

Not Approved's Plan for Getting Your Loan Approved

John Morillo, Mohil Patel, Cinah Pourhamidi,
David Rohweder, Oleg Sapranenko, and Alec Sudol

Department of Computer Science
The Pennsylvania State University
University Park, Pennsylvania, USA

Abstract

Throughout this report our main goal is to investigate the loan approval process and how Artificial Intelligence is present in it, as well as create our own model that will be accurate, transparent, and contain very little bias. We will focus on how we can answer three questions: (1) How can we make our model fairly accurate? (2) How can we make sure that our model is transparent? and (3) How can we keep bias out of our model? After answering these questions, we will be able to understand how our model can be used in the future.

Keywords – Loan Approval, Neural Nets, RBES, AI, Transparency, Bias

I. Introduction

Our task in this report is to come up with a solution to the topic of Artificial Intelligence (AI) and Capitalism. When trying to decide which avenue to take when coming up with a solution, we decided to focus on loans and the loan approval process. We were

particularly interested in automating this process and making it as accurate and transparent as possible. Part of the reason for this is many banks use AI to help them with their loan approvals; however, what they are lacking is transparency. This can cause people who are highly qualified for a loan to be rejected and not know why they were rejected. This problem is relevant because people may be rejected due to issues with the model, whereas, if a human was doing this process they would be approved. Any person applying for a loan would want a system that is the most accurate and least biased when applying for their loan, so as to maximize their chances of getting the loan. With many of us going into the real world after graduation, we all wish to see a loan approval model that is most like the one described above, so that we have the highest chance of getting approved. Thus, the problem that we are trying to solve is making the loan approval application process for housing loans to be more

accurate, more transparent, and less biased for the users of this system.

Our main objective of this paper is to answer 3 major questions. To answer these questions, we will be using a dataset from Kaggle called Loan Data Set [1]. The dataset is composed of 13 columns that includes the features Loan_ID, Gender, Married, Dependents, Education, Self_Employed, ApplicantIncome, CoapplicantIncome, LoanAmount, Loan_Amount_Term, Credit_History, Property_Area, and Loan_Status. It includes a total of 614 records that we will be using to train our model. Our first main question is how can we make a model that is fairly accurate on approving home loans. Then we will attempt to answer the question of transparency, and how can our model be transparent about what it is doing and its selection process. And finally, we will try to combat bias in our model and attempt to make a model that contains the least amount of bias.

Next, we will address how our paper is laid out and what to expect in each section. Section I is the introduction that focuses on the background of why we chose this subject and why it is relevant to us and the world. It also aims to lay out the problem statement

and objectives, while including some information about the data set that we used. We will also foreshadow our solution in the introduction. In Section II, we will focus on addressing the related works and background of Home Loan Approvals and AI. This section will include information on how home loans work, the process behind them, and different approaches and models that have been used previously in this topic space. Then in Section III Experimental Setup, we will address 3 main topics: our project setup, our project execution, and our project limitations. This section will be a good mix of how we were able to go about creating our model and how we set it up. Section IV Results will focus on the results of our model and display some visuals that we were able to gather when creating the model. In Section V Analysis & Discussion, we will attempt to provide an analysis of how our model compares to other models that were previously tried and mentioned in the Related Works section and discuss some ethics and how we managed to answer the transparency and bias questions. We will conclude with Section VI Conclusion and Future Work that will conclude our findings and suggest ideas that we can explore in the future. At the end a list of references will be provided that we used in our report.

A bit about our solution to our problem. We decided to implement a neural network that takes in the dataset and preprocesses it to eliminate bias that could be tied to certain features. This neural network runs and outputs whether the user got approved or rejected for the loan. Then we have another neural network that predicts the interest rate if the person got approved. Finally, a Rule Based Expert System (RBES) is used to categorize the interest rate for that individual.

II. Related work & Background

In the financial industry, there has been increasing demand for loans to be distributed from one party to another. This first started with individuals loaning funds to one another in a “handshake agreement” method where one individual agreed to loan to another by trusting that they would receive their money back with terms of interest [2]. This trust was the fundamental basis of creating and distributing loans as the quantity of loans at the time was manageable for the lenders as the amount of prospective borrowers was feasible for the process to be managed in this way.

As time progressed the needs of the economy changed and commercial banks started replacing individuals as the sources of lending money to the public. This was first initiated as banks lent money to people that did not have the proper exchange currencies within their respective markets and geographic locations [3]. This notion was furthered by the increase of popularity of commercial banks lending money to the public, as it created the need for federal regulation within the financial sector to ensure that large commercial banks were not taking advantage of their ability to provide loans to those in need as well as not crippling the economy by not leveraging loans to an extent where individuals would not be able to pay them back [3].

Further, this idea of regulation was deemed necessary as instances of loan abuse by the commercial banks became more apparent and had a direct impact on the economy. One of the most significant examples of this was the housing crisis of 2008 in which commercial banks extended offers to approve mortgages for borrowers that did not meet the criteria to be able to repay these loans so they defaulted on a lot of the processed mortgages approved by the banks,

causing the economy to take a significant negative hit [4].

In the present day, a lot of companies utilize technology to automate their loan approval process. There has also been a lot of research conducted within the field of artificial intelligence to apply certain aspects within it to automating the loan approval process. In the paper “Automated Loan Approval for Banks, by Abdulrahman Saeed Almheiri”, the author discusses the use of different models and tools such as random forest, neural networks, logistic regression models, and confusion matrices for the application of loan categorization as well as loan calculation [5]. The author then goes on to explain the selection of their features within their respective dataset and their importance within the calculations [5]. They utilized some attributes of borrowers such as gender, marital status, education, employment, and property district to identify significant values for their models to learn from within their dataset [5]. We observed the approach of this paper within our own research and findings to formulate our own approach and identify key points of improvement or variance to appropriately alter necessary constraints/requirements to fit the scope of our problem.

III. Experimental Setup

A. Project Setup

The code uses a neural network in order to evaluate the different input from our dataset. The neural net evaluates the different input categories to determine the person's eligibility to be able to comfortably afford and be able to pay back the loan that they requested. Another algorithm we use is a Ruled Based Expert System to evaluate the interest on the loan given the input predictions that the neural network made in order to assign that applicant a "riskiness" factor & choose an appropriate percent interest on the loan.

B. Project Execution

The model should be used with a csv file that is structured the same as the input file with separate columns for each of the evaluating categories for the neural network. The ideal environment is one that adheres to the training data & ideally expands upon it with more information about the applicant's status and other objective financial categories to analyze eligibility. Our model can adapt to the environment & add extra hidden layers to be able to process the new categories and

make updates to the decision-making of the neural networks. More parameters can also be taken into consideration on the RBES system and further expand the interest rate categories of the model to give maximum flexibility and fairness to each applicant's interest rates.

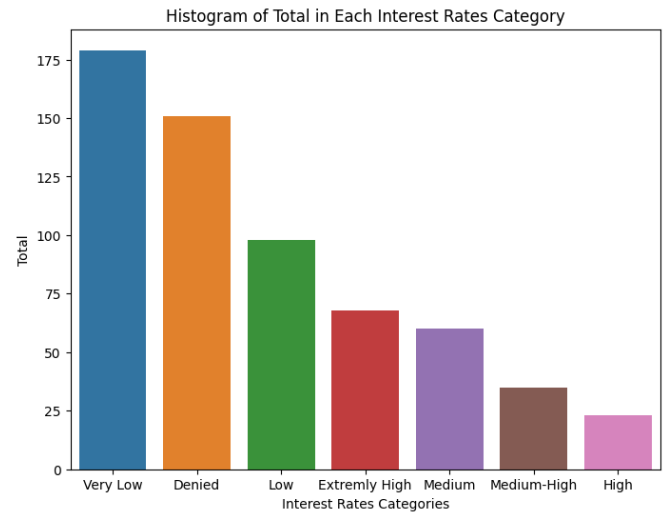
C. Project Limitations

Some limitations of our projects are unfortunately the training data. Due to the lack of availability for this type of dataset, the training is inadequate and yields less than ideal accurate percentages for test data sets.

IV. Results

In this section, we present the insightful results gathered from the `predicted_data.csv`, which shed light on the crucial factors influencing loan acceptance. By analyzing the distribution of interest rate categories, income levels, and loan durations, we gain valuable insights into the patterns and impacts that these variables have on the loan approval process. Understanding these results is essential for borrowers seeking loan approval and lenders evaluating loan applications, as it enables informed decision-making and enhances the overall

efficiency and effectiveness of the lending process.



The histogram based on the `predicted_data.csv` provides important information about the distribution of interest rate categories and their influence on loan acceptance. It reveals that the majority of loan applications in the dataset fall into the "Very Low" interest rate category, indicating a higher probability of loan approval for borrowers in this group. This emphasizes the significance of securing a lower interest rate as it signifies a stronger financial profile and increased chances of loan acceptance.

The presence of a significant number of loan applications in the "Denied" interest rate category highlights a sizable portion of applicants who displayed unfavorable characteristics resulting in loan rejection. This underscores the importance of fulfilling lender requirements, maintaining a positive

credit history, and demonstrating financial stability to enhance the likelihood of loan acceptance.

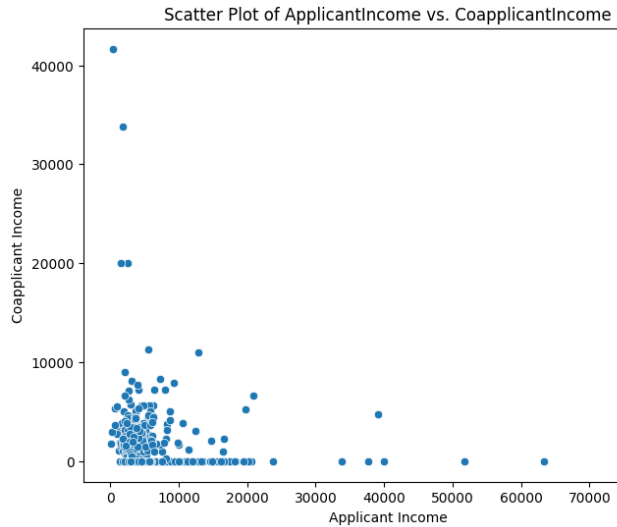
The frequencies of other interest rate categories, such as "Low," "Medium," "Medium-High," and "High," offers further insights into loan acceptance patterns. These categories represent borrowers with varying interest rates, indicating a range of credit profiles and risk levels among approved borrowers. This suggests that loan applications with different interest rates were accepted, emphasizing the need for borrowers to understand the factors considered by lenders beyond just the interest rate.

The histogram provides valuable insights into the distribution of interest rate categories and their impact on loan acceptance. It emphasizes the importance of securing a lower interest rate, meeting lender criteria, and understanding the broader factors considered in the loan approval process. This knowledge empowers borrowers to make informed decisions and take appropriate steps to improve their credit profiles, while also aiding lenders in assessing risk and making sound lending decisions.

**Confidence Intervals (95%):
(118.91683692225666,
150.58137046880785)**

The 95% confidence intervals provide a range within which we can be 95% confident that the true average loan amount for the "Low" interest rate category lies. In this case, the lower bound of the confidence interval is 118.91683692225666, and the upper bound is 150.58137046880785. This means that based on the available data, we can estimate with 95% confidence that the average loan amount for individuals in the "Low" interest rate category falls between these two values.

The calculation of the confidence intervals considers the sample mean, sample standard deviation, and sample size of the loan amounts in the "Low" interest rate category. It assumes a normal distribution of the loan amounts, which is a common assumption in statistical analysis. In this case, the relatively narrow range of the confidence intervals suggests a relatively precise estimate of the average loan amount in the "Low" interest rate category. These confidence intervals can be useful in understanding the potential range of average loan amounts for the "Low" interest rate category. They can aid in decision-making, such as setting appropriate loan terms or evaluating the fairness of loan offer.

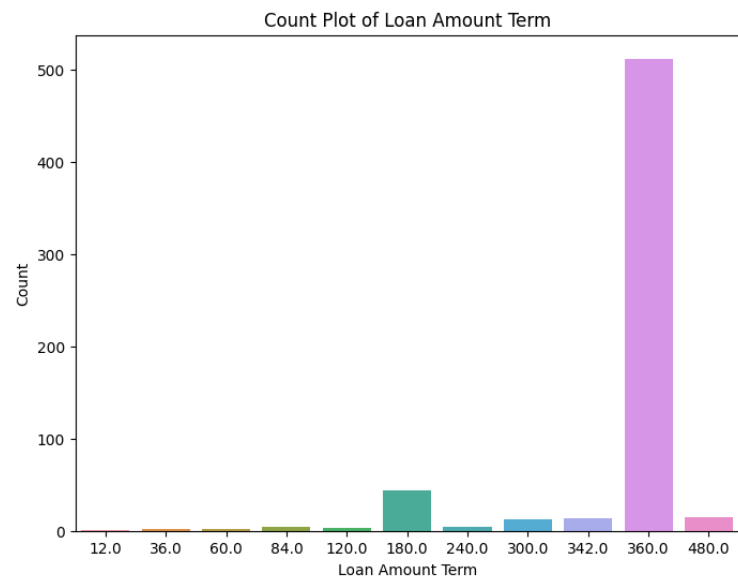


Based on the scatter plot, we can observe a significant concentration of data points and clusters within the range of 0 to 20,000 for the Applicant Income variable. Additionally, there is a notable accumulation of data points around 0 to 10,000 for the Co-applicant Income variable.

This concentration of data points in the 0 to 20,000 range for Applicant Income and 0 to 10,000 range for Co-applicant Income suggests a higher frequency of individuals with incomes within these ranges. It indicates that a considerable number of applicants and co-applicants fall into these income brackets. However, beyond these ranges, specifically for higher values of Applicant Income and Co-applicant Income, the density of data points diminishes significantly. This implies that there are relatively fewer instances where the incomes exceed these thresholds.

These observations indicate that the majority of applicants and co-applicants in the dataset have incomes in the 0 to 20,000 range for Applicant Income and the 0 to 10,000 range for Co-applicant Income. The scarcity of data points beyond these ranges suggests that higher income levels are less prevalent among the applicants.

Understanding the distribution of income levels can be valuable in assessing loan eligibility and making informed decisions related to loan acceptance. It provides insights into the income profiles of applicants and co-applicants and can aid in assessing their financial capacities and loan repayment capabilities.



The count plot of Loan Amount Term provides valuable insights into the relationship between loan acceptance and loan duration. By analyzing the distribution

of loan durations, we can gain a deeper understanding of their impact on loan acceptance rates.

The count plot reveals that loan terms of 180 and 360 months have a significantly higher number of instances compared to other durations. This suggests that borrowers who opt for these specific loan durations are more likely to have their loan applications accepted. Lenders may consider these durations as preferred options due to their perceived reliability and stability. Many borrowers seem to prefer longer loan durations, possibly because of the lower monthly payments associated with these terms. Lenders may take these preferences into account when evaluating loan applications, resulting in a higher acceptance rate for loans with these durations.

On the other hand, loan terms with shorter durations, such as 12, 36, 60, 84, 120, 240, 300, 342, and 480 months, show significantly lower counts in the plot. This suggests that loans with these durations may be less frequently accepted by lenders.

Shorter-term loans often carry higher risk, as the repayment period is shorter and monthly payments are usually higher. Lenders may exercise caution and conduct more thorough evaluations when considering applications with these durations, focusing on factors

such as creditworthiness and repayment capacity.

The findings derived from the analysis of the `predicted_data.csv` provide significant insights into the intricate relationship between various factors and loan acceptance. Additionally, examining income levels highlights the significance of meeting lender requirements and demonstrating financial stability to enhance loan acceptance prospects. The exploration of loan durations unveils the preferences of both borrowers and lenders, with longer-term loans being more frequently accepted due to their perceived reliability and lower monthly payments. By considering these factors in the loan approval process, borrowers can better position themselves to meet lender criteria, improve their financial profiles, and increase their likelihood of loan acceptance. Lenders, on the other hand, can leverage these insights to make more informed decisions, manage risk effectively, and ensure efficient allocation of resources. All in all, the findings from the `predicted_data.csv` provide actionable insights that empower both borrowers and lenders in navigating the loan approval landscape. By understanding the significance of interest rates, income levels, and loan durations, stakeholders can make

informed decisions, improve loan eligibility, and foster a more transparent and equitable lending environment.

V. Analysis & Discussion

In this section, we will delve into the intricacies of the machine learning model utilized for loan approval prediction. To achieve the most accurate model, two primary training approaches were considered: a brute-force method and a randomly optimal method, known as the `lazy_train` function. Notably, the `lazy_train` function differs from the brute-force method as it only has one hidden layer.

The model's architecture follows an input layer with nodes equivalent to the feature size, connected to a hidden layer with 16 neurons and a tanh activation function. The output layer has a single neuron with a sigmoid activation function for binary classification, mapping the result to either 1 (approved) or 0 (denied).

The training process consists of 1000 epochs, a batch size of 128, and the choice of Adagrad or simply “adam” as learning functions. The model's accuracy is evaluated against the testing dataset, and the best model is determined based on the highest accuracy achieved during training. The primary function, `train_model()`, implements a loop to optimize each layers node quantity

and after the most optimal batch size and then trains the model using the most optimized configuration. The `lazy_train()` function, on the other hand, relies on a more simplified architecture and training process. Both approaches resulted in similar results the highest the full brute-force approach was able to achieve was 80.12% accuracy in the data but took over full day to train.

The `lazy_train()` method however achieved 78.86% accuracy only taking several seconds to preprocess and train the model. The `main.py` script is the router and brain behind the application and determines whether to train the model, make predictions, or analyze the data. In the `preprocessing.py` script, the raw data is loaded, preprocessed, and split into training and testing datasets (80/20 split). The `approval_predictions.py` script loads the trained model, preprocesses new data, and generates predictions based on the model's output. By evaluating and comparing the two training approaches, we can determine the best model for predicting loan approvals. The chosen model should exhibit high accuracy and low loss in its evaluation against the testing dataset, proving its reliability and correctness compared to alternative models. This comprehensive analysis allows us to confidently deploy the

most suitable model for loan approval prediction.

In addition to predicting loan approvals with binary outputs of 1 or 0, the code provided also calculates the interest rate for each customer using a Rule-Based Expert System. This calculation leverages the raw output values from the prediction model to determine the appropriate interest rate range for each user.

The algorithm begins by loading the new data with the predicted loan approvals and corresponding raw predictions. The interest rates for each customer are calculated using the `rbes_range()` function. This function iterates over the raw predictions and computes the interest rate by subtracting the prediction value from 1 and multiplying the result by 100. Depending on the calculated interest rate, the customer is assigned to a specific rate range, such as "Very Low," "Low," "Medium," "Medium-High," "High," "Extremely High," or "Denied." Once the interest rate ranges have been calculated for each customer, the `new_data` dataframe is updated to include these values. The 'Raw_Predictions' column is dropped, as it is no longer necessary, and the 'Interest_Rates' column is added to store the assigned rate ranges. Finally, the updated dataframe is saved as a new CSV file,

'predicted_data.csv,' preserving the loan approval predictions and corresponding interest rate ranges for further analysis or decision-making.

VI. Conclusion & Future Work

A. Future Work

While working on our project and doing our report, we ran into a couple of new ideas that we thought would be interesting to investigate in the future. One topic that we thought would be interesting to investigate would be adding Explainable AI to our model. We thought that if we could get explanations to users as to why they were rejected, then users would have a better understanding of what they would need to improve on in the future, so that they would get accepted for their loan. This could also be used if someone was trying to get an appeal and they would be able to use examples that were provided to them to argue as to why they should get the loan approved. Next, we thought we should attempt to try and run our model with different datasets. We would need a way to preprocess that would only use the features that are relevant to our model and discard the other features. In theory, by training our data on more data sets, this will cause our model to have better accuracy. This would

be the case because the weights would be updated more times and more data would allow for more inclusive training. However, we would have to make sure that the data that we are adding to our dataset is not biased and does not contain any categories that could be used for bias. To tie this in with Explainable AI, if we have more data in our training dataset, when the model rejects a certain individual, then it would have more examples to compare them to and give more accurate explanations as to why someone was rejected. This would also allow for the rejected individuals to be given anonymous examples of similar individuals who were approved for loans.

B. Prescriptive Analysis

Based on our data, we believe that there is a certain course of action that we should take when choosing our future work. We thought that we should first consider improving our dataset as opposed to adding Explainable AI. Our accuracy for the first neural network was not where we wanted it to be, and we believe that we should first focus on improving the accuracy. To do that we will need to investigate different ways to improve our model. We can start off by adding more data sets and train the model that way. This would give it more data to work with and theoretically improve

accuracy. Then, if we notice that our accuracy is about the same, then we will try to consider changing the inner workings of our neural network. We will try to change the number of hidden nodes and hidden layers, to improve accuracy. Afterwards, we will compare our findings and arrive at the best possible combination of hidden node and hidden layers to improve accuracy.

C. Conclusion

To conclude this paper, we started by attempting to solve the issues in home loans that dealt with accuracy, transparency, and bias. To do this we created a 2 step model which included a neural network that predicted the approval of a loan and a second neural network that predicted the interest rate for that individual. We also combatted transparency and bias by including how our network works and getting rid of bias categories. We believe that our impact was that we were able to create a model that was accurate, transparent, and the least bias that we could make it.

D. Post-Thoughts

Utopia:

In a utopian environment, our loan approval model would work to service more users at a more efficient rate. This would allow more

applications to be reviewed in a shorter period of time, increasing turnover rates and therefore ensuring that people receive their funds needed earlier. This model would also strive to eliminate sources of bias from individual commercial banks which would encourage users of all backgrounds and attribute types to apply without the worry of being excluded due to sources of bias.

Dystopia:

In a dystopian setting, our loan approval model would work to create bias towards a specific category within its data. The model would learn to approve a certain subset of users that is based on attribute values that pertain less to financial metrics (such as race, marital status, geographic location, etc.). By having the model learn from these values, it will be encouraged to learn how to base its assessment of approval to favor certain categories of the population which would be detrimental to its success and function within its scope of the economy.

VII. References

- [1] [Loan Data Set | Kaggle](#)
- [2] <https://www.koho.ca/learn/history-of-lending/>
- [3] <https://portal.ct.gov/DOB/Consumer/Consumer-Education/ABCs-of-Banking---Banks-and-Our-Economy#:~:text=The%20first%20American%20banks%20appeared,issued%20notes%20for%20money%20deposited.>
- [4] https://www.aei.org/special-features/government-mortgage-complex/?gclid=Cj0KCQjw3a2iBhCFARIsAD4jQB33w3YJ2PSB1a9mjaqcbLqRH7mNgYJrjF-5tsiz4hTSizUACF8H7jIaAkD0EALw_wcB
- [5] <https://scholarworks.rit.edu/cgi/viewcontent.cgi?article=12535&context=theses>