



Enhancing Credit Risk Classification Using LightGBM with Deep Feature Synthesis

Sarah Rosdiana Tambunan¹, Junita Amalia², Kristina Margaret Sitorus³,
Yehezchiel Abed Rafles Sibuea⁴, Lucas Ronaldi Hutabarat⁵

^{1,2,3,4,5} Information System Study Program, Institut Teknologi Del, Toba, Sumatera Utara, Indonesia
Email: ¹sarah.tambunan@del.ac.id, ²junita.amalia@del.ac.id, ³kristinasitorus190@gmail.com,
⁴yehezchiel@gmail.com, ⁵lucastakdung@gmail.com

Abstract

In the digital financial services era, Peer-to-Peer (P2P) lending has emerged as a significant innovation in fintech. However, credit risk remains a major concern due to the potential for payment failures, which can cause losses for platforms and investors. This research explores the impact of Deep Feature Synthesis (DFS) on credit risk classification and evaluates the performance of the Light Gradient Boosting Machine (LightGBM) algorithm with and without DFS. The data used in this study was sourced from Kaggle, a peer-to-peer lending company based in San Francisco, California, United States. The dataset contains 74 attributes, with a total of 887,379 rows. DFS automatically generates new attributes, while LightGBM is used for selecting the most important features, aiming to optimize credit risk predictions and simplify the model's complexity. The results of credit risk classification models using DFS and without it. Findings reveal that DFS enhances the accuracy of the credit risk classification, achieving a 0.99 accuracy rate compared to 0.97 without DFS, achieving a recall and F1-score of 0.94 and 0.96 with DFS and 0.68 and 0.81 without DFS. These results suggest that DFS is an effective feature engineering technique for boosting credit risk model performance. This research contributes significantly to the P2P lending industry by demonstrating that combining DFS with LightGBM can improve credit risk management, making it a valuable approach for financial platforms.

Keywords: Credit Risk, P2P Lending, Feature Engineering, Feature Selection, LightGBM.

1. INTRODUCTION

In the current era, lending services are a form of financial service agreed upon between the fund provider and the borrower. The development of this system follows the advancement of digital technology. One of the leading innovations in the financial services sector besides banks, available through fintech, is Peer-to-Peer (P2P) lending. Fintech, as the application of information technology in the financial sector, has brought many innovations. This service provides a financing channel or lends money to individuals or businesses, enabling broader access to the funding world for those who previously had difficulty obtaining financing. A



form of financial. service provided by fintech is P2P lending, where this service offers financial loan facilities conducted online. The mechanism in P2P lending offers ease in borrowing without requiring documents, collateral, and without involving intermediaries. There are many types of P2P lending services, such as unsecured personal loans, business loans, student loans, and secured business loans, which provide new financing alternatives for those in need of funds [1].

Credit risk, according to Connors (2010), refers to the uncertainty regarding whether the involved party will meet their financial commitments. Credit risk models are generally divided into two main types: qualitative and quantitative. Qualitative models evaluate elements like reputation, financial leverage, income stability, collateral, business cycles, and interest rates. In contrast, quantitative models aim to generate credit scores to predict the probability of default or classify borrowers according to their default risk level [2]. This situation may arise due to issues in personal credit from customers or borrowers, including default history and other qualitative and quantitative factors reflecting credit risk or default. Therefore, risk management in the P2P industry sets credit risk as a primary focus to safeguard the interests of investors and companies and reduce potential default risks that could lead to economic losses [3].

Credit risk is necessary when lenders, such as banks or other financial institutions, need to consider the possibility of default based on the information regarding the loan applied for [3]. Credit risk evaluation is used by lenders to assess the borrower's ability to meet their financial obligations. This is important to ensure that lenders can minimize losses due to defaults and make prudent lending decisions. Credit risk is used in various situations, including when providing personal loans, mortgages, business loans, and other forms of credit [4]. To manage this risk, lenders use various credit analysis methods to assess the creditworthiness of prospective borrowers. One common approach is credit risk modelling, which involves using historical data. Lenders use credit risk results to evaluate a borrower's ability to fulfil their financial obligations [5]. This is important to ensure that lenders can minimize losses due to defaults and make prudent lending decisions. However, issues of default or late payments are the primary reasons for the emergence of credit risk, meaning there is a distinction between good loans and bad loans. This arises due to a large attribute space that can influence this. Therefore, the selection of the most optimal attributes must be conducted [6]. This includes the process of identifying the credit risk faced when extending credit to customers. The goal of credit risk analysis is to reduce credit risk and ensure that banks or financial institutions can make the right decisions in extending credit [3]. By providing information that someone has a good loan background and bad loan background, it will impact the decision to grant a loan or not. In this process, distinguishing between good loans and bad loans becomes crucial. Good loans refer to loans given to borrowers with a low risk of default,

while bad loans refer to loans with a high risk of default. This identification is important as it affects investor confidence and the operational sustainability of the P2P lending platform. By accurately classifying borrowers as good loans or bad loans, financial institutions can reduce potential financial losses caused by problematic loans. Furthermore, this separation also helps in better risk management, allowing companies to establish appropriate risk mitigation strategies and ensure financial stability. Therefore, effective credit analysis must be able to explain and identify good loans and bad loans to make smarter credit decisions and maintain the health of the loan portfolio [7].

Machine learning models, particularly deep learning in models, are frequently regarded as "black boxes," which makes it challenging make to comprehend how they arrive at their decisions. This means that even though these models can generate accurate predictions or decisions, their workings and the reasons behind these decisions are not transparent or easily interpret able by humans [8]. Therefore, to understand and confirm the decisions generated by the model, such as "good loans" and "bad loans," a deeper explanation of how the model makes these decisions is required. This research is expected to provide reasons why the model categorizes a loan as "good loan" or "bad loan." By understanding the reasons behind these decisions, researchers can improve the model's transparency and accuracy, as well as reduce the uncertainty and distrust that may arise from using a "black box" model. The focus of this research, which aims to obtain feature importance analysis that can be used to classify the model with accurate credit predictions. In this context, feature engineering using in Deep Feature Synthesis and feature selection by the model, namely LightGBM, will be involved in this research to select the most informative attributes (feature importance) in predicting credit risk from the many available attributes.

Feature engineering is a machine learning technique that involves generating new attributes from data before it is input into the model [9]. These attributes are created by utilizing available data and improving data quality to enhance model accuracy. This plays a crucial role in improving the performance of machine learning models. The problem often encountered with available datasets is their high complexity, with many irrelevant attributes and unidentified patterns. With this, DFS is expected to help find and create more relevant attributes that can improve the accuracy of machine learning models. DFS can also generate time-based, entity-based, relationship-based, pattern-based, and dependency-based attributes, which are very useful in customer analysis and behaviour prediction. Deep Feature Synthesis has the ability to automatically generate deep attributes for the datasets without human intervention [10]. Previous research conducted Feature Engineering using Deep Feature Synthesis (DFS) on relational datasets, where new attributes were automatically generated. In earlier previous research on Deep Feature Synthesis (DFS) within the automatic feature engineering framework

sought to tackle issues in loan fraud detection, including attribute dimension explosion, low interpretability, prolonged training times, and low detection accuracy. The problem faced was the explosion of attribute dimensions and low interpretability in detecting car loan fraud. To address this issue, abstract and uninterpretable attribute compression was performed to limit the depth of the DFS algorithm. The algorithm used was Deep Feature Synthesis (DFS). Research shows that using this DFS method can increase detection accuracy by 23%, reduce the number of attributes used by 92.5%, and reduce model training time by 54.3% compared to traditional methods. Jian Yang et al. suggested that future research can be conducted on techniques that can reduce complexity or the number of features. Researchers can select a subset of the most relevant and informative attributes from the original attributes before and after DFS [11].

Therefore, performance analysis of DFS in generating attributes to improve model predictions will be conducted. Deep Feature Synthesis (DFS) automatically generates attributes for relational datasets by following relationships in the data. Next, the use of feature selection aims to improve classification performance by optimizing the attributes to be used. In this research, feature selection, which is generally the process on selecting relevant attributes from the data set to be used in machine learning models, will be performed. Feature selection is important as it helps improve machine learning model performance, reduce model complexity, and save computation time [12]. In this research, feature selection will provide information related to good loan and bad loan borrowers, where feature selection is performed by the Light Gradient Boosting Machine (LGBM) model, influencing the decision to grant a loan or not. In the research [13], the main problem is how to select the most informative attributes from the datasets to build an accurate classification model. This research conducted feature selection using three methods is Chi-Square and Information Gain, and Gain Ratio, to identify the most informative attributes from the data. This study also evaluated five models: Bayesian, Naive Bayes, Support Vector Machine (SVM), Decision Tree (C5.0), and Random Forest (RF), with each model assessed using accuracy, F- measure, true positive rate, and true negative rate. Building on this prior research, it motivates researchers to perform classification testing on datasets using different models to accurately predict credit risk, with an emphasis on identifying borrowers who are likely to default.

The result of this research is that the selection technique with the methods used by the machine learning classification model, and the Random Forest model achieved an accuracy rate of up to 93% [13]. Furthermore, the research [14], where the problem addressed is related to classifying diabetes and non-diabetes patients. The results generated after pre-processing and optimization showed that decision trees achieved an accuracy of 76.07%, random forest 79.8%, multi-layer perceptron 77.60%, and K-nearest neighbor 78.58%. This study compared various

classification methods in attribute selection to predict diabetes with higher accuracy. This research encourages further classification using other classification methods to improve prediction [14]. LightGBM is designed for higher speed and efficiency compared to other boosting algorithms. This is achieved through histogram-based binning and leaf-wise tree growth techniques, which significantly reduce computation time. LightGBM can handle very large datasets due to its high memory and computation efficiency [15]. In the research [16], feature selection was conducted using the LightGBM algorithm model. LightGBM (Light Gradient Boosting Machine) is a highly efficient and fast machine learning algorithm, commonly employed in various tasks like classification and regression. One of LightGBM advantages is its ability to be applied in feature selection, a crucial process in building effective and efficient models. Based on this research conducted on the Kaggle Home Credit Default Risk datasets using boosting methods such as Cat-boost, XGBoost, and LightGBM, Light Gradient Boosting Machine is faster in training the model compared to XGBoost. This is because Light Gradient Boosting Machine uses.

2. METHODS

The methodology outlines the steps taken throughout this research. This description is essential as a guideline to ensure that the outcomes achieved align with the established objectives of the study.

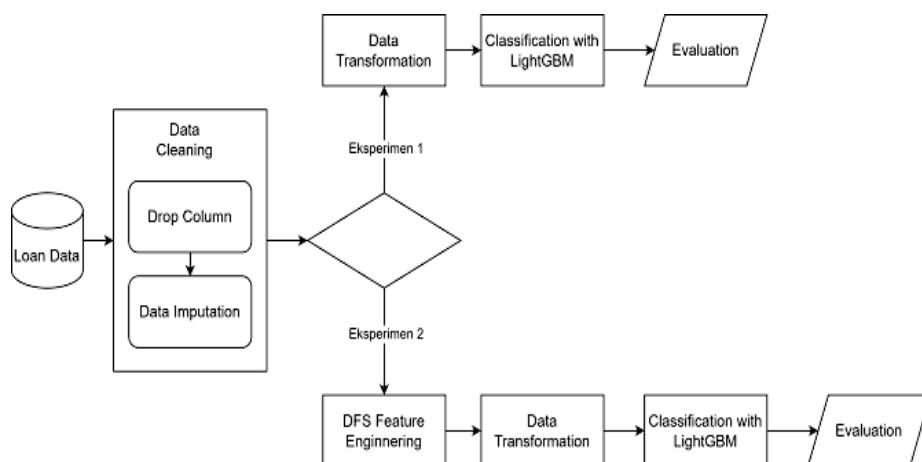


Figure 1. Research Design

According to Figure 1, this study involves several stages that need to be completed before proceeding to the subsequent stages. The stages in this research include the datasets used in this study is lc_loan, sourced from Kaggle, consisting of 887,379 rows and 74 columns. This dataset includes comprehensive information regarding

borrowers, loans, and payment statuses relevant to credit risk analysis. The datasets can be accessed via platform Kaggle lending club-loan-defaulters prediction 887,379 rows and 74 columns. It provides a broad basis for analyzing and modelling borrower behaviour, as well as predicting loan default probabilities.

2.1. Pre-processing

All data in pre-processing is a critical step in the data mining process, aimed at refining raw datasets to ensure their suitability for further analysis and modelling. This stage includes **data cleaning**, where incorrect, incomplete, or irrelevant data is identified and corrected to enhance accuracy and reliability. Data imputation follows, using the K-Nearest Neighbors method or KNN to handle missing values in numerical attributes. The rationale lies in KNN's non-parametric nature, which makes it suitable for datasets with complex and non-linear relationships. We use Mode for categorical attributes. We iterate from $K = 1$ to 10 to find the optimal value. Lastly, data transformation modifies the data format for better analysis, employing techniques such as label encoding for categorical variables, ordinal encoding for ordered data, and one-hot encoding for categorical columns, ultimately improving data integrity and the overall workflow efficiency. In the first experiment, the following steps were performed:

- 1) Data Transformation: Transforming raw data into a format suitable for machine learning models.
- 2) Classification with LightGBM: Performing classification using the LightGBM (Light Gradient Boosting Machine) algorithm, known for its efficiency in handling large datasets with high performance.
- 3) Evaluation: Evaluating the performance of the LightGBM model using appropriate evaluation metrics, such as accuracy, precision, recall, and F1-score.

In the second experiment, an additional step of DFS feature engineering was conducted before proceeding to the subsequent steps:

- 1) DFS Feature Engineering: Conducting feature engineering using the Deep Feature Synthesis (DFS) technique. This method automatically generates new features from existing raw data, which can enhance the performance of the machine learning model.
- 2) Data Transformation: Like the first experiment, transforming the engineered features into a format suitable for the model.
- 3) Classification with LightGBM: Performing classification using the LightGBM algorithm.
- 4) Evaluation: Evaluating the model's performance after feature engineering, using the same evaluation metrics.

2.2. Feature Engineering

In the implementation of Feature Engineering in credit risk cases is analyzed using DFS. Deep Feature Synthesis is a technique that allows for the automation of creating new attributes from existing ones in a datasets, which is beneficial for handling large and complex datasets. This process begins with the identification of entities and relationships within the datasets, which are then used by DFS to create attributes that combine the existing ones. DFS utilizes machine learning algorithms to explore and create combinations of attributes that may not be manually identified by human analysts, enhancing the efficiency and accuracy of predictive models. The results of applying DFS in this study show an improvement in the performance of predictive models in assessing credit risk, compared to traditional feature engineering techniques. This analysis underscores the importance of automation in the feature creation process to address the challenges in big data, as well as the potential use of DFS in improving the accuracy and efficiency of credit risk models in the financial industry. The method used in this research [11] is Deep Feature Synthesis (DFS) for automating feature engineering, and LightGBM for feature selection. DFS was chosen for its ability. It automatically generates many relevant features from previous transaction data. This has proven to be effective in detecting fraud. In this study, DFS is applied to address the issues of high dimensional, long training times, and low interpretability of features by limiting the depth of the DFS algorithm.

This research analyzes the implementation of DFS in feature engineering for credit risk cases. DFS enables the automation of creating new attributes from existing ones in a datasets, which is particularly useful for handling large and complex datasets. The process begins with the identification of entities and relationships within the datasets, which DFS then uses to create new attributes by combining existing ones. The results show that using DFS improves the performance of predictive models in assessing credit risk compared to traditional feature engineering techniques, highlighting the importance of automation in feature creation to address the challenges of big data and enhance the accuracy and efficiency of credit risk models.

In this study, DFS is employed to address issues such as high dimensionality, long training times, and low interpretability of features by limiting the depth of the DFS algorithm. The feature engineering process begins with creating an Entity Set that outlines the data entities used in the model. Relationships between entities, such as 'loan_applicant' and 'transactions,' are then established, allowing DFS to generate more complex attributes from the existing data. The results of the DFS process are displayed in a feature matrix, which can be used in a machine learning model to predict borrower eligibility.

2.3. LightGBM

LightGBM is a powerful gradient boosting algorithm that utilizes a leaf-wise approach for tree growth vertically, enhancing its efficiency in handling large datasets. Unlike traditional methods, LightGBM selects the branch of the tree that most reduces the loss for tree growth. This algorithm employs a histogram-based approach to find the best splitting candidates. LightGBM has hyperparameters that can be adjusted to optimize the algorithm's performance. The application of hyperparameters to the algorithm can be done in various ways. The parameters include crucial configurations such as the boosting type, learning rate, and the number of estimators, among others, to optimize the model's performance for binary classification tasks. These values have been carefully chosen to ensure that the model is both accurate and efficient, balancing between bias and variance. The objective is set to binary classification, with other parameters adjusted to enhance the model's capability in handling the specific characteristics of the dataset. For instance, the learning rate and the number of leaves are tuned to balance the trade-off between model accuracy and training time. Additionally, the bagging and feature fractions are set to improve generalization and prevent overfitting, ensuring the model can make accurate predictions.

3. RESULTS AND DISCUSSION

The experiments divide into two parts: feature engineering experiments with Deep Feature Synthesis, experiments without Deep Feature Synthesis, and experiments with LightGBM.

3.1 Experiment Performance

This section presents the results obtained from the two experiments conducted: with DFS and without DFS.

3.1.1 Experiments with Deep Feature Synthesis

We explore various techniques and strategies in feature creation, as well as analyze their impact on the performance of our model. Table 1 the results of these experiments provide valuable insights into the potential and effectiveness of DFS in the context of credit risk in.

Table 1. Evaluation Matrix With DFS

Data Test	Precision	Recall	F1-Score	Accuracy
0	1.00	1.00	1.00	0.99
1	0.97	0.94	0.96	0.99

For class 0, the precision, recall, and f1-score are 1.00, and accuracy is 0.99, indicating the model's consistency in identifying good loans. For class 1, the precision is 0.97, the recall is 0.94, f1-score is 0.96 and accuracy 0.99.

3.1.2 Experiments Without Deep Feature Synthesis

We focus on the use of original features and testing our credit risk model without additional features generated by DFS. We analyze the performance comparison between this approach and the approach with DFS, as well as provide an in-depth understanding of the strengths and weaknesses of each approach.

Table 2. Evaluation Matrix Without DFS

Data Test	Precision	Recall	F1-Score	Accuracy
0	0.97	1.00	0.99	0.97
1	0.99	0.68	0.81	0.97

In Table 2 for class 0 (good loans), the precision is 0.97, recall is 1.00, and F1-score is 0.99 and accuracy is 0.97. The accuracy showing the model's reliability in correctly classifying good loans.

For class 1 (bad loans), the precision is 0.99, recall is 0.68, and F1-score is 0.81. This demonstrates that while the model is not too good at identifying bad loans (perfect precision and near-perfect recall), the F1-score is lower compared to class 0.0, indicating some room for improvement in balancing precision and recall for bad loans. The overall accuracy remains high at 0.97, showing that the model performs well even without the additional features generated by DFS.

3.2 Discussion

The implementation of Deep Feature Synthesis (DFS) has been shown to enhance the performance of the LightGBM model in credit risk classification. Before applying DFS, the model's precision for "good loans" and "bad loans" is 0.97 and 0.99, respectively, with a recall of 1.00 for "good loans" and only 0.68 for "bad loans." After applying DFS, the precision for "good loans" improved to 1.00, while the precision for "bad loans" slightly decreased to 0.97. However, the most significant change was the increase in recall for "bad loans," which rose substantially to 0.94. The F1-score also demonstrated improvement after DFS was applied, reaching 1.00 for "good loans" and 0.96 for "bad loans," compared to 0.99 and 0.81 in the experiment without DFS. Overall model accuracy also increased from 0.97 to 0.99 after DFS was implemented, indicating that the model became more effective in detecting and classifying both "good loans" and "bad loans." These results suggest that DFS can enhance the quality of information used by the model, particularly in detecting "bad loans," which were previously more

challenging to identify, while LightGBM effectively leverages the abundance and quality of data to improve overall performance. Figure 2 is a visualization of the results from the experiments conducted.

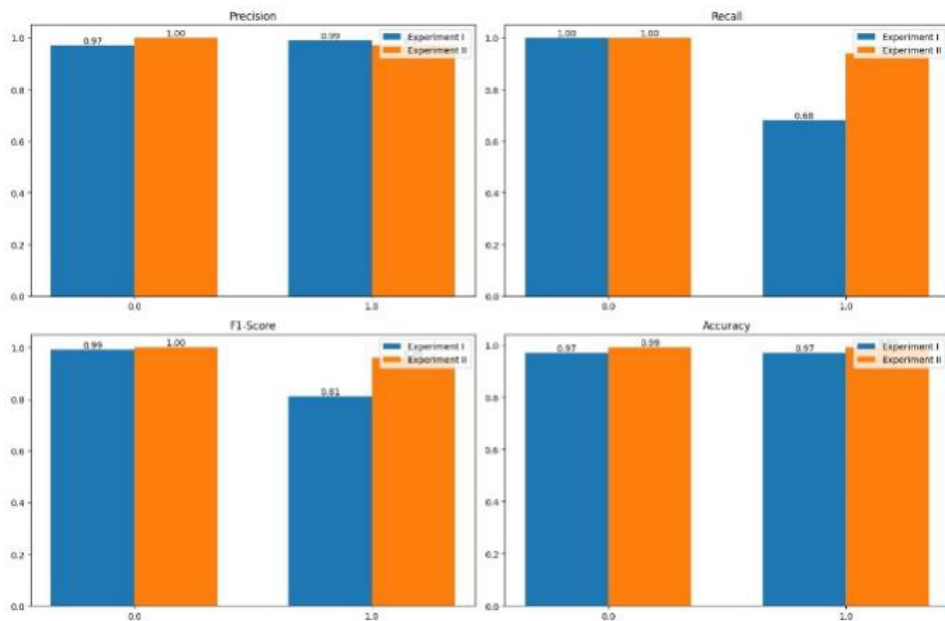


Figure 2. Comparison of evaluation matrix results

Figure 2 demonstrates that the model is highly effective in identifying good loans, with very high precision, recall, F1-score, and accuracy. However, for bad loans, although the precision remains high, the recall and F1-score indicate room for improvement in detecting all cases of bad loans. The model's accuracy remains high in both experiments, showing that the model generally performs well in predicting loan quality. During the experimental process, validating the results is crucial. By conducting the experiment without DFS, we can ensure that the performance improvement is not due to chance or overfitting but is genuinely a result of using DFS. The experiment without DFS also helps us understand how informative the original features in the dataset are. This can provide valuable insights into the quality and characteristics of the data we have and how the original features impact the model.

4. CONCLUSION

The use of Deep Feature Synthesis (DFS) in the Light Gradient Boosting Machine (LGBM) model has demonstrated a significant enhancement in performance for credit risk classification compared to models without DFS. By leveraging DFS, the model achieved exceptionally high accuracy, with precision, recall, and F1-scores

for both good and bad loan classifications indicating a strong and reliable ability to assess credit risk. In contrast, the model without DFS exhibited noticeably lower performance, highlighting the value of feature engineering in generating more informative and relevant features. These findings have practical implications for fintech platforms and credit scoring systems, as they underscore the potential of DFS to improve the reliability and accuracy of risk assessments, ultimately aiding in better decision-making processes. Future research could explore applying this method to diverse datasets to assess its generalizability across different financial contexts or investigate the integration of other advanced feature engineering techniques to further optimize model performance.

REFERENCES

- [1] J. C. Westland, T. Q. Phan, and T. Tan, "Private Information, Credit Risk and Graph Structure in P2P Lending Networks," pp. 1–31, 2018.
- [2] Y. Wang, Y. Zhang, Y. Lu, and X. Yu, "A Comparative Assessment of Credit Risk Model Based on Machine Learning ——a case study of bank loan data," *Procedia Comput. Sci.*, vol. 174, pp. 141–149, 2020, doi: 10.1016/j.procs.2020.06.069.
- [3] D. Li, S. Na, T. Ding, and C. Liu, "Credit risk management of p2p network lending," *Teh. Vjesn.*, vol. 28, no. 4, pp. 1145–1151, 2021, doi: 10.17559/TV-20200210110508.
- [4] A. Fatahuddin, P. Studi Akuntansi, and F. Ekonomi dan Bisnis, "Analisis Risiko pada Platform," pp. 209–218, 2020.
- [5] F. Doko, S. Kalajdziski, and I. Mishkovski, "Credit Risk Model Based on Central Bank Credit Registry Data," *J. Risk Financ. Manag.*, vol. 14, no. 3, 2021, doi: 10.3390/jrfm14030138.
- [6] S. Kokate and M. S. R. Chetty, "Credit risk assessment of loan defaulters in commercial banks using voting classifier ensemble learner machine learning model," *Int. J. Saf. Secur. Eng.*, vol. 11, no. 5, pp. 565–572, 2021, doi: 10.18280/IJSSE.110508.
- [7] C. Guan, H. Suryanto, A. Mahidadia, M. Bain, and P. Compton, "Responsible Credit Risk Assessment with Machine Learning and Knowledge Acquisition," *Human-Centric Intell. Syst.*, vol. 3, no. 3, pp. 232–243, 2023, doi: 10.1007/s44230-023-00035-1.
- [8] S. Lessmann, B. Baesens, H. V. Seow, and L. C. Thomas, "Benchmarking state-of-the-art classification algorithms for credit scoring: An update of research," *Eur. J. Oper. Res.*, vol. 247, no. 1, pp. 124–136, 2015, doi: 10.1016/j.ejor.2015.05.030.
- [9] V. Hassija *et al.*, "Interpreting Black-Box Models: A Review on Explainable Artificial Intelligence," *Cognit. Comput.*, vol. 16, no. 1, pp. 45–74, 2024, doi: 10.1007/s12559-023-10179-8.
- [10] Z. Liu, Y. Wang, F. Feng, Y. Liu, Z. Li, and Y. Shan, "A DDoS Detection

- Method Based on Feature Engineering and Machine Learning in Software-Defined Networks,” *Sensors*, vol. 23, no. 13, 2023, doi: 10.3390/s23136176.
- [11] J. M. Kanter and K. Veeramachaneni, “Deep feature synthesis: Towards automating data science endeavors,” *Proc. 2015 IEEE Int. Conf. Data Sci. Adv. Anal. DSAA 2015*, no. c, 2015, doi: 10.1109/DSAA.2015.7344858.
- [12] J. Yang *et al.*, “Automatic Feature Engineering-Based Optimization Method for Car Loan Fraud Detection,” *Discret. Dyn. Nat. Soc.*, vol. 2021, 2021, doi: 10.1155/2021/6077540.
- [13] N. Thomas Rincy and R. Gupta, “Feature Selection Techniques and its Importance in Machine Learning: A Survey,” *2020 IEEE Int. Students’ Conf. Electr. Electron. Comput. Sci. SCEECs 2020*, 2020, doi: 10.1109/SCEECs48394.2020.189.
- [14] S. K. Trivedi, “A study on credit scoring modeling with different feature selection and machine learning approaches,” *Technol. Soc.*, vol. 63, no. September, p. 101413, 2020, doi: 10.1016/j.techsoc.2020.101413.
- [15] R. Saxena, S. K. Sharma, M. Gupta, and G. C. Sampada, “A Novel Approach for Feature Selection and Classification of Diabetes Mellitus: Machine Learning Methods,” *Comput. Intell. Neurosci.*, vol. 2022, 2022, doi: 10.1155/2022/3820360.
- [16] S. B. Coşkun and M. Turanlı, “Credit risk analysis using boosting methods,” *J. Appl. Math. Stat. Informatics*, vol. 19, no. 1, pp. 5–18, 2023, doi: 10.2478/jamsi-2023-0001.