

Model Efficacy in Credit Risk Predictions

Foodies: Bosia N’dri, Zanderz McCluer, Sailesh Pulukuri, Reagan Todd

AI/Data science track

Introduction

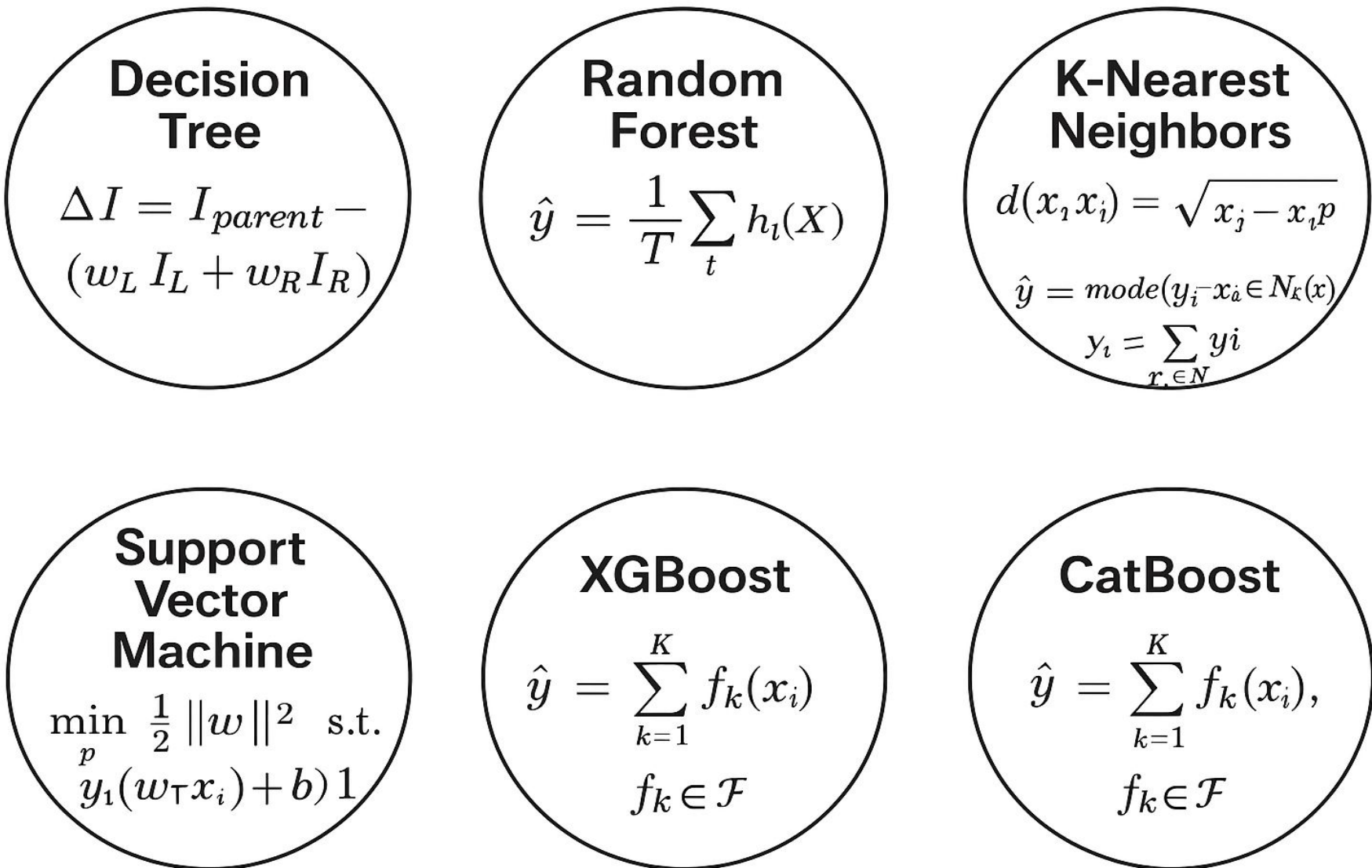
For the past several decades, credit risk estimation has been utilized to evaluate whether or not a borrower would pay back the lender^{1,2}. While the lenders want to ensure they get their money and resources returned, borrowers want to ensure that they can pay for daily goods, purchase homes or vehicles, or make other large purchases. These estimations follow the three C’s, which include character (a borrower’s reputation for repayments), capacity (a borrower’s ability to repay loans based upon income), and collateral (assets a borrower may have that can be pledged in the case of default).³ Characteristics that are prohibited from being negatively included in the evaluation of credit are demographic factors such as race, gender, or age.¹

Historically used models

Logistic Regression:

$$P(Y = 1 | X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p)}}$$

Models that would increase ML use:



There have been various ML techniques used in prior literature, including k-NN, Tree-related methods, Boosting, Bagging and deep learning methods that include Neural Networks (artificial, recurrent, or convolutional).^{4,5} To adopt a model in practice, there must be human agency involvement, transparency of the model, and accountability for the model, which are each key components to explainable artificial intelligence.⁵

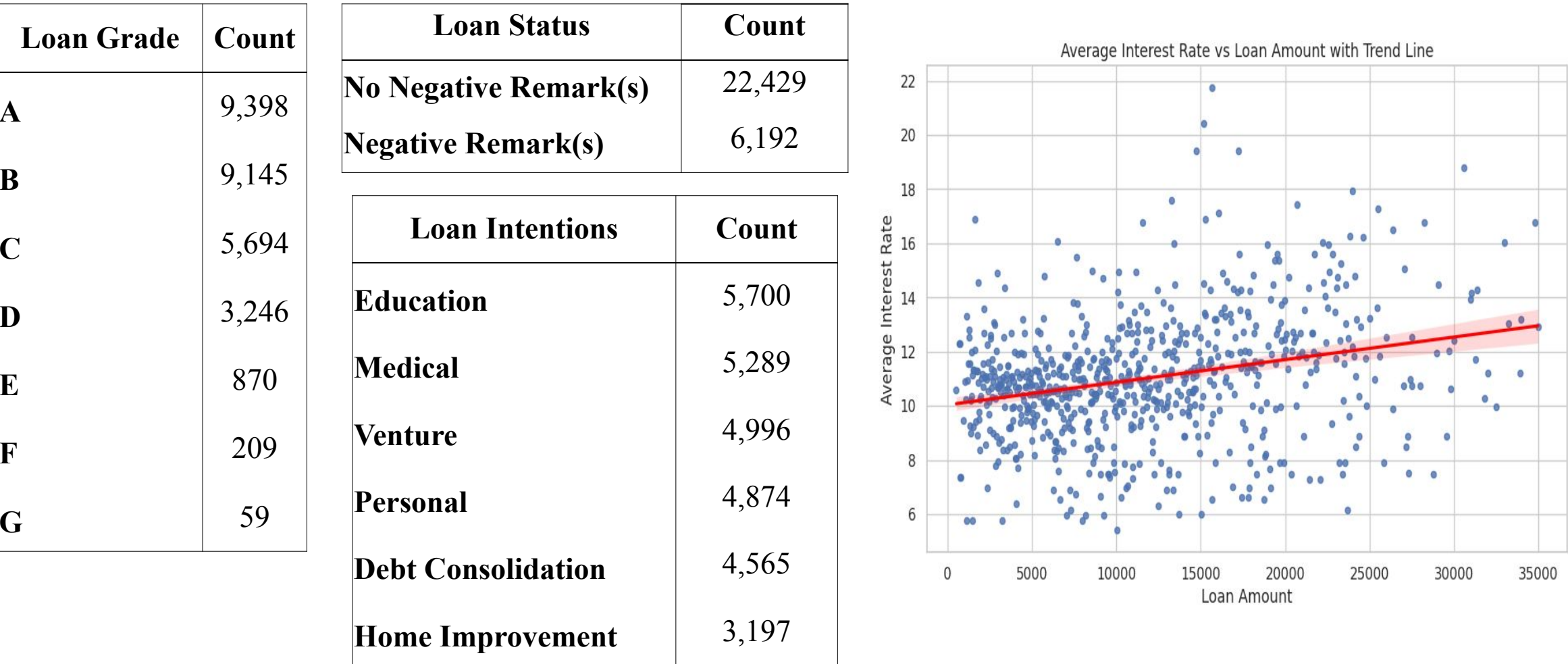
Problem Statement

How can machine learning improve the credit risk prediction process?

The purpose of this work is to employ various predictive models to illustrate how ML and AI techniques can work to create a more robust credit risk prediction model. Comparing the various ML models to a standard Logistic Regression will help understand the impact of ML onto credit risk prediction. By evaluating these models, we hope to outline the important factors contributing to a person’s credit risk as well as provide a good model for lenders to utilize to assess their borrowers.

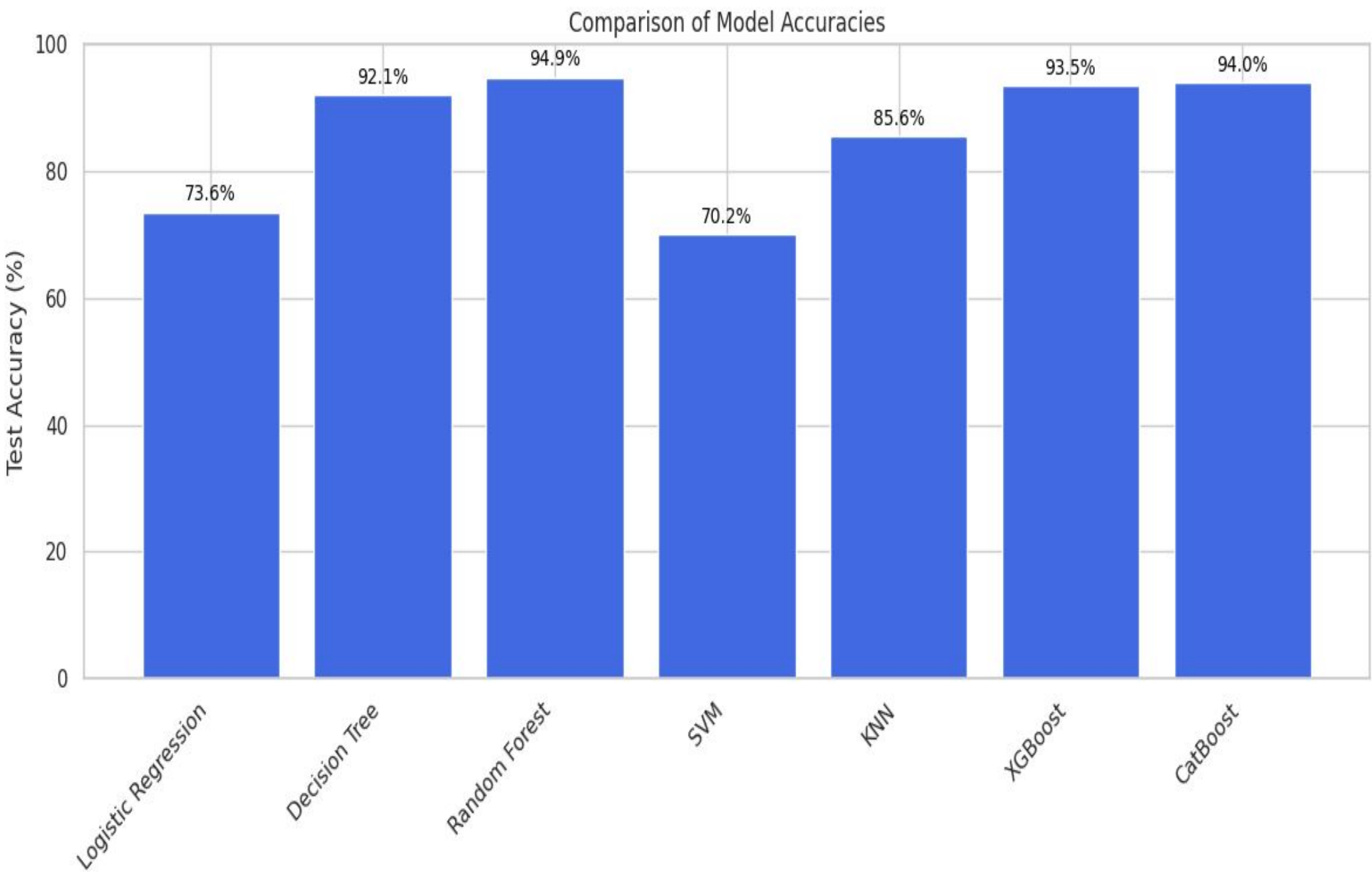
Results

From the data, most individuals had a moderate to positive credit history and requested smaller sized loans. The most variation in the data came from a person's income, employment length, and credit history length. From the data, it was also apparent that interest rate and loan amount have a positive relationship; as the loan size increases, so does the interest rate.



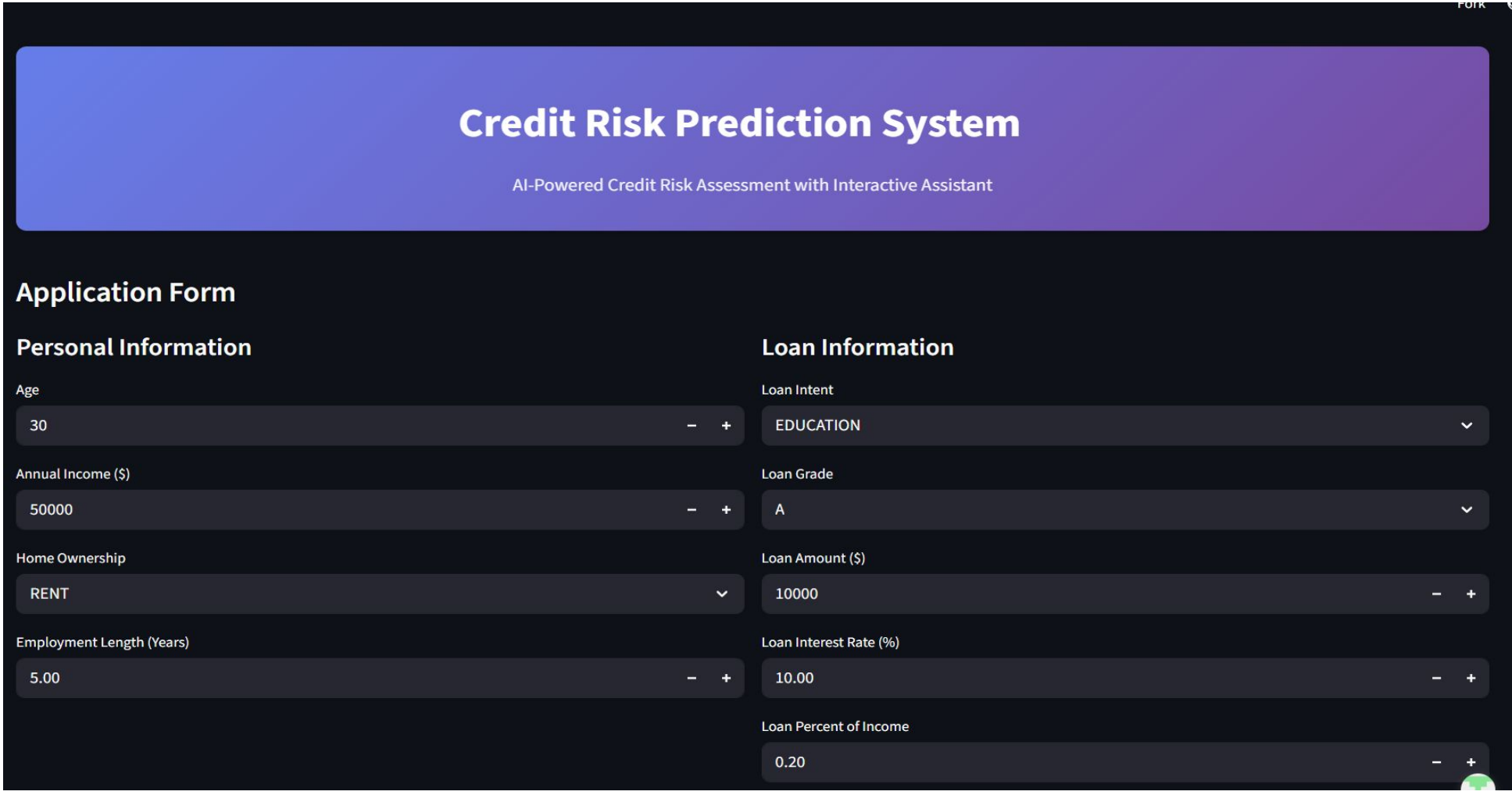
Characteristics of the Study Sample			
	Mean (SD)	Min	Max
Age	28 (6)	20	84
Income	66,435 (51,522)	4,000	2,039,784
Employment Length	5 (4)	0	41
Loan Amount	9,651 (6,318)	500	35,000
Interest Rate	11% (3%)	5%	23%
Loan Percent Income	0.2 (0.1)	0.0	0.8
Credit History Length	6 (4)	2	30

The Random Forest revealed an accuracy of 95%, which was the highest accuracy revealed within the seven models. Visualized in a confusion matrix, this model identified true positives and true negatives with high precision, showing only a small number of false predictions (120 false positives and 340 false negatives). These results suggest the model effectively distinguishes between classes, with a slight tendency to underpredict positive cases.



Methodology

The dataset from OpenML (32,569 records) was processed in Google Colab, where 4,011 missing values and outliers (e.g., age > 90 years) were removed. Categorical variables such as home ownership, loan intent, and default history were one-hot encoded, and derived ratios like loan-to-income and interest-to-loan-amount were created. Data was split 80 % train / 20 % test. Seven ML models—Logistic Regression, Decision Tree, Random Forest, SVM, KNN, XGBoost, CatBoost—were evaluated using accuracy, precision, recall, F1-score, and confusion matrix. The Random Forest model achieved the best results (Train = 0.869, Test = 0.868) and was deployed via a Streamlit web app with an interactive chatbot assistant for real-time credit-risk interpretation.



Conclusion

Among the models evaluated for predicting credit risk, Random Forest demonstrated the highest predictive accuracy and reliability, effectively distinguishing between high- and low-risk borrowers. Its performance suggests that lenders can leverage this model to make more informed lending decisions, reduce default rates, and optimize resource allocation. While simpler models like Logistic Regression offer interpretability, the chosen model balances both accuracy and practical utility in real-world credit risk assessment.

This work would suggest that ML has the ability to enhance credit risk predictions, while keeping the lender and borrowers best interests in mind. Future direction of this work should be put on testing other datasets and disseminating the results to lenders for use.

References

- FRB: Report to the Congress on Credit Scoring and Its Effects on the Availability and Affordability of Credit. August 2007. Accessed October 16, 2025. <https://www.federalreserve.gov/boarddocs/rptcongress/creditscore/demographics.html>
- Thomas LC. *Consumer Credit Models: Pricing, Profit and Portfolios*. OUP Oxford; 2009.
- The Three C’s of Credit (Lesson 9A). Accessed October 16, 2025. <https://www.federalreserveeducation.org/en/teaching-resources/personal-finance/managing-credit/the-three-c-s-of-credit-lesson-9a>
- Shi S, Tse R, Luo W, D’Addona S, Pau G. Machine learning-driven credit risk: a systemic review. *Neural Comput Appl*. 2022;34(17):14327-14339. doi:10.1007/s00521-022-07472-2
- Bussmann N, Giudici P, Marinelli D, Papenbrock J. Explainable Machine Learning in Credit Risk Management. *Comput Econ*. 2021;57(1):203-216. doi:10.1007/s10614-020-10042-0

Instructions

- **Please do not change the size or font of the template**, but you are allowed to change the color or size of the font and/or the background color.
- You can change the section headings or delete them in order to best present your research. For example, you may decide to delete the summary section.
- You can change the size of each section box for pictures/graphs, etc. However, all section headings need to be in alignment with each other. For example, you can move the credits section further down on the template and expand the conclusion section.
- Please list the person that is presenting the poster as the first author.
- Do NOT Resize the template. Template is sized accurately for printing. Please include logos for all involved institutions in the header **and delete the second page before printing.**

Important Information!!!

- In addition to this PowerPoint file, **please fill out an order request form and attach it to the email when submitting your file for printing.**
- Please submit this file along with your order request form at least 48 hours prior to the date that you need the poster by.
- **Note that if you submit your request within 24-48 hours of needing the poster, you run into the risk of not having it completed by the time you need it.**
- **Each poster requires its own order form.**