# PROJECT REPORT

## ON

## "Knowledge Representation and Insights Generation from Structured Datasets"

is submitted to

Intel Unnati Industrial Training 2024

**Prof. Ram Meghe Institute of Technology and Research, Badnera, Amravati.**

By

**Mr. Krushna Mohod**          **Mr. Sudhanshu Atalkar**

Under the Guidance of

**Dr. R. R Karwa**                              **Intel Industry Expert**
**Prof A. U. Chaudhari**

**Prof. Ram Meghe Institute of Technology and Research, Badnera, Amravati.**

**(An Autonomous Institute & NAAC Accredited)**

**2024-2025**

**PROF. RAM MEGHE INSTITUTE OF TECHNOLOGY AND RESEARCH, BADNERA, AMRAVATI.**



## CERTIFICATE

This is to certify that

**Mr. Krushna Mohod          Mr. Sudhanshu Atalkar**

has satisfactorily completed the project work towards the **Intel Unnati Industrial Training 2024** in discipline on the topic entitled **"Knowledge Representation and Insights Generation from Structured Datasets",** during the academic year 2024-2025 under my supervision. The student was trained by Intel experts.



Date:                          Dr. R. R. Karwa                          Prof A. U. Chaudhari

                               **Mentor**                               **Mentor**

# ACKNOWLEDGEMENT

The successful completion of this project would not have been possible without the support and guidance of numerous individuals and organizations. We extend our heartfelt gratitude to all those who contributed to this endeavor.
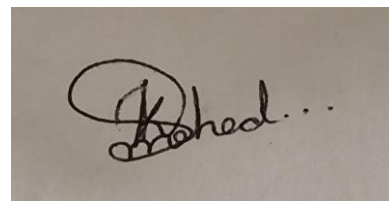
First and foremost, we are deeply thankful to **Intel Corporation** for providing us with the opportunity to undertake this project as part of the **Intel Unnati Industrial Training program**. We owe a debt of gratitude to the organizers and mentors of the **Intel Unnati Industrial Training Program 2024**. Their structured approach and industry-relevant training have significantly enhanced our understanding and project management skills.

Our sincere appreciation goes to **Prof. Ram Meghe Institute of Technology and Research, Badnera, Amravati**, for providing the essential infrastructure and academic environment that made this project possible. Our heartfelt thanks go to our faculty mentors, **Dr. R. R. Karwa** and **Prof. A. U. Chaudhari**, for their constant encouragement, constructive criticism, and steadfast support during the project's development. We would also like to express our special thanks to **Dr. Y. A. Dhumale** for their invaluable guidance, which was fundamental to the project's success.
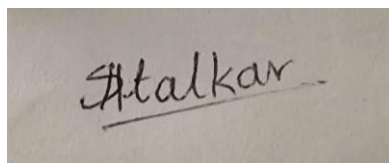
We extend our sincere gratitude to our friend Pranav Bhatkar and **Yash Atalkar** for his their contributions to our project.

This project has been an incredible learning experience, and we are profoundly grateful for the opportunity to have worked on it. The knowledge and skills we've gained will undoubtedly serve us well in our future endeavors.

**Name of student(s):**



**Mr. Krushna Mohod**



**Mr. Sudhanshu Atalkar**

# Table of Contents

# Chapter 1

# Introduction

In today's big data era, organizations from all industries produce vast volumes of data on a regular basis. When this data is processed and analyzed well, insightful conclusions can be drawn that greatly improve decision-making. Finding useful ideas and effectively conveying this knowledge are the difficult parts. Using the dataset of loans and credits from the International Bank for Reconstruction and Development (IBRD), this study develops an AI-based solution to handle this difficulty. The IBRD dataset is a great resource for knowledge representation and insight production since it offers a thorough perspective of a variety of loan and credit activities.

This project's main goal is to develop a solid solution that can process structured datasets, spot trends, and produce insightful analysis. This research intends to identify patterns and linkages in loan performance, repayment behavior, and other important financial variables by utilizing the IBRD dataset. Preparing the data, utilizing machine learning algorithms to identify patterns, creating actionable insights, and representing knowledge through visualizations are all components of the solution. This project opens the door for better informed decision-making procedures in the financial industry by improving our comprehension of financial data and showcasing the usefulness of AI in data analytics.

## 1.1 Problem Statement(PS:12)

In the era of big data and digital transformation, the ability to extract meaningful insights from structured datasets has become a cornerstone of data-driven decision-making across various sectors. This project, "Knowledge Representation and Insights Generation from Structured Datasets," addresses the critical challenge of transforming raw structured data into actionable knowledge. The significance of this endeavor is underscored by the exponential growth of data generation and collection in fields ranging from business and healthcare to scientific research and public policy.

Structured datasets, characterized by their organized format and predefined data models, offer a rich source of information. However, the sheer volume and complexity of these datasets often obscure the valuable patterns and relationships they contain. The challenge lies not merely in storing or accessing this data, but in developing sophisticated methods to represent, analyze,

and interpret it in ways that reveal its full potential for informing strategic decisions and driving innovation.

## 1.2. Motivation

The Structured International Bank for Reconstruction and Development (IBRD) dataset presents a unique opportunity for knowledge representation and insight generation due to several compelling reasons:

- The IBRD dataset likely contains a wealth of information on various aspects of international development, including loan details, project descriptions, borrower characteristics, and economic indicators. By effectively representing this knowledge, we can create a comprehensive understanding of the World Bank's lending activities and their impact.
- Extracting insights from the structured IBRD data can empower data-driven decision making within the World Bank and other development institutions. Analyzing historical trends in lending patterns, project success rates, and economic impacts can inform future strategies for allocating resources and maximizing development effectiveness.
- Knowledge representation and insight generation can enhance transparency and accountability in the development sector. By making the knowledge within the IBRD dataset more accessible and interpretable, stakeholders can gain a clearer picture of the World Bank's activities and their alignment with development goals.
- The structured format of the IBRD data facilitates comparative analysis across different countries, sectors, and lending periods. Identifying patterns and trends through insights generation can inform best practices and guide future development efforts.
- Insight generation can be used to uncover potential risks associated with specific projects or lending approaches. Additionally, it can help identify promising opportunities for development intervention based on historical data and economic indicators.
- Effective knowledge representation can enable knowledge discovery by allowing researchers and analysts to explore hidden patterns and relationships within the data. This can foster collaboration and accelerate progress towards achieving development goals.

The structured IBRD dataset offers a valuable resource for knowledge representation and insight generation. By harnessing the power of these techniques, we can gain a deeper understanding of international development efforts, promote data-driven decision making, and ultimately contribute to a more impactful and sustainable development agenda.

**1.3 Objective**

The primary objective of this project is to develop an AI-driven system capable of efficiently representing knowledge and generating insights from structured datasets. This system aims to go beyond traditional data analysis methods by incorporating advanced techniques in knowledge representation, pattern recognition, and insight generation. The envisioned solution should be adaptable to various types of structured data, capable of identifying complex patterns, and able to present findings in a clear, actionable format.

- To develop robust methods to clean, standardize, and prepare diverse structured datasets for analysis. This step is crucial for ensuring data quality and compatibility across different sources and formats.

- To explore and implement advanced knowledge representation techniques, such as semantic networks, ontologies, or graph-based models. These frameworks aim to capture not just the data itself, but also the relationships and context that give it meaning.

- To incorporate machine learning and data mining algorithms to identify significant patterns, trends, and anomalies within the datasets. This includes supervised and unsupervised learning techniques, as well as machine learning models for complex pattern recognition.

- To develop algorithms that can translate identified patterns into meaningful insights. This involves not only detecting correlations but also inferring causality, predicting future trends, and generating hypotheses for further investigation.

- To ensure that the insights generated are not only accurate but also interpretable and explainable. This is crucial for building trust in the system's outputs and facilitating their integration into decision-making processes.

- To design the system to efficiently handle large-scale datasets and complex analyses, potentially leveraging distributed computing and parallel processing techniques.

- To create intuitive interfaces and visualization tools that allow non-technical users to interact with the system, explore insights, and customize analyses according to their specific needs.

This project aims to make a significant contribution to the field of data analytics and knowledge representation. By developing an AI system capable of efficiently extracting and representing knowledge from structured datasets, it seeks to enhance decision-making processes across various domains. The potential impact of this research extends from improving business strategies and healthcare outcomes to advancing scientific research and informing public policy. As we continue

to navigate an increasingly data-rich world, the ability to transform raw data into actionable insights becomes not just an advantage, but a necessity for progress and innovation.

**Chapter 2**

**Literature Review**

The quest for effective knowledge representation and insightful analysis boasts a rich history. This section delves into the key contributions of past researchers, providing a foundation for understanding the current state-of-the-art.

This section examines how researchers have tackled the challenge of capturing and structuring knowledge in a machine-understandable format. It explores early symbolic logic approaches, the development of ontologies, the emergence of probabilistic graphical models, and the recent advancements in deep learning architectures for knowledge representation. This section investigates various techniques for extracting meaningful insights from structured knowledge. It explores statistical methods like hypothesis testing and clustering algorithms. It delves into how natural language processing (NLP) allows machines to understand the semantics of text data and extract relationships between entities. Finally, it examines advanced machine learning algorithms, including decision trees, support vector machines, and deep neural networks, which are increasingly used for complex data analysis and insight generation. Following is past researchers work:

Famili et al. [1] offered an extensive examination of data preprocessing, emphasizing the challenges posed by real-world data. These challenges necessitated meticulous comprehension and resolution prior to initiating any data analysis. The study elaborated on two principal objectives of data preprocessing: addressing data issues and preparing for subsequent analysis.

Habib and Okayli [2] conducted a thorough investigation into the influence of different data preprocessing methods on the predictive accuracy of machine learning models for estimating the compressive strength of concrete. The study involved the development of ten regression models across nine unique preprocessing scenarios, which encompassed techniques such as normalization, standardization, principal component analysis (PCA), and incorporation of polynomial features. The analysis utilized a comprehensive dataset containing both normal and high-strength concrete performances.

Rao et al. [3] extensively discussed mechanisms such as data missingness, strategies for handling missing data, encoding categorical features, discretization, outlier detection, and feature scaling to enhance the effectiveness of predictive modeling. The study presented comprehensive arguments outlining the benefits and drawbacks of prevalent data preprocessing techniques across different

distributions of variables within house price datasets.

Singh & Singh [4] offered a thorough analysis of machine learning (ML) visualization techniques, resources, and methodologies. It explored how data visualization fit into the visual analytics methodology for both customers and researchers. The study presented an analysis of various types of charts available for data visualization and discussed guidelines for their use, considering the unique circumstances of each use case. Additionally, it examined some of the latest and most innovative visualization tools at that time. The research delved into visualization challenges across different domains, acknowledging the unique characteristics of each ML model and its visualization strategies.

Waskom et al. [5] described seaborn as a library designed for creating statistical graphics in Python. It provided a high-level interface to matplotlib and tightly integrated with pandas data structures. Functions within the seaborn library offered a dataset-oriented API that facilitated the translation of data-related questions into graphical representations capable of answering them. Given a dataset and plot specifications, seaborn automatically mapped data values to visual attributes like color, size, or style. It internally computed statistical transformations and adorned the plot with descriptive axis labels and a legend.

Lu et al. [6] proposed a novel approach named multi-modal knowledge representation learning (MMKRL) to leverage multi-source knowledge including structured, textual, and visual data. Instead of straightforwardly integrating multi-modal knowledge into a unified space with structured knowledge, the study introduced a component alignment scheme. This scheme was combined with translation methods to achieve multi-modal knowledge representation learning (KRL).

Müllner [7] proposed a new approach named KAVA to address a dataset resulting from a clinical trial on a medication for treating the eye disease Uveitis, aiming to facilitate exploration and analysis of the dataset. The approach was designed and developed through a user-centered design process that involved a domain expert, combined with problem-driven visualization research. The final approach was validated through a qualitative task-oriented user study involving five visualization experts. The results indicated that the approach effectively supported both analysis and exploration of the dataset.

Ng et al. [8] proposed a Python-based tool called MUFASA (Medical Information Data Uses For AI Semantic Analysis). This tool employs cutting-edge techniques such as the Sentence Transformer library, clustering algorithms, and visualization methods. MUFASA utilizes

unsolicited medical information (MI) data alongside AI technology to improve efficiency and deliver actionable insights in medical affairs, specifically targeting healthcare providers (HCPs).

Ma & Sun [9] offered an overview of typical machine learning tasks and methodologies, contrasting them with traditional statistical and econometric approaches commonly utilized by marketing researchers. The authors contended that machine learning techniques excel in handling large-scale and unstructured datasets, offering flexible model architectures that yield robust predictive capabilities. However, they acknowledged potential drawbacks such as reduced model transparency and interpretability. The paper examined notable industry trends driven by AI and reviewed emerging academic literature in marketing that employs machine learning methods. Lastly, Ma & Sun proposed a unified conceptual framework and outlined a comprehensive research agenda.

Chatzimparmpas [10] outlined the background of the topic and introduced a categorization for visualization techniques aimed at enhancing trust in interactive machine learning (ML). They also discussed insights and future research directions. Their contribution includes categorizing trust across various aspects of interactive ML, which improves upon previous research. The study analyzed results from different angles: providing statistical summaries, summarizing key findings, conducting topic analyses, and exploring datasets used in individual papers, all facilitated through an interactive web-based survey tool. This survey aims to benefit visualization researchers interested in improving trust in ML models, as well as researchers and practitioners from various fields seeking effective visualization techniques to confidently interpret and convey meaningful insights from their data. Table 2 summarized past researchers' work.

Table 2 Summarization of Past Works

| Author | Methodology | Identified Gap |
|---|---|---|
| Famili et al. [1] | Examined data preprocessing challenges and objectives. | Need for more practical guidelines on addressing specific real-world data issues during preprocessing. |
| Habib and Okayli [2] | Investigated various data preprocessing methods' impact on predictive accuracy in concrete strength estimation. | Lack of standardized benchmarks for evaluating and comparing multiple preprocessing techniques in concrete strength prediction models. |
| Rao et al. [3] | Discussed mechanisms for enhancing predictive modeling through data preprocessing techniques. | Need for practical guidance on selecting and applying preprocessing techniques based on variable distributions in house price datasets. |
| Singh & Singh [4] | Analyzed ML visualization | Need for practical strategies to |

| | techniques and challenges across different domains. | improve ML model visualization clarity and effectiveness across diverse domains. |
|---|---|---|
| Waskom et al. [5] | Described seaborn's role in creating statistical graphics and integrating with pandas. | Lack of comprehensive evaluation on seaborn's integration with complex data structures and its practical impact on visualization efficiency. |
| Lu et al. [6] | Introduced MMKRL for multi-modal knowledge representation learning. | Need for practical assessment of MMKRL's scalability and effectiveness in integrating diverse data modalities. |
| Müllner [7] | Proposed KAVA for exploring and analyzing clinical trial datasets. | Need for practical validation of KAVA across diverse clinical trial datasets to assess its applicability and reliability. |
| Ng et al. [8] | Developed MUFASA for AI semantic analysis of medical information. | Need for practical evaluation of MUFASA's performance in handling real-world medical datasets and supporting healthcare decision-making. |
| Ma & Sun [9] | Compared machine learning methods with traditional approaches in marketing research. | Need for practical methods to enhance interpretability of ML models in marketing research and effectively apply flexible model architectures. |
| Chatzimparmpas [10] | Introduced categorization for visualization techniques in interactive ML to enhance trust. | Need for practical refinement and application of visualization techniques to improve trust and interpretability in interactive ML models. |

By examining past research efforts, the review aims to gain valuable insights into the evolution of these fields, identify successful strategies, and understand the limitations of existing approaches. This knowledge is instrumental in evaluating current research trends and exploring potential avenues for future advancements.

# Chapter 3

## Dataset Description

The IBRD (International Bank for Reconstruction and Development) Statement of Loans - Historical Data is a comprehensive dataset detailing the historical loan activities of the IBRD. As part of the World Bank Group, the IBRD plays a pivotal role in providing financial and technical assistance to developing countries worldwide. This dataset is crucial for understanding the scope and impact of international finance and development efforts spearheaded by the IBRD over time.

The dataset contains detailed records of IBRD loans, including loan amounts, interest rates, borrower countries, and other relevant financial data. It provides insights into the financial commitments and support provided by the IBRD to various countries and projects aimed at fostering development.

The data is presented in a tabular format, covering a wide time range that encapsulates the historical lending activities of the IBRD. Each record in the dataset corresponds to a specific loan and includes multiple attributes describing the loan's details.

The dataset includes extensive information on IBRD loans, offering a detailed view of international financial assistance over the years.

- Historical Trends: It allows for the analysis of historical lending trends, including changes in loan amounts, interest rates, and regional distribution.
- Diverse Loan Categories: The dataset encompasses various types of loans, reflecting the IBRD's multifaceted approach to supporting development projects across different sectors and regions.

| Data Element | Definition |
|---|---|
| **End of Period** | End of Period Date represents the date as of which balances are shown in the report. |
| **Loan Number / Credit Number** | For IBRD loans and IDA credits or grants a loan number consists of the organization prefix (IBRD/IDA) and a five-character label that uniquely identifies the loan within the organization. In IDA, all grant labels start with the letter 'H'. |
| **Region** | World Bank Region to which the country and loan belong. Country lending is grouped into regions based on the current World Bank administrative (rather than geographic) region where project implementation takes place. The "Other"" Region is used for loans to the IFC. |
| **Country** | Country to which a loan has been issued. Loans to the IFC are included under the country "World". |
| **Country Code** | Country Code according to the World Bank country list. This might be different from the ISO country code. |

| | |
|---|---|
| **Borrower** | The representative of the borrower to which the Bank loan is made. |
| **Guarantor** | The Guarantor guarantees repayment to the Bank if the borrower does not repay. |

## Data Dictionary for IBRD Statement of Loans and IDA Statement of Credits and Grants

This Data Dictionary provides descriptions of the data elements for the IBRD Statement of Loans and IDA Statement of Credits and Grants published on the World Bank Finances website.

| | |
|---|---|
| **Guarantor Country Code** | Country Code of the Guarantor according to the World Bank country list. This might be different from the ISO country code. |
| **Loan Type** | A     type of loan/loan instrument for which distinctive accounting and/or other actions need to be performed.<br><br>Loan Type Descriptions:<br>B     Loan –Co-financing lending product that includes Contingency and Regular loans and guarantees<br> Pool loan- Currency Pooled Loans<br> FSL - Fixed Spread Loans (includes both fixed spread loans and IBRD flexible loans that have either fixed spread or variable     spread terms)<br>  IFC loan – single currency loans to the IFC<br>  Non  Pool -  Original IBRD lending product, prior to currency pooled loans.<br>  Sngl crncy -  - Single Currency Loans<br>  SCP USD - Single Currency Pooled Loans - USD<br>  SCP - DEM - Single Currency Pooled Loans - EUR<br>  SCP  JPY - Single Currency Pooled Loans – JPY |
| **Loan Status  /  Credit Status** | Status of the loan.<br> Loan Status descriptions:<br>  APPROVED - Loan has been approved by the Bank<br>  SIGNED - Loan has been signed by both parties<br>  EFFECTIVE - Loan has been made effective in accordance with the terms of the legal agreement    DISBURSING - Loan is disbursing<br>  DISBURSED - Loan has no undisbursed balance<br>  REPAID - Loan has been fully repaid<br><br>  CANCELLED - Entire loan principal has been cancelled<br>  TERMINATED -  Unsigned loan that has been cancelled in full |

| | |
|---|---|
| **Currency    of Commitment** | The currency in which a borrower's loan, credit or grant is denominated. |
| **Project Name** | Short descriptive project name. |
| **Project ID** | A Bank project is referenced by a project ID (Pxxxxxxx).  More than one loan, credit, or grant may be associated with one Project ID. |

| | |
|---|---|
| **Original Principal Amount** | The original US dollar amount of the loan that is committed and approved. |
| **Cancelled Amount** | The portion of the undisbursed balance which has been cancelled (i.e. no longer available for future disbursement). Cancellations include terminations (where approved loan agreements were never signed). |
| **Undisbursed Amount** | The amount of a loan commitment that is still available to be drawn down. These currency amounts have been converted to US dollars at the exchange rates applicable at the end of period date. |
| **Disbursed Amount** | The amount that has been disbursed from a loan commitment in equivalent US dollars, calculated at the exchange rate on the value date of the individual disbursements. |
| **Repaid to IBRD** | Total principal amounts paid or prepaid to IBRD in US dollars, calculated at the exchange rate on the value date of the individual repayments. |
| **Repaid to IDA** | Total principal amounts paid or prepaid to IDA in US dollars, calculated at the exchange rate on the value date of the individual repayments. Repaid to IDA amounts include amounts written off under the Multilateral Debt Relief Initiative (MDRI). |
| **Due to IBRD** | Where the exchange adjustment is shown separately, this is the amount disbursed and outstanding expressed as a stock of debt in historical US Dollars. Where the exchange adjustment is not shown separately, this is the amount due and outstanding as of the End of Period date. |
| **Due to IDA** | Amount due and outstanding as of the End of Period date. |
| **Exchange Adjustment** | The increase (decrease) in value of disbursed and outstanding amount due to exchange rate fluctuations. This amount added to "Due to IBRD" yields "Borrower's Obligation"; includes exchange adjustments on the amounts Due to 3rd parties. |
| **Borrower's Obligation** | The Borrower Obligation is the outstanding balance for the loan as of the end of period date in US dollars equivalent.<br><br>The Borrower's Obligation includes the amounts outstanding Due to 3rd parties. |
| **Sold 3rd Party** | Portion of loan sold to a third party. |
| **Repaid 3rd Party** | Amount repaid to a third party. |
| **Due 3rd Party** | Amount due to a third party. |
| **Loans Held** | The sum of the disbursed and outstanding amounts (net of repayments, i.e. Due to IBRD/IDA) plus undisbursed available amounts. |
| **First Repayment Date** | The date on which principal repayment starts. |
| **Last Repayment Date** | The date specified in the loan/credit agreement (amended for any partial prepayments) on which the last principal installment must be repaid by the Borrower. |
| **Agreement Signing Date** | The date the borrower and the Bank sign the loan agreement. |
| **Board Approval Date** | The date the World Bank approves the loan. |

| Effective Date | The date on which a legal agreement becomes effective, or is expected to become effective. |
|---|---|
| Close Date | The date specified in the legal agreement (or extension) after which the Bank may, by notice to the borrower, terminate the right to make withdrawals from the loan account. |
| Last Disbursement Date | The date on which the last disbursement was made (prior to the end of period date). |

**Column Datatype**

- Loan Number: (String) Unique identifier for each loan.
- Region: (String) Geographic region of the borrower.
- Country Code: (String) ISO code of the borrower country.
- Country: (String) Name of the borrower country.
- Borrower: (String) Entity receiving the loan.
- Guarantor Country Code: (String) ISO code of the guarantor country.
- Guarantor: (String) Name of the guarantor.
- Loan Type: (String) Type of loan provided.
- Loan Status: (String) Current status of the loan (e.g., active, closed).
- Project ID: (String) Identifier for the associated project.
- Project Name: (String) Name of the associated project.
- Interest Rate: (Numeric) The interest rate applied to the loan.
- Original Principal Amount: (Numeric) The initial amount of the loan.
- Cancelled Amount: (Numeric) Amount of the loan that was cancelled.
- Undisbursed Amount: (Numeric) Amount of the loan that has not been disbursed yet.
- Disbursed Amount: (Numeric) Amount of the loan that has been disbursed.
- Repaid to IBRD: (Numeric) Amount repaid to the IBRD.
- Due to IBRD: (Numeric) Amount still due to the IBRD.
- Exchange Adjustment: (Numeric) Adjustment amount due to exchange rate changes.
- Borrower's Obligation: (Numeric) Total obligation of the borrower.
- Sold 3rd Party: (Numeric) Amount sold to third parties.
- Repaid 3rd Party: (Numeric) Amount repaid to third parties.
- Due 3rd Party: (Numeric) Amount still due to third parties.
- Loans Held: (Numeric) Total amount of loans held.
- End of Period: (Date) End date of the loan period.
- First Repayment Date: (Date) Date of the first repayment.
- Last Repayment Date: (Date) Date of the last repayment.
- Agreement Signing Date: (Date) Date when the loan agreement was signed.
- Board Approval Date: (Date) Date when the loan was approved by the board.
- Effective Date (Most Recent): (Date) Most recent effective date of the loan.
- Closed Date (Most Recent): (Date) Most recent closing date of the loan.
- Currency of Commitment: (String) Currency in which the loan was committed.
- Last Disbursement Date: (Date) Date of the last disbursement.

The dataset is sourced from the World Bank, a reputable international financial institution dedicated to providing financial and technical assistance to developing countries. The World Bank plays a crucial role in collecting, maintaining, and disseminating data related to global development finance.

The dataset can be accessed at [World Bank Open Finances](https://finances.worldbank.org/Loans-and-Credits/IBRD-Statement-Of-Loans-Historical-Data/zucq-nrc3/about_data). The World Bank ensures the accuracy and availability of this data, supporting transparency and research in international finance and development.

**Reason of selection of dataset**

This dataset was chosen for its rich and detailed content, making it ideal for various academic purposes, including: The dataset's complexity, with numerous columns and diverse data types, provides an excellent opportunity to practice and refine data preprocessing skills. The structured format and comprehensive nature of the data make it suitable for representing knowledge in the context of international finance. The historical and detailed nature of the dataset allows for in-depth analysis and the generation of valuable insights into global development finance trends and patterns. The dataset's complexity, including a wide range of numerical, categorical, and date fields, along with its high-quality, well-maintained data, makes it an excellent resource for academic research and analysis.

# Chapter 4

## Methodologies

Knowledge representation and insight generation from structured datasets are pivotal in transforming raw data into actionable intelligence. With the exponential growth of data across various domains, the ability to efficiently process, analyze, and visualize structured data has become essential. The proposed methodology aims to provide a comprehensive framework for handling structured datasets, from initial data preprocessing to advanced visualization techniques and insight generation. By leveraging robust algorithms and state-of-the-art tools, this methodology seeks to uncover hidden patterns, trends, and correlations within the data, thereby facilitating informed decision-making. Figure 4.1 illustrates the component of methodology.
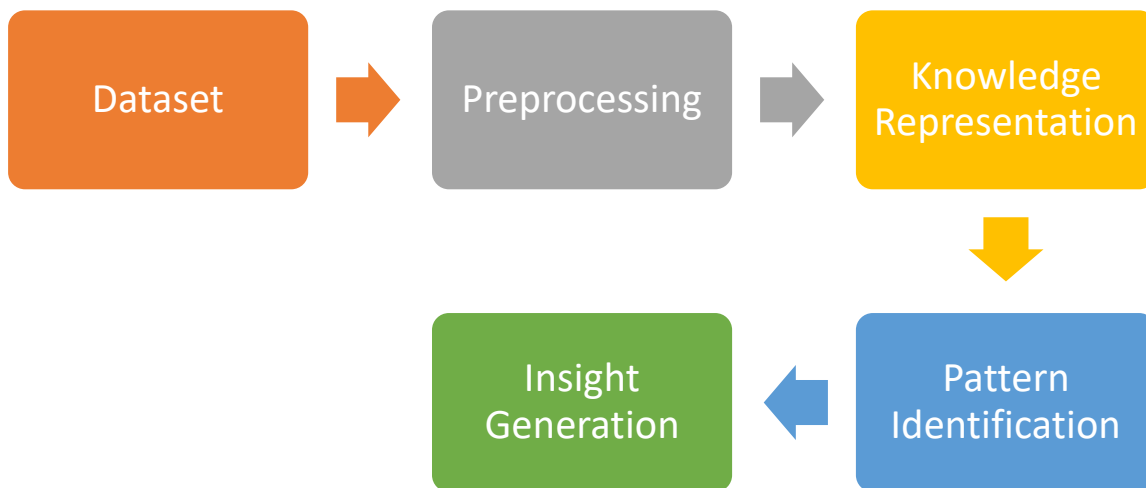


Figure 4.1. Proposed Methodology

## 4.1. Data Preprocessing

The preprocessing stage is a critical component of our project, aimed at transforming raw data into a clean and analyzable format. Effective preprocessing ensures data quality and reliability, thereby enabling accurate and insightful analysis. The main goals of this preprocessing stage include data cleaning, handling missing values, and detecting outliers.

**Data Loading and Initial Analysis:-**
**Method:** We utilized cuDF to load the CSV file containing the IBRD Statement of Loans - Historical Data.

```python
import cudf
df = cudf.read_csv("IBRD_Statement_of_Loans.csv")
```

**Result:** The initial analysis provided insights into the shape of the dataset and the distribution of missing values.

```python
shape = df.shape
missing_values = df.isnull().sum()
missing_percentage = (missing_values / len(df))  100
```

Shape of the DataFrame: (1323796, 33)

```
Column Name                 | Missing Values  | Percentage Missing
----------------------------+-----------------+------
Country Code                | 319             | 0.02 %
Borrower                    | 9122            | 0.69 %
Guarantor Country Code      | 49050           | 3.71 %
Guarantor                   | 75399           | 5.70 %
Interest Rate               | 30736           | 2.32 %
Currency of Commitment      | 1323796         | 100.00%
Project ID                  | 42              | 0.00 %
Project Name                | 159114          | 12.02%
First Repayment Date        | 4077            | 0.31 %
Last Repayment Date         | 3918            | 0.30 %
Agreement Signing Date      | 19244           | 1.45 %
Board Approval Date         | 2               | 0.00 %
Effective Date (Most Recent) | 10234          | 0.77 %
Closed Date (Most Recent)   | 1250            | 0.09 %
Last Disbursement Date      | 540402          | 40.82%
```
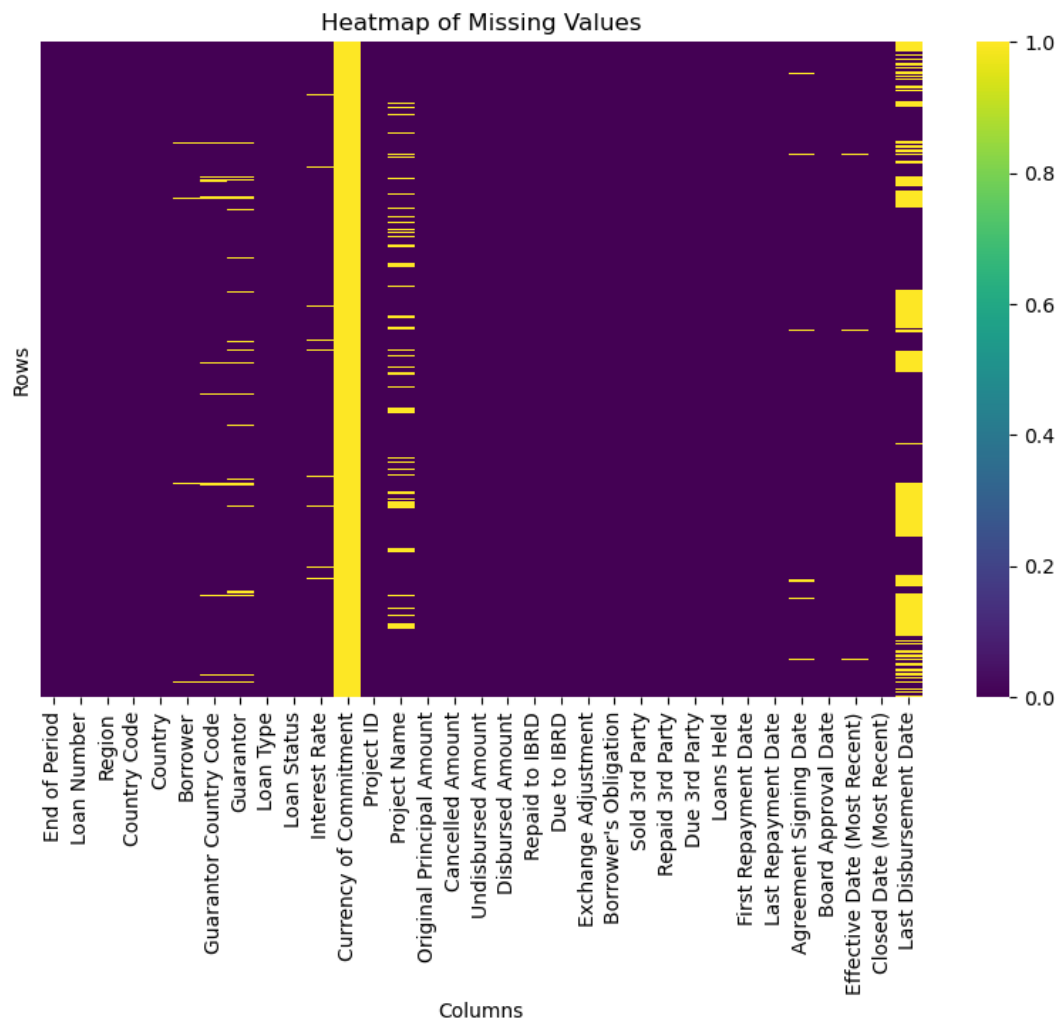
Figure 4.2. Heatmap of missing values

**Reason:** This step is crucial for understanding the dataset's structure and identifying potential issues that need to be addressed in subsequent preprocessing steps.

**Column Categorization**

Method: We employed the GPU-based column classifier to categorize columns into datetime, numeric, categorical, and ID columns.

```python
from src.data_process.GPU_columnclassifier import
categorize_columns_gpu

column_types = categorize_columns_gpu(df)
datetime_columns = column_types['datetime_columns']
numeric_columns = column_types['numeric_columns']
categorical_columns = column_types['categorical_columns']
id_columns = column_types['id_columns']
```

```
```

**Result:** The classifier identified the following types of columns: datetime, numeric, categorical, and ID.

**Datetime Columns:** End of Period First Repayment Date,Last Repayment Date,Agreement Signing Date,Board Approval Date,Effective Date (Most Recent),Closed Date (Most Recent),Last Disbursement Date,

**Numerical Columns:** Original Principal Amount,Cancelled Amount,Undisbursed Amount,Disbursed Amount,Repaid to IBRD,Due to IBRD,Exchange Adjustment,Borrower's Obligation,Sold 3rd Party,Repaid 3rd Party,Due 3rd Party,Loans Held, Interest Rate

**Categorical Columns:** ,Region,Country Code,Country,Borrower,Guarantor Country Code,Guarantor,Loan Type,Loan Status,Currency of Commitment,Project Name

**ID columns:** Loan Number, Project ID

**Reason:** Proper column categorization is essential for applying the correct preprocessing techniques to different types of data.

**Data Cleaning**
**Method:** We used the `drop_col_row` function to drop columns with high percentages of missing values and rows with missing datetime data.

```python
def drop_col_row(df, datetime_columns):
    threshold = 0.4  len(df)
    for column in df.columns:
        if df[column].isna().sum() >= threshold:
            df.drop(column, axis=1, inplace=True)
    df.dropna(subset=datetime_columns, inplace=True)
    df.drop_duplicates(inplace=True)
    df.reset_index(drop=True, inplace=True)
    return df

df = drop_col_row(df, datetime_columns)
```

**Columns Removed:** Currency of Commitment and Last Disbursement Date

**Result:** Columns with more than 40% missing values and rows with missing datetime values were removed.

**Reason:** This step is vital to ensure the remaining data is of high quality and free from excessive missing values that could skew analysis.

**DateTime Processing**

**Method:** The `process_datetime_columns` function was used to standardize datetime formats.

```python
def process_datetime_columns(df, datetime_columns):
    for column in datetime_columns:
        df[column]              =              cudf.to_datetime(df[column],
format='%d/%m/%Y')
    return df


df = process_datetime_columns(df, datetime_columns)
```

```
10/31/2020 12:00:00 AM was formatted to 31/10/2020
```

**Result:** All datetime columns were standardized to a consistent format.

**Reason:** Consistent datetime formatting is crucial for time-based analysis and ensures accurate temporal comparisons.

**Missing Value Imputation**

Method: The `kde_impute` function utilized Gaussian Kernel Density Estimation (KDE) to impute missing numeric values.

```python
from scipy.stats import gaussian_kde


def kde_impute(df, numeric_columns):
    for col in numeric_columns:
        data = df[col].dropna().to_pandas().values
        kde = gaussian_kde(data)
        missing_indices = df[col].isna()
        imputed_values                              =
kde.resample(missing_indices.sum()).flatten()
        df[col][missing_indices] = imputed_values
    return df


df = kde_impute(df, numeric_columns)
```

**Result:** Missing numeric values were imputed using KDE.

**Reason:** KDE provides a more robust method for imputation by considering the distribution of the data, which is preferable over simpler methods like mean or median imputation.

**Categorical and ID Column Handling**

Method: The `fill_empty_values` function was used to fill missing values in categorical and ID columns with 'Unknown'.

```python
def fill_empty_values(df, categorical_columns, id_columns):
    columns_to_fill = categorical_columns + id_columns
    df[columns_to_fill] = df[columns_to_fill].fillna('Unknown')
    return df


df = fill_empty_values(df, categorical_columns, id_columns)
```

**Result:** Missing values in categorical and ID columns were filled.

**Reason:** Handling missing categorical data ensures that all records are complete, which is necessary for categorical analysis and model training.

**Loan Duration Calculation**

**Method:** The `calculate_loan_duration` function calculated the duration of each loan in days, months, and years.

```python
def calculate_loan_duration(df):
    df['loan duration (days)'] = (df['end of period'] -
df['agreement signing date']).dt.days
    df['loan duration (months)'] = ((df['end of period'].dt.year
- df['agreement signing date'].dt.year)  12 +
                                    (df['end of
period'].dt.month - df['agreement signing date'].dt.month))
    df['loan duration (years)'] = df['loan duration (days)'] /
365.25
    return df


df = calculate_loan_duration(df)
```

**Result:** New columns for loan duration were added.

**Reason:** These derived features are crucial for understanding the time dynamics of loans and for subsequent temporal analysis.

**Outlier Handling**

**Method:** The `cap_outliers_iqr_with_zeros_cudf` function was used to cap outliers using the Interquartile Range (IQR) method.

```python
def cap_outliers_iqr_with_zeros_cudf(df, numerical_columns):
    for col in numerical_columns:
        Q1 = df[col].quantile(0.25)
        Q3 = df[col].quantile(0.75)
        IQR = Q3 - Q1
        lower_bound = Q1 - 1.5  IQR
        upper_bound = Q3 + 1.5  IQR
        df[col] = df[col].clip(lower=lower_bound,
upper=upper_bound)
    return df


df = cap_outliers_iqr_with_zeros_cudf(df, numeric_columns)
```

**Result:** Outliers were identified and capped.

**Reason:** Outlier handling is crucial to prevent skewed results and maintain data integrity.

### Categorical Encoding

Method: We used the `encode_categorical_columns` function to handle the encoding of categorical columns. For columns with fewer than 20 unique values, one-hot encoding was applied. For columns with 20 or more unique values, label encoding was used.

```python
from sklearn.preprocessing import LabelEncoder

def encode_categorical_columns(df, categorical_columns):
    encoded_df = df.copy()

    for col in categorical_columns:
        unique_count = df[col].nunique()

        if unique_count < 20:
             One-hot encoding
            dummies = cudf.get_dummies(df[col], prefix=col,
prefix_sep='_', dtype='int8')
            dummies = dummies.rename(columns={c:
f"{col}_{c.split('_')[-1]}" for c in dummies.columns})
```

```python
            encoded_df = cudf.concat([encoded_df, dummies],
axis=1)
        else:
             Label encoding
            le = LabelEncoder()
            encoded_df[f"{col}_encoded"] =
le.fit_transform(df[col].astype('str'))

     Reorder columns to move encoded columns to the end
    original_cols = df.columns.tolist()
    new_cols = [c for c in encoded_df.columns if c not in
original_cols]
    encoded_df = encoded_df[original_cols + new_cols]

    return encoded_df

df = encode_categorical_columns(df, categorical_columns)
```

**Result:** Categorical columns were appropriately encoded, with new encoded columns added to the DataFrame.

**Reason:** Proper encoding of categorical variables is essential for machine learning models to process these features effectively, ensuring that the model can interpret and use the categorical information during training.

**Final Processing and Output**

**Method:** Final steps included dropping unnecessary columns and saving the processed data.
```python
df.to_csv("processed_IBRD_Statement_of_Loans.csv", index=False)
```

**Result:** The final preprocessed dataset was saved in a clean format.

Reason: These final steps ensure that the data is ready for subsequent analysis, providing a clean and well-structured dataset for further processing and insights generation.-

**Memory Management:**

Here's a comprehensive explanation of the Memory Management technique used in the data preprocessing pipeline for the IBRD Statement of Loans - Historical Data project:

Memory Management in Data Preprocessing

**Method:** In our preprocessing pipeline, we employ a straightforward memory management technique using Python's 'del' statement. This simple yet effective approach allows us to explicitly remove variables from memory when they are no longer needed. The basic syntax we use is:

```python

del variable_name

```

For example, after processing a large DataFrame and creating a new, cleaned version, we might use:

```python

del original_df

```

This command removes the reference to the original DataFrame, allowing Python's garbage collector to free up the associated memory.

**Result:** By using the 'del' statement, we achieve immediate memory relief. When a variable is deleted, Python removes the reference to the object it was pointing to. If there are no other references to that object, Python's garbage collector can then reclaim the memory occupied by that object. This is particularly beneficial when working with large data structures like DataFrames in our IBRD loan data analysis.

The immediate effect is a reduction in the Python interpreter's memory usage. This can be especially noticeable when dealing with large datasets or when running the preprocessing pipeline on systems with limited resources.

**Reason:** We chose this straightforward memory management approach for several reasons:

1. **Simplicity:** The 'del' statement is easy to understand and implement, making our code more maintainable.

2. **Immediate effect:** It allows us to free up memory at specific points in our workflow, rather than relying solely on Python's automatic garbage collection.

3. **Large dataset handling:** The IBRD Statement of Loans - Historical Data is a substantial dataset. By actively managing memory, we can process larger chunks of data without running into memory constraints.

4. **Performance**: Freeing up memory can potentially improve overall performance, especially when running multiple preprocessing steps sequentially.

**Implementation:**

In our preprocessing workflow, we typically use 'del' statements:

1. After creating cleaned or transformed versions of DataFrames

2. When intermediate results are no longer needed

3. At the end of major processing steps

For example, we might delete original DataFrames after cleaning, intermediate calculation results after they've been used, or large temporary variables created during complex transformations.

**Considerations:**

While using 'del' statements offers memory benefits, we've been mindful of the following considerations:

1. Code readability: We ensure that the use of 'del' doesn't obscure the main logic of our preprocessing steps.

2. Timing: We're careful to delete variables only when we're certain they won't be needed again.

3. Limitations: 'del' only removes the reference, not the object itself. If other references exist, the memory won't be freed immediately.

4. Debugging: Excessive use of 'del' can make debugging more challenging, so we use it judiciously.

By incorporating this simple memory management technique, we've been able to enhance the efficiency of our preprocessing pipeline for the IBRD Statement of Loans - Historical Data. It allows us to handle the large dataset more effectively, maximizing the use of available system resources without significantly complicating our code structure. This approach contributes to the overall robustness and scalability of our data analysis process, ensuring that we can handle the extensive loan data efficiently throughout our preprocessing stages.

## 4.2. Knowledge Representation

Data visualization plays a pivotal role in our data analysis project focusing on the IBRD Statement of Loans - Historical Data. It transforms raw data into a visual context, making it easier to identify patterns, trends, and insights that might be missed in text-based data. The main goals of our visualization process are to enhance pattern identification and insight generation, ultimately aiding in better decision-making and understanding of the data.

**Visualization Framework**
For creating visualizations, we utilize a combination of `matplotlib`, `seaborn`, and `pandas`. These libraries offer powerful tools for data manipulation and visual representation:

- **Matplotlib:** A comprehensive library for creating static, animated, and interactive visualizations in Python.
- **Seaborn:** Built on top of `matplotlib`, it provides a high-level interface for drawing attractive statistical graphics.
- **Pandas:** Used for data manipulation and analysis, making it easy to preprocess data before visualization.

We use the `Agg` backend for `matplotlib`, which is particularly useful for environments where a display server is not available, such as in automated scripts or cloud-based environments. This ensures that our visualizations can be generated and saved as image files without requiring a graphical display.

**Plot Types and Their Purposes**

a) **Correlation Matrix**
- Method: `sns.heatmap` with `pandas.DataFrame.corr()`
- Insights: Shows the strength and direction of relationships between numerical variables.
- Usefulness: Ideal for identifying potential predictors in regression analysis or finding highly correlated variables that might indicate redundancy.
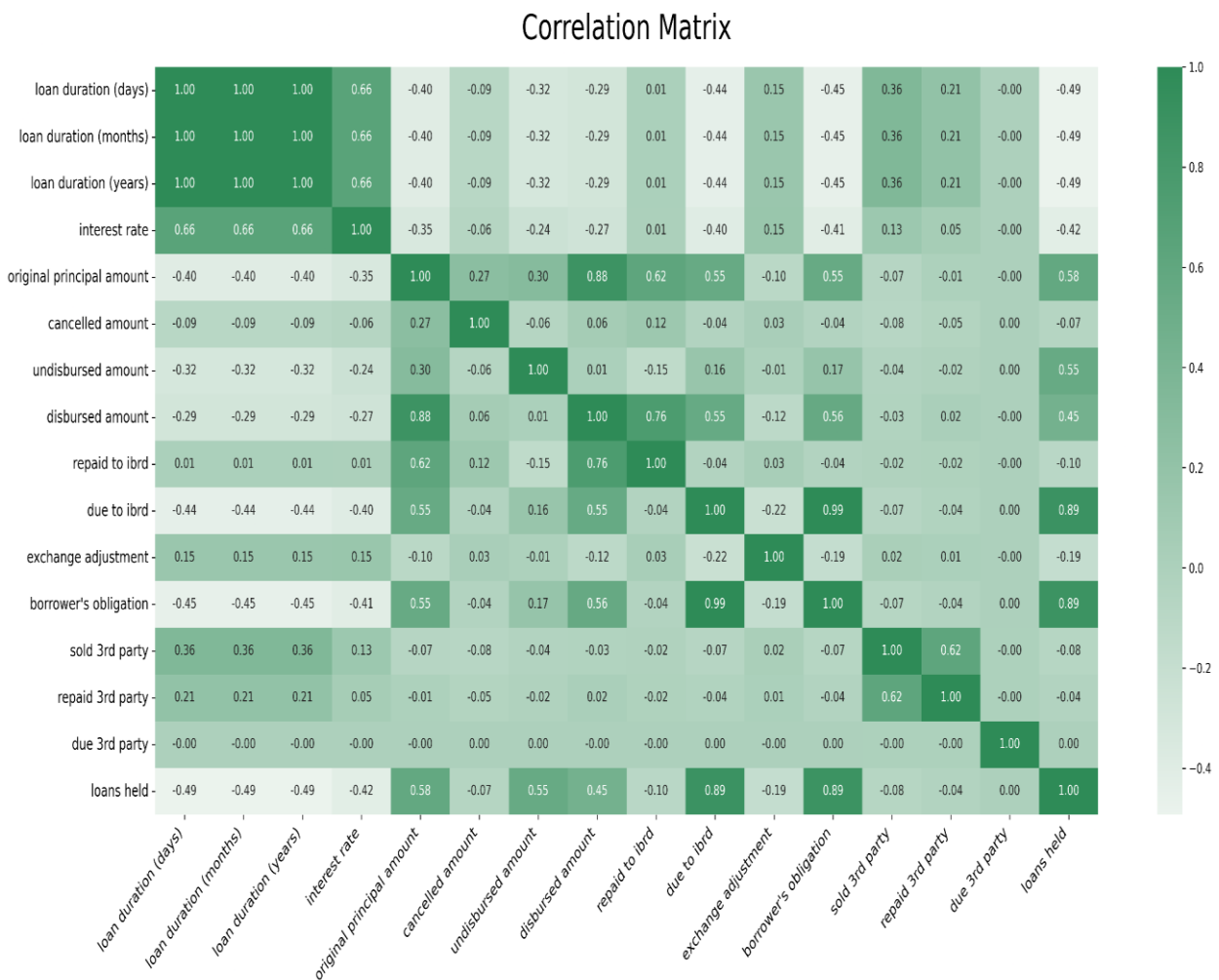


Figure 4.2.1 Correlation between all columns

b) **Pie Chart**

- Method: `plt.pie` with `pandas.Series.value_counts()`
- Insights: Displays the relative proportions of categories within a single variable.
- Usefulness: Useful for showing the composition of a dataset, such as the distribution of loan types or regions.
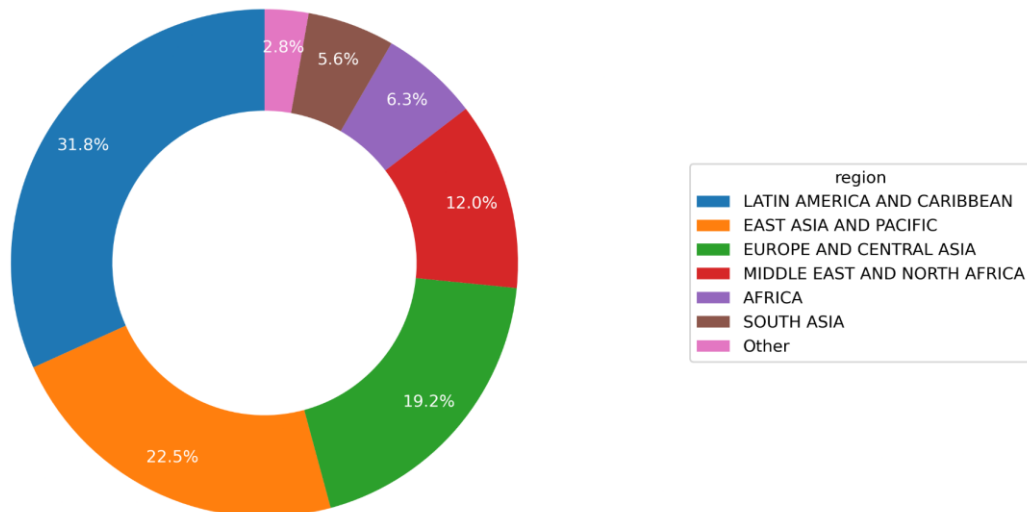


Figure 4.2.2  Region count Piechart

 c) **Box Plot**

- Method: `sns.boxplot`
- Insights: Visualizes the distribution, central tendency, and variability of a dataset across different categories.
- Usefulness: Effective for detecting outliers and comparing distributions between groups.
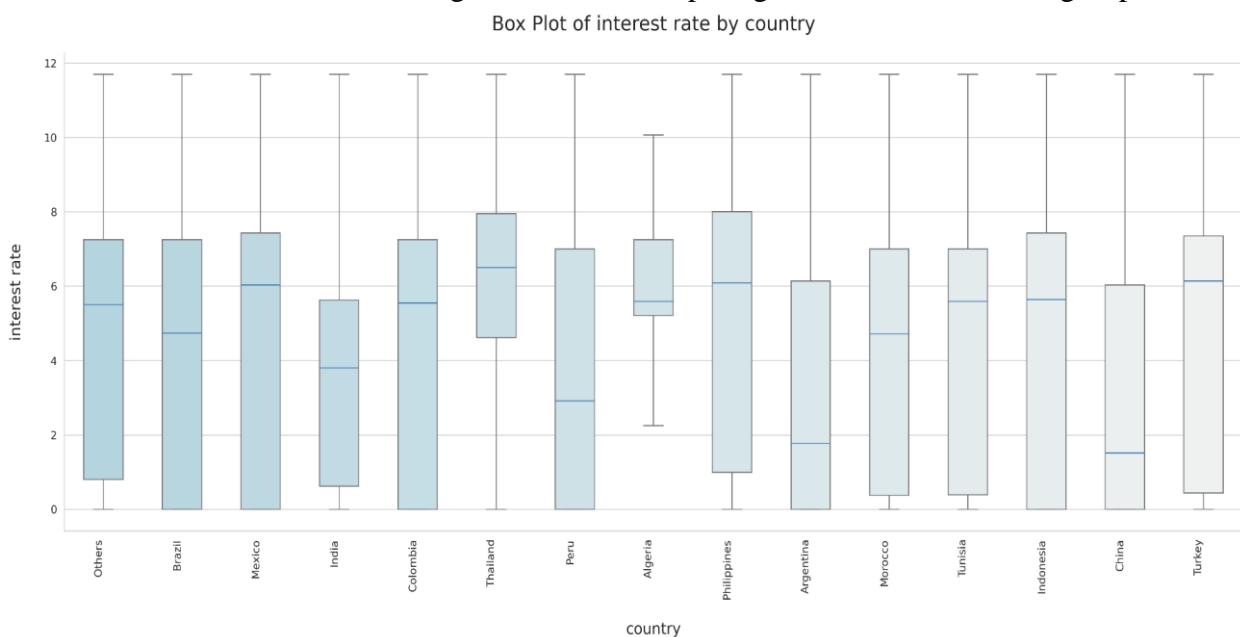


Figure 4.2.3 Box plot of interest rate by country

### d) Line Plot
- Method: `sns.lineplot`
- Insights: Tracks changes over time or ordered categories.
- Usefulness: Best for time series data to show trends and patterns over a period.



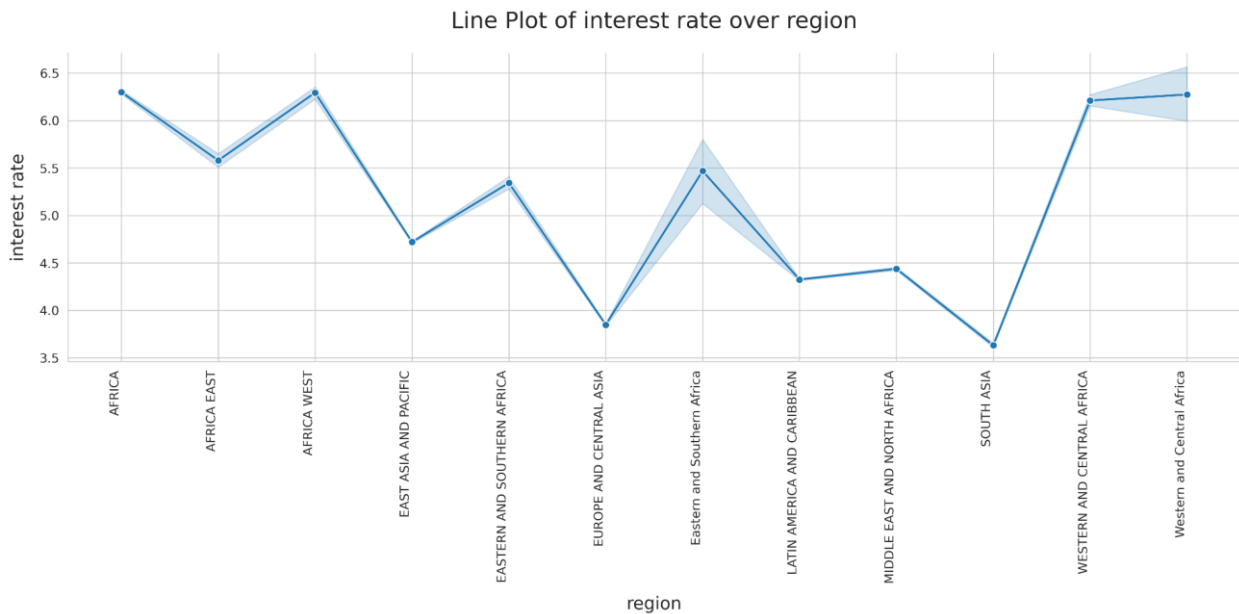Figure 4.2.4 Line plot of interest rate over region

### e) Histogram
- Method: `plt.bar` for categorical data and `plt.hist` for numerical data
- Insights: Illustrates the frequency distribution of a single variable.
- Usefulness: Helps in understanding the distribution shape, central tendency, and variability of data.
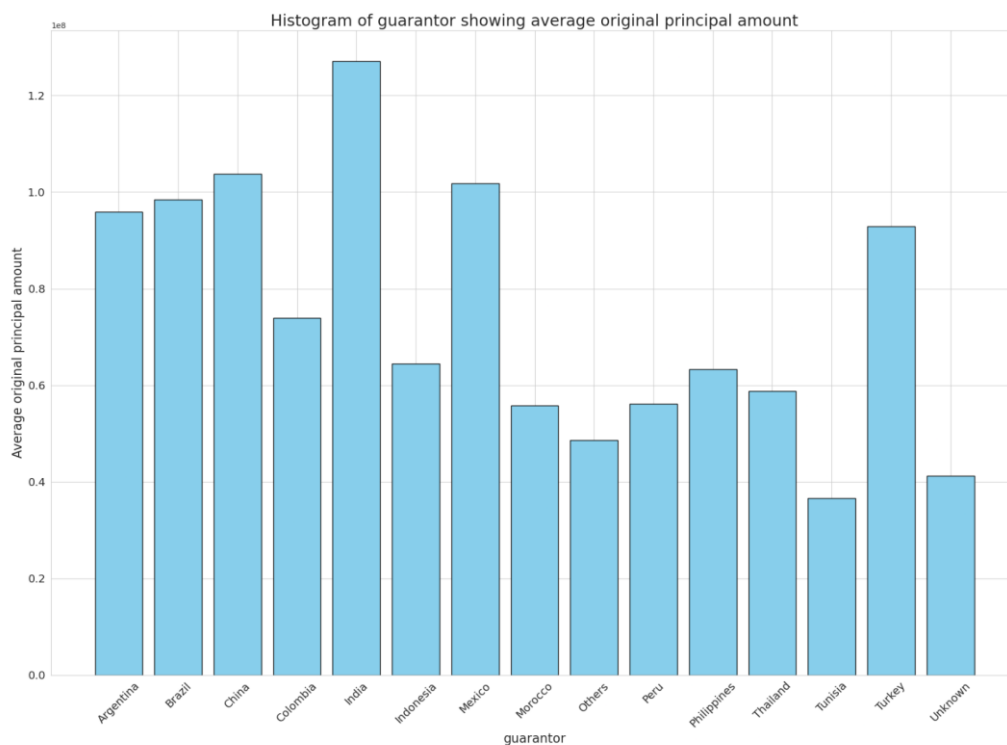


Figure 4.2.5 histogram of guarantor showing average original principal amount

**f) Scatterplot**

- Method: `sns.scatterplot`
- Insights: Displays the relationship between two numerical variables.
- Usefulness: Excellent for identifying correlations, clusters, and trends between variables.
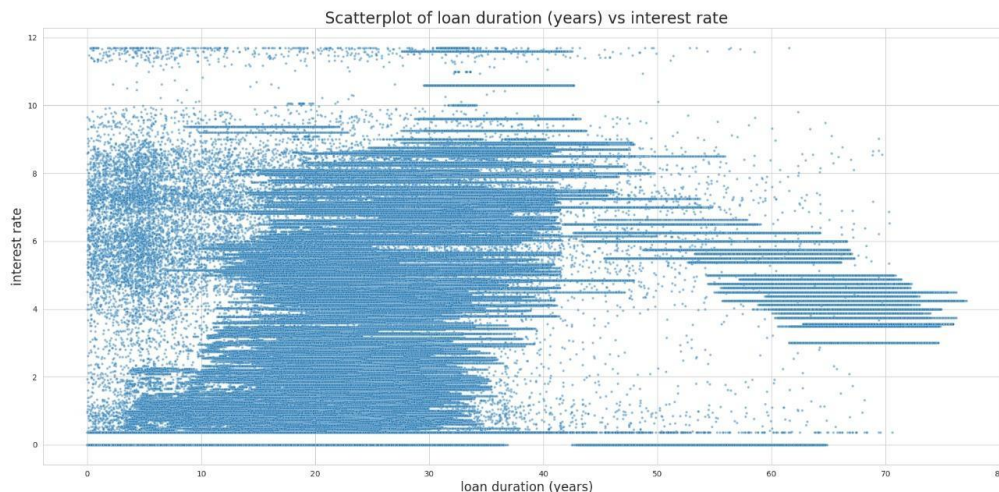


Figure 4.2.6 Scatterplot of loan duration(Years) vs interest rate

**Data Preprocessing for Visualization**

The `process_x_axis` function handles categorical data by limiting the number of categories to a manageable number, grouping less frequent categories into an "Others" category. This preprocessing enhances the readability and interpretability of plots, ensuring that visualizations remain clear and informative without being overwhelmed by too many categories.

**Customization and Styling**

We use various customization and styling options to improve the effectiveness of our visualizations:

- Color Palettes: `seaborn` palettes like `light_palette` are used to maintain a consistent and visually appealing color scheme.
- Figure Sizes: Adjusting `figsize` ensures that plots are appropriately scaled for readability and presentation.
- Styling Elements: Titles, labels, and legends are carefully styled to enhance clarity and aesthetics.

These choices contribute to creating professional, easy-to-understand visualizations that effectively communicate the underlying data insights.

**Handling Large Datasets**
To manage large datasets, we employ techniques such as limiting categories in pie charts and aggregating less frequent categories into an "Others" category. This approach maintains clarity and performance, preventing visual clutter and ensuring that key insights are not lost in the noise.

**Plot Saving and Organization**
The `save_plot` function organizes output by saving plots into a structured directory. This function ensures that visualizations are systematically stored in a `visuals_processed/visual_images` directory, making it easy to retrieve and reference specific plots.

**Dynamic Plot Generation**
The `generate_plot` function allows for flexible plot creation based on user input. This dynamic approach enables interactive analysis, where users can generate different types of plots by specifying the plot type and relevant axes, facilitating a more tailored data exploration experience.

**Error Handling and Edge Cases**
Special case management and error handling are incorporated to ensure robust visualization:
- Invalid Data Types: Checks are in place to ensure that appropriate data types are used for specific plots (e.g., numerical data for scatterplots).
- Missing Data: Techniques like grouping infrequent categories handle edge cases, maintaining the integrity of visualizations.

These considerations are crucial for creating reliable and accurate visualizations, even in the presence of data inconsistencies.
Our visualization methodology leverages powerful libraries and robust preprocessing to create insightful and accessible visualizations. These visualizations are integral to identifying patterns and generating insights from the IBRD Statement of Loans - Historical Data, ultimately supporting better data-driven decisions and deeper understanding of the data.

## 4.3. Pattern Detection

The data preprocessing step is essential for ensuring that the dataset is in the right format for analysis and modeling. In the provided code, the preprocessing involves encoding categorical variables using LabelEncoder from scikit-learn.

```
# Preprocess data
categorical_columns = df.select_dtypes(include=['object']).columns
label_encoders = {}

for col in categorical_columns:
    le = LabelEncoder()
    df[col] = le.fit_transform(df[col].astype(str))
    label_encoders[col] = le
```

This snippet selects columns with categorical data types (object) and applies label encoding,

converting categorical values into numerical values. This preprocessing step is crucial for machine learning models that require numerical input.

**Machine Learning Models Used**

The provided code loads three pre-trained Random Forest models to generate insights on loan status, interest rate, and disbursed amount.

**Code Snippet: Loading Models**

```
# Load models
loan_status_model = load('loan status_rf_model.joblib')
interest_rate_model = load('interest rate_rf_model.joblib')
disbursed_amount_model = load('disbursed amount_rf_model.joblib')
```

The models are loaded using joblib, which is efficient for handling large numpy arrays and other data structures.

**Key Insights Generated**

The code generates insights related to regional distribution, loan status, interest rate, and disbursed amount. Each insight involves predictions using the respective models and visualizations to convey the insights effectively.

**Regional Distribution Insight** This insight calculates the total disbursed amount for each region and visualizes it using a bar plot.

**Code Snippet: Regional Distribution Insight**

```python
def regional_distribution_insight(df, label_encoders,
output_file_path):
    regional_distribution = df.groupby('region')['disbursed
amount'].sum().sort_values(ascending=False)

    # Plotting
    fig, ax = plt.subplots(figsize=(12, 6))
    regional_distribution.plot(kind='bar', ax=ax)
    ax.set_title('Regional Distribution of Loans (Disbursed
Amount)')
    ax.set_xlabel('Region')
    ax.set_ylabel('Total Disbursed Amount')
```
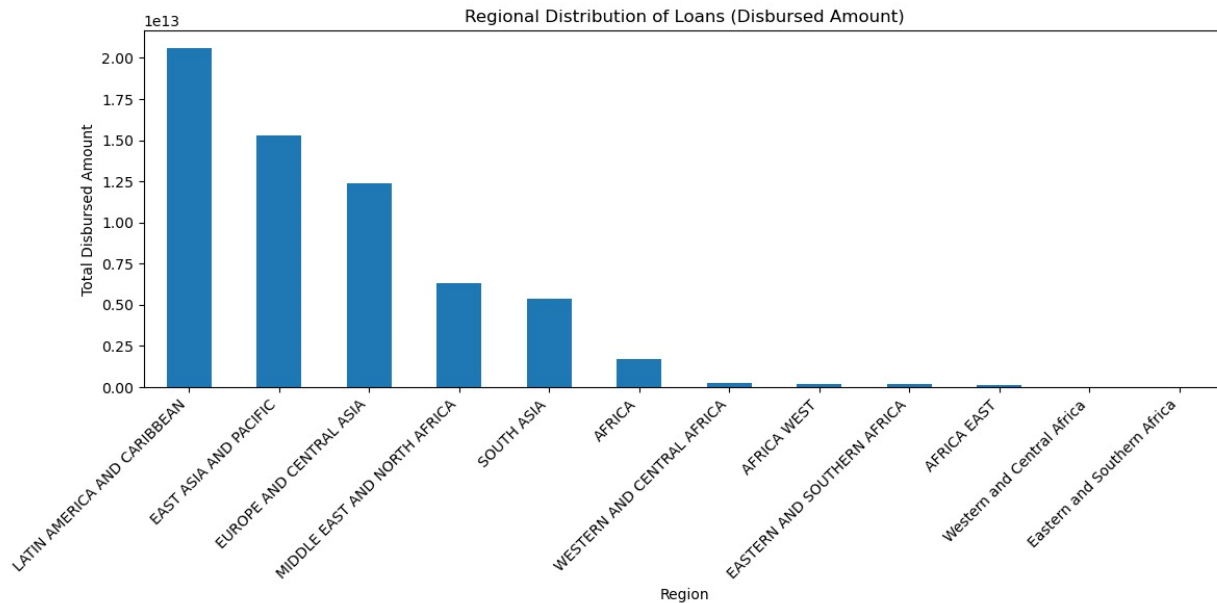
Figure 4.3.1 Regional Distribution of loans

1. The group by method aggregates the disbursed amounts by region, and the resulting data is plotted to show the distribution of loans across regions.

**Loan Status Insight** This insight uses the loan_status_model to predict loan statuses and visualizes the distribution of these statuses using a pie chart.
**Code Snippet: Loan Status Insight**

```
def loan_status_insight(df, model, label_encoders, output_file_path):
X = df[model.feature_names_in_]
predicted_status = model.predict(X)
predicted_status = label_encoders['loan
status'].inverse_transform(predicted_status)
status_distribution = pd.Series(predicted_status).value_counts()

# Plotting
fig, ax = plt.subplots(figsize=(12, 5))
wedges, texts, autotexts = ax.pie(status_distribution,
labels=status_distribution.index,
autopct='%1.1f%%',
startangle=90,
textprops=dict(color="w"))
```
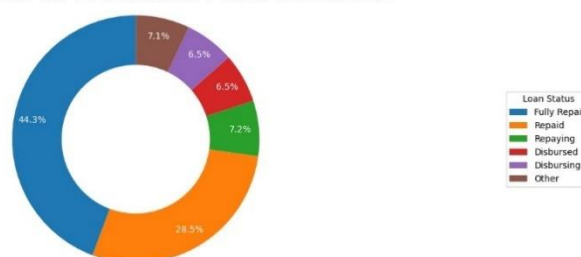


Fig 4.3.2 Pie chart of Distribution of predicted Loan status

2. The model predicts loan statuses, which are then inverse-transformed back to their original labels. The distribution of these statuses is plotted as a pie chart to visualize the proportions of each status.

**Interest Rate Insight** This insight predicts interest rates using the interest_rate_model and visualizes the relationship between loan duration and interest rate using a hexbin plot, scatter plot, and linear regression line. It also calculates the Pearson correlation coefficient.

```python
def interest_rate_insight(df, model, output_file_path):
    X = df[model.feature_names_in_]
    predicted_rates = model.predict(X)

    # Plotting
    fig = plt.figure(figsize=(12, 10))
    gs = fig.add_gridspec(3, 3)
    ax_main = fig.add_subplot(gs[1:, :2])
    ax_right = fig.add_subplot(gs[1:, 2], sharey=ax_main)
    ax_top = fig.add_subplot(gs[0, :2], sharex=ax_main)

    hb = ax_main.hexbin(df['loan duration (years)'],
predicted_rates, gridsize=20, cmap='YlOrRd')
    ax_main.scatter(df['loan duration (years)'],
predicted_rates, alpha=0.3, color='blue', s=5)

    # Linear regression
    z = np.polyfit(df['loan duration (years)'], predicted_rates,
1)
    p = np.poly1d(z)
    ax_main.plot(df['loan duration (years)'], p(df['loan
duration (years)']), "r--", alpha=0.8)

    # Statistical analysis
    correlation = stats.pearsonr(df['loan duration (years)'],
predicted_rates)[0]
```
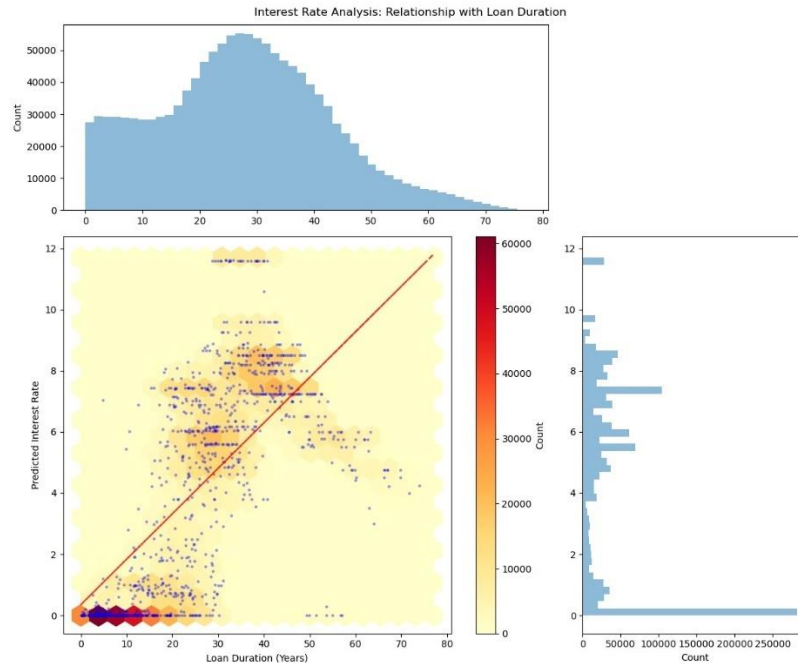
Fig 4.3.3 Interest rate analysis

The hexbin plot shows the density of data points, while the scatter plot overlays individual data points. The linear regression line indicates the trend between loan duration and interest rate, and the Pearson correlation coefficient quantifies the strength and direction of the linear relationship.

**Disbursed Amount Insight** This insight predicts disbursed amounts using the disbursed_amount_model and visualizes the relationship between the original principal amount and predicted disbursed amount using similar visualization techniques as the interest rate insight.

```
def disbursed_amount_insight(df, model, output_file_path):
    X = df[model.feature_names_in_]
    predicted_amounts = model.predict(X)

    # Plotting
    fig = plt.figure(figsize=(12, 10))
    gs = fig.add_gridspec(3, 3)
    ax_main = fig.add_subplot(gs[1:, :2])
    ax_right = fig.add_subplot(gs[1:, 2], sharey=ax_main)
    ax_top = fig.add_subplot(gs[0, :2], sharex=ax_main)

    hb  =  ax_main.hexbin(df['original  principal  amount'],
predicted_amounts, gridsize=20, cmap='YlOrRd')
    ax_main.scatter(df['original      principal      amount'],
predicted_amounts, alpha=0.3, color='blue', s=5)

    # Linear regression
    z    =    np.polyfit(df['original     principal    amount'],
predicted_amounts, 1)
    p = np.poly1d(z)
    ax_main.plot(df['original principal amount'], p(df['original
principal amount']), "r--", alpha=0.8)
```

```
    # Statistical analysis
    correlation = stats.pearsonr(df['original principal amount'],
predicted_amounts)[0]
```

The hexbin plot, scatter plot, and linear regression line visualize the relationship between the original principal amount and the predicted disbursed amount. The Pearson correlation coefficient provides a measure of the linear correlation between these variables.

**Visualization Methods**

The code employs various visualization techniques to represent insights effectively:

- **Bar Plot** for regional distribution of loans.
- **Pie Chart** for distribution of loan statuses.
- **Hexbin Plot** and **Scatter Plot** for relationships between variables such as loan duration vs. interest rate and original principal amount vs. predicted disbursed amount.
- **Linear Regression Line** to highlight trends in scatter plots.
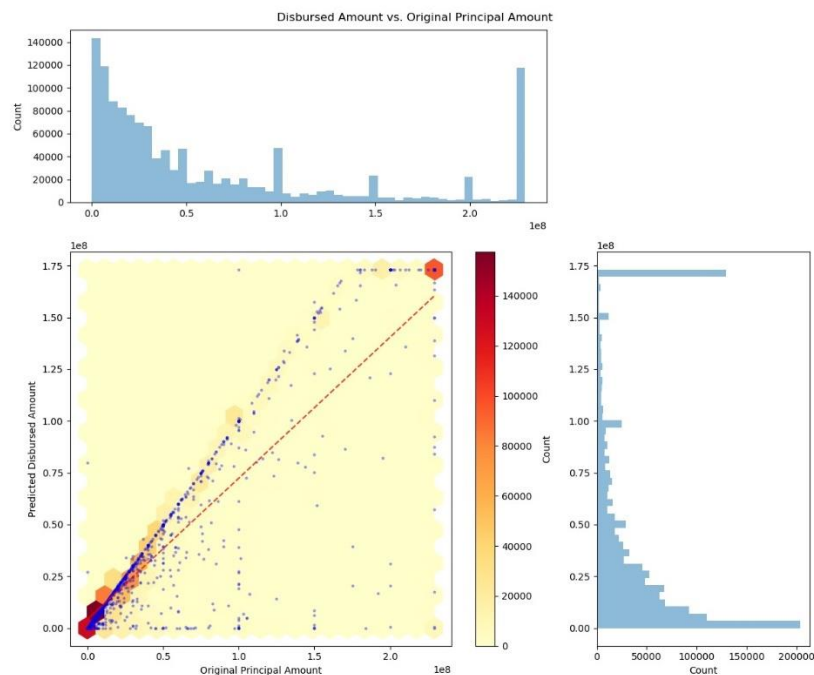


Fig 4.3.4 Disbursed amount analysis

**Statistical Analyses Performed**

Statistical analysis in the code includes the calculation of the Pearson correlation coefficient to quantify the linear relationships between variables:

**Code Snippet: Statistical Analysis**

python
Copy code
# Statistical analysis

correlation = stats.pearsonr(df['loan duration (years)'], predicted_rates)[0]
correlation = stats.pearsonr(df['original principal amount'], predicted_amounts)[0]

The Pearson correlation coefficient measures the strength and direction of the linear relationship between two continuous variables, providing a numerical insight into the data's behavior.

**Conclusion**

The code systematically preprocesses the data, leverages pre-trained machine learning models to generate predictions, and employs various visualization techniques to convey key insights. Statistical analyses further enhance the understanding of relationships between variables, making the insights more robust and actionable.

**4.4 Insights Generation**

This report details the methodology for an advanced loan portfolio analysis project, which leverages machine learning, statistical analysis, data visualization, and large language models to generate deep insights from banking data. The project encompasses several stages: data preprocessing and feature engineering, model development and training, pattern detection and insight generation, data visualization and reporting, advanced analytics with large language models, workflow integration and automation, ethical considerations and bias mitigation, and scalability and future improvements.

**Data Preprocessing and Feature Engineering**

- **Initial Data Cleaning and Preparation**

The initial step involved loading the dataset and performing essential cleaning operations to handle missing values, remove duplicates, and ensure data consistency. Data was loaded using the pandas library, and missing values were managed by either filling them with appropriate statistical measures or dropping incomplete rows/columns, depending on the extent of the missing data.

- **Encoding Categorical Variables**

Categorical variables were encoded using a custom EncoderDecoderModel class which leverages LabelEncoder for variables with high cardinality, while variables with fewer unique values were encoded using one-hot encoding. This ensured that the categorical data was appropriately transformed for model training.

- **Feature Scaling Techniques**

Feature scaling was applied where necessary, particularly for algorithms sensitive to the scale of data. StandardScaler from scikit-learn was used to standardize features by removing the mean and scaling to unit variance.

- **Feature Selection and Dimensionality Reduction**

Feature selection was guided by domain knowledge and statistical analysis, aiming to retain the most informative features while reducing noise. Dimensionality reduction techniques like Principal Component Analysis (PCA) were considered but ultimately not employed due to the

sufficient performance of the models without it. Features highly correlated with target variables or deemed essential for predictive accuracy were prioritized.

**Model Development and Training**

- **Machine Learning Algorithms**

Several machine learning algorithms were explored, with Random Forest being the primary choice for both regression and classification tasks due to its robustness and interpretability. Random Forest Regressors were used for predicting continuous variables like loan duration, interest rate, and disbursed amount, while Random Forest Classifiers were employed for predicting categorical variables such as loan status.

- **Model Architecture and Hyperparameters**

The Random Forest models were configured with 100 estimators, and hyperparameters such as the maximum depth of the trees and the minimum samples per leaf were optimized using grid search and cross-validation techniques. Cross-validation ensured that the models were not overfitted and could generalize well to unseen data.

- **Training and Performance Metrics**

Models were trained using an 80-20 train-test split strategy. Performance metrics for regression tasks included Mean Squared Error (MSE) and R-squared (R²) scores, while classification tasks were evaluated using accuracy, precision, recall, and F1 scores. Detailed reports generated using the classification_report function provided insights into model performance across different classes.

**Pattern Detection and Insight Generation**

- **Statistical Methods**

Statistical methods, including correlation analysis and hypothesis testing, were employed to identify significant relationships and patterns within the data. These analyses provided a foundation for understanding the underlying dynamics of the loan portfolio.

- **Utilizing Trained Models for Predictions**

Trained models were utilized to generate predictions on new data, facilitating the generation of insights regarding loan performance, risk assessment, and other financial metrics. Model predictions were critically evaluated to ensure they aligned with domain knowledge and real-world expectations.

- **Ensemble Methods and Model Stacking**

To enhance prediction accuracy, ensemble methods and model stacking techniques were explored. These methods involved combining predictions from multiple models to produce a more robust and accurate final prediction.

**Data Visualization and Reporting**

- **Types of Visualizations**

Various visualizations, including hexbin plots, histograms, and bar charts, were created to convey insights effectively. These visualizations were generated programmatically using libraries like matplotlib and seaborn.

- **Generating Visualizations**

Visualizations were created using a dedicated function in the raw_insight_maker script, which handled the generation and saving of plots. These visualizations played a crucial role in presenting the analysis results clearly and intuitively.

- **Contribution to Analysis and Decision-Making**

The visualizations facilitated a better understanding of the data, highlighting key patterns and trends that informed decision-making. For instance, the regional distribution of loans helped identify high-priority areas for future investments, while the loan status distribution provided insights into the overall health of the loan portfolio.

**Advanced Analytics with Large Language Models**

- **Integration of Gemini 1.5 Pro Model**

The Gemini 1.5 Pro model was integrated to generate detailed financial analyses, leveraging its capabilities to produce narrative-based insights from raw data. This integration was facilitated through the report_generator script.

- **Prompt Engineering**

Prompt engineering involved crafting specific instructions to guide the LLM in generating coherent and insightful narratives. The predefined prompt ensured that the generated analyses were focused on key financial aspects and presented in a professional tone.

**Challenges and Considerations**

Challenges included ensuring the accuracy and relevance of the generated content. These were addressed by iteratively refining the prompts and incorporating feedback loops to improve the quality of the outputs.

**Workflow Integration and Automation**

**Cohesive Workflow Integration**

The various components of the project were integrated into a cohesive workflow, automating the processes of data preprocessing, model training, visualization, and insight generation. This automation streamlined the analysis process and ensured consistency across different stages of the project.

**Automation Techniques**

Automation techniques included the use of scripts to handle repetitive tasks and the deployment of models to automatically generate predictions and insights. This significantly reduced the manual effort required and improved the efficiency of the analysis.

**Compilation and Presentation of Results**

The results were compiled into comprehensive reports, which included both narrative insights and visualizations. These reports provided a holistic view of the loan portfolio, facilitating informed decision-making.

**Ethical Considerations and Bias Mitigation**

- **Ensuring Fairness and Mitigating Bias**

Steps were taken to ensure the fairness of the models and mitigate potential biases. This involved careful selection of training data, rigorous evaluation of model performance across different demographic groups, and ongoing monitoring to detect and address any emerging biases.

- **Data Privacy and Security**

Data privacy and security were maintained throughout the process by adhering to best practices in data handling and storage. Sensitive information was anonymized, and access to the data was restricted to authorized personnel only.

**Scalability and Future Improvements**

- **Scalability of Methodology**

The current methodology is scalable to larger datasets and more complex financial products. The use of robust algorithms and efficient data processing techniques ensures that the analysis can handle increased data volumes without compromising performance.

- **Potential Areas for Improvement**

Future improvements could include the integration of more advanced machine learning techniques, such as deep learning, to capture more complex patterns in the data. Additionally, expanding the scope of the analysis to include more diverse financial products would provide a more comprehensive view of the loan portfolio.

## Conclusion

This methodology represents a significant advancement in financial analytics, offering a comprehensive and automated approach to loan portfolio analysis. By leveraging machine learning, statistical analysis, data visualization, and large language models, this project provides deep insights that can inform strategic decision-making and enhance loan portfolio management. The potential impact of this methodology is substantial, paving the way for more informed and effective financial strategies.

**Chapter 5**

## Results and Discussions

### 5.1 Presentation of Results

This section presents the results of the KnowledgeInsight Hub website project, showcasing various pages and features through a series of screenshots. The purpose of this section is to provide a detailed visual tour of the website, highlighting its design, functionalities, and how these elements align with the project's objectives. By exploring these screenshots, readers will gain a comprehensive understanding of how the website facilitates knowledge representation and insight generation from structured datasets.



Figure 5.1.1: Flow of Execution of Project

**Home Page**

We begin with the **Home Page**, which serves as the entry point for users. The screenshot of the landing page highlights the main features such as the website logo, navigation menu, hero section, and key call-to-action buttons. The design focuses on user engagement, providing clear messaging and easy access to core functionalities. This intuitive layout ensures that users can quickly understand the website's purpose and start exploring its features.



Figure 5.1.2 Home page

## Visualization Page

Next, the **Visualization Page** showcases the various types of visualizations available and their interactivity features. The screenshot focuses on the diversity of charts and graphs, user interactivity, and customization options. These visual tools are essential for enhancing users' understanding of data patterns, making complex information more accessible and actionable.

**Data Report**

**Before Processing**

Shape of the DataFrame: (1323796, 33)

```
Column Name                | Missing Values | Percentage Missing
---------------------------+----------------+------
Country Code               | 319            | 0.02 %
Borrower                   | 9122           | 0.69 %
Guarantor Country Code     | 49050          | 3.71 %
Guarantor                  | 75399          | 5.70 %
Interest Rate              | 30736          | 2.32 %
Currency of Commitment     | 1323796        | 100.00%
Project ID                 | 42             | 0.00 %
Project Name               | 159114         | 12.02%
First Repayment Date       | 4077           | 0.31 %
Last Repayment Date        | 3918           | 0.30 %
Agreement Signing Date     | 19244          | 1.45 %
Board Approval Date        | 2              | 0.00 %
Effective Date (Most Recent) | 10234        | 0.77 %
Closed Date (Most Recent)  | 1250           | 0.09 %
Last Disbursement Date     | 540402         | 40.82%
```
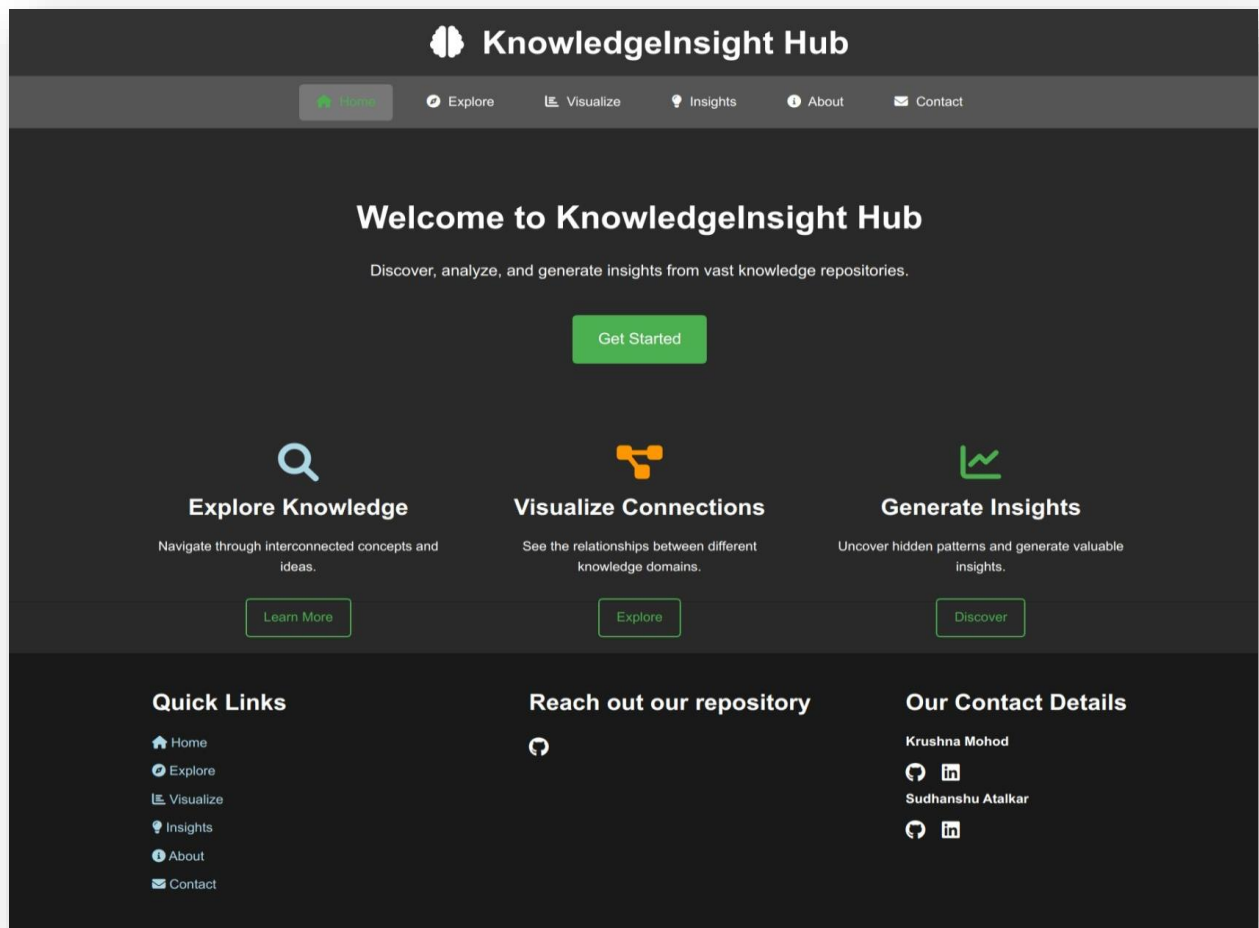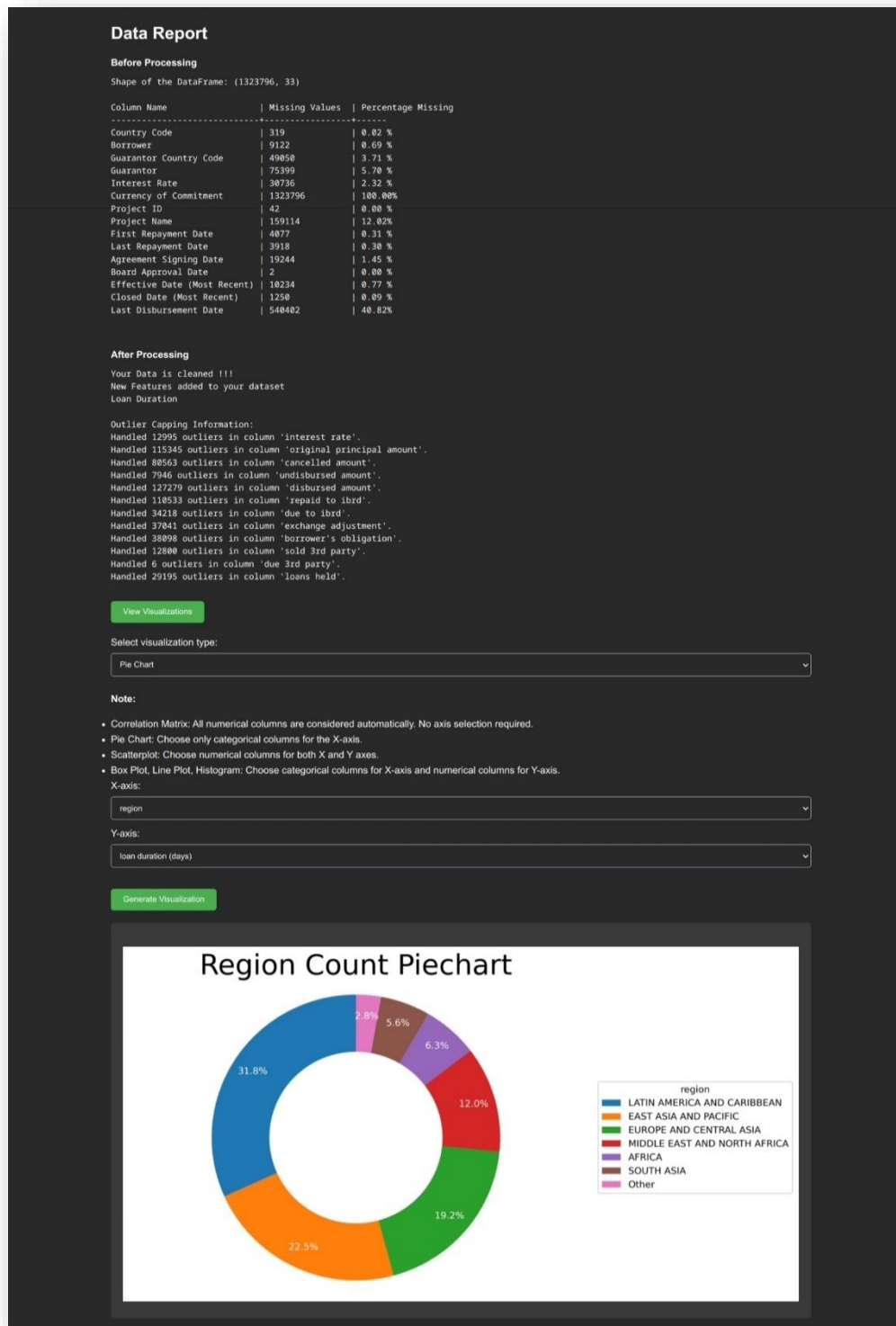
**After Processing**

Your Data is cleaned !!!
New Features added to your dataset
Loan Duration

Outlier Capping Information:
Handled 12995 outliers in column 'interest rate'.
Handled 115345 outliers in column 'original principal amount'.
Handled 80563 outliers in column 'cancelled amount'.
Handled 7946 outliers in column 'undisbursed amount'.
Handled 127279 outliers in column 'disbursed amount'.
Handled 110533 outliers in column 'repaid to ibrd'.
Handled 34218 outliers in column 'due to ibrd'.
Handled 37041 outliers in column 'exchange adjustment'.
Handled 38098 outliers in column 'borrower's obligation'.
Handled 12800 outliers in column 'sold 3rd party'.
Handled 6 outliers in column 'due 3rd party'.
Handled 29195 outliers in column 'loans held'.

[ View Visualizations ]

Select visualization type:

[ Pie Chart ]

**Note:**

- Correlation Matrix: All numerical columns are considered automatically. No axis selection required.
- Pie Chart: Choose only categorical columns for the X-axis.
- Scatterplot: Choose numerical columns for both X and Y axes.
- Box Plot, Line Plot, Histogram: Choose categorical columns for X-axis and numerical columns for Y-axis.

X-axis:

[ region ]

Y-axis:

[ loan duration (days) ]

[ Generate Visualization ]

### Region Count Piechart

31.8%
2.8%
5.6%
6.3%
12.0%
19.2%
22.5%

region
- LATIN AMERICA AND CARIBBEAN
- EAST ASIA AND PACIFIC
- EUROPE AND CENTRAL ASIA
- MIDDLE EAST AND NORTH AFRICA
- AFRICA
- SOUTH ASIA
- Other

Figure 5.1.3 Visualization page

## Insights Page

Moving to the **Insights Page**, we see the core functionality of the website in action The screenshot displays the data upload section, insight generation tools, and the results display area. Key features include an intuitive upload process, real-time insight generation, and a clear presentation of insights. These elements are crucial for enabling users to derive meaningful insights from their data, directly supporting the project's goal of providing a platform for data analysis and decision-making.



Figure 5.1.4 Insight page

**Explore Knowledge Page**

The **Explore Knowledge Page** allows users to delve deeper into the knowledge base. The screenshot highlights the search functionality, categorization of knowledge, and user interactivity. These features enable users to efficiently explore and interact with the knowledge base, supporting informed decision-making and enhancing the overall user experience.
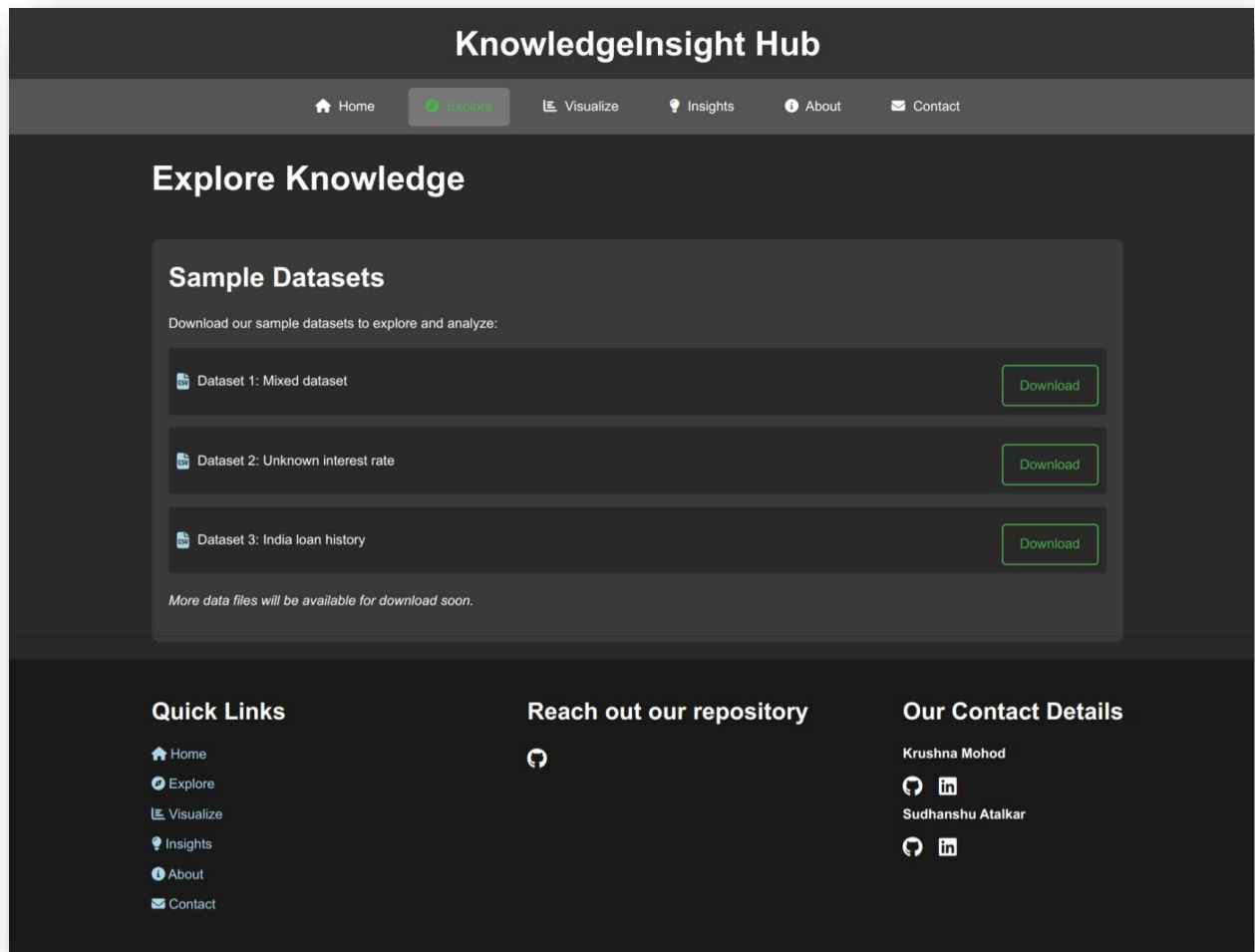


Figure 5.1.5 Explore page

**About and Contact Pages**

Finally, the About and Contact Pages provide information about the team and offer user support options. The screenshots showcase the team information, contact form, and additional resources. These sections build trust and ensure users have access to the support they need, which is crucial for maintaining a positive user experience.
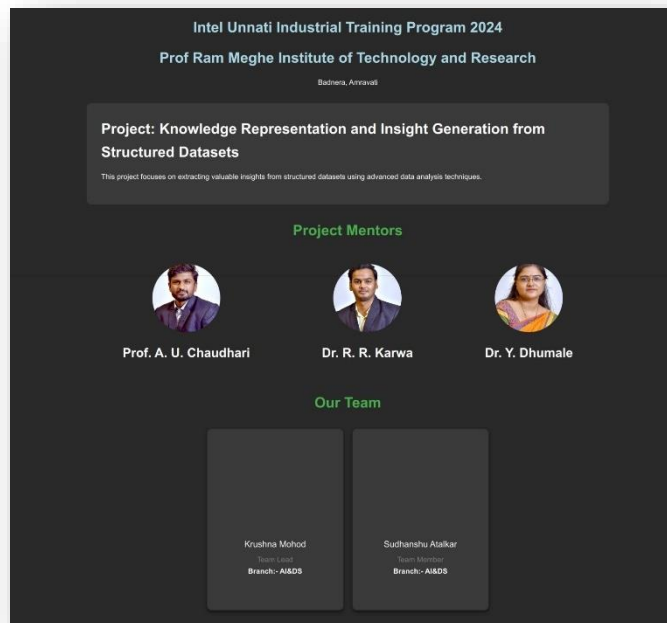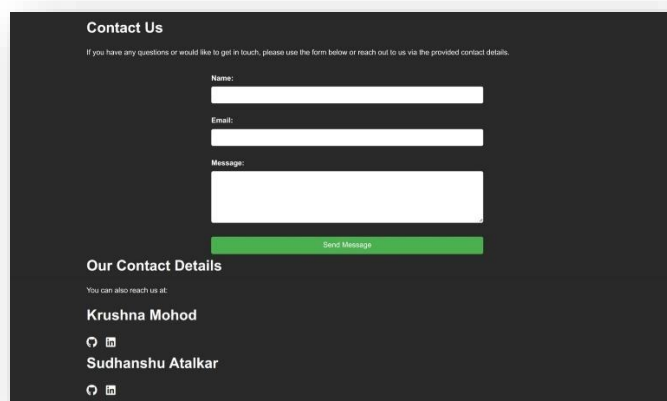
Figure 5.1.6 About Page



Figure 5.1.7 Contact page

The screenshots presented in this section illustrate the comprehensive and user-centric design of the KnowledgeInsight Hub. Each feature has been thoughtfully integrated to support the project's goal of knowledge representation and insight generation. The website's intuitive interface, robust functionalities, and visually appealing design collectively enhance the user experience and facilitate effective data analysis and decision-making.

### 5.2  Identified Patterns

Analysis Report on Loan Portfolio Visualizations

### 5.2.1. Regional Distribution of Loans (Disbursed Amount)

**Key Patterns and Trends:**

- **Highest Disbursements**: Latin America and the Caribbean, followed by East Asia and Pacific, and Europe and Central Asia.
- **Significant Drop**: Middle East and North Africa, South Asia, and Africa.
- **Minimal Disbursements**: Specific regions within Africa, such as Western and Central Africa, and Eastern and Southern Africa.

**Potential Causes:**

- **Economic Stability and Growth**: Higher disbursements in regions with stable economies.
- **Political Stability and Governance**: More stable regions likely to receive higher loan amounts.
- **Historical Ties and Trade Relations**: Influence of past relationships with international banks.

**Implications:**

- **Investment Opportunities**: Regions with higher disbursements indicate better investment prospects and lower perceived risks.
- **Targeted Financial Strategies**: Need for strategies to support underrepresented regions.
- **Risk Assessment Models**: Should consider regional economic conditions for optimal performance.

**Anomalies and Outliers:**

- **Minimal Disbursements in Africa**: Could indicate lack of demand or significant risk factors.
- **Investigation Needed**: To understand underlying issues such as political instability or poor creditworthiness.

**Broader Context:**

- **Global Economic Trends**: Emerging markets attract substantial investments.
- **Economic Inequalities**: Discrepancies in disbursements may reflect broader economic disparities.

**Limitations:**

- **Size of Economy and Population**: Not accounted for, which could skew perceptions.
- **Loan Performance Data**: Would provide a more comprehensive risk assessment.

**Further Analysis:**

- **Loan Performance Metrics**: Such as repayment rates and defaults across regions.
- **Socio-Economic Factors**: Influencing loan demand and disbursement.

**Strategic Recommendations:**

- **Short-term**: Increase financial support and capacity-building in underrepresented regions.
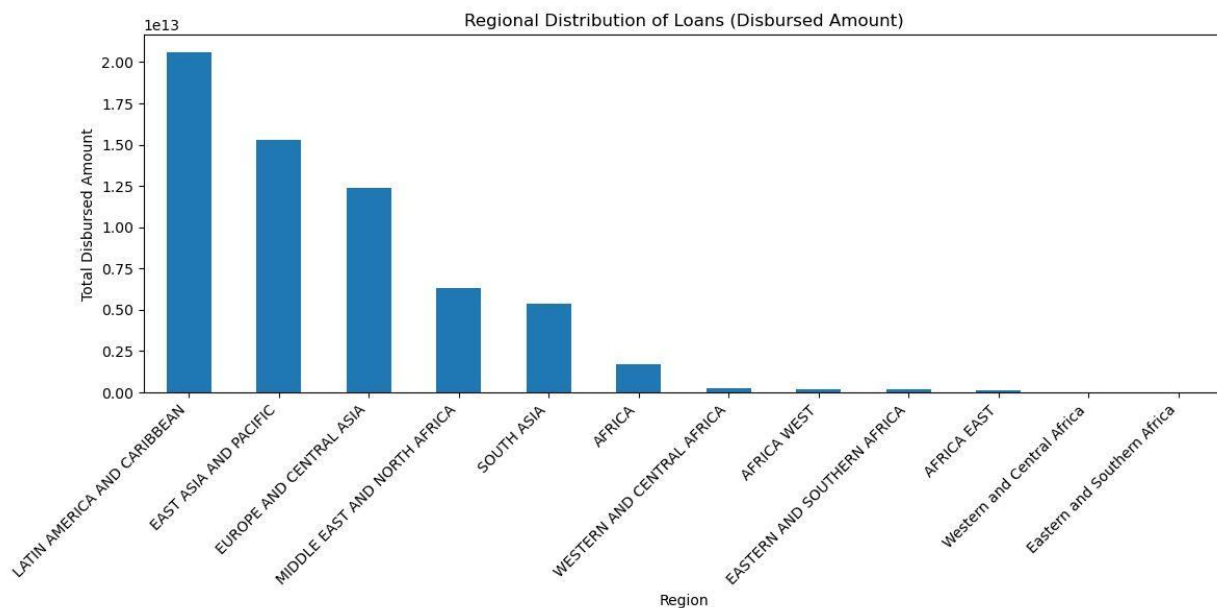- **Long-term**: Develop tailored financial products and risk mitigation strategies.



Figure 5.2.1.1 Regional distribution of loans

### 5.2.2. Distribution of Predicted Loan Statuses

**Key Patterns and Trends:**

- **Fully Repaid Loans**: Largest segment at 44.3%.
- **Ongoing Financial Engagements**: 'Repaid' (28.5%) and 'Repaying' (7.2%).
- **Active Loan Cycles**: 'Disbursed' and 'Disbursing' both at 6.5%.
- **Other Statuses**: Account for 7.1%.

**Potential Causes:**

- **Effective Loan Management**: Contributing to high repayment rates.
- **'Other' Category**: Needs investigation for specific issues affecting these loans.

**Implications:**

- **Strong Portfolio Performance**: High repayment rates support sustainable practices.
- **Understanding 'Other'**: Crucial for identifying potential risks.

**Anomalies and Outliers:**

- **'Other' Category**: Warrants detailed analysis to uncover specific issues.

**Broader Context:**

- **Positive Loan Performance**: Aligns with global economic recovery trends.
- **Continued Monitoring**: Essential for adapting to changing conditions.

**Limitations:**

- **Loan Types**: Not differentiated, which might have varying risk profiles.
- **External Economic Factors**: Not accounted for, affecting interpretation.

**Further Analysis:**

- **Segmenting Loan Statuses**: By region, sector, and type for granular insights.
- **Factors in 'Other' Category**: To refine risk assessment models.

**Strategic Recommendations:**

- **Short-term**: Review loans in the 'Other' category.
- **Long-term**: Implement continuous monitoring and adaptive strategies.
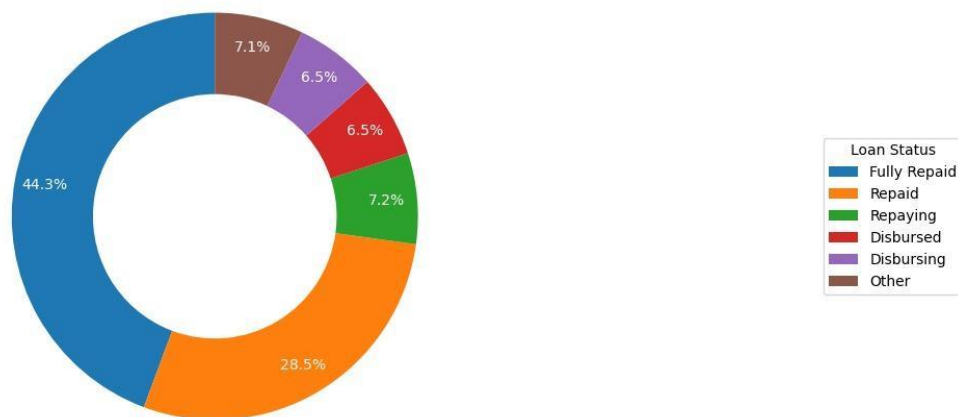
## Distribution of Predicted Loan Statuses

Figure 5.2.2.1 Distribution of Pre

### 5.2.3. Disbursed Amount vs. Original Principal Amount

**Key Patterns and Trends:**

- **Distribution Shape**: Right-skewed for both disbursed and original principal amounts.
- **Low-Value Loans**: High concentration close to zero.
- **High-Value Loans**: Noticeable spikes at specific points.
- **Positive Correlation**: Between original principal and disbursed amount.

**Potential Causes:**

- **Economic Factors**: Reflect strategy to serve low-risk borrowers.
- **Loan Policies**: Specific programs targeting high-value borrowers.

**Implications:**

- **Risk Management**: Low-value loans reduce overall risk; high-value loans need monitoring.
- **Financial Strategy**: Loan policies may need adjustment.

**Anomalies:**

- **Outliers**: High principal amounts need further investigation.

**Comparative Analysis:**

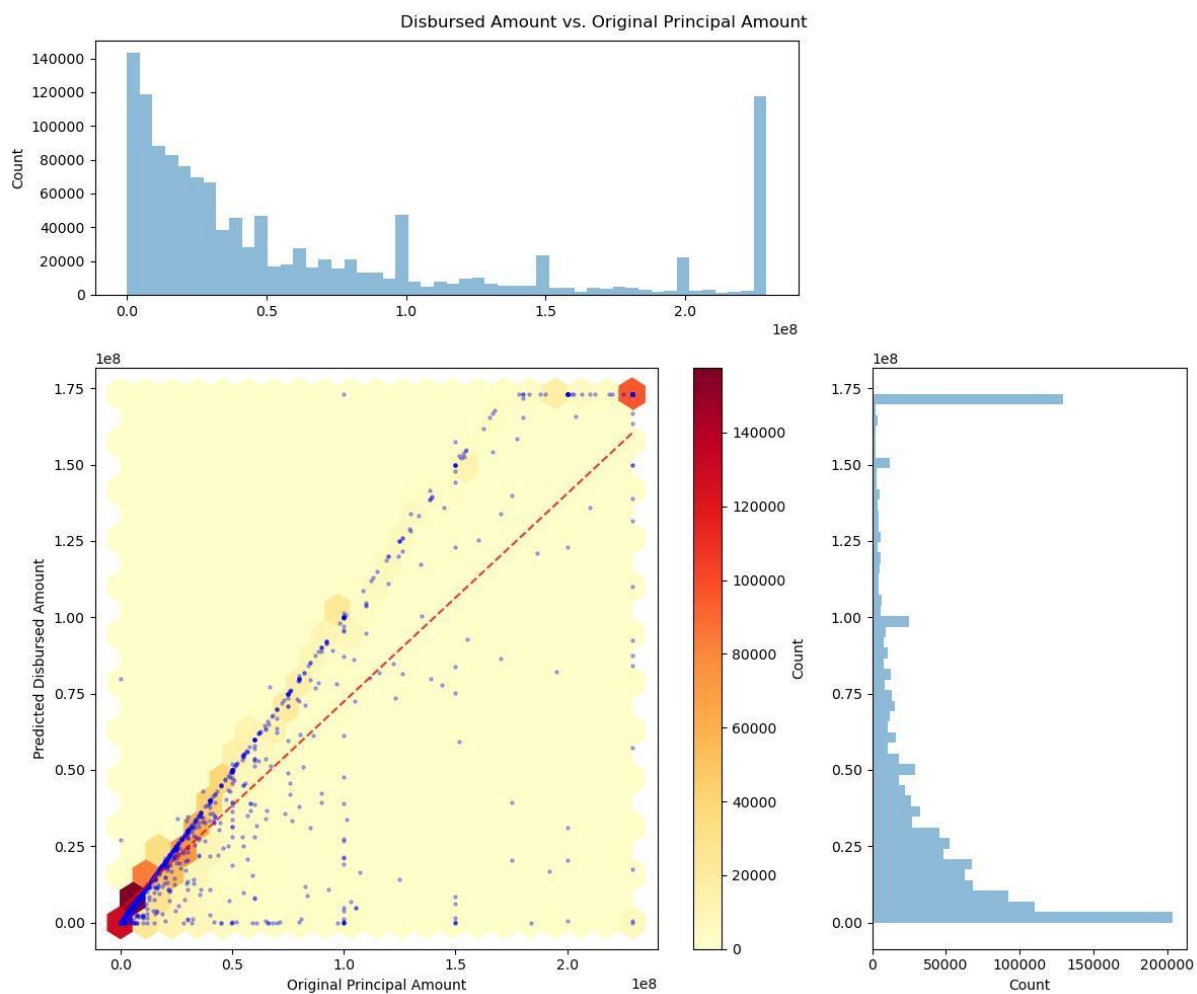- **Correlation Consistency**: Should align with other financial indicators.



Figure 5.2.3.1 Disbursed amount analysis

### 5.2.4. Interest Rate Analysis: Relationship with Loan Duration

**Key Patterns and Trends:**

- **Distribution Shape**: Bell-shaped for loan durations.
- **Positive Correlation**: Between loan duration and interest rate.
- **High-Density Regions**: At lower interest rates and shorter durations.

**Potential Causes:**

- **Economic and Policy Factors**: Longer durations come with higher interest rates.
- **Market Conditions**: Influencing interest rate trends.

**Implications:**

- **Loan Management**: Setting interest rates based on durations.
- **Risk Assessment**: Higher interest rates for longer durations indicate higher risk.

**Anomalies:**

- **Outliers**: High interest rates or durations need scrutiny.

**Comparative Analysis:**

- **Pattern Consistency**: Should match other financial indicators.
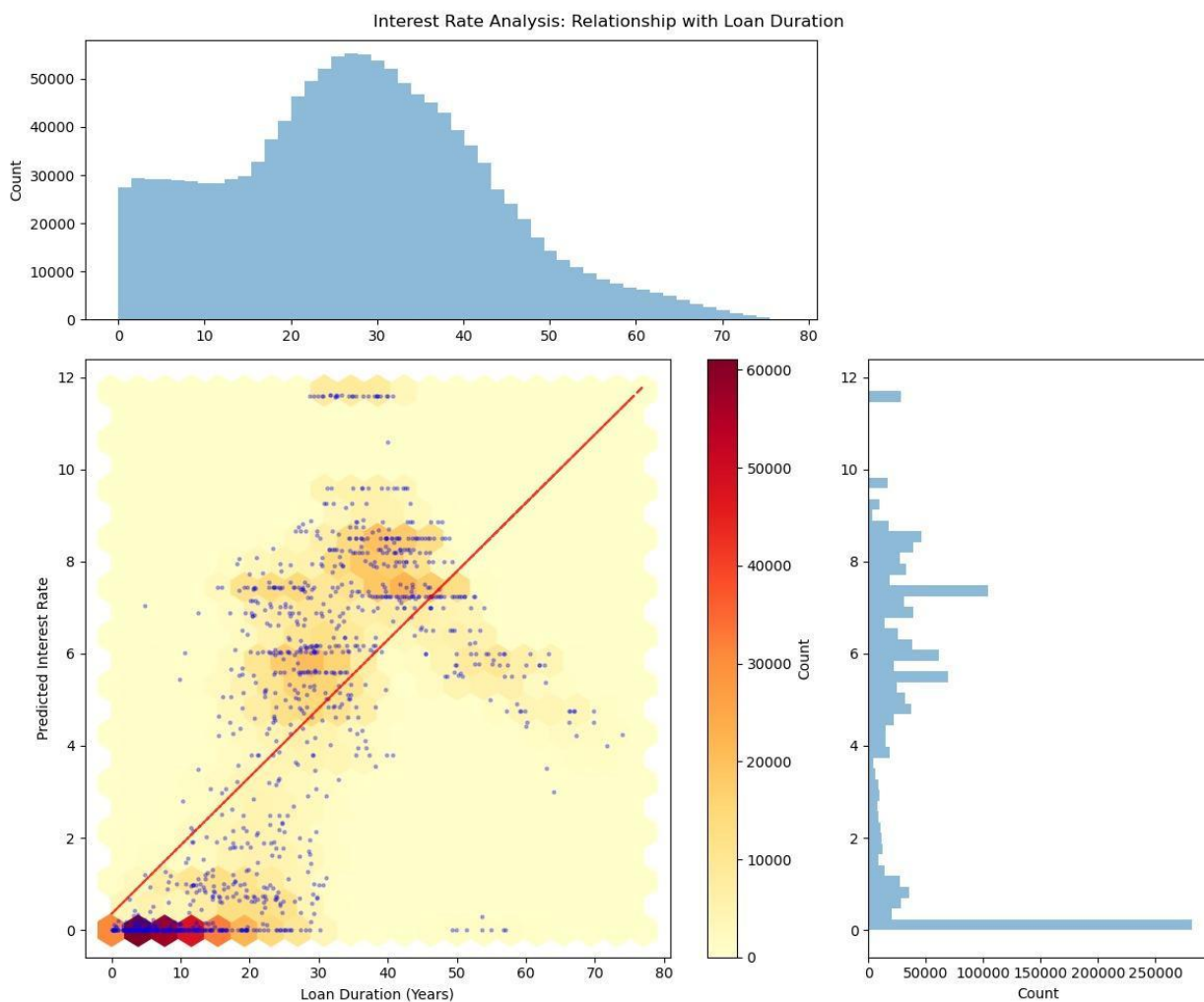


Figure 5.2.4.1 Interest rate analysis

The analysis highlights significant patterns in the loan portfolio, such as the concentration of low-value loans, the positive correlation between loan amount and disbursed amount, and the relationship between loan duration and interest rates. These insights can guide strategic decisions in loan management, risk assessment, and financial planning, ensuring the portfolio's stability and profitability in the broader context of international banking trends.

# Chapter 6

# Conclusion and future scope

The project "Knowledge Representation and Insight Generation from Structured Datasets" has demonstrated a robust and comprehensive workflow using the IBRD Bank dataset. The journey began with meticulous data preprocessing, which is critical for ensuring the accuracy and reliability of subsequent analyses. This step involved cleaning, transforming, and preparing the data, thus setting a strong foundation for the next stages.

Following preprocessing, the project focused on knowledge representation through a variety of visualizations. These visualizations were instrumental in providing clear, intuitive insights into the dataset's structure and key attributes, making complex data more understandable and accessible.

To uncover hidden patterns and relationships within the data, a Random Forest Classifier was employed. This machine learning model, known for its robustness and accuracy, successfully identified significant patterns, enhancing our understanding of the dataset.

The process of generating insights was further refined using the Gemini 1.5 Pro, a powerful tool that facilitated the extraction of actionable insights from the dataset. This step was crucial in translating raw data into meaningful information that can inform decision-making processes.

Scalability was a key consideration throughout the project. By utilizing parallel processing techniques and leveraging the concurrent.futures and CUDA libraries, the project was able to efficiently handle large datasets and computationally intensive tasks. This approach not only improved processing times but also ensured that the solution could scale to meet increasing demands.

A significant aspect of the project was the development of a user-friendly interface using HTML, CSS, and JavaScript, complemented by a Flask web application framework. This ensured a seamless user experience (UX) and a visually appealing user interface (UI). The integration of these technologies allowed for interactive data visualizations and an accessible platform for users to explore insights.

Overall, the project has showcased a seamless integration of data preprocessing, knowledge representation, pattern identification, and insight generation, supported by advanced computational techniques to ensure scalability and efficiency. The strong focus on UI and UX design, along with the use of modern web technologies, has resulted in a powerful and user-friendly application.

**Future Scope**

1. **Integration of Advanced Machine Learning Models**:
    o Incorporate more advanced machine learning models such as Gradient Boosting Machines (GBM), Deep Neural Networks (DNN), or ensemble methods to improve pattern identification and predictive accuracy.
2. **Enhanced Data Preprocessing Techniques**:
    o Develop more sophisticated data preprocessing techniques to handle missing values, outliers, and data normalization more effectively. This could include the use of advanced imputation methods and automated feature engineering.

3.  **Real-time Data Processing and Insight Generation**:
    o   Implement real-time data processing capabilities to handle streaming data, enabling the generation of insights in real-time. This would be particularly useful for applications requiring immediate decision-making.
4.  **Incorporation of Natural Language Processing (NLP)**:
    o   Integrate NLP techniques to extract and analyze textual data, providing a more comprehensive insight generation that combines both structured and unstructured data.
5.  **Collaborative Insight Generation**:
    o   Develop collaborative features that allow multiple users to work together on insight generation, sharing findings, and collaboratively refining models and visualizations.

By addressing these future scope areas, the project can evolve into a more comprehensive and versatile solution for knowledge representation and insight generation from structured datasets. This evolution will ultimately enhance its ability to support informed and effective decision-making processes across various domains and applications.

# References

[1] Famili, A., Shen, W. M., Weber, R., & Simoudis, E. (1997). Data preprocessing and intelligent data analysis. *Intelligent data analysis*, *1*(1), 3-23.

[2] Habib, M., & Okayli, M. (2024). Evaluating the sensitivity of machine learning models to data preprocessing technique in concrete compressive strength estimation. *Arabian Journal for Science and Engineering*, 1-19.

[3] Rao, VVR Maheswara, and V. Valli Kumari(2011). "An Enhanced Pre-Processing Research Framework For Web Log Data Using A Learning Algorithm." *Computer Science and Information Technology, DOI* (2011): 1-15.

[4] Singh, S. and Singh, N. (2011) Big Data Analytics. 2012 International Conference on Communication, Information & Computing Technology Mumbai India, IEEE, October 2011.

[5] Waskom, Michael L. "Seaborn: statistical data visualization." *Journal of Open Source Software* 6, no. 60 (2021): 3021.

[6] Lu, X., Wang, L., Jiang, Z., He, S., & Liu, S. (2022). MMKRL: A robust embedding approach for multi-modal knowledge graph representation learning. *Applied Intelligence*, 1-18.

[7] Müllner, L. (2021). *Knowledge-assisted visual analytics: data exploration and insight generation of health care data* (Doctoral dissertation, Wien).

[8] Ng, K. K. Y., & Zhang, P. C. (2023). Advancing medical affair capabilities and insight generation through machine learning techniques. *Journal of Pharmaceutical Policy and Practice*, *16*(1), 165.

[9] Ma, L., & Sun, B. (2020). Machine learning and AI in marketing–Connecting computing power to human insights. *International Journal of Research in Marketing*, *37*(3), 481-504.

[10] Chatzimparmpas, A., Martins, R. M., Jusufi, I., & Kerren, A. (2020). A survey of surveys on the use of visualization for interpreting machine learning models. *Information Visualization*, *19*(3), 207-233.