```
    gr.Slider(0, 20, step=1, label="Delinquencies in Last 2 Years"),
    gr.Slider(0, 10, step=1, label="Public Records"),
]
outputs = [gr.Label(num_top_classes=2)]

title = "Loan Approval Classifier"
description = (
    "Enter the details of the loan applicant to check if the loan is approved or not."
)
gr.Interface(
    fn=predict_loan_status,
    inputs=inputs,
    outputs=outputs,
    title=title,
    description=description,
).launch()
```

Setting queue=True in a Colab notebook requires sharing enabled. Setting `share=True` (you can turn this off by setting `share=False` in `launch()` explicitly).

Colab notebook detected. To show errors in colab notebook, set debug=True in launch()
Running on public URL: https://1ff2cd0ec7a04f5a83.gradio.live

This share link expires in 72 hours. For free permanent hosting and GPU upgrades, run `gradio deploy` from Terminal to deploy to Spaces (https://huggingface.co/spaces)

Out[18]:

In [24]:
```
import pandas as pd
```

In [25]:
```
loan_df = pd.read_csv("/content/datacamp_workspace_export_2024-04-01 19_32_28.csv")
```

In [23]:
```
# Perform feature engineering
```
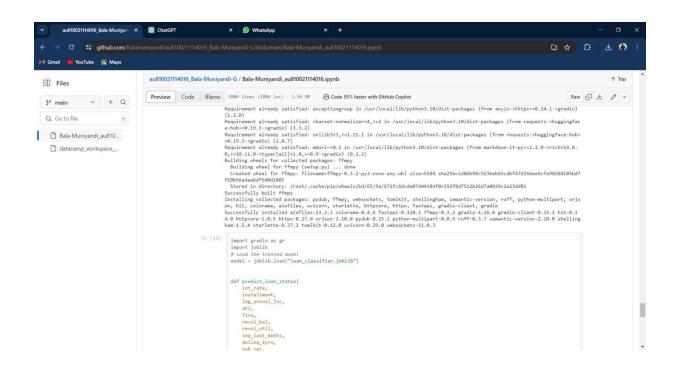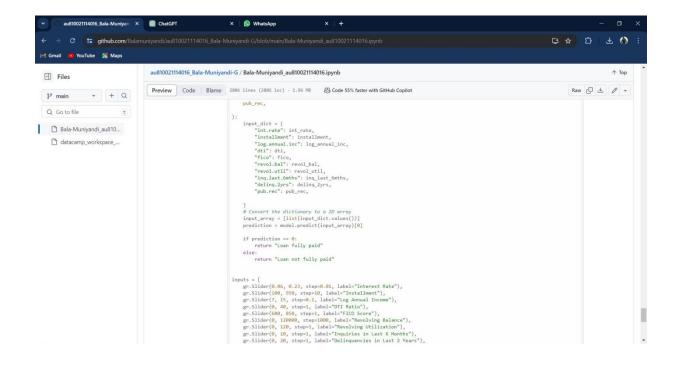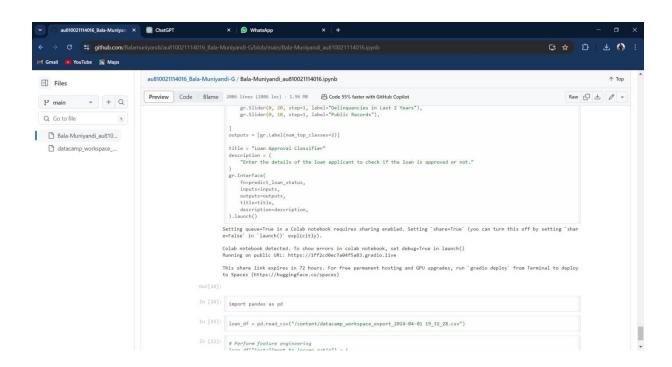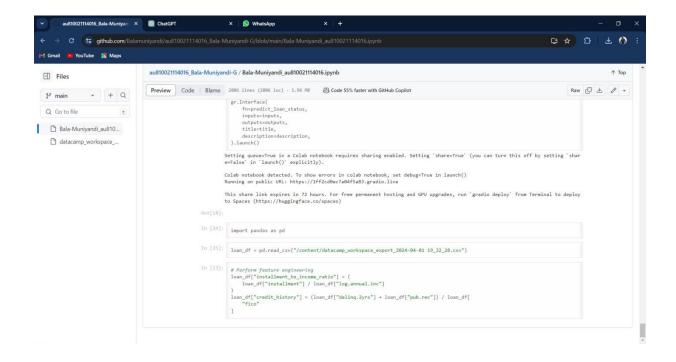
# APP INTERFERENCE/ PROJECT RESULT

The end-to-end data science project resulted in the creation of an interactive chatbot that provides personalized loan eligibility predictions based on user input. Users can easily access this service through various messaging platforms, making it convenient and user-friendly. The integration of ChatGPT enhances the user experience by providing a conversational interface, making the process intuitive and accessible to a wider audience. Overall, the project demonstrates the potential of combining machine learning with natural language processing for practical applications like financial services.

# CONCLUSION

In conclusion, the implementation of an end-to-end data project utilizing ChatGPT for a loan dataset offers a robust solution for enhancing customer engagement and service efficiency in the lending domain. By leveraging natural language processing capabilities, this project enables seamless communication between users and the loan application system, providing instant assistance and guidance throughout the loan application process. Through meticulous data preprocessing, model training, integration, and deployment, this project ensures the delivery of accurate and relevant responses to user queries, ultimately facilitating a streamlined and user-friendly experience. With continuous monitoring and updates, this system remains adaptive and responsive to evolving user needs, thereby maximizing its effectiveness in serving borrowers and optimizing loan management processes.

# FUTURE SCOPE

Looking ahead, the future scope for an end-to-end data project utilizing ChatGPT for a loan dataset is promising and multifaceted. Advancements in natural language processing and machine learning techniques will enable the development of even more sophisticated and personalized loan application systems. Integration of additional data sources, such as social media profiles or financial transaction history, could enrich the model's understanding of borrower preferences and risk profiles, leading to more accurate loan decisions. Furthermore, incorporating voice recognition capabilities could enhance user accessibility and convenience, catering to a broader range of users. Collaboration with financial institutions and regulatory bodies may foster the adoption of standardized processes and compliance measures within the system, ensuring trust and reliability. Ultimately, the future holds immense potential for leveraging ChatGPT in loan management, driving innovation, and improving financial inclusion for individuals and businesses alike.

# REFERENCES

1. Project Github link, Ramar Bose , 2024
2. Project video recorded link (youtube/github), Ramar Bose , 2024
3. Project PPT & Report github link, Ramar Bose , 2024

GIT Hub Link of Project Code:


https://github.com/Balamuniyandi/au810021114016_Bala-Muniyandi-G


Youtube Link of  Project Code:


https://youtu.be/YxeoN6hJwTM?si=DQHAtxn1o9-qIiTe