The background is a dark navy blue. On the left, there is a large, light green parallelogram and a blue parallelogram below it, both with black outlines. In the bottom left, there is a circular inset showing a close-up of a circuit board. In the top right, there is a grey, 3D-rendered circuit board pattern.

Credit Cards Payment Default Prediction

Capstone Project by Violeta Furculita

April 2024



Context and Problem Overview

Context:

- Credit cards are lines of credit extended to customers to use for everyday purchases
- Most customers make monthly payments toward the principal and interest charged, where applicable. Some customers may forgo making payments due to financial distress or intent to commit fraud
- Credit cards are a major revenue stream for financial institutions
- Even small increases in default rates can lead to significant losses
- Being able to accurately predict which customers may default allows targeted risk management

Problem:

- Conventional methods of predicting likelihood of default at underwriting or during usage may not be well equipped to predict default
- Identifying at-risk customers early allows contact and/or appropriate loss management
- However, default risk depends on complex interactions of many customer attributes




Understanding the problems

- 01 Default costs are high for both the bank and the customer. The bank may experience losses, lost revenue, collection costs. The customer may experience credit rating negative impact
- 02 The conventional approach such as FICO is not well equipped to predict default. Lost of other factors such as demographic, employment, recent behaviors, credit limits, balance levels etc factored together may be able to help.
- 03 Machine learning is well-suited to solve the problem of default prediction because it has the ability to look at vast troves of data and make predictions that can help at underwriting and during ongoing credit management.



Project objective



Analyze public credit card default data and build machine learning models that can help predict payment default rates on credit cards with reasonably high accuracy of at least 75-80% accuracy, which would be a significant improvement over existing methods.



Target audience

Financial institutions that offer credit cards either directly or through an intermediary to consumers.



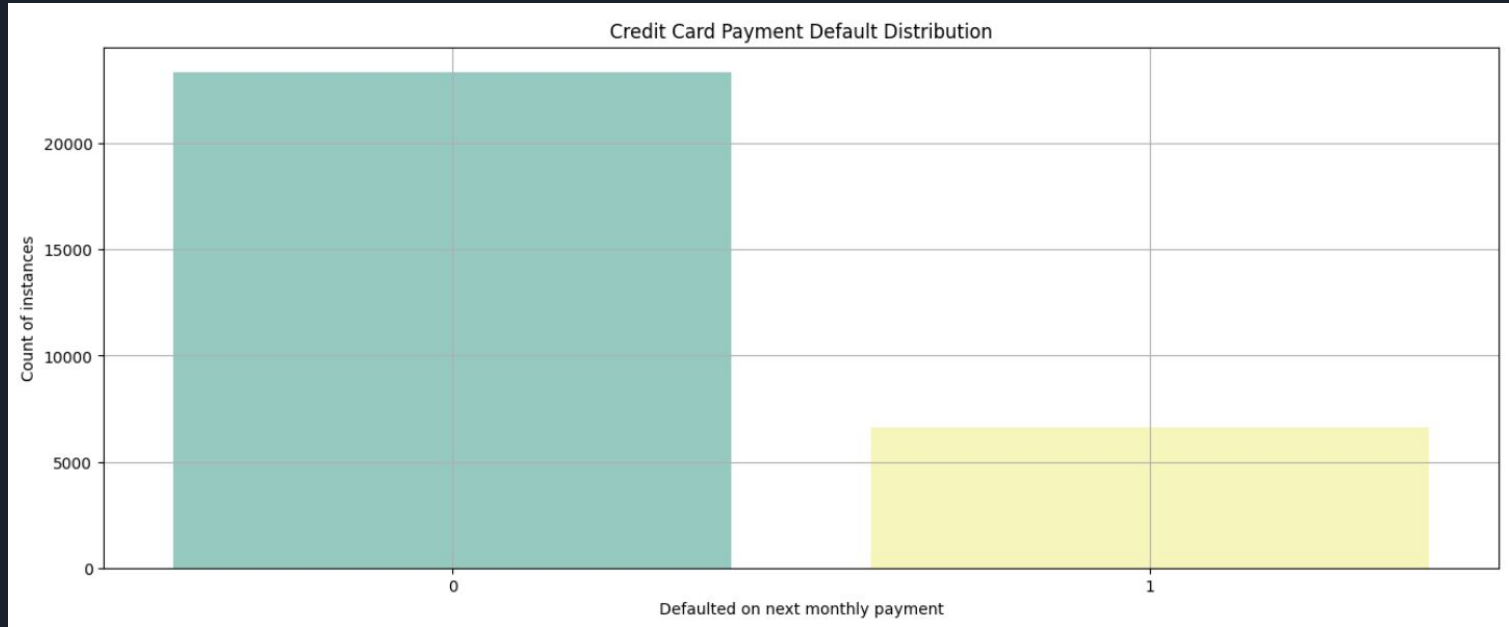


Project Methodology

1. Clean, explore and visualize the credit card default dataset to understand trends, patterns and relationships with default behaviors.
2. Engineer relevant features from the raw data for model building.
3. Evaluate various classification algorithms like Logistic Regression, Decision Trees, Neural Networks , K Nearest Neighbors and Gradient Boosting ensembles like XG Boost to identify models that perform best on this problem.
4. Tune model hyperparameters to optimize for key metrics like accuracy and train time.
5. Analyze feature importance to understand driver attributes of default risk.
6. Deploy the best performing model(s) and evaluate on test data to assess predictive power.
7. Present model performance, key findings and recommendations to stakeholders on how predictive analytics can help mitigate future credit card default rates.

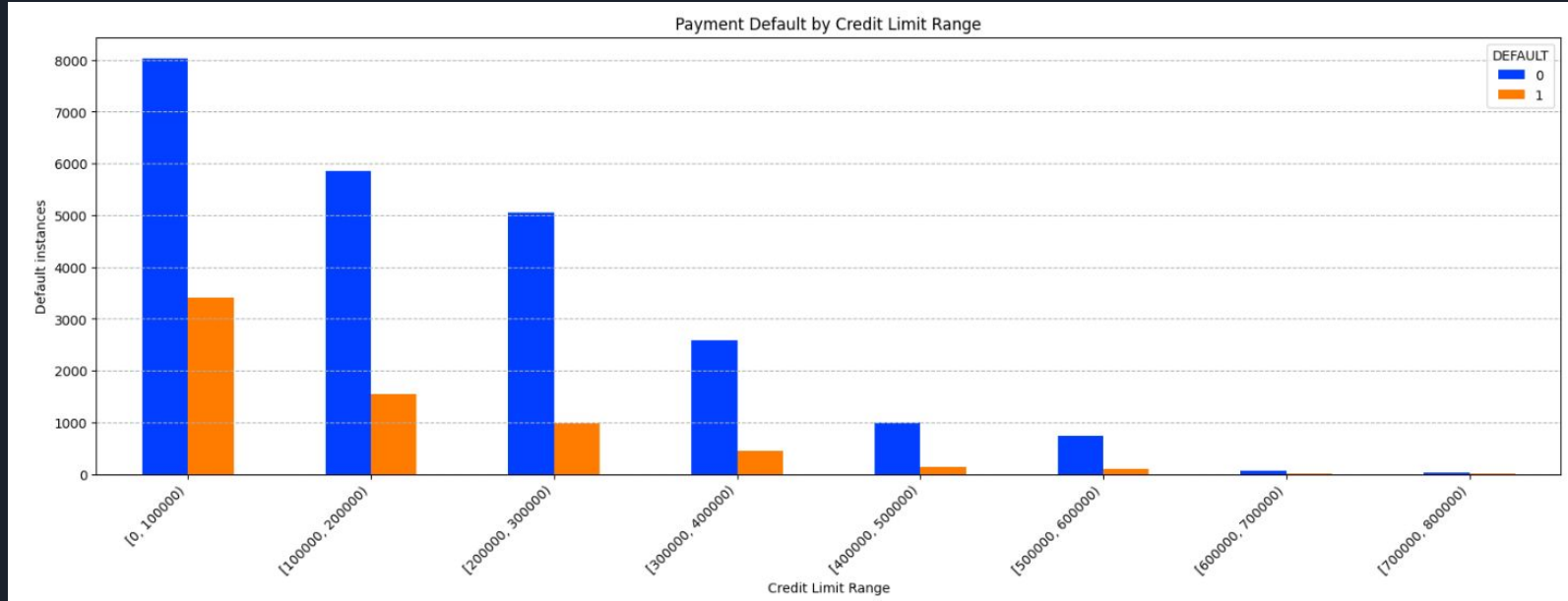
Overview of initial credit card user data

22% of customers defaulted on their credit card payment



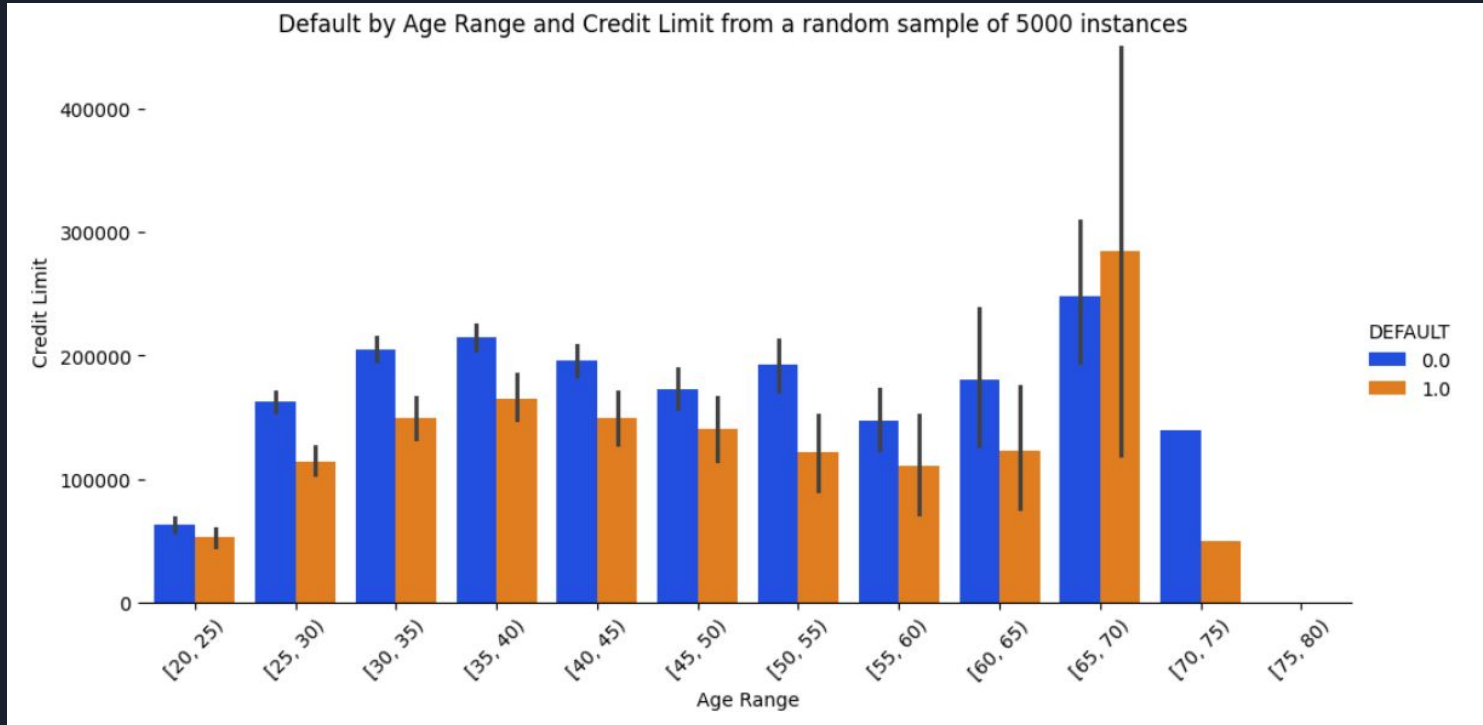
Overview of initial credit card user data

Users with credit limits under 200k default the most



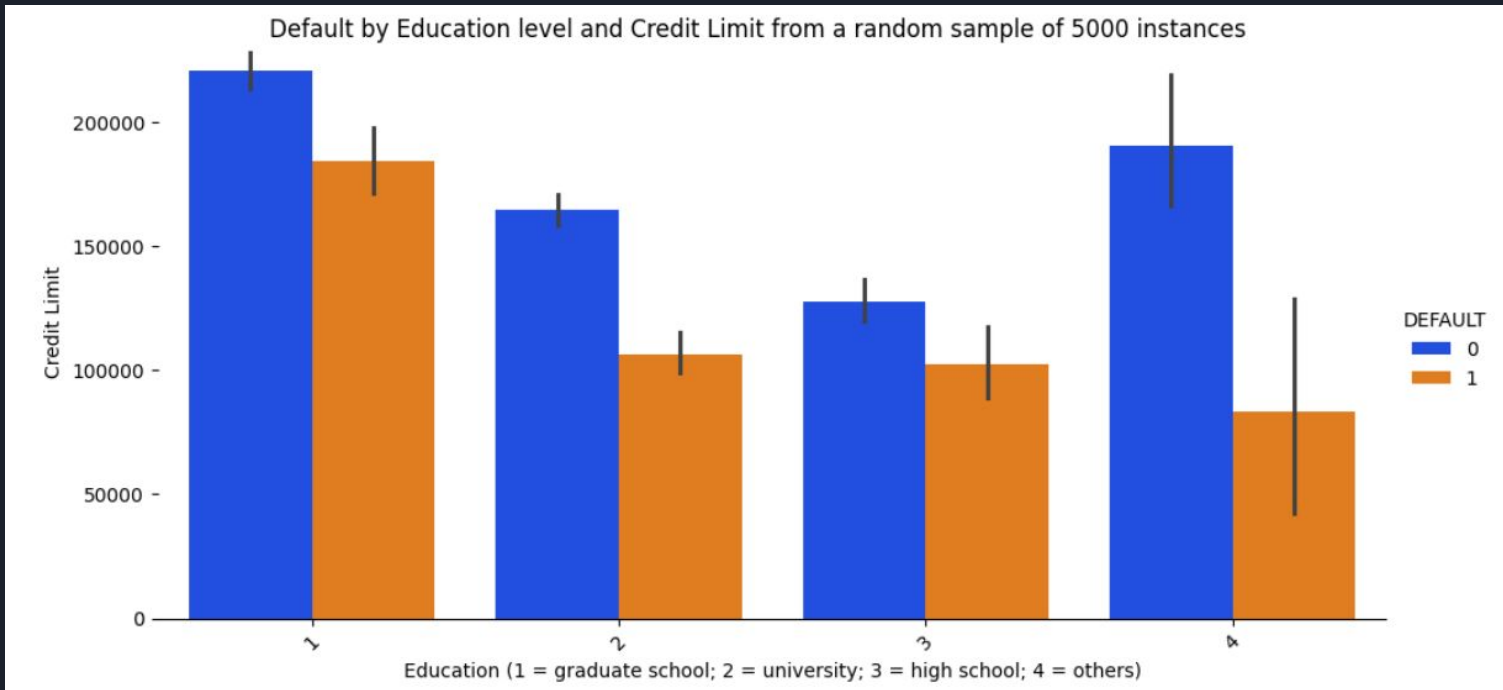
Overview of initial credit card user data

Users with ages between 65-70 and credit limit above 200k have a high default rate



Overview of initial credit card user data

Users with graduate degrees and credit limits above 100k tend to default more than users from other education groups





Preliminary Findings From User Data

- Payment default rate is high at over 20% of the user base
- Users with credit limits under 100k have the highest default rate followed by those with credit limits under 200k
- For user with ages higher than 65 and credit limits higher than 200k the default rate is higher as compared to sage group age and lower credit limits
- Most defaults happen in the population aged 25-55
- Users with credit limits higher than 100k and graduate degrees tend to default more than users with the same credit limit and different degree types



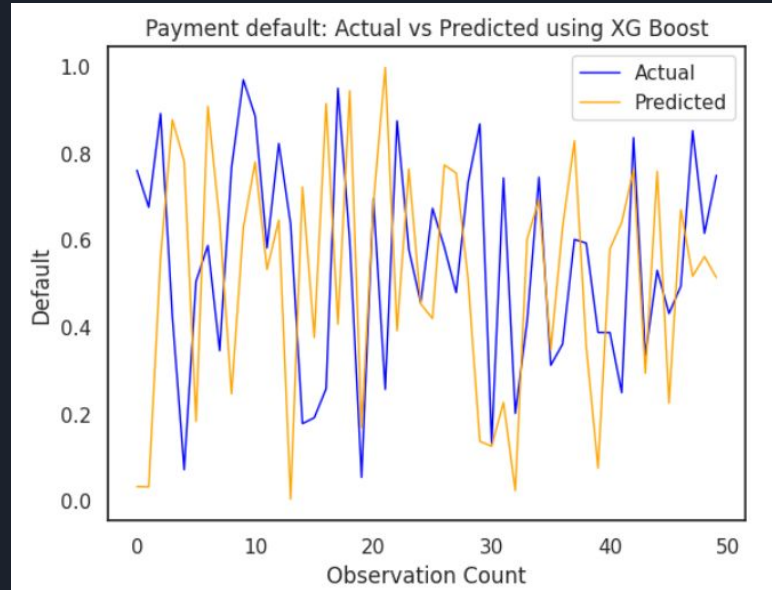
Preliminary Findings From Building Models

Using supervised learning models we can achieve robust accuracy and training times for our models that predict payment default. Specifically XG Boost and Neural Networks can produce prediction accuracy of over 82% with fast training time.

	Train Time	Train Accuracy	Test Accuracy
Logistic Regression	0.241966	0.777238	0.777238
Decision Tree	0.00915127	0.792452	0.772143
Neural Networks	0.121528	0.820048	0.820048
Support Vector Machines	111.48	0.782262	0.774095
K Nearest Neighbors	1.25893	0.777393	0.777333
XG Boost	0.313887	0.837429	0.822524

Preliminary Findings From Building Models

Randomly selecting 50 users from the test set, we showcase the predictive power of 82.25% of the XG Boost model ensemble.





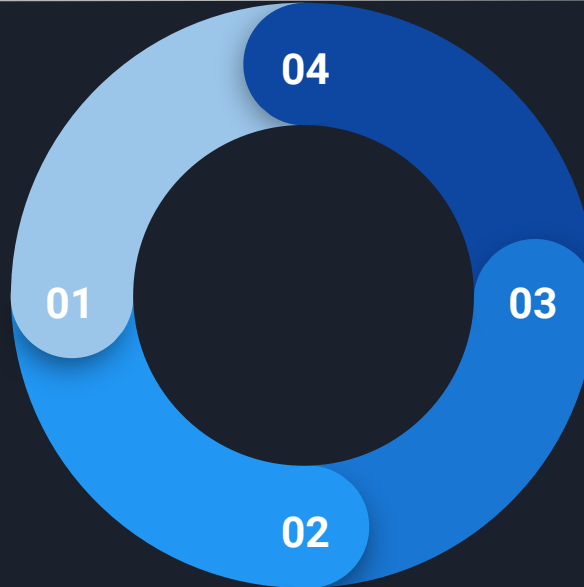
Recommendations and next steps

Fine tune models

Further model fine tuning is required to achieve higher accuracy rate on predictions

Use more data

Consider looking at more categorical data about users' behaviors such as income, assets, real time employment status...



Deploy

Select a model(s) that deliver the best accuracy and train time to deploy and be used by the credit card company

Get feedback

As more user data gets collected, ingest it in the models deployed to check whether predictive accuracy increases