

Credit rating by clustering algorithm in the Vietnam Stock Exchange market

Tam Phan Huy^{1,2,*}, Thuy Chu Quang^{1,2}



Use your smartphone to scan this QR code and download this article

ABSTRACT

This study employs the K-means clustering algorithm to develop a corporate credit rating framework tailored to the Vietnamese market. By analyzing financial data from 568 non-financial firms listed on the Ho Chi Minh City Stock Exchange and the Hanoi Stock Exchange between 2019 and 2023, the research identifies vital financial indicators, including financial health ratios, management efficiency ratios, growth ratios, and dividend payout ratios. The K-means clustering model effectively categorizes these companies into six distinct clusters, each representing different levels of financial performance and credit risk. The clusters range from A+ (very low credit risk) to C (very high credit risk), providing a clear differentiation based on financial stability and operational efficiency. This systematic approach offers valuable insights for investors, managers, and government agencies, enhancing their ability to make informed decisions. Despite some limitations, such as reliance on historical data and sensitivity to initial cluster centroids, the K-means clustering model proves to be a robust starting point for assessing the creditworthiness of companies. This research contributes to the growing body of literature on machine learning applications in credit rating by demonstrating the superiority of clustering algorithms over traditional methods. It highlights how financial health and management efficiency indicators can be integrated into a data-driven framework to enhance credit risk assessment. The results suggest that the K-means clustering approach improves the accuracy of credit ratings and promotes transparency and efficiency in the financial market. Furthermore, the proposed framework can be a foundation for developing more sophisticated models, incorporating additional financial and non-financial variables. Future research could expand on this by integrating real-time data and exploring the impact of external economic factors on credit risk. By leveraging advanced machine learning techniques, this study paves the way for more reliable and comprehensive credit rating systems, ultimately supporting the stability and growth of financial markets in emerging economies like Vietnam.

Key words: K-Means, Credit Rating, Clustering, Vietnam

¹University of Economics and Law, Ho Chi Minh City, Vietnam

²Vietnam National University Ho Chi Minh City, Vietnam.

Correspondence

Tam Phan Huy, University of Economics and Law, Ho Chi Minh City, Vietnam

Vietnam National University Ho Chi Minh City, Vietnam.

Email: tamphan.ntc@gmail.com

History

- Received: 17-5-2024
- Revised: 23-7-2024
- Accepted: 27-9-2024
- Published Online:

DOI :



Copyright

© VNUHCM Press. This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International license.



1 INTRODUCTION

In today's fiercely competitive market, all enterprises must utilize their resources efficiently. Companies with high financial leverage ratios often mobilize short-term capital through credit¹. Some surveys also indicate that most businesses utilize credit². In the banking sector, efficiency and productivity can be measured by the profits from loans extended to customers. As a result, the credit rating process, used to measure credit risk, has become an important issue in recent years³. With accurate business credit ratings, investors and financial institutions can make better investment and lending decisions. Additionally, credit ratings serve as a reference channel, increasing transparency in the market. Current credit rating methods and indicators often rely on financial statements and credit information of businesses⁴. The evaluation mainly focuses on borrowing situations, operational efficiency, debt collection ability, and as-

set utilization efficiency. Globally, credit ratings are usually performed by large and well-established credit rating agencies such as Standard & Poor's (S&P), Moody's, and Fitch Group. In Vietnam, many banks have developed and implemented their own internal credit scoring systems tailored to their specific needs and criteria. The Credit Information Centre (CIC) under the State Bank of Vietnam is a notable entity that provides credit information for customers who have borrowed from the commercial banking system. However, it does not perform business credit ratings. These internal systems and thallowsormation from CIC allow banks to better manage and assess the credit risk of their clients. Although domestic credit ratings have been implemented, they still face limitations in terms of data and tools, so only a few units perform this activity professionally and publicly. In academics, few published research works related to domestic business credit ratings have been published.

Cite this article : Huy T P, Quang T C. **Credit rating by clustering algorithm in the Vietnam Stock Exchange market.** *Sci. Tech. Dev. J. - Eco. Law Manag.* 2024; ():1-18.

39 Moreover, the increasing risks in lending highlight
 40 the necessity for robust corporate credit ratings. Cur-
 41 rently, most credit ratings are conducted internally
 42 by commercial banks, which means that external in-
 43 vestors do not have access to comprehensive credit in-
 44 formation. This lack of transparency can lead to unin-
 45 formed investment decisions and increased financial
 46 instability. Therefore, establishing a standardized and
 47 publicly accessible credit rating system is crucial for
 48 providing investors with the information they need to
 49 make well-informed decisions, ultimately promoting
 50 a more stable and transparent financial market.

51 Thus, business credit ratings in Vietnam are a fasci-
 52 nating and practical topic in the financial field. Re-
 53 search on this subject will help us better understand
 54 the credit rating process, the factors affecting this pro-
 55 cess, and the methods for evaluating business credit
 56 rankings. Furthermore, with a reasonable credit rat-
 57 ing basis, financial institutions can make decisions on
 58 granting loans or raising credit limits for businesses,
 59 and investors can gain a broader perspective on busi-
 60 nesses' financial stability, enabling them to make in-
 61 formed investment decisions.

62 Currently, most business credit risk ratings are con-
 63 ducted by experts, but this method is not immune
 64 to human risks and disagreements among experts.
 65 Therefore, applying machine learning to the business
 66 credit rating process can help reduce workload, min-
 67 imize disagreements and human risks, and increase
 68 evaluation accuracy. Through machine learning al-
 69 gorithms, we can perform calculations of financial
 70 indicators for thousands of businesses and visualize
 71 analyses automatically and quickly. In the long run,
 72 by combining theoretical foundations with compu-
 73 tational power, financial institutions with clear data
 74 structures and fast information updates will be able to
 75 proactively assess business credit ratings in real time.

76 The objective of this research is to develop a corpo-
 77 rate credit rating framework specifically tailored for
 78 the Vietnamese market, utilizing the K-means clus-
 79 tering algorithm. This framework leverages data from
 80 the financial statements of non-financial firms listed
 81 on the Ho Chi Minh City Stock Exchange and the
 82 Hanoi Stock Exchange from 2019 to 2023. By ana-
 83 lyzing key financial indicators such as financial health
 84 ratios, management efficiency ratios, growth ratios,
 85 and dividend payout ratios, the framework aims to
 86 categorize companies into distinct clusters that re-
 87 flect their credit risk levels. This systematic and data-
 88 driven approach will provide investors, lenders, and
 89 other stakeholders with a clearer understanding of

these companies' creditworthiness and financial sta- 90
 bility, thereby promoting more informed decision- 91
 making and contributing to a more transparent and 92
 efficient financial market. 93

LITERATURE REVIEW 94

Background theories 95

Credit rating through clustering is an innovative ap- 96
 proach that combines both financial theories and ma- 97
 chine learning techniques to assess the creditworthi- 98
 ness of businesses. The foundational financial theo- 99
 ries related to this topic include the Modigliani-Miller 100
 theorem, the Trade-off theory, and the Pecking Or- 101
 der theory. These theories focus on firms' capital 102
 structure, the implications of their financing choices 103
 on overall credit risk, and the foundation of machine 104
 learning and clustering algorithms^{5,6}. 105

Modigliani-Miller Theory 106

The Modigliani-Miller (M-M) theorem, proposed by 107
 Franco Modigliani and Merton Miller in 1958, is an 108
 influential financial theory that lays the groundwork 109
 for understanding the relationship between a firm's 110
 capital structure and its credit risk. The M-M theorem 111
 posits that a firm's value is independent of its capital 112
 structure under certain assumptions such as no taxes, 113
 no bankruptcy costs, and perfect capital markets⁵. In 114
 other words, the choice between debt and equity fi- 115
 nancing does not impact a firm's overall value. 116

In the context of the research topic on credit rating 117
 by clustering, the Modigliani-Miller theorem is crucial 118
 in establishing the fundamental principles of cap- 119
 ital structure and financing choices. Despite the theo- 120
 rem's assumptions not holding in the real world, it 121
 still provides a theoretical foundation that helps re- 122
 searchers and practitioners understand how different 123
 financing choices may affect a firm's credit risk. Re- 124
 searchers can identify relevant financial ratios and in- 125
 dicators that reflect a company's credit risk by exami- 126
 ning the deviations from the M-M theorem's assump- 127
 tions, such as the presence of taxes and bankruptcy 128
 costs. For example, higher leverage ratios, which 129
 represent the proportion of debt in a firm's capi- 130
 tal structure, may indicate a higher credit risk due 131
 to the increased likelihood of financial distress and 132
 bankruptcy. These financial ratios can then be used as 133
 input features for clustering algorithms, which group 134
 companies with similar financial profiles and credit 135
 risk characteristics⁷. 136

137 **Trade-off theory**

138 The Trade-off theory is a significant financial concept relevant to the research topic of clustering-based credit rating. This theory posits that companies strive to find an optimal balance between the advantages and disadvantages of debt financing to minimize their overall capital costs⁸. Debt financing's primary benefit stems from tax shields gained through interest payments, while the costs are associated with a heightened risk of financial distress and bankruptcy resulting from increased leverage.

148 In relation to credit rating through clustering, the Trade-off theory aids in pinpointing essential financial ratios and indicators that signify a company's credit risk. For example, a business with elevated leverage ratios may be more vulnerable to financial distress, while one with lower leverage ratios might possess a more stable capital structure and, consequently, reduced credit risk. Furthermore, the theory implies that companies with greater profitability and diminished bankruptcy risk will likely have superior credit ratings, as they can accommodate higher debt levels. By leveraging the insights offered by the Trade-off theory, researchers can select pertinent financial ratios, such as those about leverage, liquidity, and profitability, as input variables for clustering algorithms. Subsequently, these algorithms, including hierarchical clustering, k-means clustering, and density-based clustering, can be employed to categorize companies based on similar financial characteristics and credit risk profiles⁷.

168 **Pecking Order theory**

169 The Pecking Order theory is another crucial financial concept relevant to the research topic of credit rating using clustering methods. This theory posits that firms prioritize their financing sources based on the information asymmetry and costs associated with each option, preferring internal financing first, followed by debt, and finally equity financing⁹. The rationale behind this order is that internal financing minimizes asymmetric information problems, while equity financing is considered the most expensive due to the adverse selection issue arising from information asymmetry.

181 In clustering-based credit rating, the Pecking Order theory helps identify vital financial ratios and indicators that reflect a company's credit risk. For instance, a firm that relies heavily on debt financing, as opposed to equity financing, may have a higher credit risk due to the potential for financial distress. On the other hand, companies with a greater reliance

on internal financing and lower debt levels might exhibit lower credit risk. By incorporating the insights derived from the Pecking Order theory, researchers can choose relevant financial ratios, such as leverage, liquidity, and profitability ratios, as input features for clustering algorithms. These algorithms, including hierarchical clustering, k-means clustering, and density-based clustering, can then be utilized to group companies with similar financial characteristics and credit risk profiles¹⁰.

198 In the context of credit rating by clustering, financial theories, such as the Modigliani-Miller theorem, Trade-off theory, and Pecking Order theory, can be employed to identify relevant financial ratios and indicators that reflect a company's credit risk. Key financial ratios include leverage ratios (e.g., debt-to-equity and debt-to-assets), liquidity ratios (e.g., current and quick ratios), profitability ratios (e.g., return on assets and return on equity), and efficiency ratios (e.g., asset turnover and inventory turnover)¹¹. These financial ratios and indicators serve as the basis for clustering algorithms, which analyze patterns in large datasets to group companies with similar financial profiles and credit risk characteristics. Machine learning techniques, such as hierarchical clustering, k-means clustering, and density-based clustering, are particularly well-suited for this task.

215 Hierarchical clustering creates a tree-like structure, called a dendrogram, representing the hierarchical relationships between different clusters¹². This approach allows for a more intuitive understanding of the relationships between clusters, which can be particularly helpful for credit rating purposes. K-means clustering is a popular centroid-based clustering algorithm that partitions the dataset into a predefined number of clusters by minimizing the within-cluster sum of squared distances¹³. This technique provides a simple and efficient way to group companies based on their financial ratios, thus facilitating comparisons of credit risk across different firms.

228 Combining these machine learning techniques and financial theories allows for a more comprehensive and data-driven approach to credit rating. This could potentially improve the accuracy and reliability of credit assessments and aid investors, lenders, and other stakeholders in their decision-making process.

234 **The foundation of machine learning**

235 The foundation of machine learning lies in its ability to learn patterns and make predictions from data without explicit programming for each specific task.

Machine learning algorithms, such as clustering, classification, and regression, are designed to identify underlying structures in data, enabling more accurate and automated decision-making processes. In the context of credit rating, clustering algorithms like K-means play a crucial role in categorizing companies based on their financial profiles.

Clustering algorithms are unsupervised learning techniques that group data points based on similarity measures. K-means clustering, one of the most widely used clustering algorithms, partitions data into k distinct clusters by minimizing the within-cluster variance¹³. This algorithm operates iteratively, assigning each data point to the nearest cluster centroid and recalculating centroids until convergence. Various studies have demonstrated the effectiveness of K-means clustering in financial applications, including credit rating¹⁴.

The advantage of using machine learning, mainly clustering algorithms, in credit rating, lies in its ability to handle large datasets and uncover complex patterns that may not be evident through traditional methods. Clustering algorithms can provide a more nuanced and data-driven credit risk assessment by analyzing a comprehensive set of financial indicators. This approach enhances the objectivity, consistency, and transparency of credit ratings, addressing many of the limitations associated with traditional expert-driven methods.

Incorporating machine learning into credit rating processes aligns with the broader trend of leveraging big data and advanced analytics in financial decision-making. As financial markets become increasingly complex, the ability to process and analyze large volumes of data efficiently is crucial for maintaining accurate and reliable credit assessments. Studies have shown that machine learning models, including clustering algorithms, outperform traditional statistical methods in various aspects of credit risk prediction^{15,16}.

By combining the theoretical foundations of capital structure with the analytical power of machine learning, credit rating through clustering represents a significant advancement in credit risk assessment. This innovative approach not only improves the accuracy and reliability of credit ratings but also provides valuable insights into the financial health and stability of businesses, ultimately supporting more informed investment and lending decisions.

287 Credit Rating Methods

288 One of the earliest and most prominent methods in
289 this group of credit rating systems was developed by

Moody's Investors Service in 1909¹⁷. Moody's employed an alphabetical rating system to assess the debt repayment ability of businesses. In descending order, the ratings are Aaa, Aa, A, Baa, Ba, B, Caa, Ca, C, with Aaa being the safest and C being the most dangerous. This method uses the following primary criteria to evaluate a company's debt repayment ability:

- Debt and interest repayment capacity: This is the most crucial factor in assessing a company's ability to repay its debt. Moody's evaluates a company's capacity to repay its principal and interest based on its profitability, assets, and debt repayment history.
- Financial health: This criterion is assessed based on measurements of outstanding debt, net assets, profitability, and cash flow.
- Market and competition: Moody's assess the market in which a company operates, including its competitors, pricing power, and value creation for shareholders.
- Management and business strategy: This includes evaluations of innovation, adaptability to the business environment, and motivation to create value for shareholders.

In addition to Moody's credit rating method, Standard & Poor's (S&P) introduced its credit rating system in 1917¹⁷. They also use an alphabetical rating system to assess the creditworthiness of businesses but employ different symbols to distinguish rating levels. The S&P credit rating method uses various criteria to evaluate a company's debt repayment ability, including:

- The company's financial situation: This is the most important factor used to assess a company's debt repayment ability. It includes indicators such as debt-to-total assets ratio, return on equity, free cash flow, and financial leverage.
- Product and service diversification: A company with diversified products and services is better able to mitigate risks than one focused on a single business area.
- Market position: A company's market position is assessed by examining market share and industry competition. A company with a strong market position is better able to maintain sales and profits.
- Management and business strategy: S&P also assesses the ability of the company's leadership to manage the business and its overall business strategy.

339 • External factors: S&P considers external factors
 340 such as the impact of the economic, political,
 341 and legal environment on the company.

342 Furthermore, Fitch Ratings introduced another credit
 343 rating method in 1913¹⁷. Like the other agencies,
 344 Fitch Ratings uses an alphabetical rating system with
 345 different symbols to distinguish rating levels. The
 346 Fitch Ratings method uses various evaluation criteria
 347 to assess a company's debt repayment ability, includ-
 348 ing:

- 349 • Financial Strength: This criterion assesses a
 350 company's financial ability, including its prof-
 351 itability, cash flow management, debt repay-
 352 ment capacity, and market opportunity seizing.
- 353 • Operating Performance: This criterion evalu-
 354 ates a company's ability to achieve its long-term
 355 operational objectives, including growth, prof-
 356 itability, and cost reduction.
- 357 • Business Profile: This criterion assesses a com-
 358 pany's ability to maintain and grow its sales,
 359 profits, and market share in the industry, includ-
 360 ing strategic direction, human resource man-
 361 agement, and customer relations.
- 362 • Risk Management: This criterion evaluates a
 363 company's ability to manage and control risks
 364 in its business operations, including credit risk,
 365 market risk, capital risk, and environmental
 366 risk.

367 Globally, major credit rating agencies such as Stan-
 368 dard & Poor's (S&P), Moody's, and Fitch Ratings
 369 have established well-defined criteria for assessing the
 370 creditworthiness of companies. These criteria typ-
 371 ically include debt and interest repayment capacity,
 372 financial health, market and competition, manage-
 373 ment and business strategy, and external factors. Debt
 374 and interest repayment capacity evaluate a company's
 375 ability to repay its principal and interest based on its
 376 profitability, assets, and debt repayment history. Fi-
 377 nancial health is assessed by measuring outstanding
 378 debt, net assets, profitability, and cash flow. Mar-
 379 ket and competition consider the market in which a
 380 company operates, including its competitors, pricing
 381 power, and value creation for shareholders. Manage-
 382 ment and business strategy evaluate the company's in-
 383 novation, adaptability to the business environment,
 384 and motivation to create shareholder value. External
 385 factors consider the economic, political, and legal en-
 386 vironments affecting the company.

In Vietnam, commercial banks have developed inter-
 387 nal credit scoring systems to evaluate their clients, tai-
 388 lored to their specific needs and criteria. These inter-
 389 nal systems typically include liquidity, leverage, prof-
 390 itability, and efficiency ratios. Liquidity ratios, such
 391 as the current ratio and quick ratio, assess a com-
 392 pany's ability to meet short-term obligations. Leverage
 393 ratios, including debt-to-equity and debt-to-asset
 394 ratios, evaluate financial leverage. Profitability ratios,
 395 such as return on assets (ROA) and return on equity
 396 (ROE), measure financial performance. Efficiency ra-
 397 tios, like asset turnover and inventory turnover, gauge
 398 management efficiency.
 399

Given these established criteria, the input variables
 400 for the K-means model in this study are selected to
 401 provide a comprehensive assessment of a company's
 402 financial performance. The variables include finan-
 403 cial health ratios (quick ratio, current ratio, short-
 404 term liabilities to equity, short-term liabilities to as-
 405 set, debt to equity, debt to asset, long-term debt to equ-
 406 ity, and long-term debt to asset), management effi-
 407 ciency ratios (ROA, asset turnover, accounts receiv-
 408 able turnover, and payment period turnover), growth
 409 ratios (sales growth rate and EBIT growth rate), and
 410 the dividend payout ratio. These variables are es-
 411 sential for labeling the clusters obtained from the K-
 412 means algorithm and developing a robust credit rating
 413 system.
 414

By incorporating these financial variables as inputs
 415 for the K-means model, this study aims to create a
 416 comprehensive credit rating system that accurately re-
 417 flects various aspects of a company's financial perfor-
 418 mance and credit risk profile. The identified clusters
 419 will provide meaningful and reliable credit ratings for
 420 various stakeholders in the financial sector, ultimately
 421 promoting a more transparent and efficient financial
 422 market.
 423

Despite the widespread use of traditional credit rat-
 424 ing methods, these approaches have notable areas for
 425 improvement. Traditional methods often rely heavily
 426 on expert judgment, which can introduce subjectiv-
 427 ity and potential biases into the credit rating process.
 428 This subjectivity can lead to consistency in ratings, es-
 429 pecially when different experts assess the same com-
 430 pany. Additionally, traditional methods may need to
 431 efficiently handle large datasets or rapidly changing fi-
 432 nancial environments, making it difficult to provide
 433 timely and accurate credit ratings. They also need to
 434 improve in their ability to uncover complex patterns
 435 and relationships within financial data, as they often
 436 focus on a narrow set of financial indicators and his-
 437 torical performance.
 438

Machine learning techniques, particularly clustering algorithms like K-means, offer solutions to these limitations. Machine learning models can quickly process vast amounts of data and identify intricate patterns and relationships that human analysts may miss. By leveraging data-driven insights, machine learning can enhance the objectivity and consistency of credit ratings. Clustering algorithms, specifically, can group companies based on a comprehensive set of financial indicators, providing a more nuanced understanding of their credit risk profiles. This approach reduces the reliance on subjective expert judgment and improves the transparency and accuracy of the credit rating process.

Clustering Algorithm

This study employs the k-means algorithm as the primary machine learning technique to achieve the research objective. As discussed earlier, the k-means algorithm offers several advantages, including simplicity, computational efficiency, scalability, and proven effectiveness in various applications, particularly in finance and credit risk assessment. By utilizing k-means as the chosen machine learning algorithm, this research aims to effectively uncover patterns and groupings within the dataset, facilitating a deeper understanding of the relationships between financial and non-financial variables and credit ratings. Ultimately, the application of the k-means algorithm in this study is expected to contribute to improved credit rating prediction accuracy, providing valuable insights to support informed decision-making in the credit assessment process.

The k-means algorithm was chosen for this research topic on credit rating prediction for several reasons. First, the simplicity and computational efficiency of the k-means algorithm make it an attractive choice for researchers¹⁰. The algorithm's straightforward nature allows for rapid prototyping and experimentation, enabling researchers to quickly assess its potential utility in predicting credit ratings. Second, k-means has been proven effective in various applications, including finance and credit risk assessment. Its ability to identify patterns and groupings in data makes it suitable for uncovering distinct credit risk categories based on financial and non-financial variables. This feature can enhance the understanding of the underlying relationships between variables and credit risk, ultimately leading to better prediction accuracy.

Third, k-means is capable of handling large datasets efficiently¹³. As credit rating prediction often involves the analysis of large amounts of data from

numerous companies, the algorithm's scalability is a critical factor. K-means can process large datasets quickly, making it suitable for this research context. Lastly, k-means has been successfully applied in previous credit rating research, showing promising results in comparison to other techniques^{14,18}. Its previous success in the field adds credibility to its use in the current research topic and suggests that it may provide valuable insights into credit rating prediction. To summarize, the k-means algorithm's simplicity, effectiveness in various applications, scalability, and successful application in previous credit rating research make it a suitable choice for the current research topic. Its ability to efficiently handle large datasets and identify underlying patterns can contribute to improved credit rating prediction accuracy. The k-means algorithm is an unsupervised machine learning technique widely employed for clustering and partitioning datasets into meaningful groups^{10,13}. It aims to identify underlying structures and patterns in the data based on similarity among data points. The algorithm's simplicity, computational efficiency, and effectiveness in various applications make it a popular choice for researchers and practitioners¹⁰.

The k-means algorithm operates by initializing a predetermined number of centroids (k), representing the centers of each cluster. These centroids are generally initialized randomly within the dataset's feature space¹³. The algorithm then iteratively assigns each data point to the nearest centroid, based on a distance metric, such as Euclidean distance¹⁰. Once all data points are assigned to their respective centroids, the centroids are recalculated to represent the meaning of all data points within each cluster. This process is repeated until convergence is reached, i.e., the centroids' positions stabilize, or a predefined number of iterations have been completed¹³.

By partitioning the dataset into distinct groups, the k-means algorithm facilitates the identification of relationships between variables and allows researchers to uncover hidden patterns within the data¹⁰. In the context of credit rating prediction, the k-means algorithm can be applied to cluster companies based on their financial and non-financial characteristics, providing insights into the factors that drive credit risk and potentially contributing to improved prediction accuracy.

To evaluate the performance of the k-means algorithm in credit rating prediction, various performance metrics can be utilized. One standard method is the silhouette score, which measures the clustering quality by computing the average distance between observations within the same cluster and comparing it

545 to the average distance to the nearest neighboring
546 cluster¹⁹. A higher silhouette score indicates better-
547 defined clusters and implies that the algorithm has ef-
548 fectively identified distinct risk categories in the con-
549 text of credit rating prediction.

550 The elbow method is a popular technique to deter-
551 mine the optimal number of clusters (k) in k-means
552 clustering. It involves plotting the variance explained
553 or within-cluster sum of squared distances (WSS) as a
554 function of the number of clusters and identifying the
555 "elbow point," where adding more clusters does not
556 significantly reduce the WSS²⁰. The rationale behind
557 the elbow method is that as the number of clusters in-
558 creases, the WSS decreases since each additional clus-
559 ter can capture a portion of the remaining variance.
560 However, at some point, adding more clusters will
561 not lead to a substantial decrease in the WSS, and the
562 curve will begin to flatten. The elbow point represents
563 the number of clusters at which the diminishing re-
564 turns in variance reduction are no longer worth the
565 added complexity of having more clusters²¹. To im-
566 plement the elbow method, researchers can perform
567 k-means clustering for a range of cluster values (e.g., k
568 = 1 to k = 10) and compute the WSS for each value of
569 k. By visualizing the WSS values on a line chart, the el-
570 bow point can be identified, representing the optimal
571 number of clusters for the dataset.

572 In conclusion, employing the elbow method and sil-
573 houette score in this research provides a robust ap-
574 proach to determining the optimal number of clus-
575 ters for the k-means algorithm in credit rating predic-
576 tion. The elbow method allows us to identify the point
577 where adding more clusters does not significantly re-
578 duce the within-cluster sum of squared distances, en-
579 suring the model's simplicity without compromising
580 its explanatory power. On the other hand, the silhou-
581 ette score evaluates the quality of clustering by assess-
582 ing the cohesion within clusters and the separation
583 between them, ensuring that the chosen clusters are
584 meaningful and well-defined.

585 By combining the elbow method and silhouette score,
586 this research benefits from a comprehensive approach
587 to cluster selection, balancing the trade-off between
588 model complexity and prediction accuracy. These
589 techniques enhance the reliability and validity of the
590 credit rating predictions derived from the k-means
591 algorithm. It contributes to a better understanding
592 of the underlying relationships between variables and
593 credit risk. Ultimately, this approach can lead to more
594 accurate credit rating predictions, benefiting both fi-
595 nancial institutions and companies in their decision-
596 making processes.

Previous studies

597
598 In recent years, the application of machine learn-
599 ing techniques for predicting corporate credit ratings
600 has become an increasingly popular research topic.
601 A wide range of studies have explored various al-
602 gorithms, input variables, and methodologies to im-
603 prove the accuracy and reliability of credit rating pre-
604 dictions.

605 Early research laid the groundwork for using machine
606 learning in credit rating prediction. Huang et al.¹⁴
607 compared support vector machines (SVMs) to tra-
608 ditional statistical methods like linear discriminant
609 analysis and logistic regression, while Altman and
610 Sabato²² explored hybrid models that combined lo-
611 gistic regression with SVM. Both studies found that
612 machine-learning approaches outperformed conven-
613 tional methods in accuracy and robustness.

614 Subsequent research has built upon these initial find-
615 ings. Kim and Kang¹⁵, for example, investigated the
616 performance of decision trees, artificial neural net-
617 works (ANNs), and logistic regression in predicting
618 Korean firms' credit ratings. Their study demon-
619 strated that ANNs provided superior accuracy com-
620 pared to the other methods. Similarly, other stud-
621 ies have compared various machine learning algo-
622 rithms, such as logistic regression, decision trees, ran-
623 dom forests, SVMs, ANNs, and k-nearest neighbors
624 (KNN), to identify the best-performing models for
625 credit rating prediction²³⁻²⁵.

626 In terms of input variables, most studies have utilized
627 financial ratios related to liquidity, leverage, prof-
628 itability, and efficiency^{16,26}. However, some research
629 has also explored the incorporation of industry-
630 specific variables, such as asset turnover and net profit
631 margin as well as non-financial data like macroeco-
632 nomic indicators and textual information from news
633 articles²⁷. These studies have found that the inclu-
634 sion of industry-specific and non-financial variables
635 can improve the accuracy of credit rating prediction
636 models.

637 The performance of machine learning models in
638 credit rating prediction has been assessed using var-
639 ious evaluation metrics, such as accuracy, precision,
640 recall, and F1 score. Overall, the literature suggests
641 that machine learning algorithms can effectively pre-
642 dict corporate credit ratings using financial ratios
643 as input variables, and that incorporating industry-
644 specific and non-financial variables may further en-
645 hance the accuracy of these models^{14,16,22,25,27,28}.

646 In summary, the growing body of literature on pre-
647 dicting corporate credit ratings using machine learn-
648 ing models has demonstrated the potential of these

649 approaches in providing more accurate and reli- 700
 650 able predictions compared to traditional statistical 701
 651 methods. Researchers have explored various algo- 702
 652 rithms, input variables, and methodologies, and have 703
 653 found that a combination of financial ratios, industry- 704
 654 specific variables, and non-financial data can lead 705
 655 to improved performance in credit rating prediction. 706
 656 Future research may further refine these models and 707
 657 explore the potential of emerging machine learning 708
 658 techniques in this area. 709

659 **Research Gaps**

660 Despite the extensive research conducted on credit 710
 661 rating and risk assessment using machine learning 711
 662 techniques, several gaps remain that this study aims 712
 663 to address. Previous studies have predominantly fo- 713
 664 cused on well-established markets and large corpora- 714
 665 tions, leaving a significant gap in understanding the 715
 666 credit risk dynamics within emerging markets such 716
 667 as Vietnam. For instance, research by Huang et al. 14 717
 668 and Altman and Sabato 22 primarily explored the use 718
 669 of support vector machines (SVMs) and logistic re- 719
 670 gression in more developed markets, thereby limiting 720
 671 the applicability of their findings to the Vietnamese 721
 672 context. 722

673 Furthermore, while studies by Kim and Kang 15 and 723
 674 Barboza et al. 16 have shown the efficacy of machine 724
 675 learning models such as artificial neural networks 725
 676 (ANNs) and decision trees in credit rating prediction, 726
 677 they often neglect the specific financial indicators rel- 727
 678 evant to smaller firms and emerging economies. This 728
 679 study bridges this gap by incorporating a comprehen- 729
 680 sive set of financial ratios specifically tailored to non- 730
 681 financial firms listed on the Ho Chi Minh City Stock 731
 682 Exchange and the Hanoi Stock Exchange. 732

683 Additionally, the existing literature, including works 733
 684 by Abdou and Pointon 23 and Galindo and Tamayo 24, 734
 685 has largely overlooked the practical implementation 735
 686 challenges and the need for a standardized and pub- 736
 687 licly accessible credit rating framework in emerging 737
 688 markets. This study addresses this issue by proposing 738
 689 a robust credit rating system based on the K-means 739
 690 clustering algorithm, which enhances prediction ac- 740
 691 curacy but also provides a transparent and systematic 741
 692 approach to credit risk assessment. 742

693 Moreover, while the integration of non-financial data 743
 694 and industry-specific variables has been explored to 744
 695 some extent 26,27, there is still a lack of research focus- 745
 696 ing on the unique financial environments of emerg- 746
 697 ing markets. This study fills this void by analyzing key 747
 698 financial indicators such as liquidity ratios, leverage 748
 699 ratios, profitability ratios, and efficiency ratios, which 749

are crucial for assessing the creditworthiness of com- 700
 panies in Vietnam. 701

In conclusion, this research contributes to the existing 702
 body of knowledge by addressing these critical gaps 703
 and providing a nuanced understanding of credit risk 704
 assessment in the Vietnamese market. By leveraging 705
 machine learning techniques and a detailed set of fi- 706
 nancial indicators, this study offers a practical tool 707
 for financial institutions, investors, and policymakers 708
 to make informed decisions, ultimately promoting a 709
 more transparent and efficient financial market. 710

711 **METHODOLOGY**

712 **Data**

713 In this study, we focus on non-financial firms listed on 713
 both the Ho Chi Minh City Stock Exchange and the 714
 Hanoi Stock Exchange from 2019 to 2023. The ini- 715
 tial dataset comprised data collected from 692 firms. 716
 Upon inspection, observations with missing values 717
 or duplicates were identified and subsequently elim- 718
 inated from the dataset. Consequently, the refined 719
 dataset encompassed 568 firms, resulting in 2,567 720
 unique observations. The yearly distribution of com- 721
 panies within the dataset is as follows: 510 compa- 722
 nies in 2018, 525 companies in 2019, 534 companies 723
 in 2020, 532 companies in 2021, and 466 companies in 724
 2022. This comprehensive dataset offers a solid foun- 725
 dation for investigating the credit rating prediction 726
 of these non-financial firms using machine learning 727
 techniques. 728

729 **Input Variables**

730 The input data for the K-means model in this study 730
 comprises a comprehensive set of financial variables, 731
 which can be broadly categorized into four groups: fi- 732
 nancial health ratios, management efficiency ratios, 733
 growth ratios, and dividend payout ratio . These vari- 734
 ables provide a detailed assessment of a company’s fi- 735
 nancial performance and are essential criteria for la- 736
 beling the clusters obtained from the K-means algo- 737
 rithm as described in Table 1. 738

739 Financial health ratios include the quick ratio, current 739
 ratio, short-term liability on equity, short-term liabil- 740
 ity on the asset, long-term debt on equity, long-term 741
 debt on the asset, debt on equity, and debt on asset. 742
 These ratios offer insights into a company’s liquidity, 743
 solvency, and overall financial stability, capturing the 744
 its ability to meet its short-term and long-term obli- 745
 gations. 746

747 Management resource management comprise ROA, 747
 asset turnover, account receivable turnover, and pay- 748
 ment period turnover. These ratios evaluate a com- 749
 pany’s ability to generate returns from its assets and 750

751 the efficiency with which it manages its operations.
 752 Efficient management of resources is a critical factor
 753 in assessing a company’s creditworthiness, as it re-
 754 flects the firm’s capacity to generate profits and meet
 755 its financial commitments.
 756 Growth ratios, including sales and EBIT growth rates,
 757 capture a company’s ability to expand its operations
 758 and increase its earnings. Companies with strong
 759 growth potential are generally considered less risky,
 760 as their expanding revenue base allows them to ser-
 761 vice their debts better..
 762 Lastly, the dividend payout ratio is important in deter-
 763 mining a company’s financial health and credit risk.
 764 This ratio measures the proportion of earnings paid
 765 out to shareholders as dividends, providing insights
 766 into a firm’s ability to retain earnings for future growth
 767 and its commitment to returning value to sharehold-
 768 ers.
 769 By incorporating these financial variables as inputs for
 770 the K-means model, this study aims to develop a com-
 771 prehensive credit rating system that accurately reflects
 772 the various aspects of a company’s financial perfor-
 773 mance and credit risk profile. The identified clusters
 774 will be labeled based on their unique combination of
 775 these financial variables, providing a meaningful and
 776 reliable credit rating system for various stakeholders
 777 in the financial sector.

778 **RESULTS & DISCUSSION**

779 The elbow method graph displays a sharp decline in
 780 the SSE (sum of squared errors) from 900 to 400 as the
 781 number of clusters (k) increases from 1 to 5. After this
 782 point, the SSE continues to decrease, albeit at a slower
 783 rate, reaching around 300 at k=7.5. Beyond this point,
 784 the SSE exhibits a more gradual decline, decreasing to
 785 approximately 200 by the time k reaches 18.

786 Figure 1 suggests that the optimal value for k is around
 787 6 clusters, as the most significant reduction in SSE oc-
 788 curs up to that point. Beyond k=6, the SSE decreases
 789 at a diminished rate, indicating that adding more clus-
 790 ters does not contribute substantially to the reduction
 791 of the within-cluster sum of squared distances. There-
 792 fore, selecting k=6 strikes a reasonable balance be-
 793 tween model simplicity and its ability to capture the
 794 underlying patterns in the data, making it a suitable
 795 choice for credit rating prediction using the k-means
 796 algorithm.

797 According to Figure 2, upon analyzing the silhou-
 798 ette scores, we observe a gradual decline from 0.28 to
 799 approximately 0.25 as the number of clusters (k) in-
 800 creases from 1 to 5. The silhouette score remains re-
 801 latively stable, fluctuating around 0.25, as k increases
 802 from 5 to 8. However, beyond k=8, the silhouette

score experiences a sharp drop, decreasing to 0.2 as
 k continues to increase up to 20.

803
 804
 805 Considering the results from both the elbow method
 806 and silhouette score analyses, we can conclude that se-
 807 lecting k=6 is an appropriate choice for our credit rat-
 808 ing prediction model. With the elbow method reveal-
 809 ing a significant drop in SSE at k=6 and the silhou-
 810 ette score maintaining a relatively stable level around
 811 k=5 to k=8, it is reasonable to proceed with fitting
 812 the k-means model using k=6. This choice balances
 813 the trade-off between model complexity and perfor-
 814 mance, thus allowing us to effectively uncover the
 815 underlying relationships between variables and credit
 816 risk in our dataset.

817 In the three-dimensional space depicted in Figure 3,
 818 it is evident that the k-means clustering algorithm ef-
 819 fectively partitions the data into distinct clusters with
 820 clear convergence. To further assess the differences
 821 between these six clusters, it is necessary to examine
 822 additional graphical representations or employ de-
 823 scriptive statistical methods, as discussed below. By
 824 doing so, we can better understand the criteria that
 825 set each cluster apart and solidify our confidence in
 826 the effectiveness of using k=6 in the k-means cluster-
 827 ing algorithm for credit rating prediction.

Table 2: Number of observations for each cluster with K=6

Cluster	Observations
0	213
1	208
2	623
3	369
4	463
5	691

828 Table 2 displayed above provides a comprehensive
 829 overview of the distribution of observations within
 830 the six clusters generated by the k-means clustering
 831 algorithm. The different number of observations in
 832 each cluster suggests that the dataset comprises di-
 833 verse patterns and relationships, which have been suc-
 834 cessfully captured by the algorithm. Cluster 0 con-
 835 tains 213 observations, indicating a group of com-
 836 panies with certain shared characteristics. Similarly,
 837 Cluster 1 comprises 208 observations, revealing an-
 838 other set of companies with distinct features. Cluster
 839 2, the largest group with 623 observations, represents
 840 a significant portion of the dataset and highlights a
 841 more prevalent pattern among the companies. Clus-
 842 ter 3, consisting of 369 observations, and Cluster 4,

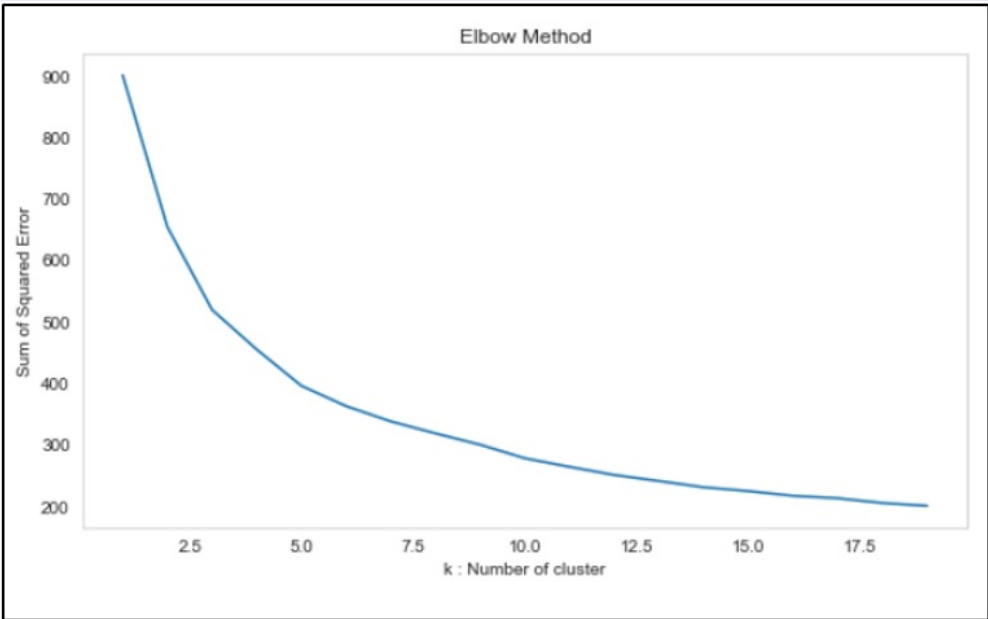


Figure 1: Sum of Squared Error by number of clusters (Source: Author's Calculation)

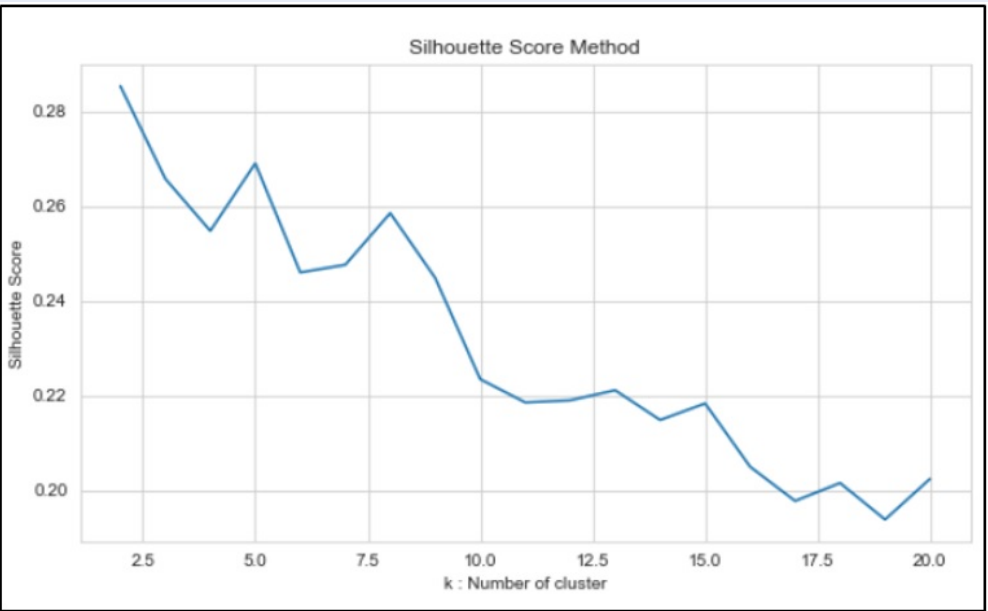


Figure 2: Silhouette Score by number of clusters (Source: Author's Calculation)

Table 1: Credit Rating Criteria and Measurement Methods

Criteria Group	Criteria	Measurement Method	Referenced Standards
Financial Health Ratios	Quick Ratio	Quick Assets / Current Liabilities	Standard & Poor's, Moody's, Fitch Ratings
	Current Ratio	Current Assets / Current Liabilities	Standard & Poor's, Moody's, Fitch Ratings
	Short-term Liabilities to Equity	Short-term Liabilities / Equity	Internal Standards of Vietnamese Commercial Banks
	Short-term Liabilities to Asset	Short-term Liabilities / Total Assets	Internal Standards of Vietnamese Commercial Banks
	Debt to Equity	Total Debt / Equity	Standard & Poor's, Moody's, Fitch Ratings
	Debt to Asset	Total Debt / Total Assets	Standard & Poor's, Moody's, Fitch Ratings
	Long-term Debt to Equity	Long-term Debt / Equity	Internal Standards of Vietnamese Commercial Banks
	Long-term Debt to Asset	Long-term Debt / Total Assets	Internal Standards of Vietnamese Commercial Banks
Management Efficiency Ratios	Return on Assets (ROA)	Net Income / Total Assets	Standard & Poor's, Moody's, Fitch Ratings
	Asset Turnover	Net Sales / Average Total Assets	Standard & Poor's, Moody's, Fitch Ratings
	Accounts Receivable Turnover	Net Credit Sales / Average Accounts Receivable	Internal Standards of Vietnamese Commercial Banks
	Payment Period Turnover	Number of Days in Period / Payables Turnover	Internal Standards of Vietnamese Commercial Banks
Growth Ratios	Sales Growth Rate	(Current Year Sales - Previous Year Sales) / Previous Year Sales	Standard & Poor's, Moody's, Fitch Ratings
	EBIT Growth Rate	(Current Year EBIT - Previous Year EBIT) / Previous Year EBIT	Standard & Poor's, Moody's, Fitch Ratings
Dividend Payout Ratio	Dividend Payout Ratio	Dividends / Net Income	Standard & Poor's, Moody's, Fitch Ratings

Source: by authors

843 with 463 observations, illustrate additional variations
 844 within the dataset. Lastly, Cluster 5 encompasses 691
 845 observations, making it the second-largest group and
 846 pointing to another common pattern among the companies.
 847
 848 These varying cluster sizes demonstrate the k-means
 849 algorithm's effectiveness in identifying and segregating
 850 diverse patterns within the dataset. The k-means
 851 clustering algorithm with k=6 has resulted in the formation
 852 of six distinct clusters, which the author proposes
 853 to use as the basis for a new credit rating system.
 854 This system is outlined in the Table 3 and consists of

the following credit ratings.

855
 856 The K-means clustering algorithm applied in this
 857 study identified six distinct clusters (0, 1, 2, 3, 4, 5),
 858 each representing different levels of financial performance
 859 and credit risk. These clusters provide valuable
 860 insights into the financial health and creditworthiness
 861 of the companies analyzed, which can be understood
 862 through theoretical, empirical, and practical lenses.

- 863 • Cluster 0 (C): Companies in Cluster 0 exhibit
 864 significant liquidity challenges and lower management
 865 efficiency. The high levels of both

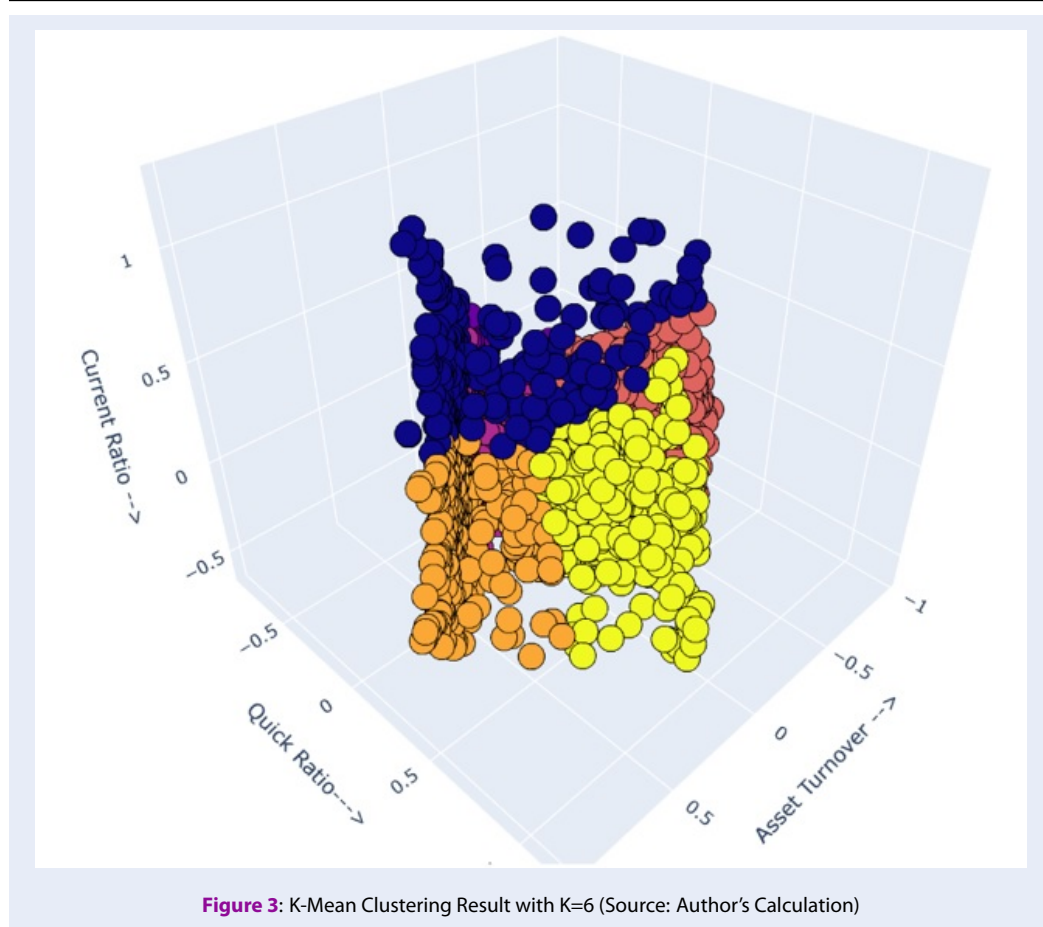


Table 3: Suggested label for credit scoring.

Label	Description
A+	Very good (very low credit risk)
A	Good (low credit risk)
B+	Fairly good (credit risk in the middle range from fair to good)
	Average (medium credit risk)
C+	Poor (high credit risk)
C	Very poor (very high credit risk)

Source: Author's Suggested

866 short-term and long-term debt indicate a sub-
 867 stantial credit risk. Theoretically, this aligns
 868 with the Pecking Order Theory⁹, suggesting that
 869 companies facing financial distress are more re-
 870 liant on debt. Empirically, the observed low re-
 871 turn on assets (ROA) and subpar growth rates
 872 support categorizing these companies as high-
 873 risk. Practically, investors and financial institu-
 874 tions should approach these firms with caution,
 875 considering their high likelihood of financial in-

stability.

- Cluster 1 (A+): This cluster is characterized
 877 by outstanding liquidity, low indebtedness, and
 878 strong financial health, positioning these com-
 879 panies as very low credit risk. The Trade-off
 880 Theory supports the high creditworthiness of
 881 firms with optimal leverage, which is evident in
 882 this cluster. Empirically, the high ROA and ef-
 883 ficient management practices confirm the the-
 884 oretical expectations. Practically, companies in
 885

886 this cluster are attractive investment opportuni- 939
 887 ties due to their financial stability and low risk 940
 888 of default. 941

- 889 • Cluster 2 (A): Companies in Cluster 2 also dis- 942
 890 play robust financial health with above-average 943
 891 management efficiency and growth potential. 944
 892 However, their liquidity is not as strong as that 945
 893 in Cluster 1. This finding is consistent with 946
 894 the Modigliani-Miller Theorem, which suggests 947
 895 that firm value is independent of capital struc- 948
 896 ture under certain conditions⁵. Empirically, the 949
 897 strong ROA and EBIT growth rate validate the 950
 898 theoretical foundation. Practically, these firms 951
 899 are still considered low-risk and are suitable 952
 900 candidates for investment, albeit with slightly 953
 901 higher caution than Cluster 1. 954
- 902 • Cluster 3 (B+): This cluster includes compa- 955
 903 nies with mixed financial health and manage- 956
 904 ment efficiency. While they have reasonable 957
 905 liquidity, their high debt levels increase credit 958
 906 risk. The theoretical backing from the Trade- 959
 907 off Theory indicates that these firms balance the 960
 908 benefits of debt with the risk of financial dis- 961
 909 tress. Empirically, the average ROA and above- 962
 910 average growth rates provide a nuanced under- 963
 911 standing of their creditworthiness. Practically, 964
 912 these companies offer moderate investment po- 965
 913 tential but require a thorough risk assessment. 966
- 914 • Cluster 4 (B): Firms in Cluster 4 show weaker 967
 915 financial health and lower management effi- 968
 916 ciency, coupled with higher debt ratios. The 969
 917 Pecking Order Theory again explains the re- 970
 918 liance on debt due to financial constraints. Em- 971
 919 pirically, their low ROA and mixed growth rates 972
 920 indicate medium credit risk. Practically, while 973
 921 investment in these firms is riskier, potential re- 974
 922 turns could be balanced against the higher risk, 975
 923 making them suitable for risk-tolerant investors. 976
- 924 • Cluster 5 (C+): Companies in this cluster have 977
 925 better financial health than those in Cluster 0 978
 926 but still face significant credit risk due to lower 979
 927 management efficiency and growth rates. The 980
 928 theoretical implications align with the Trade-off 981
 929 Theory, indicating an ongoing struggle to main- 982
 930 tain financial stability. Empirically, the find- 983
 931 ings of moderate ROA and low dividend payout 984
 932 ratios reinforce their classification. Practically, 985
 933 these firms are higher-risk investments, and in- 986
 934 vestors should be cautious. 987

935 This proposed credit rating system aims categorizes 988
 936 companies based on their credit risk levels, as deter- 989
 937 mined by the k-means clustering analysis. By assign- 990
 938 ing specific credit ratings to each cluster, the author

has established a comprehensive framework to assess 939
 the creditworthiness of companies. The ratings range 940
 from A+ for those exhibiting shallow credit risk to C 941
 for companies with very high credit risk. 942
 The suggested credit rating system provides a valuable 943
 tool for investors, financial institutions, and regula- 944
 tors to make informed decisions and assess the credit 945
 risk of different companies effectively. By leverag- 946
 ing the insights from the k-means clustering analysis, 947
 the proposed system captures the underlying relation- 948
 ships between financial and non-financial variables, 949
 contributing to determining credit risk levels. 950
 The k-means clustering algorithm with k=6 has suc- 951
 cessfully grouped the data into six distinct clusters, 952
 each with different characteristics regarding financial 953
 health, management efficiency, growth potential, and 954
 dividend payout capacity. These clusters offer valu- 955
 able insights into the various credit risk profiles and 956
 can aid in developing a credit rating system (see Ap- 957
 pendix 1 & 2). 958
 Upon examination of the clusters, it is evident that 959
 companies in Cluster 1 exhibit outstanding liquid- 960
 ity and low indebtedness, indicating strong financial 961
 health. However, they have lower growth rates and 962
 dividend payout ratios than the average. Cluster 2 963
 companies, on the other hand, demonstrate above- 964
 average management efficiency and growth potential 965
 but have average liquidity and lower dividend payout 966
 ratios. 967
 Clusters 3 and 4 present a more mixed picture, with 968
 companies in these groups showing weaker financial 969
 health and management efficiency, alongside varied 970
 growth potential. Both clusters have lower dividend 971
 payout ratios compared to the average. Companies in 972
 Cluster 5 display better financial health, average man- 973
 agement efficiency, and higher growth rates, but their 974
 dividend payout ratios remain low. Finally, Cluster 0 975
 companies face liquidity challenges and lower man- 976
 agement efficiency, along with average growth rates 977
 and below-average dividend payout ratios. 978
 These findings suggest that companies within each 979
 cluster share common financial and operational char- 980
 acteristics, which can help inform credit risk assess- 981
 ment and decision-making. It is crucial to note that 982
 further research, including the evaluation of addi- 983
 tional graphs and the application of descriptive statis- 984
 tical methods, is necessary to validate the differences 985
 between clusters and refine the proposed credit rat- 986
 ing system. Moreover, it is essential to consider exter- 987
 nal factors, such as market conditions and industry- 988
 specific risks, to ensure a comprehensive and accurate 989
 credit risk assessment. 990

991 Upon revisiting the clusters with the new naming con- 1041
 992 vention, the author proposed the following credit rat- 1042
 993 ing suggestions: Cluster 1 as A+, Cluster 2 as A, Clus- 1043
 994 ter 3 as B+, Cluster 4 as B, Cluster 5 as C+, and Clus- 1044
 995 ter 0 as C. This rating system aligns with the compa- 1045
 996 nies' observed financial and operational characteris- 1046
 997 tics within each cluster. 1047
 998 Companies in Cluster A+ (Cluster 1) demonstrate 1048
 999 exceptional financial health, while those in Cluster 1049
 1000 A (Cluster 2) exhibit above-average management ef- 1050
 1001 ficiency and growth potential. Cluster B+ (Cluster 1051
 1002 3) and Cluster B (Cluster 4) include companies with 1052
 1003 varying financial health and management efficiency. 1053
 1004 Companies in Cluster C+ (Cluster 5) display better 1054
 1005 financial health and higher growth rates, but lower 1055
 1006 dividend payout ratios. Finally, Cluster C (Cluster 0) 1056
 1007 comprises companies facing liquidity challenges and 1057
 1008 lower management efficiency. The suggested credit 1058
 1009 rating system appears to be a logical classification 1059
 1010 based on the distinct characteristics observed in each 1060
 1011 cluster. 1061

1012 **CONCLUSIONS &**
 1013 **RECOMMENDATIONS**

1014 **Conclusions**

1015 In conclusion, this study has made a significant con- 1062
 1016 tribution to the development of a credit rating system 1063
 1017 based on companies' financial and operational char- 1064
 1018 acteristics using the K-means clustering algorithm. 1065
 1019 The research objectives were successfully met, with 1066
 1020 the K-means model effectively clustering the compa- 1067
 1021 nies into six distinct groups, each exhibiting unique 1068
 1022 financial and operational attributes. The author has 1069
 1023 suggested a credit rating system consisting of A+, A, 1070
 1024 B+, B, C+, and C labels, representing varying levels of 1071
 1025 credit risk. 1072

1026 The findings of this study provide valuable insights 1073
 1027 into the financial and operational features that distin- 1074
 1028 guish companies with different credit risk profiles. By 1075
 1029 identifying these characteristics, the proposed credit 1076
 1030 rating system offers a practical tool for assessing credit 1077
 1031 risk, which various stakeholders, including financial 1078
 1032 institutions, credit rating agencies, and investors can 1079
 1033 use. 1080

1034 Furthermore, this research has demonstrated the po- 1081
 1035 tential of clustering techniques, notably the K-means 1082
 1036 algorithm, for addressing complex financial problems 1083
 1037 such as credit risk assessment. The methodology em- 1084
 1038 ployed in this study can serve as a foundation for fu- 1085
 1039 ture research endeavors that aim to improve and re- 1086
 1040 fine credit rating systems. 1087

The practical application of the K-means clustering 1041
 model developed in this study can significantly en- 1042
 hance credit rating processes within various financial 1043
 institutions. Commercial banks can implement this 1044
 model to improve their internal credit scoring sys- 1045
 tems, allowing for more accurate risk management 1046
 and loan pricing strategies by better segmenting cor- 1047
 porate clients based on credit risk. Credit rating 1048
 agencies in Vietnam can utilize this model to sup- 1049
 plement traditional credit rating methods, providing 1050
 a data-driven approach that complements expert as- 1051
 sessments. Additionally, government and regulatory 1052
 bodies, such as the State Bank of Vietnam, can use the 1053
 model to monitor and evaluate the financial health of 1054
 businesses within the economy, facilitating more in- 1055
 formed policymaking. 1056

To ensure the credibility and usability of the model, 1057
 the results should be published and disseminated in 1058
 a transparent manner. This can be achieved through 1059
 periodic reports that detail the credit ratings of com- 1060
 panies segmented by the identified clusters, making 1061
 these reports accessible to investors, financial institu- 1062
 tions, and other stakeholders. Furthermore, develop- 1063
 ing an online platform where stakeholders can access 1064
 real-time credit ratings and updates will provide de- 1065
 tailed insights into rated companies' financial health 1066
 and risk profiles. 1067

Several factors underscore the reliability of the K- 1068
 means clustering model in assessing credit risk. The 1069
 model is grounded in quantitative data, utilizing com- 1070
 prehensive financial indicators to ensure robust credit 1071
 ratings. Using the elbow method and silhouette scores 1072
 to determine the optimal number of clusters enhances 1073
 the model's robustness and validity. Additionally, 1074
 the clustering results align with established financial 1075
 theories, providing empirical support for the model's 1076
 conclusions. To maintain continuous reliability, it is 1077
 essential to periodically update the model with new 1078
 data and refine the input variables based on evolving 1079
 market conditions and financial environments. Reg- 1080
 ular validation against actual financial outcomes will 1081
 enhance the model's accuracy and credibility. 1082

1083 **Recommendations**

Overall, this study's findings contribute to the existing 1084
 body of knowledge on credit risk assessment and of- 1085
 fer a foundation for the development of more accurate 1086
 and reliable credit rating systems. By addressing the 1087
 identified limitations and recommendations, future 1088
 research can continue to advance our understanding 1089
 of credit risk and support improved decision-making 1090
 processes in the financial sector. 1091

1092 For investors, focusing on companies categorized in
 1093 clusters A+ and A, as they demonstrate robust fi-
 1094 nancial health, efficient management, and promising
 1095 growth potential. These companies will likely offer
 1096 higher returns on investment and lower credit risk.
 1097 Additionally, investors should consider diversifying
 1098 their portfolio by including companies from clusters
 1099 B+ and B, as they may present moderate risk and po-
 1100 tential for growth. However, investors should cau-
 1101 tiously approach investments in clusters C+ and C due
 1102 to their relatively weaker financial health and manage-
 1103 ment efficiency.

1104 Managers of companies within clusters B+, B, C+, and
 1105 C should improve their financial health and manage-
 1106 ment efficiency. This may include enhancing liquidity
 1107 management, reducing debt levels, optimizing work-
 1108 ing capital, and implementing cost control measures.
 1109 Furthermore, managers should focus on sustainable
 1110 growth strategies and aim for higher operational effi-
 1111 ciency to increase profitability and competitiveness.

1112 Government agencies can utilize the clustering results
 1113 to understand the financial landscape better and iden-
 1114 tify potential areas of concern. This information can
 1115 be used to develop targeted policies and regulations to
 1116 promote a healthier financial environment for com-
 1117 panies. Additionally, government agencies can sup-
 1118 port and incentivize companies in lower-ranked clus-
 1119 ters to improve their financial stability and promote
 1120 growth. This might include offering tax incentives,
 1121 providing access to low-interest loans, or facilitat-
 1122 ing collaboration between companies and relevant stake-
 1123 holders to foster innovation and technological ad-
 1124 vancements.

1125 For Credit Rating Agencies, adopting the K-means
 1126 clustering algorithm can lead to more accurate and re-
 1127 liable credit ratings. The algorithm's ability to handle
 1128 large datasets efficiently and its robustness in identifi-
 1129 ing distinct credit risk profiles can improve the overall
 1130 quality of credit assessments. Credit Rating Agencies
 1131 can integrate this algorithm into their existing frame-
 1132 works to complement expert evaluations, thereby en-
 1133 hancing the transparency and credibility of their rat-
 1134 ings. Several policies and solutions should be consid-
 1135 ered to help Credit Rating Agencies achieve more ac-
 1136 curate and reliable credit ratings using the K-means
 1137 clustering algorithm. Firstly, Credit Rating Agencies
 1138 should invest in advanced data analytics infrastruc-
 1139 ture to support the implementation of machine learn-
 1140 ing models. This includes acquiring the necessary
 1141 hardware, software, and skilled personnel to manage
 1142 and analyze large datasets. Additionally, staff train-
 1143 ing and development programs should be established

to ensure they are proficient in the latest data anal- 1144
 ysis and machine learning techniques. Financial in- 1145
 stitutions should collaborate with credit rating agen- 1146
 cies to share relevant financial data, enhancing the 1147
 robustness of the clustering models. This collabora- 1148
 tion can be facilitated through standardized data- 1149
 sharing agreements that protect the confidentiality 1150
 and integrity of sensitive information. Moreover, fi- 1151
 nancial institutions should consider integrating these 1152
 advanced credit rating models into their risk man- 1153
 agement and loan pricing strategies to optimize their 1154
 credit assessment processes. 1155

Government and regulatory bodies play a crucial role 1156
 in fostering an environment conducive to adopting 1157
 such advanced technologies. They should establish 1158
 guidelines and regulations that encourage using data- 1159
 driven credit rating methods while ensuring data pri- 1160
 vacy and security. Incentives, such as tax breaks 1161
 or grants, could be provided to CRAs and financial 1162
 institutions that invest in these technologies. Fur- 1163
 thermore, regulatory bodies should promote trans- 1164
 parency and standardization in credit rating practices 1165
 to enhance the comparability and reliability of credit 1166
 ratings across the market. 1167

However, it is important to acknowledge that the pro- 1168
 posed credit rating system may have limitations, and 1169
 further research is needed to ensure its robustness and 1170
 accuracy. Additional validation, incorporation of ex- 1171
 ternal factors, longitudinal analysis, and comparison 1172
 with other methods are recommended to enhance the 1173
 credit rating system's comprehensiveness and predic- 1174
 tive power. While the K-means clustering model pro- 1175
 vides valuable insights, there are certain limitations to 1176
 consider. First, the analysis is based on a set of finan- 1177
 cial ratios, which may not capture all aspects of a com- 1178
 pany's performance. Second, the model is sensitive to 1179
 the initial cluster centroids, which can affect the re- 1180
 sults. Finally, the model relies on historical data, and 1181
 thus may not accurately predict future performance 1182
 or account for external factors such as economic or 1183
 industry changes. 1184

FUNDING 1185

The research is funded by the University of Economics 1186
 and Law, Vietnam National University, Ho Chi Minh 1187
 City, Vietnam. 1188

ABBREVIATIONS 1189

SVM: Support Vector Machine 1190
 LDA: Linear Discriminant Analysis 1191
 LR: Logistic Regression 1192
 HOSE: Ho Chi Minh City Stock Exchange 1193
 HNX: Hanoi Stock Exchange 1194

1195 IPO: Initial Public Offering
 1196 ML: Machine Learning
 1197 ROA: Return on Assets
 1198 ROE: Return on Equity
 1199 EPS: Earnings Per Share
 1200 DPR: Dividend Payout Ratio
 1201 CR: Current Ratio
 1202 QR: Quick Ratio
 1203 DER: Debt to Equity Ratio
 1204 GPR: Gross Profit Ratio
 1205 NPM: Net Profit Margin
 1206 ATO: Asset Turnover Ratio

1207 **CONFLICT OF INTEREST**

1208 The authors declare that they have no conflicts of inter-
 1209 est

1210 **AUTHORS' CONTRIBUTION**

1211 **Tam Phan Huy:** research ideas, data processing, data
 1212 collecting, methodology, results interpreting, conclu-
 1213 sion and implication writing.

1214 **Thuy Chu Quang:** coordinator, data collecting,
 1215 methodology, data visualizing, results interpreting,
 1216 conclusion and implication writing, table and figure
 1217 editing.

1218 **APPENDIXES**

1219 Figures 4 and 5

1220 **REFERENCES**

1221 1. Chung KJ, Chang SL, Yang WD. The optimal cycle time for ex-
 1222 ponentially deteriorating products under trade credit financ-
 1223 ing. *The Engineering Economist*. 2001;46(3):232-42;Available
 1224 from: <https://doi.org/10.1080/00137910108967575>.
 1225 2. Scherr FC. Credit-granting decisions under risk. *The Engineer-*
 1226 *ing Economist*. 1992;37(3):245-62;Available from: <https://doi.org/10.1080/00137919208903072>.
 1227 3. Yilmaz MK, Kucukcolak A. Effects of Basel II standards on small-
 1228 medium size enterprises: evidence from the Istanbul Stock
 1229 Exchange. *Am J Finance Account*. 2009;1(4):408-31;Available
 1230 from: <https://doi.org/10.1504/AJFA.2009.031776>.
 1231 4. Yamanaka S. Credit scoring method using estimated forward
 1232 financial statements based on purchase order information.
 1233 *JSIAM Lett*. 2019;11:33-6;Available from: <https://doi.org/10.14495/jsiaml.11.33>.
 1234 5. Modigliani F, Miller MH. The cost of capital, corpora-
 1235 tion finance and the theory of investment. *Am Econ Rev*.
 1236 1958;48(3):261-97;.
 1237 6. Myers SC, Majluf NS. Corporate financing and investment de-
 1238 cisions when firms have information that investors do not
 1239 have. *J Financ Econ*. 1984;13(2):187-221;Available from: [https://doi.org/10.1016/0304-405X\(84\)90023-0](https://doi.org/10.1016/0304-405X(84)90023-0).
 1240 7. Xu R, Wunsch DC. Clustering algorithms in biomedical
 1241 research: a review. *IEEE Rev Biomed Eng*. 2010;3:120-
 1242 54;Available from: <https://doi.org/10.1109/RBME.2010.2083647>.
 1243 8. Kraus A, Litzenberger RH. A state-preference model of opti-
 1244 mal financial leverage. *J Finance*. 1973;28(4):911-22;Available
 1245 from: <https://doi.org/10.1111/j.1540-6261.1973.tb01415.x>.
 1246 9. Myers SC. Capital structure puzzle. 1984;Available from: <https://doi.org/10.3386/w1393>.
 1247 10. Jain AK, Murty MN, Flynn PJ. Data clustering: a review. *ACM*
 1248 *Comput Surv*. 1999;31(3):264-323;Available from: <https://doi.org/10.1145/331499.331504>.
 1249 11. Gitman LJ, Juchau R, Flanagan J. Principles of managerial fi-
 1250 nance. Pearson Higher Education AU; 2015;.
 1251 12. Murtagh F, Legendre P. Ward's hierarchical agglomerative
 1252 clustering method: which algorithms implement Ward's crite-
 1253 rion? *J Classif*. 2014;31:274-95;Available from: <https://doi.org/10.1007/s00357-014-9161-z>.
 1254 13. MacQueen J. Some methods for classification and analysis of
 1255 multivariate observations. *Proc Fifth Berkeley Symp Math Stat*
 1256 *Probab*. 1967;1(14):281-97;.
 1257 14. Huang Z, Chen H, Hsu CJ, Chen WH, Wu S. Credit rat-
 1258 ing analysis with support vector machines and neural
 1259 networks: a market comparative study. *Decis Support*
 1260 *Syst*. 2004;37(4):543-58;Available from: [https://doi.org/10.1016/S0167-9236\(03\)00086-1](https://doi.org/10.1016/S0167-9236(03)00086-1).
 1261 15. Kim MJ, Kang DK. Ensemble with neural networks for
 1262 bankruptcy prediction. *Expert Syst Appl*. 2010;37(4):3373-
 1263 9;Available from: <https://doi.org/10.1016/j.eswa.2009.10.012>.
 1264 16. Barboza F, Kimura H, Altman E. Machine learning models
 1265 and bankruptcy prediction. *Expert Syst Appl*. 2017;83:405-
 1266 17;Available from: <https://doi.org/10.1016/j.eswa.2017.04.006>.
 1267 17. Cantor R, Packer F. Determinants and impact of sovereign
 1268 credit ratings. *Econ Policy Rev*. 1996;2(2);Available from: <https://doi.org/10.1111/j.1468-036X.1996.tb00040.x>.
 1269 18. Vellido A, Lisboa PJ, Vaughan J. Neural networks
 1270 in business: a survey of applications (1992-1998).
 1271 *Expert Syst Appl*. 1999;17(1):51-70;Available from:
 1272 [https://doi.org/10.1016/S0957-4174\(99\)00016-0](https://doi.org/10.1016/S0957-4174(99)00016-0).
 1273 19. Rousseeuw PJ. Silhouettes: a graphical aid to the interpreta-
 1274 tion and validation of cluster analysis. *J Comput Appl Math*.
 1275 1987;20:53-65;Available from: [https://doi.org/10.1016/0377-0427\(87\)90125-7](https://doi.org/10.1016/0377-0427(87)90125-7).
 1276 20. Kodinariya TM, Makwana PR. Review on determining number
 1277 of clusters in K-means clustering. *Int J*. 2013;1(6):90-5;.
 1278 21. Ketchen DJ, Shook CL. The application of cluster analy-
 1279 sis in strategic management research: analysis and cri-
 1280 tique. *Strateg Manag J*. 1996;17(6):441-58;Available from:
 1281 [https://doi.org/10.1002/\(SICI\)1097-0266\(199606\)17:6<441::AID-SMJ819>3.0.CO;2-G](https://doi.org/10.1002/(SICI)1097-0266(199606)17:6<441::AID-SMJ819>3.0.CO;2-G).
 1282 22. Altman EI, Sabato G. Modelling credit risk for SMEs: Evidence
 1283 from the US market. *Abacus*. 2007;43(3):332-57;Available
 1284 from: <https://doi.org/10.1111/j.1467-6281.2007.00234.x>.
 1285 23. Abdou HA, Pointon J. Credit scoring, statistical techniques
 1286 and evaluation criteria: a review of the literature. *Intell Syst*
 1287 *Account Finance Manag*. 2011;18(2-3):59-88;Available from:
 1288 <https://doi.org/10.1002/issaf.325>.
 1289 24. Galindo J, Tamayo P. Credit risk assessment using statistical
 1290 and machine learning: basic methodology and risk modeling
 1291 applications. *Comput Econ*. 2000;15:107-43;Available from:
 1292 <https://doi.org/10.1023/A:1008699112516>.
 1293 25. Min JH, Lee YC. Bankruptcy prediction using support vector
 1294 machine with optimal choice of kernel function parameters.
 1295 *Expert Syst Appl*. 2005;28(4):603-14;Available from: <https://doi.org/10.1016/j.eswa.2004.12.008>.
 1296 26. Kovalerchuk B, Vityaev E. Data mining for financial applica-
 1297 tions. *Data Min Knowl Discov Handb*. 2005;1203-24;Available
 1298 from: https://doi.org/10.1007/0-387-25465-X_57.
 1299 27. Yu L, Wang S, Lai KK. A novel nonlinear ensemble forecast-
 1300 ing model incorporating GLAR and ANN for foreign exchange
 1301 rates. *Comput Oper Res*. 2005;32(10):2523-41;Available from:
 1302 <https://doi.org/10.1016/j.cor.2004.06.024>.
 1303 28. Oreski S, Oreski G. Genetic algorithm-based heuristic for fea-
 1304 ture selection in credit risk assessment. *Expert Syst Appl*.
 1305 2014;41(4):2052-64;Available from: <https://doi.org/10.1016/j.eswa.2013.09.004>.
 1306 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318

Current Ratio

Class	count	mean	std	min	25%	50%	75%	max
0	213.0	1.249014	0.857532	0.13	0.75	1.05	1.420	5.73
1	208.0	6.958048	5.258556	1.75	3.58	5.36	8.055	29.41
2	623.0	2.768780	2.117755	0.37	1.57	2.21	3.320	17.78
3	369.0	1.235474	0.307692	0.19	1.09	1.22	1.420	2.41
4	463.0	1.149264	0.245689	0.25	1.05	1.15	1.275	2.05
5	691.0	1.799190	0.727213	0.27	1.30	1.67	2.195	5.34

Quick Ratio

Class	count	mean	std	min	25%	50%	75%	max
0	213.0	0.785399	0.533792	0.01	0.4680	0.690	0.98	3.35
1	208.0	3.918048	2.302878	0.50	2.2275	3.215	5.15	10.80
2	623.0	1.378963	0.991135	0.06	0.7380	1.130	1.72	9.35
3	369.0	0.695962	0.322262	0.03	0.4580	0.680	0.88	1.62
4	463.0	0.680130	0.271519	0.08	0.3950	0.590	0.79	1.38
5	691.0	1.056715	0.558862	0.06	0.6780	1.020	1.36	4.81

Short-term Liabilities to Equity

Class	count	mean	std	min	25%	50%	75%	max
0	213.0	1.088480	0.751633	0.055429	0.423898	0.813385	1.413927	3.962566
1	208.0	0.147958	0.086888	0.016252	0.080886	0.124993	0.204092	0.427761
2	623.0	0.429208	0.263983	0.030697	0.224366	0.375488	0.584218	1.595643
3	369.0	1.942869	1.119583	0.668347	1.153385	1.618127	2.474522	5.987988
4	463.0	2.415175	1.108245	0.881542	1.566512	2.177750	2.948786	5.888186
5	691.0	0.753455	0.418643	0.088264	0.444888	0.661971	0.968715	2.613387

Short-term Liabilities to Asset

Class	count	mean	std	min	25%	50%	75%	max
0	213.0	0.254304	0.114913	0.022833	0.152942	0.243483	0.348461	0.473744
1	208.0	0.113156	0.057034	0.015992	0.069363	0.102083	0.156146	0.292874
2	623.0	0.256920	0.128000	0.029783	0.166092	0.258835	0.346322	0.615169
3	369.0	0.571548	0.118188	0.283094	0.487540	0.568901	0.654615	0.851655
4	463.0	0.622120	0.106821	0.383086	0.542269	0.619415	0.708412	0.849648
5	691.0	0.358073	0.116947	0.068836	0.275002	0.354378	0.441235	0.693855

Debt on Equity

Class	count	mean	std	min	25%	50%	75%	max
0	213.0	1.833243	1.264268	0.499222	1.037094	1.508820	2.086185	8.842545
1	208.0	0.052819	0.093467	0.000000	0.000000	0.003652	0.045736	0.611728
2	623.0	0.160101	0.180212	0.000000	0.000000	0.112836	0.253899	1.060719
3	369.0	1.039450	0.724339	0.000000	0.550661	0.869829	1.472228	4.585333
4	463.0	1.431520	0.706615	0.044589	0.929022	1.321126	1.777834	4.216249
5	691.0	0.527490	0.253585	0.000000	0.127207	0.284606	0.489233	1.468681

Debt on Asset

Class	count	mean	std	min	25%	50%	75%	max
0	213.0	0.470990	0.097547	0.280033	0.396494	0.462748	0.532574	0.756367
1	208.0	0.037748	0.062159	0.000000	0.000000	0.003311	0.053012	0.363810
2	623.0	0.092953	0.093425	0.000000	0.000000	0.072108	0.152994	0.465239
3	369.0	0.307641	0.148741	0.000000	0.210302	0.319588	0.401976	0.752839
4	463.0	0.384554	0.136564	0.012194	0.287844	0.386639	0.476051	0.734588
5	691.0	0.158812	0.104418	0.000000	0.073650	0.158558	0.236323	0.411688

Long-term Debt on Equity

Class	count	mean	std	min	25%	50%	75%	max
0	213.0	1.378774	1.089121	0.312612	0.657698	1.115807	1.719013	8.308940
1	208.0	0.024746	0.078434	0.000000	0.000000	0.000000	0.005613	0.502721
2	623.0	0.064562	0.127442	0.000000	0.000000	0.000000	0.077193	0.907584
3	369.0	0.180735	0.284595	0.000000	0.000000	0.048108	0.207556	1.726118
4	463.0	0.220851	0.308953	0.000000	0.006356	0.088834	0.334022	1.947064
5	691.0	0.099784	0.152441	0.000000	0.000000	0.023061	0.147575	1.250674

Sale Growth Rate

Class	count	mean	std	min	25%	50%	75%	max
0	213.0	9.832394	37.351898	-89.50	-9.980	8.86	21.57	199.04
1	208.0	-0.322188	50.075472	-100.00	-25.690	1.76	19.20	263.28
2	623.0	7.047864	30.093048	-103.92	-4.085	5.91	14.75	233.34
3	369.0	9.326775	31.488429	-70.50	-5.940	6.38	22.21	295.72
4	463.0	13.750886	41.013789	-100.00	-7.710	10.51	30.07	355.93
5	691.0	11.198321	44.949804	-94.29	-13.440	6.89	24.31	273.91

EBIT Growth Rate

Class	count	mean	std	min	25%	50%	75%	max
0	213.0	3.633366	88.108567	-467.45	-19.9100	7.390	34.830	461.12
1	208.0	-18.689837	120.598959	-469.95	-73.1975	-14.955	25.895	452.34
2	623.0	10.914039	49.254717	-75.96	-11.3600	3.410	24.810	499.31
3	369.0	11.974946	49.813616	-90.30	-12.6200	5.130	23.980	409.64
4	463.0	7.386981	93.228668	-509.60	-23.4750	2.260	37.320	493.23
5	691.0	8.294373	102.852295	-503.73	-35.8700	1.080	34.950	499.78

Dividend Payout Ratio

Class	count	mean	std	min	25%	50%	75%	max
0	213.0	21.095587	28.631616	0.00	0.000	0.00	41.18	98.22
1	208.0	3.743125	10.739039	0.00	0.000	0.00	0.00	53.06
2	623.0	65.336950	17.747126	79.83	50.935	64.76	78.92	99.95
3	369.0	67.929431	16.023563	34.22	54.500	65.54	81.70	98.65
4	463.0	4.308294	10.143728	0.00	0.000	0.00	0.00	47.96
5	691.0	4.762171	10.453051	0.00	0.000	0.00	0.00	39.21

ROA

Class	count	mean	std	min	25%	50%	75%	max
0	213.0	3.169953	3.844279	-10.76	1.1600	2.910	4.9600	23.83
1	208.0	5.362644	8.905155	-36.97	0.8675	4.095	8.4175	42.58
2	623.0	10.521621	7.270686	1.03	5.7850	8.490	13.1500	54.15
3	369.0	4.818650	3.905776	0.20	2.3200	4.100	6.5200	19.31
4	463.0	2.744881	4.047027	-11.91	0.4650	2.060	4.4850	21.41
5	691.0	5.369190	7.389523	-51.72	1.1900	4.120	8.2800	38.88

Account Receivable Turnover

Class	count	mean	std	min	25%	50%	75%	max
0	213.0	14.907887	26.277373	0.30	3.640	6.400	13.1300	172.36
1	208.0	13.059418	31.652645	0.09	2.245	5.325	9.1175	271.35
2	623.0	30.899530	62.258880	0.00	5.780	10.600	21.0450	317.17
3	369.0	9.213461	10.687214	0.35	2.810	5.390	10.5600	65.81
4	463.0	11.661102	27.644446	0.26	2.905	5.730	9.7900	317.17
5	691.0	11.661102	28.895539	0.01	2.720	5.100	9.5450	317.17

Payment Period Turnover

Class	count	mean	std	min	25%	50%	75%	max
0	213.0	8.827624	21.698545	0.30	3.380	5.610	8.9400	305.39
1	208.0	30.644279	62.679121	0.18	6.835	11.785	27.7975	305.39
2	623.0	23.950842	36.679793	0.00	7.680	12.910	23.9900	305.39
3	369.0	11.539946	12.649028	0.00	4.590	8.270	13.9000	92.27
4	463.0	13.267970	19.632760	0.00	4.040	7.820	13.6000	180.51
5	691.0	12.422208	16.710830	0.02	4.625	7.320	12.9500	151.61

Figure 4: Descriptive Statistics Of Clusters By Variables

	Financial health factors	Management efficiency	Growth potential	Dividend payout capacity
Cluster 1	Outstanding liquidity, more than twice the average level. Short-term payables over equity are about 15% of the average. Total debt over total assets or total debt over total equity are both very low, about 15% of the average. Long-term and short-term debt ratios are approximately equal. More assets than equity.	ROA is approximately around 75% of the average, the number of collections per year is approximately the average turnover, but the number of disbursements is double.	Relatively low revenue growth rate, and low pre-tax profit and interest rate growth rate, about 3 times the average.	Dividend payout ratio is about 6-7 times lower than the average
Cluster 2	Short-term liquidity is approximately average, not dependent on inventory, short-term payables over equity is about 42%, about 50% of the average. Total debt over equity is about 20% and over assets about 37% compared to the average. Short-term debt is higher than long-term debt, but the difference is not significant. Assets are twice the equity.	Outstanding ROA, twice the average, high accounts receivable turnover, double the average, and payment ability is 2/3 of the average.	Positive revenue growth rate, below average but not too much, EBIT growth rate is 20% higher than the average	Dividend payout ratio is about 2 times lower than the average
Cluster 3	Liquidity is available but less than half of the average, loss of short-term liquidity if inventory value is excluded. Short-term payables over equity are approximately twice the average. Total debt over equity or assets are both about 1.25 times the average. Assets are three times the equity.	ROA is around 70% of the average. The cash collection cycle is lower than the payment cycle, and both are below 75% of the average.	Revenue growth rate is at an average level, but EBIT growth rate is about 1.25 to 1.5 times the average.	Dividend payout ratio is about 2 times lower than the average.
Cluster 4	Short-term liquidity is available but very low, more than half lower than the average, dependent on inventory value. Short-term payables over equity are approximately 2 to 2.4 times the average. Total debt over	ROA is low, ranging from 50% to 75% of the average. The number of collections and disbursements is about once a month, at an average level.	Revenue growth rate is 1.6 times the average, but EBIT growth rate is about 50% lower than the average.	Dividend payout ratio is about 6-7 times lower than the average.
Cluster 5	equity or assets is about 1.3 times the average. Assets are three times the equity. Short-term liquidity is available and 50% higher than the average, not dependent on inventory value. Short-term payables over equity are approximately 75%. Total debt over equity is below 50% of the average, with short-term debt being three times the long-term debt. Assets are twice the equity.	ROA is approximately average. The number of collections and disbursements is about once a month, at an average level.	Revenue growth rate is 1.5 times the average, but EBIT growth rate is about 50% lower than the average.	Dividend payout ratio is about 6-7 times lower than the average.
Cluster 0	Short-term liquidity is available but less than 50% of the average, losing liquidity when inventory value is excluded. Short-term payables are approximately equal to equity, equivalent to the average. Total debt over equity or total assets is more than twice the average, with most of the debt being long-term.	ROA is 80% lower than the average. The number of collections for accounts receivable is approximately twice that of short-term payables, about 50% of the average.	Revenue growth rate is equivalent to the average, but EBIT growth rate is about 50% lower than the average.	Dividend payout ratio is about 50% lower than the average.

Figure 5: Descriptive By Clusters

Xếp hạng tín dụng bằng thuật toán phân cụm tại thị trường Chứng khoán Việt Nam

Phan Huy Tâm^{1,2,*}, Chu Quang Thuý^{1,2}



Use your smartphone to scan this QR code and download this article

¹Trường Đại học Kinh tế - Luật, Tp. Hồ Chí Minh, Việt Nam

²Đại học Quốc gia Tp. Hồ Chí Minh, Tp. Hồ Chí Minh, Việt Nam.

Liên hệ

Phan Huy Tâm, Trường Đại học Kinh tế - Luật, Tp. Hồ Chí Minh, Việt Nam

Đại học Quốc gia Tp. Hồ Chí Minh, Tp. Hồ Chí Minh, Việt Nam.

Email: tamphan.ntc@gmail.com

Lịch sử

- Ngày nhận: 17-5-2024
- Ngày sửa đổi: 23-7-2024
- Ngày chấp nhận: 27-9-2024
- Ngày đăng:

DOI:



Bản quyền

© ĐHQG Tp.HCM. Đây là bài báo công bố mở được phát hành theo các điều khoản của the Creative Commons Attribution 4.0 International license.



TÓM TẮT

Nghiên cứu này áp dụng thuật toán phân cụm K-means để phát triển khung xếp hạng tín dụng doanh nghiệp cho thị trường Việt Nam. Bằng cách phân tích dữ liệu tài chính từ 568 công ty phi tài chính niêm yết tại thị trường Chứng khoán Thành phố Hồ Chí Minh (HOSE) và thị trường Giao dịch Chứng khoán Hà Nội (HNX) trong giai đoạn từ 2019 đến 2023, nghiên cứu xác định các chỉ số tài chính quan trọng bao gồm tỷ lệ sức khỏe tài chính, tỷ lệ hiệu quả quản lý, tỷ lệ tăng trưởng và tỷ lệ chi trả cổ tức. Mô hình phân cụm K-means cho thấy tính hiệu quả trong phân loại các doanh nghiệp này thành sáu cụm khác nhau, mỗi cụm đại diện cho các mức độ hiệu suất tài chính và rủi ro tín dụng khác nhau. Các cụm này được xếp từ A+ (rủi ro tín dụng rất thấp) đến C (rủi ro tín dụng rất cao), cung cấp sự phân biệt rõ ràng dựa trên sự ổn định tài chính và hiệu quả hoạt động. Cách tiếp cận hệ thống này mang lại những hiểu biết có giá trị cho các nhà đầu tư, nhà quản lý và các cơ quan chính phủ, nâng cao khả năng đưa ra quyết định thông minh. Mặc dù có một số hạn chế như phụ thuộc vào dữ liệu lịch sử và độ nhạy cảm đối với các tâm cụm ban đầu, mô hình phân cụm K-means chứng minh là một điểm khởi đầu mạnh mẽ để đánh giá độ tín nhiệm của các công ty. Nghiên cứu này đóng góp vào tài liệu ngày càng tăng về các ứng dụng học máy trong xếp hạng tín dụng bằng cách chứng minh sự vượt trội của các thuật toán phân cụm so với các phương pháp truyền thống. Nghiên cứu nêu bật cách các chỉ số sức khỏe tài chính và hiệu quả quản lý có thể được tích hợp vào một khung dữ liệu để nâng cao đánh giá rủi ro tín dụng. Kết quả gợi ý rằng cách tiếp cận phân cụm K-means không chỉ cải thiện độ chính xác của xếp hạng tín dụng mà còn thúc đẩy tính minh bạch và hiệu quả trong thị trường tài chính. Hơn nữa, khung đề xuất có thể đóng vai trò là nền tảng để phát triển các mô hình phức tạp hơn, tích hợp thêm các biến tài chính và phi tài chính. Nghiên cứu trong tương lai có thể mở rộng điều này bằng cách tích hợp dữ liệu theo thời gian thực và khám phá tác động của các yếu tố kinh tế bên ngoài đối với rủi ro tín dụng. Bằng cách tận dụng các kỹ thuật học máy tiên tiến, nghiên cứu này mở đường cho các hệ thống xếp hạng tín dụng đáng tin cậy và toàn diện hơn, hỗ trợ sự ổn định và phát triển của các thị trường tài chính tại các nền kinh tế đang nổi như Việt Nam.

Từ khóa: K-Means, Xếp hạng tín dụng, Phân cụm, Việt Nam

Trích dẫn bài báo này: Tâm P.H, Thuý C.Q. Xếp hạng tín dụng bằng thuật toán phân cụm tại thị trường Chứng khoán Việt Nam. *Sci. Tech. Dev. J. - Eco. Law Manag.* 2024; ():1-1.