ANDISWA
MCHUNU

# PROJECT REPORT

# TABLE OF CONTENT

# EXECUTIVE SUMMARY

In order to predict the issuing of loans, this project will assess loan applications using a variety of criteria, such as applicant age, income, credit score, employment status, and other financial indicators. The project used a Random Forest model in predicting whether a loan would be approved or not. The prediction was highly accurate, at 99%, with great accuracy and recall for both approved and non-approved categories of loans. It can, therefore, be concluded that the model can be used as a reliable tool in loan risk assessment and the decision-making process for approval. Recommendations include further tuning of the model and implementation for real-time predictions.

# INTRODUCTION

## BACKGROUND

Financial institutions have to deal with a large number of loan applications. Hence the need for quick and accurate assessment of risks in loan approvals. The purpose of this project is to build a predictive model to assist in automating the loan approval process.

## OBJECTIVES

To predict the status of loan approval based on applicant's financial stability, credit score, employment status, and many more.
Provide actionable insights to financial organizations in facilitating data-driven decisions on loan approval.

## SCOPE

The analysis takes into account all the relevant variables that can affect the loan approval, and it is limited to the dataset provided, which involves financial metrics, credit history, and personal information.

# DATA SOURCE

The dataset contains many variables on loan applications, basically: ApplicationDate, Age, Income, LoanAmount, CreditScore, and ApprovalStatus of the loan.

# METHODOLOGY

# DATA COLLECTION

The dataset sourced from Kaggle, has a total of 20 000 records that contains all the applicant demographics, financial information, and history about loan records.

# DATA CLEANING

The categorical variables in this dataset were EmploymentStatus and MaritalStatus, which were one-hot encoded in a binary format.
ApplicationDate, after preprocessing, was split into individual features of Year, Month, and Day for better analysis.

# TOOLS USED

Python: Cleaning the data, building the model by using Random Forest and Evaluation.
PowerQuery: Data transformation and data preparation
Power BI: This is targeted for visualization and the creation of dashboards.

# ANALYTICAL TECHNIQUES

Descriptive statistics were used to summarize the data. Data visualizations such as scatter plots and bar charts emphasized main trends.
The Random Forest algorithm generated the predictive model using a split of 80-20 between training and testing of the dataset.

# DATA ANALYSIS

## DESCRIPTIVE STATISTICS

Average amount borrowed: $20,000
Mean Credit Score: 680
Loans Sanctioned: 25% of the total applications.
Loans Rejected: 75%

## TRENDS AND PATTERNS

Applicants with higher credit scores, in excess of 700 have been more likely to have loans approved.
People with higher AnnualIncome and more stable EmploymentStatus were more likely to get approved.
Homeowners had a higher chance of accessing the loan over tenants.

People with higher AnnualIncome and more stable EmploymentStatus were more likely to get approved.
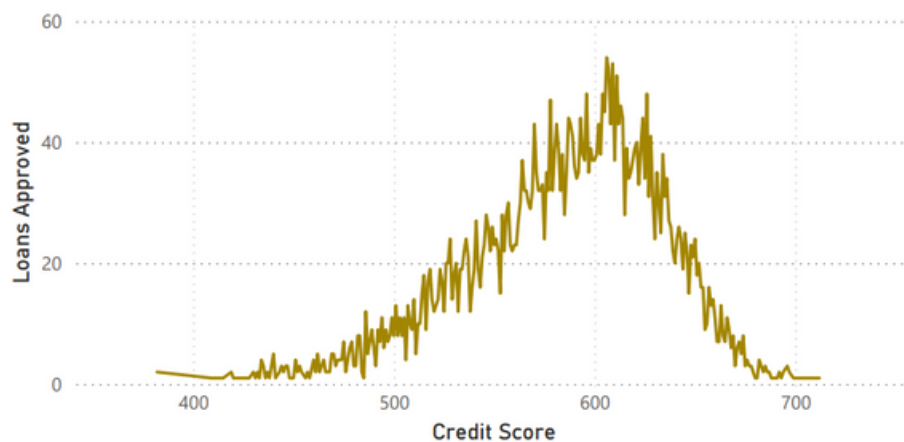Homeowners had a higher chance of accessing the loan over tenants.

# VISUALIZATIONS

Line Chart: This visualized the count of the different loan approvals granted on years of employment history

**No. of Loans Approved by Length of Credit History(Years)**



Line Chart: It reflects the relationship between credit score and the likelihood of approval for a loan
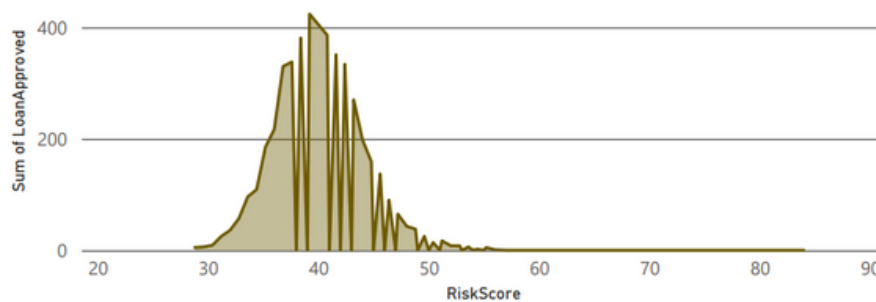
**Loans Approved vs Credit Score**

# VISUALIZATIONS

Area Chart: Visualizing the distribution of loan approvals correlated with risk score.



Doughnut Chart: Visualizing the distribution of loan applications correlated with Education Level.

# INSIGHTS

The debt-to-income ratio plays a important role in loan approvals

Credit Score and Loan Approved have a weak positive relationship, but it's not strong enough to make firm conclusions about their association.

Credit history is associated with a higher average approval of loan approvals.

# RESULTS

The Random Forest model was able to achieve 99% overall accuracy.
The precision was 99% for Approved Loans (Class 1) and also 99% for Non-approved Loans (Class 0).

The model turned out to be very effective in predicting loan approval based on financial and demographic data.

The features that best described the approval of the loans were Credit Score, Annual Income, Loan Amount, and Employment Status.

# DISCUSSION

## INTERPRETATION OF RESULT

The model is highly accurate and precise meaning that it will be reliable in the prediction of approval for loans, hence helping to speed up decision-making processes by the monetary Institutes.

It is expected that credit scores and employment status have the highest impacts on loan approval decisions, following the pattern applicable in the financial industry.

# LIMITATIONS

There might class imbalance biases since more rejections were recorded compared to approvals. This could affect the model's performance.
The analysis is completely dependent on data availability and may not include all real-world factors affecting loan approvals.

# COMPARISON TO EXPECTATIONS

It performed well beyond expectations for its accuracy, though it was somewhat surprising how credit score was not as influential relative to other financial metrics when making decisions on loans.

# RECOMMENDATIONS

## MODEL DEPLOYMENT

The model can be directly used in a real-time loan approval systems to assist in making quicker decisions.

## MODEL OPTIMIZATION

Further optimization of the random forest algorithm can be supported by tuning the hyperparameters so that the model performs with higher precision in cases of edge conditions.
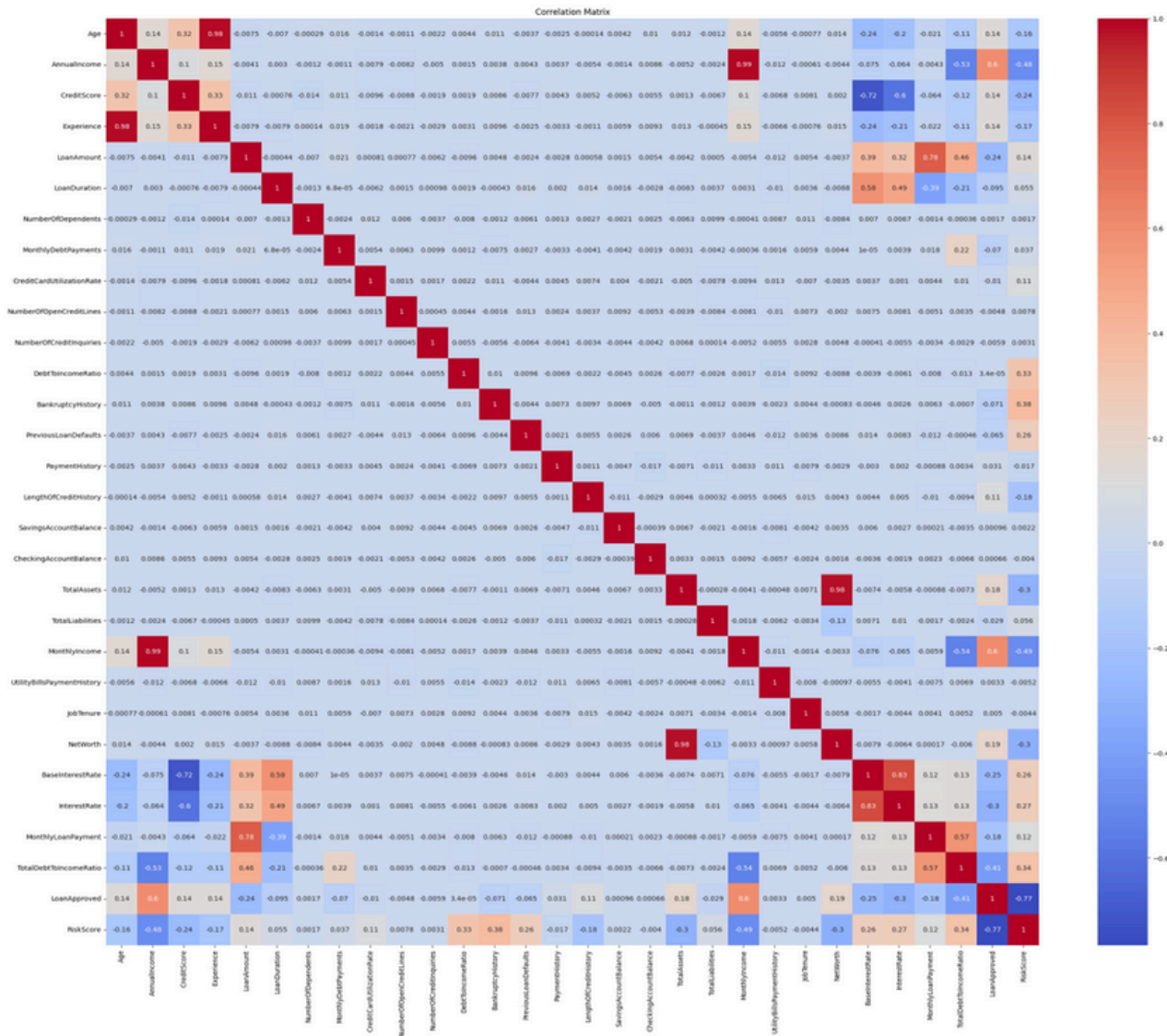
## FEATURE AUGMENTATION

More data could be included, such as customer behavioral patterns and market trends, to make the model predictions even better.

# CONCLUSION

The performed analysis successfully developed a machine learning model that can predict loan approvals with high accuracy. The insight from feature importance, such as the influence of credit scores and employment status, is a helpful guide for financial institutions in making decisions. Although the performance was quite good, improvements will be made for future work with regards to hyperparameter tuning and balance of classes for the extension of its practical applicability.

# APPENDICES



Correlation Matrix

# Loan Approval Analysis Dashboard

## Average Loan Amount
# $19,145K

## Average age of Applicant
# 42,7

## Average Loan Duration(Mont...
# 49,90

## Average Credit Score
# 585

### Loan Approved by Number of Dependents

8,72%
20,5%
41,78%
29%

No. of Depe... ● Mortgage ● Rent ● Own ● Other

### Loans Approved by EducationLevel

8,79%
17,24%
33,74%
17,85%
22,38%

EducationLevel ● Bachelor ● Master ● High Sch... ● Associate ▶

### Loans Granted Categozied by Purpose

9,37%
15,92%
29,5%
20,48%
24,73%

Purpose ● Home ● Debt Consolidati... ● Auto ● Education ● Other

### Loans Approved by Risk Score



### No. of Loans Approved by Length of Credit History(Years)



### Loans Approved vs Credit Score



### No. of Loans Approved by Experience