# Introduction

The global financial landscape is intricately linked with the banking sector's ability to assess and manage risks associated with lending. In this context, the project delves into the predictive analysis of loan defaulters using historical data obtained from a renowned German bank. This dataset encapsulates crucial information pertaining to customers' financial behaviors, employment details, and personal demographics, aiming to construct a robust machine learning model capable of foreseeing potential loan defaults.

In this project, we delve into a credit dataset aiming to analyze various factors influencing credit outcomes. The dataset encompasses several features such as 'months_loan_duration', 'amount', 'percent_of_income', 'years_at_residence', 'age', 'existing_loans_count', and 'dependents'. Our primary goal is to understand the relationships between these factors and the target variable while employing machine learning models to predict credit outcomes.

# Methods and Materials

The initial phase of this project involved an extensive Exploratory Data Analysis (EDA) to comprehend the dataset's structure and inherent patterns. Utilizing Python's Pandas and Matplotlib libraries, descriptive statistics were computed to gain insights into the central tendencies and distributions of key features. Visualizations, such as histograms for continuous variables like 'months_loan_duration' and 'amount', and bar plots for categorical variables like 'purpose' and 'credit_history', offered a comprehensive overview of the data distribution. These visualizations, alongside correlation matrices and pair plots, unveiled potential relationships between variables, aiding in feature selection and understanding potential predictive attributes.

### *Machine Learning Models*

Following the EDA, the dataset underwent meticulous pre-processing steps. Missing values in numeric columns were handled by imputing them with the mean of respective features. Furthermore, categorical variables were transformed using one-hot encoding to convert them into a format suitable for machine learning models. The resulting dataset was then split into features and target variables, with the former encompassing 'months_loan_duration', 'amount', 'percent_of_income', 'years_at_residence', 'age', 'existing_loans_count', and 'dependents'.

Subsequently, the prepared dataset was subjected to several machine learning models to predict credit outcomes. The selected models encompassed diverse methodologies, including Logistic Regression, Random Forest Classifier, Decision Tree Classifier, Support Vector Machine (SVM), and Gradient Boosting Classifier. Each model was trained on the dataset and evaluated based on appropriate evaluation metrics (e.g., accuracy, precision, recall, F1-score, ROC-AUC, etc.) to assess its predictive performance. These models were employed to ascertain the most suitable algorithm for predicting credit outcomes based on the provided dataset.

# Results

Our data preprocessing phase involved handling missing values by filling them with the mean of respective columns. Subsequently, we utilized various machine learning models, including Linear Regression, Decision Tree, Support Vector Machine (SVM), and K-Nearest Neighbors (KNN), to predict credit outcomes. The mean squared errors for these models were as follows:

Linear Regression: 0.5967

Decision Tree: 0.9629

SVM: 0.5339

KNN: 0.6921

## *Discussion*

Our analysis revealed that the Support Vector Machine (SVM) model exhibited the lowest mean squared error, suggesting better performance in predicting credit outcomes compared to other models. However, it's crucial to interpret these results cautiously, considering potential limitations. Further examination into feature importance, hyperparameter tuning, and additional feature engineering could enhance model performance.

We acknowledge limitations in this study, including the necessity for more robust feature selection and exploring other advanced algorithms for improved predictions. Additionally, the dataset's size and scope might have influenced model performance.

# Conclusions

In conclusion, our project embarked on understanding credit outcomes through machine learning models. While SVM displayed promising performance, there remains room for refinement in feature engineering and model selection to achieve more accurate predictions. This study highlights the complexity of credit prediction and emphasizes the need for continuous exploration and refinement in the field.This report encapsulates our journey from data exploration to employing various models, presenting insights, limitations, and prospects for future research in credit prediction using machine learning.