

Data Analytic Methods for Institutional Discrimination Detection in Finance Applications

Carlton Wilson
cwillso39@gmu.edu
George Mason University

Abstract—Banks are slowly recovering from the COVID 19 pandemic. This is the effect of The American Rescue Plan of 2021. As a result, consumers are spending and putting money back into the economy. As we begin to see light at the end of the tunnel, this is good news for the banking industry and especially beneficial to the working-class people. Based on our research, we've discovered that there are unorthodox approaches and factors that can contribute to help mapping and identifying the risk when calculating the credit score.

Index Terms—creditworthiness, institutional, discrimination, styling, insert (key words)

I. INTRODUCTION

To understand creditworthiness, first, we need to understand the meaning of this concept. Creditworthiness is based on how the borrower handled debt and credit. Creditworthiness is how a lender decides if the person or company who requests for money can repay the loan that will be borrowed. The first step to get a loan is to complete and fill an application. Based on the applicant's current credit score, the lender takes into consideration how likely the applicant will repay the obligations of the debt on time. Once the lender determines if the applicant

[1], [2] is deemed not a risk and is worthy of credit, the decision will be made if eligible to get the loan. There are some very important factors other than credit score that determine the approval of a loan application status. In other words, proving to the creditors and lenders that timely payments will be made will establish trust with any applicant. Throughout this project, we will attempt to examine other factors and find some patterns using datasets to determine who can establish creditworthiness, without the use of institutional discriminations [1].

The problem we have discovered and will discuss in this project is that there more than 45 million Americans, who don't have a credit score. [3] They are unable to borrow money or use credit cards because they do not appear in the credit scoring system. [4] In another way, one in five Americans have no traditional credit score and are not eligible. We understand that it is a major socio-economic problem which needs immediate attention. [5]

Just because there is no established credit score, many times people will go for alternate options, such as increased interest rate loans such as payday loans, title loans etc. and even they are not available in many cases. [6] This intern could lead them into more and more troubles and could make them default

some payments [7]. [6]–[8]. Machine Learning approach for finance [3] [4] [9]–[20] [5] [7], [10] [6] [6] and Social Media [7], [8], [21]–[29] improve prediction results.

II. LITERATURE REVIEW

Through various references, it is clear that the current credit lending system is broken due to various reasons. [30], [31] Within their research, Kumar Arun, Garg Ishan, and Kaur Sanmeet [32] demonstrated that it is indeed possible to use machine learning to predict loan approval odds. Machine Learning can determine new relationships that a person would never think to test. This is the legality of ethics of using machine learning. Another such attempt is done by Mridul Bhandari, using IBM Watson technology suite. [33]

Similarly, there are other studies that have identified other factors other than credit score that could foster the inclusion of information that will impact creditworthiness [34].

Through various references, it is clear that the current credit lending system is broken due to various reasons. [30], [31] Within their research, Kumar Arun, Garg Ishan, and Kaur Sanmeet [32] demonstrated that it is indeed possible to use machine learning to predict loan approval odds. Machine Learning can determine new relationships that a person would never think to test. This is the legality of ethics of using machine learning. Another such attempt is done by Mridul Bhandari, using IBM Watson technology suite. [33]

Similarly, there are other studies that have identified other factors other than credit score that could foster the inclusion of information that will impact creditworthiness [34]. This can include data points such as payment of utility bills, rent, and personal habits. [35]

AI will build on our existing system's dual goals of pricing financial services based on the true risk the individual consumer poses while aiming to prevent discrimination (e.g., race, gender, DNA, marital status, etc.). Currently, there are not enough sources of standardized information to base decisions and too little credit being made available. Those conditions allowed rampant discrimination by loan officers who could simply deny people because they "didn't look credit worthy." [35], [36]

It is important to encourage and develop a great sense of money management. Managing how applicants spend their money will ease the stress down the road. The fundamental practice of paying bills on time and saving money will eventually help. The path to financial freedom is not always

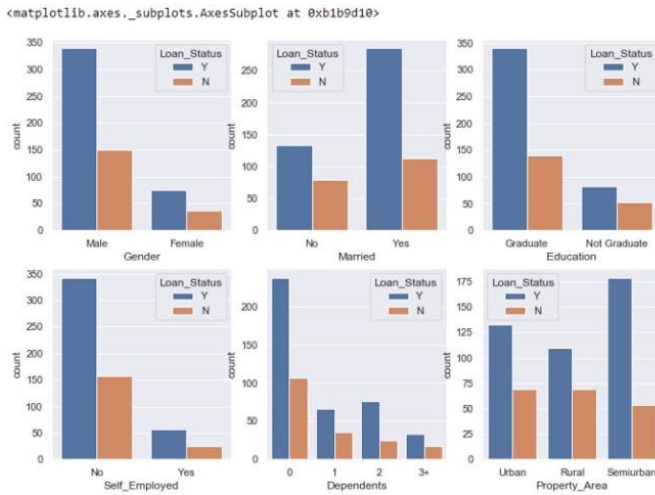


Fig. 1. Findings from the bivariate analysis from loan prediction practice problem [37]

an easy task. In this research, we have identified that loans are widely used around the world, and we have revealed that there are many reasons why people want to get money from banks or any other financial companies that offer loan services. [38].

III. DATASETS

After the text edit has been completed, the paper is ready for the template. Duplicate the template file by using the Save As command, and use the naming convention prescribed by your conference for the name of your paper. In this newly created file, highlight all of the contents and import your prepared text file. You are now ready to style your paper; use the scroll down window on the left of the MS Word Formatting toolbar.

A. BigML.com. The dataset title: Loan Risk Data [39]

The dataset is taken from BigML.com. The dataset title: Loan Risk Data [39] link: <https://bigml.com/user/bigml/gallery/dataset/4f89c38f155268645900033#info> The dataset is about loan risk data. It is having around 1000 records which shows the creditworthiness of applications and contains 21 attributes.

B. Lending Club Loan Dataset 2007_2011

The other dataset we use in our project is “Lending Club Loan Dataset 2007_2011” which is a big data set; it contains around 39,000 rows and 111 columns [40]. Link: <https://www.kaggle.com/imspars/h/lending-club-loan-dataset-2007-2011?select=loan.csv>

C. Data on loan delinquency

The dataset is about loan delinquency, data has around 50,000 loans data and 19 attributes. The size of the dataset is 4.3 MB. [41]

Link: <https://bigml.com/user/bigml/gallery/dataset/4f8b5eae155268645900033#info>

Fig. 2. ATTRIBUTES FROM LOAN RISK DATA

Sl No	Attribute Name	Data type	Description
1	checking_status	Categorical	Status of the loan (status can be in process, grace, repayment, forbearance, etc..)
2	duration	Numeric	Measure of the bond with sensitivity of price, or other debt to change in interest rates.
3	credit_history	Categorical	Records of how a person maintained their credit history in the past.
4	purpose	Categorical	Purpose of the loan
5	credit_amount	Numeric	Amount the customer promises to repay
6	savings_status	Categorical	Status of the savings account
7	employment	Categorical	The customer's employment
8	installment_commitment	Numeric	Includes all the terms and conditions as per amount
9	personal_status	Categorical	Personal Status of the customer

IV. PURPOSED APPROACH

Like any other data science project, the approach we are planning to use includes multiple steps. The steps we are planning to follow are described in Figure ??.

The goal is already defined for this project, which is to find the impact of social and economic factors on Creditworthiness. While we already found one dataset [37], [42] to start with the research, we will continue the research for more data sources, which could help us investigate the problem. Subsequently a data clean up activity is planned and then normalization and grouping of data is also planned. In the next stage, we will be looking for patterns and derive the required knowledge to address the topic under consideration. Finding insights and visualizing the same will be done at this stage. Next step of using machine learning is a bit ambitious for us, with which we will try to find clusters within the dataset(s) under consideration to gain necessary wisdom to solve the problem under consideration.

V. FRAMEWORK

Through the research done so far, we identified that the below attributes have a significant impact on determining the credit worthiness of a person.

10	other_parties	Categorical	All other parties included in the loan agreement
11	residence_since	Numeric	Dates of since when a person is living at a particular residence.
12	property_magnitude	Categorical	Type and importance of the loan
13	age	Numeric	Age of the customer
14	other_payment_plans	Categorical	Other payment plans included in the bank
15	housing	Categorical	Housing status of the customer
16	existing_credits	Numeric	Available information/history of the customer
17	job	Categorical	A customer's basic job information
18	num_dependents	Numeric	Number of dependents included in the loan
19	own_telephone	Categorical	If customers have a contact number
20	foreign_worker	Categorical	If the customer is a foreign worker.
21	class	Categorical	Two different classes of loan-good/bad

- Marital Status
- Employment
- Income
- Property Type
- Rent and/or Utility payments
- Purchase history

In the next phase we are looking to find patterns among them and the weightage of those attributes on determining creditworthiness.

In typical situations, lenders use the credit score as the main factor to determine if people were eligible for a loan. Since 45 million people do not have a credit score, we will use the AI, ML to create a new framework called “qualification score”. The qualification score is a calculation of multiple factors such as social, educational, and financial factors to name a few. Based on this information, then the lender can use the “qualification score” instead of the credit score to determine creditworthiness for those who do not have one.

We will use AI to adopt a new system to find a proper way to rank each element based on the datasets we have. As a result, a will inherit and then calculate the “qualification score”. [36], [43]

VI. FINDINGS

VII. FEATURE WALK THROUGH

VIII. CONCLUSION

Ultimately, measuring creditworthiness without institutional discriminations should be the law. Fairness is the act of treating an individual equally or in a way that is right or reasonable. This is what we learn from our life experience. However, in life, there are instances where there are misinterpretations and different views of the meaning of the term, fairness.

All applications should be given an equal opportunity and accommodations to help gain access to the same lender model. With the help of the most powerful tool available in the world, artificial intelligence and bank institutions are collaborating to provide new and alternative approaches to help increase credit scores. With the help of alternative data source, Artificial Intelligence and Machine Learning will make the decision-making process for the lender much faster and provide an insight on whom will repay their loans. [44]–[46]

REFERENCES

- [1] D. Valatsas, “Coronavirus recovery: why non-us central banks must protect countries from rates cycle.” <https://www.scmp.com/comment/opinion/article/3127424/coronavirus-recovery-why-non-us-central-banks-must-protect>, 29 March 2021.
- [2] A. Iyyengar, “40% of financial services use ai for credit risk management. want to know why?” <https://blog.aspiresys.com/banking-and-finance/40-financial-services-use-ai-for-credit-risk-management/>, 11 August 2020.
- [3] “45 million americans have no credit score.” <https://www.cnbc.com/2015/05/05/credit-invisible-26-million-have-no-credit-score.html>, 16 March 2021.
- [4] Nina, “Credit cards are for transactions; not to borrow money.” <https://www.queercents.com/credit-cards-are-for-transactions-not-to-borrow-money/>, 14 August 2009.
- [5] M. Backman, “1 in 5 americans have no credit score. here’s why that’s a problem.” <https://www.fool.com/the-ascent/credit-cards/articles/1-5-americans-have-no-credit-score-heres-why-problem/>, 5 February 2021.
- [6] M. Tatham, “How to qualify for new credit with no credit score.” <https://www.experian.com/blogs/ask-experian/how-to-qualify-for-new-credit-with-no-credit-score/>, 18 December 2018.
- [7] A. Randazzo and C. Young, “What caused the meltdown: A financial crisis faq.” <https://reason.org/faq/what-caused-the-meltdown-a-fin/>, 25 January 2010.
- [8] E. Hermes, “How to assess the creditworthiness of a customer,” 14 February 2021.
- [9] M. Heidari and S. Rafatirad, “Using transfer learning approach to implement convolutional neural network model to recommend airline tickets by using online reviews,” in *2020 15th International Workshop on Semantic and Social Media Adaptation and Personalization (SMA)*, pp. 1–6, 2020.
- [10] J. Wang, H. Song, and X. Zhou, “A collaborative filtering recommendation algorithm based on biclustering,” in *2015 IEEE International Conference on Cyber Technology in Automation, Control, and Intelligent Systems (CYBER)*, pp. 803–807, 2015.
- [11] M. Heidari and S. Rafatirad, “Bidirectional transformer based on online text-based information to implement convolutional neural network model for secure business investment,” in *IEEE 2020 International Symposium on Technology and Society (ISTAS20)*, ISTAS20 2020, 2020.
- [12] S. Chen, S. Owusu, and L. Zhou, “Social network based recommendation systems: A short survey,” in *2013 International Conference on Social Computing*, pp. 882–885, 2013.
- [13] S. Chen, S. Owusu, and L. Zhou, “Social network based recommendation systems: A short survey,” in *2013 International Conference on Social Computing*, pp. 882–885, 2013.

- [14] S. Lin, C. Liu, and Z.-K. Zhang, "Multi-tasking link prediction on coupled networks via the factor graph model," in *IECON 2017 - 43rd Annual Conference of the IEEE Industrial Electronics Society*, pp. 5570–5574, 2017.
- [15] M. Heidari and S. Rafatirad, "Semantic convolutional neural network model for safe business investment by using bert," in *IEEE 2020 Seventh International Conference on Social Networks Analysis, Management and Security, SNAMS 2020*, 2020.
- [16] Y. Chu, F. Huang, H. Wang, G. Li, and X. Song, "Short-term recommendation with recurrent neural networks," in *2017 IEEE International Conference on Mechatronics and Automation (ICMA)*, pp. 927–932, 2017.
- [17] M. Heidari, S. Zad, and S. Rafatirad, "Ensemble of supervised and unsupervised learning models to predict a profitable business decision," in *IEEE 2021 International IOT, Electronics and Mechatronics Conference, IEMTRONICS 2021*, 2021.
- [18] C. Yang, X. Chen, T. Song, B. Jiang, and Q. Liu, "A hybrid recommendation algorithm based on heuristic similarity and trust measure," in *2018 17th IEEE International Conference On Trust, Security And Privacy In Computing And Communications/ 12th IEEE International Conference On Big Data Science And Engineering (TrustCom/BigDataSE)*, pp. 1413–1418, 2018.
- [19] S. Ji and J. Liu, "Interpersonal ties and the social link recommendation problem," in *2019 6th International Conference on Systems and Informatics (ICSAI)*, pp. 456–462, 2019.
- [20] M. Heidari, S. Zad, B. Berlin, and S. Rafatirad, "Ontology creation model based on attention mechanism for a specific business domain," in *IEEE 2021 International IOT, Electronics and Mechatronics Conference, IEMTRONICS 2021*, 2021.
- [21] M. Heidari, J. H. J. Jones, and O. Uzuner, "Deep contextualized word embedding for text-based online user profiling to detect social bots on twitter," in *IEEE 2020 International Conference on Data Mining Workshops (ICDMW)*, ICDMW 2020, 2020.
- [22] M. Heidari and J. H. Jones, "Using bert to extract topic-independent sentiment features for social media bot detection," in *2020 11th IEEE Annual Ubiquitous Computing, Electronics Mobile Communication Conference (UEMCON)*, pp. 0542–0547, 2020.
- [23] A. Gatzoura, J. Vinagre, A. M. Jorge, and M. Sánchez-Marrè, "A hybrid recommender system for improving automatic playlist continuation," *IEEE Transactions on Knowledge and Data Engineering*, vol. 33, no. 5, pp. 1819–1830, 2021.
- [24] Z. Liao, Y. Song, Y. Huang, L.-w. He, and Q. He, "Task trail: An effective segmentation of user search behavior," *IEEE Transactions on Knowledge and Data Engineering*, vol. 26, no. 12, pp. 3090–3102, 2014.
- [25] M. Heidari, J. H. J. Jones, and O. Uzuner, "An empirical study of machine learning algorithms for social media bot detection," in *IEEE 2021 International IOT, Electronics and Mechatronics Conference, IEMTRONICS 2021*, 2021.
- [26] C.-Y. Chi, Y.-S. Wu, W.-r. Chu, D. C. Wu, J. Y.-j. Hsu, and R. T.-H. Tsai, "The power of words: Enhancing music mood estimation with textual input of lyrics," in *2009 3rd International Conference on Affective Computing and Intelligent Interaction and Workshops*, pp. 1–6, 2009.
- [27] A. Gatzoura, J. Vinagre, A. M. Jorge, and M. Sánchez-Marrè, "A hybrid recommender system for improving automatic playlist continuation," *IEEE Transactions on Knowledge and Data Engineering*, vol. 33, no. 5, pp. 1819–1830, 2021.
- [28] H. Yang, C. He, H. Zhu, and W. Song, "Prediction of slant path rain attenuation based on artificial neural network," in *2000 IEEE International Symposium on Circuits and Systems (ISCAS)*, vol. 1, pp. 152–155 vol.1, 2000.
- [29] S. Zad, M. Heidari, J. H. J. Jones, and O. Uzuner, "A survey on concept-level sentiment analysis techniques of textual data," in *IEEE 2021 World AI IoT Congress, AllIoT2021*, 2021.
- [30] N. Campisi, "From inherent racial bias to incorrect data—the problems with current credit scoring models." <https://www.forbes.com/advisor/credit-cards/from-inherent-racial-bias-to-incorrect-data-the-problems-with-current-credit-scoring-models/>, 26 February 2021.
- [31] S. Wu, "10 problems with credit in the united states." <https://bloom.co/blog/10-problems-with-credit-in-the-united-states/>, 25 September 2017.
- [32] K. S. Kumar Arun, Garg Ishan, "Loan approval prediction based on machine learning approach," *IOSR J. Comput. Eng.*, vol. 18, pp. 18–21, 2016.
- [33] GitHub-Mridulrb, "Predict loan eligibility using ibm watson studio." <https://github.com/mridulrb/Predict-loan-eligibility-using-IBM-Watson-Studio>, 2020.
- [34] C. Bieber, "5 factors besides your credit that affect personal loan approval." <https://www.fool.com/the-ascent/personal-loans/articles/5-factors-besides-your-credit-that-affect-personal-loan-approval/>, 16 March 2021.
- [35] "New research shows that your looks, creditworthiness may go hand in hand." <https://phys.org/news/2009-03-creditworthiness.html>, 12 March 2009.
- [36] A. Klein, "Reducing bias in ai-based financial services." <https://www.brookings.edu/research/reducing-bias-in-ai-based-financial-services/>, 10 July 2020.
- [37] A. Smith, "7 fundamental steps to complete a data analytics project." <https://blog.dataiku.com/2019/07/04/fundamental-steps-data-project-success>, 4 July 2019.
- [38] C. Bai, B. Shi, F. Liu, and J. Sarkis, "Banking credit worthiness: Evaluating the complex relationships," *Omega*, vol. 83, pp. 26–38, Mar. 2019.
- [39] bigml.com, "Loan risk data," 14 February 2021.
- [40] S. Gupta, "Lending club loan dataset 2007_2011." <https://www.kaggle.com/imsparsh/lending-club-loan-dataset-2007-2011?select=loan.csv>, 29 May 2020.
- [41] bigml.com, "Loanstats' dataset," 2021.
- [42] P. Meena, "Step-by-step guide to exploratory data analysis in python." <https://towardsdatascience.com/an-introduction-to-exploratory-data-analysis-in-python-9a76f04628b8>, 14 August 2020.
- [43] P. Crosman, "Can ai be programmed to make fair lending decisions?." <https://www.americanbanker.com/news/can-ai-be-programmed-to-make-fair-lending-decisions-27>, 27 September 2016.
- [44] A. Dobrin D.S.W., "It's not fair! but what is fairness?." <https://www.psychologytoday.com/us/blog/am-i-right/201205/its-not-fair-what-is-fairness>, 11 May 2012.
- [45] R. Chawla, "How ai supports financial institutions for deciding creditworthiness." <https://www.entrepreneur.com/article/310262>, 11 March 2018.
- [46] The New York State Senate, "Unlawful discriminatory practices in relation to credit." <https://www.nysenate.gov/legislation/laws/EXC/296-A>, 2021.