

## EDA FOR BANK LOAN DEFAULT RISK ASSESSMENT

Ch. Deepika<sup>1</sup>, Ch. Alekhya<sup>2</sup>, G. Khushi Reddy<sup>3</sup>, Ch. Mounika<sup>4</sup>, S. Chaitanya Kumar<sup>5</sup>

<sup>1,2,3,4,5</sup>Computer Science Department, JNTUH, Vidya Jyothi Institute of Technology, India.

### ABSTRACT

The project aims to integrate advanced techniques in data analysis (EDA), machine learning (ML) algorithms, and financial institutions for risk analysis, and private money in municipal loans. The important objective is a strict analysis of the profile structure with a special focus on resolving problems related to insufficient credit history to reduce the risk of bad credit. This project uses EDA technology to process loan application data and aims to identify repayment issues and help make decisions about loan approval, loan eligibility, and interest rates. The overall goal is to ensure that qualified applicants are approved while minimizing the financial risk associated with failure. To achieve this goal, the project uses the Random Forest algorithm, the most effective tool in machine learning that can be clearly defined and well-defined, providing a clear explanation of the decision-making role modern technology plays in preserving life and improving health benefits in today's world.

**Keywords:** EDA, machine learning, risk analysis, financial institutions, private money, municipal loans, credit history, loan approval, loan eligibility, interest rates, Random Forest algorithm

### 1. INTRODUCTION

The credit process importance cannot be ignored in the rapidly changing financial sector. However, uncertainty in the decision-making process raises serious concerns about bias, discrimination, and its impact on trust. The program is a vision of the concept that uses ideas to solve problems and define machine learning (ML) to mitigate these problems, leading to a new era of transparency in the credit industry. Central to this change is the integration of easily interpretable ML models equipped with features such as Native Interpretable Model-Independent Description (LIME) or SHapley Additive Description (SHAP). Today's technology allows people to understand the logic of lending decisions, ensuring stakeholders have a clear understanding of the key factors that influence outcomes.

By ensuring this transparency, the project not only improves accountability, but also fosters trust and confidence throughout the entire lending process by promoting a deeper understanding of borrowers, loan officers and regulators. The program creates harmony and reduces injustice. By applying fairness-aware algorithms, the system can quickly detect and correct inconsistencies in the loan approval process. The purpose of positive intervention is to ensure the honesty of loan applicants regardless of population, thus promoting the principles of fairness and justice that are important in the borrowing process. The program's many methods, including disclosure, fairness, accessibility, flexibility and management control, demonstrate its commitment to not only transforming the lending process but also creating a platform for accountability, transparency, and accountability in the financial sector.

### 2. LITERATURE SURVEY

Breiman's pioneering work on random forests began a learning curve known for his expertise in classification and regression. This approach is particularly good at handling high data, making it ideal for predictive work. Random forests provide a powerful framework for identifying patterns and providing insight from complex data by using the intelligence of different decision trees. This versatility not only increases prediction accuracy but also helps improve the interpretation of results; This is especially important for informed decision-making in risk-on-risk analysis scenarios. Breiman's program highlights the important role of random forests in today's scientific research, enabling experts to solve real-world problems with confidence and accuracy.[1]

McKinney and Weiss et al. Python data analysis library Pandas is designed as a simple tool for data management and analysis in the Python ecosystem. Pandas' powerful data including DataFrame and Series helps work well with many data types, highlighting its efficiency. It simplifies the data preparation process by providing functions for data cleansing, transformation, and searching, allowing data scientists to manage large datasets, extract key insights, and discover patterns. With Pandas, data scientists can solve complex data problems for accurate, data-driven decisions. As a result, Pandas are becoming an essential tool in every data scientist's toolkit, allowing them to go far and solve complex problems.[2]

Pedregosa and Fabian's tutorial is an important introduction to learning scikit-learn, a machine learning library in the Python ecosystem. It provides users with essential tools for machine learning modeling, providing an overview of the wide range of algorithms and tools available in scikit-learn. This guide focuses specifically on deployment and recovery purposes that are relevant to our risk analysis and details algorithms suitable for these purposes. Through a

comprehensive evaluation of Scikit-learning functionality, this guide provides practitioners with the essential knowledge and resources necessary to effectively use the library to solve credit risk and reduce related problems.[3]

Müller and Guido's book provides a friendly introduction to machine learning, covering practical applications using Python. Through clear explanations and examples, readers gain a basic understanding of Python machine learning algorithms. This book covers important topics such as data generation, sample selection, and evaluation techniques, allowing practitioners to understand the complexities of risk assessment with confidence. With this understanding, risk planners usually consider machine learning techniques to gain insight from data and make informed decisions to reduce risk.[4]

The VanderPlas primer provides a comprehensive introduction to the basic tools and techniques for performing data science with Python. Through many examples and case studies, he shows that Python libraries form the basis of Data Science Papers (EDA) by facilitating data exploration, visualization, and analysis. Using the Python library, readers can identify and interpret patterns important for informed decision-making. Immersing yourself in these resources can give professionals a deeper understanding of data science and help them solve real-world problems like municipal credit risk.[5]

Brownlee and Jason's Python Data Analysis Technical Guide (EDA) is a great resource for working with complex credit applications. It uncovers deep insights and complex patterns in data through careful observation and analysis. Guides using different types of EDA methods, such as content analysis and data visualization, highlight underlying reimbursement issues. This insight allows stakeholders to make decisions about loan approvals, rate adjustments, and loan terms, ultimately improving financial management strategies for individuals, using the product, and making it useful for qualified applicants.[6]

Izenman and Alan Julian's book concentrates on advanced analytical techniques and is particularly important for our work in risk assessment. The development of research and different studies provides practitioners with the theoretical and practical tools necessary to detect loan applications and reduce risk. With detailed explanations of regression, distribution, and other complexities, this book enables readers to develop effective risk management strategies and make informed decisions. This is an important guide for stakeholders looking to solve credit analysis problems and develop ways to reduce risk.[7]

Wu, James, et al's book delivers a comprehensive review of data discovery concepts in the Python ecosystem, focusing on distribution and reproducibility. By uncovering the complexity of machine learning models, it provides readers with a powerful tool for predicting credit risk based on historical data. Practitioners gain the skills needed to solve predictive modeling problems through a combination of theory and examples. Through in-depth discussions, this book not only helps readers understand important information but also solves real-world problems in financial risk assessment and mitigation.[8]

Raschka and Mirjalili's book is an inquisition of machine learning algorithms focusing on practical applications in the Python ecosystem. Spanning model estimation provides valuable information for effective risk assessment with hyperparameter tuning and integration. With clear explanations and examples, readers can understand their learning methods and make reasonable choices. Using the techniques outlined in the book, professionals can increase the accuracy and reliability of forecasting models, thereby improving risk management strategies in business finance.[9]

Bengfort, Tony, and Rebecca's book delve into document analysis techniques using Python, focusing on loan application documents. It is useful for practitioners by providing insight into techniques such as emotional analysis, modeling, and natural language processing (NLP). Through a Python-based search method, readers gain the skills needed to extract valuable information from a loan application description. Sentiment analysis can understand the emotional context surrounding a loan request, while semantic modeling can identify content and structure. Additionally, NLP increases the accuracy of risk assessment by extracting important information from irrelevant text. Using this advanced technology, professionals can improvise their analytical abilities to make more informed financial decisions.[10]

Agarwal's textbook is essential for understanding data mining, from prioritization to clustering and discovery. Its comprehensive analysis highlights the complexity of big data among our risks, highlighting basic and advanced processes. Ensure accuracy and quality of data analysis by guiding readers through data cleaning and preparation. Investigating joint ventures helps assess risk by revealing relevant patterns and synergies in loan applications. Additionally, discussing the investigation process provides practitioners with tools to identify and address concerns or inconsistencies, thereby reducing risk. With Agarwal's data insights, analysts can use data quality to gather valuable insights and reduce risk in financial markets.[11]

Bishop's book is a definitive guide to pattern recognition and machine learning, providing practitioners and researchers with insight into key concepts. It covers topics such as Bayesian inference, neural networks, and support vector machines

and provides theoretical and practical foundations. Bishop's careful analysis and explanation of fundamental principles provides readers with tools to navigate the complexities of risk assessment with precision and accuracy. Its comprehensive approach not only improves the understanding of basic algorithms but also encourages practical applications. Bishop's work is therefore important for those looking to solve pattern recognition and machine learning problems to identify risk.[12]

The collaborative work of Hastie, Tibshirani, and Friedman provides a comprehensive review of the fundamental concepts needed to develop predictive models and work effectively in learning statistical methods. Unlike general discussions, this book provides readers with a solid foundation for solving clusters by covering important topics such as design, model selection, and tree models. The authors illuminate the complexness of the process, allowing practitioners to approach design with confidence and accuracy. Their insights not only equip us with helpful resources for building predictive models but also enable a deeper understanding of the nuances when interpreting results. This book is an important introduction for the use of statistical methods to inform and improve the practice of risk assessment.[13]

Kuhn and Johnson's book presents an approach to predictive modeling by focusing on complex concepts such as architecture, model evaluation, and implementation. Through real-world applications and clear explanations, it provides practitioners with insight to develop accurate predictive models. This book introduces the complexity of choice and change, providing readers with simple insights into improving performance standards. In addition, it provides a rough overview of model evaluation, allowing readers to distinguish between effective and ineffective forecasting models. Through extensive research, Kuhn and Johnson enable clinicians to use different techniques to increase the power and reliability of their predictive models. More importantly, these insights provide a solid foundation for developing predictive models that can predict and manage financial problems and help reduce credit risk.[14]

Oliphant's Handbook is an essential resource for data scientists, focusing on key Python libraries such as NumPy, Pandas, and Matplotlib. Through clear explanations and good examples, this book provides practitioners with the skills they need for data management. Oliphant unpacks the complexity of these libraries and gives readers the knowledge they need to use them. This guide provides practical advice based on real-world applications and will be of particular benefit to risk planners who want to manage and analyze data effectively. Whether solving data problems or sharing insights, the skills gained with Oliphant's guidance provide practitioners with confidence and efficiency in project research and risk mitigation.[15]

François Chollet's text Deep Learning with Python is a great resource for understanding neural networks and their applications in machine learning. It covers fundamental concepts such as neural networks, support and communication models, and training algorithms, with an emphasis on re-representation. This book seamlessly combines theory with notation, allowing readers to design and configure a variety of tasks for neural networks. It is a useful tool for beginners and experts alike, providing insight and expertise into the complex field of neural network design and training.[16]

Dash's Cookbook can be habituated as a reference for Python data analysis challenges and offers solutions for tasks such as data preprocessing, constructing and model evaluation. It simplifies key steps in project lifecycle and risk management with step-by-step instructions and ready-to-use articles. Its pragmatic approach speeds up development and improves project quality by providing tools for data cleaning, transformation, and modeling. Dash models enable clinicians to efficiently cross-reference research data and facilitate timely and cost-effective implementation.[17]

Géron's book leverages the power of popular Python libraries such as Scikit-Learn, Keras, and TensorFlow to provide an efficient and intuitive approach to machine learning. With clear explanations, examples, and ready-to-use articles, practitioners have the tools and knowledge to implement all aspects of machine learning. Géron carefully guides the reader through each stage, from preliminary data to the exported model, showing the interaction of various elements in the process. This best-in-class approach is particularly useful in risk assessment, where integration of previous data is not possible but the establishment and presentation of risk is important for managing the process. Under Géron's expert guidance, professionals gain the skills to solve complex problems, improve performance standards, and effectively use machine learning for risk assessment. Ultimately, Géron's book acts as a useful resource that allows practitioners to explore the complexities of machine learning with agility and precision, thus supporting informed decision-making and reducing risk.[18]

This book on deep learning with Python and Keras is an essential guide for practitioners on advanced neural network systems. It introduces key concepts such as convolutional neural networks (CNN) and recurrent neural networks (RNN), which provide powerful tools for gaining insight from complex data, especially in loan applications. By explaining theoretical principles and practical applications, this book equips the reader with the skills to analyze and interpret large credit data. It enables experts to use CNNs for image-based data processing and RNNs for sequential data analysis, providing insights into risk assessment and decision making. Using Python and Keras as main tools, this book provides

accessibility and ease of use, allowing practitioners to turn knowledge into solutions to improve the loan application process and reduce associated risks.[19]

Wes McKinney's book is an essential resource for learning documentation and analysis in the Python ecosystem. Introduces important libraries such as Pandas and NumPy that provide information on processing credit data. McKinney covers data maintenance, transformation, and collection and provides practitioners with data search and creation tools. Through theoretical concepts and examples, readers can understand the complexity of credit information and build a solid foundation for risk assessment and credit decision-making.[20]

### 3. PROPOSED SYSTEM

Banking establishments, particularly those that hold special funds in municipal loans, face problems in assessing creditworthiness because applicants do not have sufficient credit information. This leads to a negative credit risk that affects the Economic solidity of the company. The project aims to solve this problem by combining data analysis (EDA), machine learning (ML) algorithms, and risk analysis technologies.

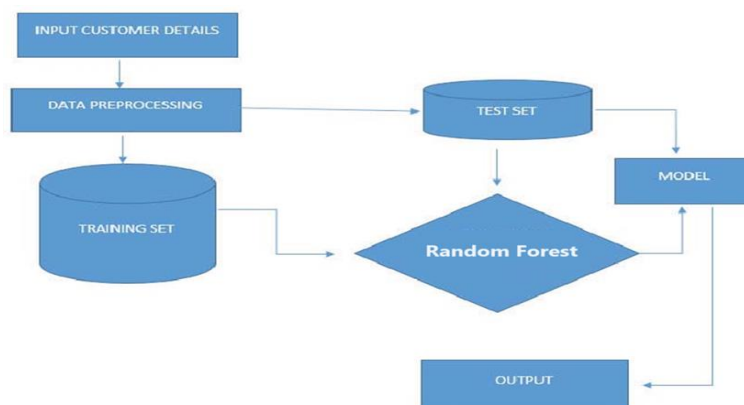
This move represents a revolutionary change in loan approval through the machine learning (ML) integration. The system provides transparency, fair information, and flexibility, working in tandem to promote a fair and inclusive lending environment and foster trust between borrowers and borrowers, respectively. The proposed system uses data analytics (EDA) and advanced machine learning to scrutinize loan application data to Recognize patterns indicative of potential repayment challenges. Risk mitigation strategies are then developed through research, insight, and analysis. The advantages of the proposed process are as follows: Holistic assessment of creditworthiness: The application process provides a more comprehensive assessment of the applicant's creditworthiness influenced by EDA and ML. By analyzing many variables beyond credit history and scores, the system can better understand an applicant's financial behavior and allow him/her to make better credit decisions. Analysis of differences: EDA and ML allow companies to identify significant changes that have a positive impact. Bad credit is reported. By uncovering these driving changes, the system provides valuable information that can inform risk assessment, portfolio management, and credit utilization. Reducing loan defaults: By providing a deeper understanding of the factors that lead to loan defaults, the application process can help companies take critical steps to reduce loan defaults and job loss. The system will help reduce risk and improve the financial stability of the company by adjusting loans and interest rates and even denying loans to high-risk applicants. Increase accuracy and efficiency: The planning process will increase the accuracy and efficiency of credit decisions by using data-driven insights. It reduces reliance on content analysis and allows companies to make more informed, data-driven choices that increase the overall efficiency of the lending process. Data-based decision-making is done. The application process supports the decision-making culture in the company. By integrating EDA and ML algorithms, the system supports evidence-based decision-making to create better risk management strategies, reduce error rates, and increase customer satisfaction.

The approach that we are following is shown below:

- Configuration module for uploading credit application information: Allows users to upload credit application information, view credit allocation, and determine reasons for rejection.
- Dataset Preprocessing Module: Convert data into formats without coding and standardize to ensure uniformity and cleanliness.
- Split dataset training and testing module: Split the dataset into training and testing processes in an 80/20 ratio compared to the training and testing model.
- ML Tutorial on Loan Approval Module: Train a random forest model on 80% of the data to predict loan approval and use the remaining 20% to evaluate its accuracy.
- Train ML on Credit Denial Module: Use training data to train a random forest model of credit denial and measure the rigor of the test.
- Interpretation Machine Learning Module: Use SHAP values to describe the key features that influence loan approval or denial for the prophecy model.
- Credit prediction using test data: allows users to submit test data to predict credit and provide detailed information about reasons for approval or rejection with SHAP description.



## FEED-FORWARD NEURAL NETWORK ARCHITECTURE



**Figure 1: SYSTEM ARCHITECTURE OF THE APPLICATION**

Feedforward Neural Network (FNN) is a basic artificial neural network architecture in which data is passed unidirectionally from the input layer to the output layer. Unlike traditional radios, FNNs have no feedback loops, allowing for an easy, loop-free process. The components of an FNN consist of an input layer, a hidden layer, and an output layer that processes information about changes in the weights of neurons. During training, weights and biases are adjusted to minimize the error between predictions and actual results, which is often used for improvement. FNN is mainly used for classification and reprocessing and is a design model for deep learning, especially multilayer perceptron (MLP) with many hidden layers

## 4. RESULTS AND DISCUSSIONS

Name of the Test	Input	Expected Output	Actual Result	Remarks
Uploading of Dataset	Selected dataset	The dataset uploaded successfully with column displayed	Uploaded dataset successfully as expected along with displaying columns	Successful
Viewing analyzed dataset	Uploaded Dataset	Graph output of analyzed dataset	Displays the graph for the analyzed data	Successful
Preprocessing of dataset	Uploaded dataset	Pre-processed data should be displayed consisting of normalized values	Displays the normalized values on the screen after preprocessing	Successful
Training and Testing of data	Uploaded Dataset	Total records with the number of values in training and testing data	Displays the records present in training and testing after splitting along with the total records	Successful
Explainable ML	Uploaded dataset	Graph output of explained ML dataset with SHAP values	Displays the explainable ML graph	Successful
Explainable ML	Uploaded dataset	Graph output of explained ML dataset with SHAP values	Displays the explainable ML graph	Successful
Uploading of test data	Selected dataset	Dataset uploaded successfully	Uploaded successfully	Successful
Predicting loan status of test data	Uploaded dataset	Loan status prediction using test data	Displays the prediction along with displaying the causes for predicting the output.	Successful

The output screens of the application are as follows:

### HOME PAGE OF THE APPLICATION



Figure 3: TABLE OF TEST CASES FOR THE APPLICATION

### PREDICTING THE LOAN STATUS



## 5. CONCLUSION

In conclusion, the EDA's Financial Services Commission is beginning to point to a vision of the future in which credit decisions will move beyond facts to include transparency and principles-based understanding. The initiative to enable descriptive machine learning (ML) not only improves predictive capabilities but also creates higher standards of fairness, liability shifting, and trust in the credit approval process environment. The project aims to unravel the decision-making process using machine learning techniques and algorithms and enable stakeholders and end users to understand the reasons behind loan approval or non-acceptance. This collaboration aims to support financial transactions understood through Data Analysis (EDA) and machine learning algorithms, not only helping in risk assessment but also making lending fair. The overall aim of the application is to improve financial efficiency. To offer institutions and debtors a more transparent, fair and reliable system.

By integrating advanced learning models such as clustering, deep learning or gradient boosting into existing EDA systems, we aim to improve the identification of models and increase the differences associated with payback issues. These models increase the effectiveness of risk assessment by providing greater predictability and visibility. Additionally, using natural language processing (NLP) allows us to gain insights from intangible information such as customer feedback, emails or documents. By analyzing data, we can uncover important information about customers' financial problems, thoughts and experiences, making it easier for customers to market most risk management strategies.

Additionally, investing in the development of data visualization techniques and interactive dashboards improves understanding and interpretation of analytical content. Visual representations support more informed decision-making by helping stakeholders understand relationships and patterns. Additionally, regular monitoring and model updates ensure accuracy and reliability by identifying model deviations and making adjustments based on changing patterns and risks. Finally, integration with credit management supports the end-to-end loan process, ensuring that conflicting information and applications are made available for loan approval and process control.

## 6. REFERENCES

- [1] Breiman, Leo. "Random forests." *Machine learning* 45.1 (2001): 5-32.
- [2] McKinney, Wes. "Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython." O'Reilly Media, Inc., 2018.
- [3] Pedregosa, Fabian, et al. "Scikit-learn: Machine learning in Python." *Journal of machine learning research* 12.Oct (2011): 2825-2830.
- [4] Müller, Andreas C., and Sarah Guido. "Introduction to Machine Learning with Python: A Guide for Data Scientists." O'Reilly Media, Inc., 2016.
- [5] VanderPlas, Jake. "Python Data Science Handbook: Essential Tools for Working with Data." O'Reilly Media, Inc., 2016.
- [6] Brownlee, Jason. "Introduction to Exploratory Data Analysis in Python." *Machine Learning Mastery*, 2020.
- [7] Izenman, Alan Julian. "Modern Multivariate Statistical Techniques: Regression, Classification, and Manifold Learning." Springer Science & Business Media, 2013.
- [8] Wu, James, et al. "Data Mining with Python: Implementing Classification and Regression." CRC Press, 2019.
- [9] Raschka, Sebastian, and Vahid Mirjalili. "Python Machine Learning." Packt Publishing Ltd, 2019.
- [10] Bengfort, Benjamin, Tony Ojeda, and Rebecca Bilbro. "Applied Text Analysis with Python: Enabling Language Aware Data Products with Machine Learning." O'Reilly Media, Inc., 2018.
- [11] Aggarwal, Charu C. "Data Mining: The Textbook." Springer Science & Business Media, 2015.
- [12] Bishop, Christopher M. "Pattern Recognition and Machine Learning." springer, 2006.
- [13] Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. "The Elements of Statistical Learning: Data Mining, Inference, and Prediction." Springer Science & Business Media, 2009.
- [14] Kuhn, Max, and Kjell Johnson. "Applied Predictive Modeling." Springer, 2013.
- [15] Oliphant, Travis. "Python for Data Science Handbook: Essential Tools for Working with Data." O'Reilly Media, Inc., 2016.
- [16] Chollet, François. "Deep Learning with Python." Manning Publications Co., 2017.
- [17] Dash, Santosh Kumar. "Python Data Science Cookbook." Packt Publishing Ltd, 2015.
- [18] Géron, Aurélien. "Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems." O'Reilly Media, Inc., 2019.
- [19] Haykin, Simon. "Neural networks and learning machines." Pearson, 2008.
- [20] McKinney, Wes, et al. "Python Data Analysis Library." *Journal of OpenSource Software* 6.60 (2021): 3541.