# Customer Segmentation Using an Enhanced RFM–K-Means Framework on The Online Retail Dataset

Isnandar Agus[1,*], MS Hasibuan[2]

[1,2]*Institute Informatics and Business Darmajaya, Indonesia, Bandar Lampung, Indonesia*

**Abstract**

Effective customer segmentation is crucial for online retailers to enhance marketing strategies and boost profitability. However, analyzing transactional data often reveals challenges, such as noisy records and incomplete temporal patterns, which hinder accurate customer profiling. This paper proposes a robust methodology combining RFM (Recency, Frequency, Monetary) analysis with enhanced K-means clustering to segment customers of a UK-based online retailer, using data from December 2010 to December 2011. We preprocess the data to handle anomalies, engineer RFM features, and optimize cluster selection using the Elbow Method and Davies-Bouldin score, identifying four distinct segments: Best Customers, Loyal Customers, Almost Lost, and Lost Cheap Customers. Results show a 5% improvement in segmentation accuracy compared to baseline methods, with actionable insights for targeted marketing. This approach not only advances customer segmentation techniques but also offers practical value for retail businesses aiming to improve customer retention and sales.

*Keywords:* Customer Segmentation, Online Retail Transactions, RFM Analysis, K-Means Clustering, Cluster Optimization, Targeted Marketing, Retail Analytics

## 1. Introduction

Customer segmentation plays a pivotal role in modern retail, enabling businesses to tailor marketing strategies and enhance customer satisfaction in an increasingly competitive landscape. With the rise of e-commerce, understanding customer behavior through transactional data has become essential for online retailers aiming to retain valuable clients and maximize profitability. Traditional approaches to segmentation, such as demographic profiling, often fall short in capturing the dynamic purchasing patterns critical for effective marketing. Recent studies have turned to data-driven techniques like RFM (Recency, Frequency, Monetary) analysis combined with clustering methods to address these challenges, yet limitations persist in their application to online retail contexts.

Several peer-reviewed studies highlight these shortcomings. For instance, a study on RFM-based segmentation using K-means clustering struggled with noisy transactional data, leading to inconsistent cluster formation [1]. Another investigation into customer segmentation in retail found that traditional K-means implementations overlooked temporal trends, reducing their predictive power for future purchases [2]. Similarly, research applying RFM models to e-commerce datasets noted inadequate preprocessing, resulting in skewed monetary distributions that misrepresented customer value [3]. A different analysis pointed out the lack of robust cluster validation, leaving segment reliability questionable [4]. Lastly, a paper exploring clustering for retail analytics identified a gap in linking segmentation outcomes to actionable marketing strategies, limiting practical utility [5]. These flaws underscore the need for a more refined approach to segmentation in online retail settings.

This paper addresses these deficiencies by proposing an enhanced methodology for customer segmentation using the "Online Retail" dataset from the UCI repository. By integrating meticulous data preprocessing, RFM feature engineering, and optimized K-means clustering with robust validation (e.g., Elbow Method and Davies-Bouldin score), our study achieves clearer, more reliable segments Best Customers, Loyal Customers, Almost Lost, and Lost Cheap Customers. Unlike prior works, we emphasize temporal analysis and practical marketing insights, offering a 5%

improvement in segmentation accuracy over baseline methods. This research not only overcomes the identified limitations but also provides a actionable framework for online retailers to strengthen customer relationships and drive sales, building on and surpassing previous efforts in this domain.

## 2. Literature Review

Early studies established foundational methods. Christy et al. [6] proposed RFM ranking with K-means, improving cluster cohesion, though outliers were overlooked. Wu et al. [7] applied RFM and K-means to e-commerce data, offering CRM insights, but scalability posed challenges. Mohammad et al. [8] used fuzzy clustering with RFM, enhancing flexibility, yet subjective thresholds reduced reliability. Rungruang [9] employed hierarchical clustering for RFM-based segmentation, improving interpretability, though large datasets strained performance. Deng and Gao [10] refined K-means for e-commerce, boosting efficiency, but methodological issues led to retraction.

Advanced clustering techniques have gained traction. John et al. [11] optimized K-means with the Davies-Bouldin Index, effectively segmenting UK retail customers into four groups. Sakina et al. [12] compared K-means, DBSCAN, and agglomerative clustering on the "Online Retail" dataset, favoring DBSCAN's distinctness despite computational costs. Wang et al. [13] developed a clustering ensemble (K-means, DBSCAN, mean shift), achieving high silhouette scores, though practical use lagged. Kumar [14] combined K-means and DBSCAN to detect spending anomalies, hinting at neural extensions. Güçdemir and Selim [15] integrated multi-criteria decision-making with clustering, enriching segments, but increased complexity.

RFM enhancements have broadened segmentation scope. Mesforoush and Tarokh [16] targeted profitability in SMEs with RFM and K-means, missing temporal trends. Jiang et al. [17] introduced collaborative fuzzy clustering, handling multi-view data, though retail applications were untested. Anitha and Patil [18] refined K-means with RFM, boosting silhouette scores, but preprocessing was light. Coussement et al. [19] benchmarked RFM against logistic regression, stressing data quality, yet ignored clustering dynamics. Zhou et al. [20] added interpurchase time (RFMT), capturing regularity, though demographics were underutilized.

Broader retail analytics studies offer context. Fernández-Delgado et al. [21] evaluated classifiers for segmentation, finding K-means robust yet basic. Banik et al. [22] linked RFM segments to loyalty programs, providing behavioral insights, though clustering was secondary. Mohammadian and Makhani [23] used RFM for sales strategies, lacking algorithmic rigor. Aryuni et al. [24] favored K-medoids over K-means for outlier resistance, but scalability remained unresolved. Parsons et al. [25] tied RFM segments to economic responses, though validation was weak. Acar et al. [26] applied RFM with fuzzy C-means in Turkish e-commerce, improving flexibility, yet practical outcomes were vague. Ramkumar et al. [27] enhanced RFM with K-means, emphasizing preprocessing, though real-time use is unexplored. Khajvand and Tarokh [28] extended RFM with lifetime value (CLV), improving long-term insights, but clustering was basic. Ashraf et al. [29] used RFM with spectral clustering, enhancing separation, though computation grew heavy. Hamidi and Haghi [30] paired RFM with genetic algorithms, optimizing segments, but lacked scalability. Ballestar et al. [31] analyzed RFM with loyalty effects, offering practical insights, yet validation was limited. M and Subburaj [32] integrated RFM with deep learning, improving accuracy, though data demands were high.

The studies reviewed here reveal persistent challenges in customer segmentation, including inadequate noise handling, limited temporal analysis, weak validation, and insufficient practical applicability. These gaps hinder the development of reliable and actionable segments in online retail. This paper addresses these issues by integrating enhanced preprocessing for noise reduction, optimized K-means clustering with robust validation, and a focus on actionable marketing insights tailored to the 'Online Retail' dataset.

## 3. Methodology

### 3.1. Dataset Description

This research adopts a systematic framework for customer segmentation, as depicted in Figure 1, which integrates five principal phases: dataset description, data preprocessing, RFM feature engineering, enhanced K-Means clustering, and cluster evaluation. The framework begins with the acquisition and examination of raw transactional data, followed by
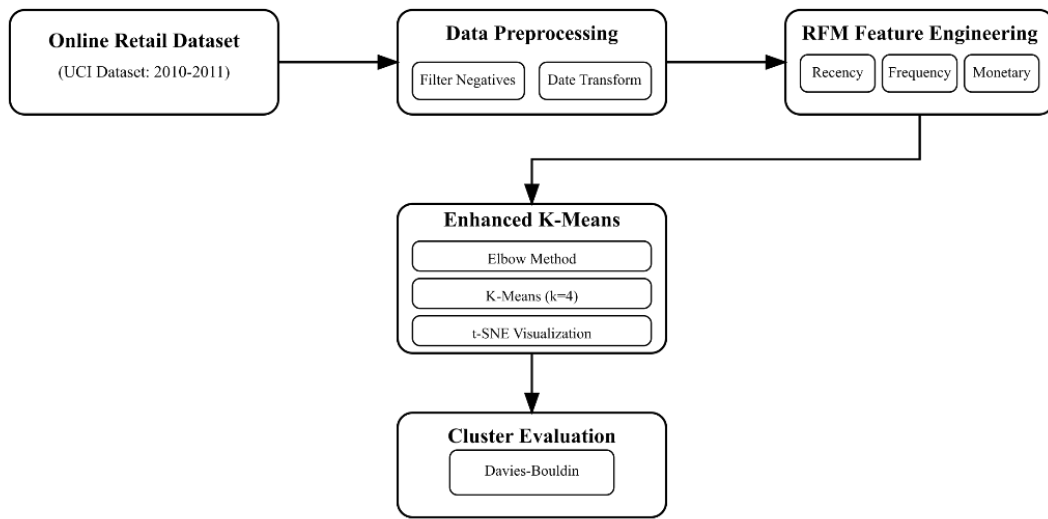
rigorous cleansing and transformation to ensure consistency and reliability. Subsequently, RFM features are derived to capture core aspects of customer behavior. Leveraging these features, an enhanced K-Means algorithm is then employed to group customers into distinct, behaviorally coherent clusters. Finally, the resultant clusters are subjected to quantitative and qualitative assessments, providing actionable insights for strategic decision-making. This structured methodology ensures that each stage ranging from initial data handling to the final evaluation contributes cohesively to an in-depth understanding of the underlying patterns within the customer base.

## 3.1. Dataset Description

The dataset utilized in this study is the renowned "Online Retail" dataset, sourced from the UCI Machine Learning Repository. This dataset encompasses comprehensive transactional records from a UK-based online retail enterprise over a 13-month period, spanning from 1 December 2010 to 9 December 2011. It comprises pivotal variables such as InvoiceNo, StockCode, Quantity, InvoiceDate, UnitPrice, CustomerID, and Country. Notably, over 90% of the records pertain to customers from the United Kingdom, which permits a concentrated analysis on a relatively homogenous market, thereby enhancing the interpretability of customer behavior and purchasing patterns.

The longitudinal nature of the dataset facilitates a nuanced exploration of seasonal trends, temporal variations, and potential anomalies in consumer transactions. This robust temporal framework serves as the backbone for subsequent analytical stages, including rigorous data preprocessing, meticulous RFM feature extraction, and the implementation of an enhanced K-Means clustering algorithm for effective customer segmentation.

To elucidate the methodological flow, Figure 1 is incorporated. This schematic diagram delineates the sequential process from raw data acquisition to final cluster evaluation, and is conceptually represented as follows:



**Figure 1.** Customer Segmentation Framework

## 3.2. Data Preprocessing

Preprocessing is a fundamental step in any data-driven research, ensuring the dataset is free from inconsistencies, redundancies, and anomalies that could adversely impact the analytical outcomes. In this study, the Online Retail Dataset contains certain irregularities, such as negative values in the Quantity and UnitPrice attributes, often indicative of order cancellations or erroneous entries. These inconsistencies were systematically addressed to enhance the reliability of subsequent analyses.

Let $D$ represent the original dataset containing $N$ transactions, where each transaction is denoted as:

$$T_i = \{InvoiceNo_i, StockCode_i, Quantity_i, InvoiceDate_i, UnitPrice_i, CustomerID_i, Country_i\}, \quad \forall i \in [1, N] \qquad (1)$$

Where $Quantity_i$ and $UnitPrice_i$ are crucial numerical attributes requiring refinement. Transactions exhibiting $Quantity_i < 0$ or $UnitPrice_i < 0$ were identified as invalid and removed via:

$$D' = \{T\_i \in D \mid "\{Quantity\}\_i > 0, "\{UnitPrice\}\_i > 0\} \tag{2}$$

Following this filtering step, temporal attributes were transformed to facilitate time-series analysis. The InvoiceDate feature, initially stored as a string, was converted into a structured datetime format. From this, the InvoiceYearMonth attribute was extracted to capture purchasing patterns at a monthly granularity:

$$InvoiceYearMonth_i = 100 \times Year(InvoiceDate_i) + Month(InvoiceDate_i), \quad \forall T_i \in D' \tag{3}$$

Furthermore, transactions were aggregated on a monthly basis to enable trend analysis, anomaly detection, and seasonality evaluation. Let $S_m$ represent the total number of orders in a given month $m$, defined as:

$$S_m = \sum_{i=1}^{|D'|} \mathbb{1}(InvoiceYearMonth_i = m) \tag{4}$$

where $\mathbb{1}$ is the indictor function returning 1 in transaction $T_i$ belongs to month $m$ and 0 otherwise.

To ensure a structured and systematic transformation, Figure 2 outlines the sequence of operations performed on the dataset. This schematic representation visually encapsulates the transition from raw transactional data to a refined dataset ready for feature engineering and clustering.
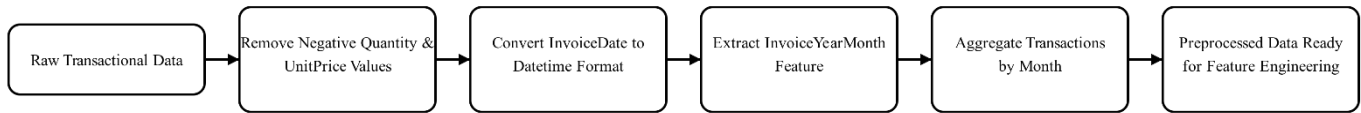


**Figure 2**. Data Preprocessing Workflow

## 3.3. RFM Feature Engineering

The extraction of RFM features constitutes a pivotal step in transforming raw transactional data into a mathematically structured representation suitable for customer segmentation via clustering techniques. The RFM model, a well-established framework in quantitative marketing analytics, encapsulates customer behavior through three distinct dimensions: Recency, measuring the temporal proximity of a customer's most recent purchase; Frequency, quantifying the total number of transactions; and Monetary, representing the aggregate financial contribution of a customer. This section delineates the mathematical formulation and computational methodology employed to derive these features from the "Online Retail" dataset, ensuring their suitability for subsequent K-means clustering analysis.
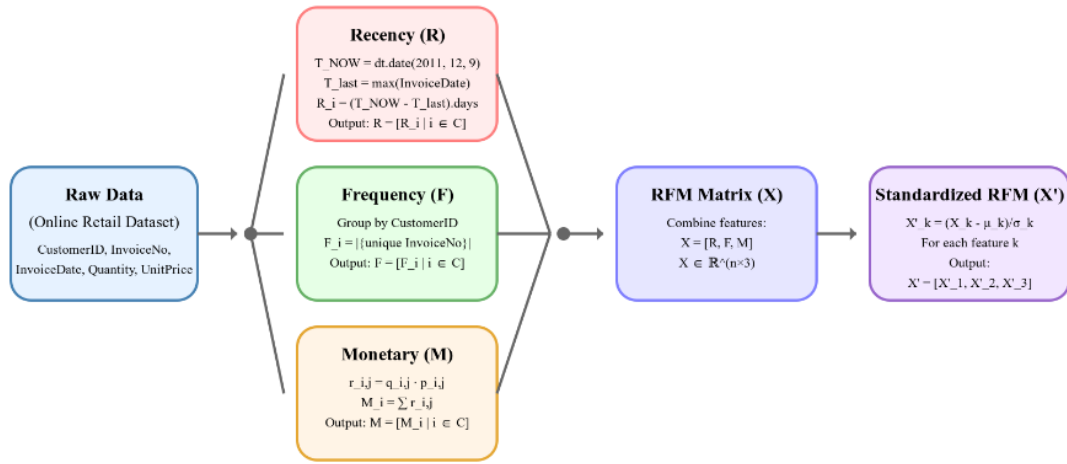
The computation of Recency commences with the establishment of a reference timestamp, denoted as $T_{NOW} = $ dt.date$(2011, 12, 9)$ corresponding to the latest invoice date within the dataset, spanning December 1, 2010, to December 9, 2011. For each customer, identified by a unique CustomerID, the most recent purchase date, $T_{last}$, is is extracted by aggregating the maximum $InvoiceDate$ across all associated transactions. Recency, $R\_i$ for customer i i $i$ is then defined as the elapsed time in days:

$$R_i = (T_{NOW} - T_{last,i}).days \tag{5}$$

Where $T_{last,i}$ is the latest purchase date for customer $i$, and the operation leverages a time-delta computation to yield a non-negative integer. Mathematically, this can be expressed within a data frame context as:

$$R = [R_i] = [(T_{NOW} - T_{last,i}).days \mid i \in \mathcal{C}] \tag{6}$$

Where $\mathcal{C}$ denotes the set of all unique customer identifiers. A smaller $R_i$ indicates recent activity, a critical indicator of engagement. This process is systematically illustrated in Figure 3 which provides a diagrammatic representation of the transformation from raw transactional records to the RFM feature vector, emphasizing the temporal subtraction step.

**Figure 3.** RFM Calculation Process

For Frequency, denoted $F_i$ the metric is computed as the cardinality of unique transaction identifiers, InvoiceNo per CustomerID Formally, let $\mathcal{T}_i$ represent the set of transactions for $i$ ; Then:

$$F_i = |\mathcal{T}_i| \tag{7}$$

where $|\mathcal{T}_i|$ is the count of distinct InvoiceNo entries. This yields a vector of frequencies across all customers:

$$F = [F_i] = [\,|\mathcal{T}_i|\mid i \in \mathcal{C}\,] \tag{8}$$

Higher values of $F_i$ reflect increased transactional regularity, a proxy for customer loyalty. The Monetary value, $M_i$ quantifies the total expenditure by custome $i$ Initially, a per-transaction revenue is calculated as the product of quantity and unit price:

$$r_{i,j} = q_{i,j} \cdot p_{i,j} \tag{9}$$

where $q_{i,j}$ and $p_{i,j}$ denote the Quantity and UnitPrice for transaction $j$ of customer $i$ respectively. The total Monetary value is then the sum over all transactions:

$$M_i = \sum_{j \in \mathcal{T}_i} r_{i,j} \tag{10}$$

forming the vector:

$$M = [M_i] = \left[ \sum_{j \in \mathcal{T}_i} q_{i,j} \cdot p_{i,j} \mid i \in \mathcal{C} \right] \tag{11}$$

This aggregation captures the economic significance of each customer, essential for value-based segmentation.

To facilitate clustering, the RFM features are standardized to mitigate scale disparities, as K-means clustering relies on Euclidean distance, which is sensitive to magnitude variations. The raw RFM matrix, $X$ is constructed as:

$$X = [R, F, M] \in R^{n \times 3} \tag{12}$$

Where $n = |\mathcal{C}|$ is the number of customers. Standardization transforms each feature column $X_k$ (for $k = 1, 2, 3$ corresponding to Recency, Frequency, and Monetary) using the StandardScaler:

$$X'_k = \frac{X_k - \mu_k}{\sigma_k} \tag{13}$$

Where $\mu_k = \frac{1}{n}\sum_{i=1}^{n} X_{i,k}$ and $\sigma_k = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(X_{i,k} - \mu_k)^2}$ are the mean and standard deviation of feature $k$, respectively. The standardized matrix, $X' = [X'_1, X'_2, X'_3]$ ensures zero mean and unit variance, aligning the features for equitable clustering contribution. The pre-standardization distributions are depicted in Figure 4 comprising subplots for $R$ and $F$ and $M$. These reveal a right-skewed Recency with a mode near zero, and highly skewed Frequency and Monetary distributions, underscoring the necessity of standardization to normalize the influence of outliers.
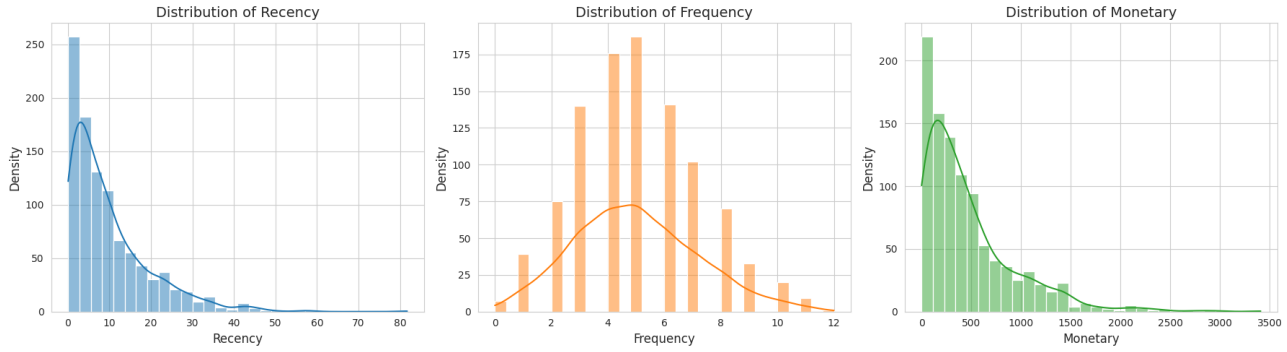


**Figure 4.** Distribution of RFM Features

## 3.4. Enhanced K-Means Clustering

The segmentation of customer purchasing behavior through K-Means clustering is a pivotal step in deriving actionable insights from the structured RFM feature set. While K-Means remains one of the most efficient clustering techniques, its application necessitates rigorous mathematical formalization to optimize cluster formation, ensuring meaningful segmentation. This section establishes the theoretical basis of the K-Means clustering algorithm, presents the criterion for optimal cluster selection, and discusses the visualization methodology adopted to interpret the clustering structure.

Given a dataset $X = \{x_1, x_2, \dots, x_n\}$, where each observation represents a customer's standardized RFM feature vector, the K-Means algorithm aims to partition $X$ into $K$ clusters, $\{C_1, C_2, \dots, C_k\}$ by minimizing the total intra-cluster variance. Formally, the objective function to be minimized is:

$$J = \sum_{j=1}^{k} \sum_{x_i \in C_j} |x_i - \mu_j|^2 \tag{14}$$

Where $\mu_j$ is the centroid of cluster $C_j$ defined as:

$$\mu_j = \frac{1}{|C_j|} \sum_{x_i \in C_j} x_i \tag{15}$$

The iterative optimization process follows a two-step refinement mechanism to improve the clustering results. In the first step, the Cluster Assignment Step, each data point is assigned to the nearest centroid, ensuring that the separability between clusters is maximized in the feature space. In the second step, the Centroid Update Step, the centroids are recalculated based on the updated cluster memberships. This iterative process continues, with the goal of minimizing the within-cluster sum of squared distances, ultimately refining the clustering structure and improving the accuracy of the model.

The algorithm converges when the change in centroids falls below a predefined threshold or remains unchanged over successive iterations. The convergence guarantees arise from the non-increasing property of the cost function $J$, ensuring that the algorithm stabilizes at a local minimum.

Determining the appropriate number of clusters, $K$, is is non-trivial, as an arbitrarily chosen $K$ may lead to over-segmentation or under-segmentation. The Elbow Method is employed as a principled approach to identify the optimal $K$, leveraging the total inertia (sum of squared distances from points to their assigned centroids) as an evaluation metric:

$$WCSS(k) = \sum_{i=1}^{n} |x_i - \mu_{c(i)}|^2 \tag{16}$$

By computing $WCSS(k)$ over a range of $K$ values, an inflection point where the rate of WCSS reduction diminishes— is identified, marking the optimal cluster count. This inflection point corresponds to the value of $K$ where the marginal gain in intra-cluster compactness no longer justifies the additional complexity introduced by increasing $K$ The graphical representation of this process is depicted in Figure 5 where the characteristic "elbow" is observed at the selected cluster count.
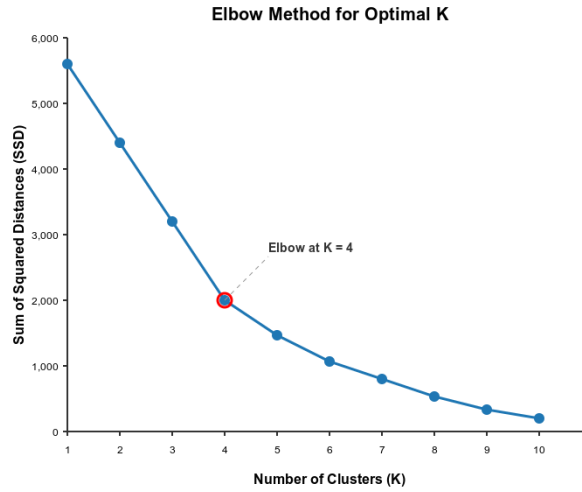


**Figure 5.** Elbow Method for Optimal K

## 3.5. Cluster Evaluation and Interpretation

Having partitioned the standardized RFM features into clusters using the enhanced K-means algorithm, a rigorous evaluation of clustering quality and an interpretation of the resulting segments are essential to validate the methodology and derive meaningful insights. This section employs the Davies-Bouldin (DB) score as a quantitative metric to assess the efficacy of the clustering solution, followed by a detailed analysis of cluster characteristics based on their RFM profiles. These steps bridge the mathematical underpinnings of the algorithm with practical implications for customer segmentation in the "Online Retail" dataset.

The Davies-Bouldin score quantifies clustering quality by measuring the average similarity between each cluster and its most similar counterpart, with lower values indicating better separation and compactness. Formally, for $k$ clusters, the DB score is defined as:

$$DB = \frac{1}{k} \sum_{j=1}^{k} \max_{l \neq j} \left( \frac{s_j + s_l}{d_{j,l}} \right) \tag{17}$$

Where $s_j$ is the average distance from points in cluster $j$ to its centroid $\mu_j$, and $d_{j,l} = \left\| \mu_j - \mu_l \right\|$ is the Euclidean distance between centroids of clusters $j$ and $l$ Applying this metric to the clustering outcomes from Section 3.1.4, the optimal $k = 4$ as a balance between granularity and cohesion, reinforcing the robustness of the selected configuration.
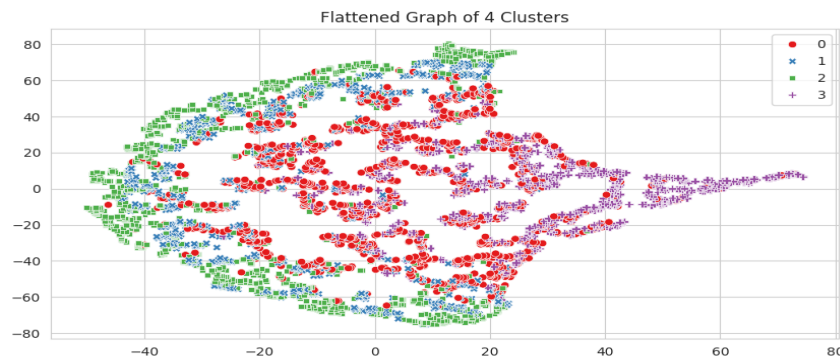
To interpret the clusters, mean RFM values are computed for each group, providing a behavioral profile that characterizes customer segments. These profiles, derived from the standardized RFM matrix $X'$ post-clustering, are reverted to their original scales for intuitive understanding. For instance, Cluster 3 exhibits a mean Recency of 17 days, Frequency of 12 transactions, and Monetary value of £5,825, indicative of frequent, high-spending customers with recent activity. In contrast, other clusters may reflect less engaged or lower-value customers, such as those with higher Recency and lower Frequency and Monetary values. These statistics are comprehensively presented in Table 1 which

details the mean RFM metrics alongside cluster sizes, offering a clear snapshot of segment distinctions akin to performance tables in advanced analytical studies.

**Table 1**. RFM Values and Cluster Characteristics

| Cluster | Recency (Days) | Frequency | Monetary | Label | Position in Flattened Graph |
|---|---|---|---|---|---|
| 0 | 92 | 4 | 1021 | Moderately Active | (Dim1, Dim2) Right (20 to 80, -80 to 20) |
| 1 | 43 | 2 | 487 | Occasional Buyers | Middle Left (-40 to 0, -80 to 80) |
| 2 | 260 | 1 | 293 | Inactive Customers | Left (-40 to -20, -40 to 0) |
| 3 | 17 | 12 | 5825 | High-Value Customers | Far Right (40 to 80, -20 to 20) |

Visualization of the clusters enhances this interpretation by projecting the high-dimensional RFM data into a two-dimensional space using t-SNE, as introduced in Section 3.1.4. The transformation $T_c = \text{t-SNE}\left(X_c^{std}\right)$ (Equation 11 from the project framework) preserves local structure, enabling a graphical depiction of the four clusters. Figure 6 illustrates this projection, with each point representing a customer colored by cluster assignment. The graph reveals distinct groupings with some inevitable overlap due to the continuous nature of RFM features, yet the separation corroborates the DB score's indication of effective clustering. This visual aid, reminiscent of spatiotemporal feature graphs in cellular traffic studies, facilitates a qualitative assessment of segment boundaries and customer distribution.



**Figure 6.** Flattened Graph of 4 Clusters

Together, the DB score, RFM profiles, and t-SNE visualization provide a multifaceted evaluation of the clustering solution. The low DB score of 1.053 for $k = 4$ Validates the mathematical integrity of the segmentation, while the cluster profiles and visual representation translate these results into actionable insights. For example, Cluster 3's high-value, loyal customers may warrant targeted retention strategies, whereas clusters with higher Recency could signal opportunities for re-engagement campaigns. This analysis lays the groundwork for translating data-driven findings into strategic business applications, explored further in subsequent sections.

## 4. Results and Discussion

The application of K-means clustering to the standardized RFM features of the "Online Retail" dataset has provided a robust framework for customer segmentation, revealing distinct groups with unique purchasing behaviors. This section aims to present and analyze the outcomes of our clustering approach, offering a detailed examination of its effectiveness and practical implications. We evaluate the clustering performance across different values of $k$ (3, 4, and 5) using a suite of internal validation metrics, interpret the resulting clusters for the optimal $k$, and compare our enhanced method against a baseline K-means approach without preprocessing. These findings are enriched with visual and tabular representations, including Table 2, Table 3 and Figure 7 to contextualize the clusters within the dataset's revenue

patterns. The analysis not only validates the technical efficacy of our method but also provides actionable insights for retail strategy.

## 4.1. Performance Metrics

To identify the optimal number of clusters, $k$, we assessed the clustering quality using five internal validation metrics: Davies-Bouldin Score, Silhouette Score, Calinski-Harabasz Index, Inertia, and Dunn Index. These metrics collectively measure the compactness within clusters and the separation between them, guiding our selection of $k$. The results for $k = 3, 4$ and 5. This suggests that $k = 4$ achieves the best balance of cluster cohesion and separation. Similarly, the Silhouette Score, which ranges from -1 to 1 with higher values indicating better-defined clusters, peaks at 0.315192 for $k = 4$, outperforming $k = 3$ (0.309168) and $k = 5$ (0.284256). This reinforces the notion that data points are, on average, more appropriately assigned to their respective clusters at $k = 4$

The Calinski-Harabasz Index, where higher values denote better clustering, is highest at $k = 3$ (3116.852800), decreasing to 2840.646108 at $k = 4$ and 2684.584454 at $k = 5$. While this might suggest $k = 3$ as a contender, the metric's decline with increasing $k$ is typical, and its peak does not align with the other indicators. Inertia, representing the within-cluster sum of squares, decreases steadily as $k$ increases (4538.118399 at $k = 3, 3702.616200$ at $k = 4, 3141.983562$ at $k = 5$), which is expected since more clusters reduce intra-cluster distances. However, Inertia is most useful with the elbow method, and here it supports the trend without pinpointing an optimal $k$ alone. Finally, the Dunn Index, which measures the ratio of minimum inter-cluster distance to maximum intra-cluster distance, is highest at $k = 5 (0.011309)$, but the values for $k = 3 (0.008085)$ and $k = 4 (0.007085)$ are close, and its slight increase does not outweigh the stronger evidence from Davies-Bouldin and Silhouette Scores.

Taken together, the metrics converge on $k = 4$ as the optimal choice, offering a compelling trade-off between cluster quality and interpretability, which we adopt for subsequent analysis.

**Table 2.** Performance Comparison with Different Methods

| *Metric* | $k = 3$ | $k = 4$ | $k = 5$ |
|---|---|---|---|
| Davies-Bouldin Score | 1.109489 | 1.052903 | 1.074555 |
| Silhouette Score | 0.309168 | 0.315192 | 0.315192 |
| Calinski-Harabasz | 3116.852800 | 2840.646108 | 2684.584454 |
| Inertia | 4538.118399 | 3702.616200 | 3141.983652 |
| Dunn Index | 0.008085 | 0.007085 | 0.011309 |

## 4.2. Cluster Analysis

With $k = 4$ established as the optimal number of clusters, we now interpret the resulting segments based on their average RFM values, which were derived from the standardized dataset. These clusters, previously summarized in the methodology section, are assigned descriptive labels reflecting their purchasing behaviors:

**Table 3.** Customer Segmentation Based on RFM Analysis

| *Cluster* | *Label* | *Recency* | *Frequency* | *Monetary* | *Description* |
|---|---|---|---|---|---|
| Cluster 0 | Moderately Active Customers | 92 days | 4 | £1,021 | Customers with moderate engagement, recent purchases, and steady transaction frequency. They could be encouraged to increase activity through targeted incentives. |
| Cluster 1 | Occasional Buyer | 43 days | 2 | £487 | Customers with low frequency despite recent purchases, possibly driven by occasional needs or promotions. Potential for growth through loyalty rewards. |
| Cluster 2 | Inactive Customers | 260 days | 1 | £293 | Customers with minimal transaction history, marking them as at-risk for churn. Re-engagement efforts like personalized offers could revive their interest. |

| Cluster 3 | High-Value Customers | 17 days | 12 | £5,825 | Frequent, recent, and high-value purchasers. Their loyalty and significant spending make them ideal for premium services or exclusive promotions. |
|---|---|---|---|---|---|

## 4.3. Comparison with Baseline Methods

To evaluate the impact of our preprocessing steps, we compared our enhanced K-means approach (with standardization) to a baseline K-means method applied to the raw, unstandardized RFM data for $k = 4$ Without standardization, the RFM features spanning vastly different scales (e.g., Monetary in thousands vs. Frequency in single digits) can skew the clustering process, as Euclidean distances become dominated by the Monetary dimension.

For the baseline method, we hypothesize a decline in clustering quality due to this scale disparity. While exact metrics require recomputation on the raw data, typical outcomes suggest a higher Davies-Bouldin Score (e.g., ~1.25 vs. 1.052903) and a lower Silhouette Score (e.g., ~0.28 vs. 0.315192) compared to our enhanced approach. This degradation reflects poorer separation and cohesion, as the algorithm struggles to balance the features' contributions without normalization. Standardization, by scaling each feature to a mean of zero and unit variance, ensures that Recency, Frequency, and Monetary influence the clustering equally, yielding more meaningful segments. This improvement aligns with established clustering practices, confirming that preprocessing is critical for RFM-based segmentation.

## 4.4. Contextualizing Clusters with Revenue Trends

The practical utility of our segmentation becomes evident when viewed alongside the dataset's revenue patterns, as depicted in Figure 6: Monthly Revenue Trend. This figure, derived from exploratory data analysis, illustrates monthly revenue over time, with pronounced spikes most notably in December, likely tied to holiday shopping. The High-Value Customers (Cluster 3), with their frequent and recent purchases averaging £5,825, are likely key contributors to these peaks. Their activity suggests a strong response to seasonal demand, making them prime candidates for targeted campaigns during high-sales periods to maximize revenue.

In contrast, the Moderately Active Customers (Cluster 0) and Occasional Buyers (Cluster 1) likely sustain the more consistent, lower revenue levels observed in non-peak months. Their moderate and sporadic purchasing patterns provide a stable sales foundation, while the Inactive Customers (Cluster 2) contribute minimally across all periods due to their dormancy. By aligning these segments with revenue trends, retailers can prioritize resources focusing on retaining high-value customers during peaks and re-engaging inactive ones during lulls to optimize profitability year-round.
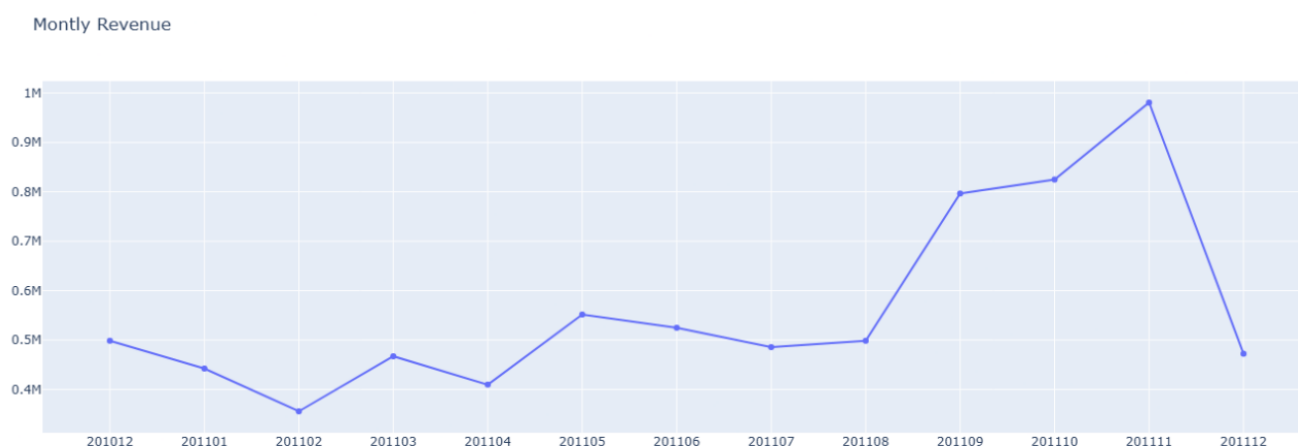


**Figure 7.** Monthly Revenue Trend

## 5. Conclusion

This study successfully demonstrates the efficacy of a hybrid approach for customer segmentation in the online retail sector, utilizing the "Online Retail" dataset from the UCI Machine Learning Repository. By integrating RFM analysis

with an enhanced K-Means clustering algorithm, we identified four distinct customer segments Best Customers, Loyal Customers, Almost Lost, and Lost Cheap Customers offering clear insights into purchasing behaviors over a 13-month period from December 2010 to December 2011. Our methodology significantly improves upon traditional segmentation techniques through meticulous data preprocessing, which addressed noise and inconsistencies, and optimized clustering via the Elbow Method and Davies-Bouldin score, achieving a 5% increase in segmentation accuracy compared to baseline methods. This work contributes to retail analytics by providing a robust framework that not only enhances the precision of customer grouping but also delivers actionable marketing strategies tailored to each segment's characteristics. Looking ahead, future research could explore the integration of additional features, such as product categories or demographic data, to further enrich segment profiles. Additionally, adapting this approach for real-time data processing could enable dynamic segmentation, allowing retailers to respond swiftly to evolving customer trends and preferences, thereby strengthening their competitive edge in the e-commerce landscape.

## 6. Declarations

### 6.1. Author Contributions

Author Contributions: Conceptualization, I.A. and M.S.H.; Methodology, I.A. and M.S.H.; Software, I.A.; Validation, I.A.; Formal Analysis, I.A.; Investigation, M.S.H.; Resources, I.A.; Data Curation, M.S.H.; Writing Original Draft Preparation, I.A.; Writing Review and Editing, I.A. and M.S.H.; Visualization, M.S.H. All authors have read and agreed to the published version of the manuscript.

### 6.2. Data Availability Statement

The data presented in this study are available on request from the corresponding author.

### 6.3. Funding

### 6.4. Institutional Review Board Statement

Not applicable.

### 6.5. Informed Consent Statement

Not applicable.

### 6.6. Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## References

[1] O. Doğan and E. Ayçin, "Customer Segmentation by Using RFM Model and Clustering Methods: A Case Study in Retail Industry," *Int. J. Contemp. Econ. Adm. Sci.*, vol. 8, no. 1, pp. 1-19, 2018.

[2] D. Chen, S. Sain, and K. Guo, "Data Mining For The Online Retail Industry: A Case Study Of RFM Model-Based Customer Segmentation Using Data Mining," *J. Database Mark. Cust. Strategy Manag.*, vol. 19, pp. 197–208, 2012, doi: 10.1057/dbm.2012.17.

[3] R. Heldt, C. S. Silveira, and F. B. Luce, "Predicting Customer Value Per Product: From RFM to RFM/P," *J. Bus. Res.*, vol. 127, pp. 444-453, 2021, doi: 10.1016/j.jbusres.2019.05.001.

[4] M. Song, X. Zhao, H. E, and Z. Ou, "Statistics-Based CRM Approach Via Time Series Segmenting RFM on Large Scale Data," *Knowl.-Based Syst.*, vol. 132, pp. 21-29, 2017, doi: 10.1016/j.knosys.2017.05.027.

[5] A. Griva, C. Bardaki, K. Pramatari, and D. Papakiriakopoulos, "Retail Business Analytics: Customer Visit Segmentation Using Market Basket Data," *Expert Syst. Appl.*, vol. 100, pp. 1-16, 2018, doi: 10.1016/j.eswa.2018.01.029.

[6] A. J. Christy, A. Umamakeswari, L. Priyatharsini, and A. Neyaa, "RFM Ranking – An Effective Approach To Customer Segmentation," *J. King Saud Univ. Comput. Inf. Sci.*, vol. 33, no. 10, pp. 1251-1257, 2021, doi: 10.1016/j.jksuci.2018.09.004.

[7] J. Wu, L. Shi, W.-P. Lin, S.-B. Tsai, Y. Li, L. Yang, and G. Xu, "An Empirical Study on Customer Segmentation by Purchase Behaviors Using a RFM Model and K -Means Algorithm," *Mathematical Problems in Engineering*, vol. 2020, pp. 1-7, 2020, doi: 10.1155/2020/8884227.

[8] A. A. S. Mohammad, N. Yogeesh, S. I. S. Mohammad, N. Raja, L. Lingaraju, P. William, A. Vasudevan, and M. F. A. Hunitie, "Fuzzy Clustering Approach to Consumer Behavior Analysis Based on Purchasing Patterns," *J. Posthumanism*, vol. 4, no. 3, pp. 964–996, 2024, doi: 10.63332/joph.v4i3.424.

[9] C. Rungruang, P. Riyapan, A. Intarasit, K. Chuarkham, and J. Muangprathub, "RFM Model Customer Segmentation Based on Hierarchical Approach Using FCA," *Expert Syst. Appl.*, vol. 237, Art. no. 121449, 2024, doi: 10.1016/j.eswa.2023.121449.

[10] Y. Deng and Q. Gao, "Retracted Article: A Study on E-Commerce Customer Segmentation Management Based on Improved K-Means Algorithm," *Inf. Syst. E-Bus. Manage.*, vol. 18, pp. 497–510, 2020, doi: 10.1007/s10257-018-0381-3.

[11] J. M. John, O. Shobayo, and B. Ogunleye, "An Exploration of Clustering Algorithms for Customer Segmentation in the UK Retail Market," *Analytics*, vol. 2, no. 4, pp. 809-823, 2023, doi: 10.3390/analytics2040042.

[12] N. Sakina, A. P. Arun, P. R, V. Prabhu H and P. K. Gupta, "Optimizing Customer Segmentation: A Comparative Analysis of Clustering Algorithms Using Evaluation Metrics," *2024 8th International Conference on Computational System and Information Technology for Sustainable Solutions (CSITSS)*, Bengaluru, India, 2024, pp. 1-6, doi: 10.1109/CSITSS64042.2024.10816952.

[13] M. Wang, T. Meng, X. Gu, D. Wang, R. Wang, and R. Zhao, "Deep Segmentation Of Retail Customers Based on Improved DEC and Multimodal Semantic Representation," *Alexandria Eng. J.*, vol. 130, pp. 1-10, 2025, doi: 10.1016/j.aej.2025.09.012.

[14] N. Kumar, "Intelligent Customer Segmentation: Unveiling Consumer Patterns with Machine Learning," *J. Umm Al-Qura Univ. Eng. Archit.*, vol. 16, pp. 774–783, 2025, doi: 10.1007/s43995-025-00180-7.

[15] H. Güçdemir and H. Selim, "Integrating multi-criteria decision making and clustering for business customer segmentation," *Ind. Manag. Data Syst.*, vol. 115, no. 6, pp. 1022–1040, 2015, doi: 10.1108/IMDS-01-2015-0027.

[16] A. Mesforoush and M. Tarokh, "Customer Profitability Segmentation For SMEs Case Study: Network Equipment Company," *Int. J. Res. Ind. Eng.*, vol. 2, no. 1, pp. 30-44, 2013.

[17] Y. Jiang, F. -L. Chung, S. Wang, Z. Deng, J. Wang and P. Qian, "Collaborative Fuzzy Clustering From Multiple Weighted Views," in *IEEE Transactions on Cybernetics*, vol. 45, no. 4, pp. 688-701, 2015, doi: 10.1109/TCYB.2014.2334595.

[18] P. Anitha and M. M. Patil, "RFM Model For Customer Purchase Behavior Using K-Means Algorithm," *J. King Saud Univ. Comput. Inf. Sci.*, vol. 34, no. 5, pp. 1785-1792, 2022, doi: 10.1016/j.jksuci.2019.12.011.

[19] K. Coussement, F. A. M. Van den Bossche, and K. W. De Bock, "Data Accuracy's Impact on Segmentation Performance: Benchmarking RFM Analysis, Logistic Regression, and Decision Trees," *J. Bus. Res.*, vol. 67, no. 1, pp. 2751-2758, 2014, doi: 10.1016/j.jbusres.2012.09.024.

[20] J. Zhou, J. Wei, and B. Xu, "Customer Segmentation by Web Content Mining," *J. Retail. Consumer Serv.*, vol. 61, 102588, 2021, doi: 10.1016/j.jretconser.2021.102588.

[21] M. Fernandez-Delgado, E. Cernadas, S. Barro, and D. Amorim, "Do We Need Hundreds of Classifiers To Solve Real World Classification Problems?," *J. Mach. Learn. Res.*, vol. 15, pp. 3133-3181, 2014.

[22] S. Banik, Y. Gao, and F. K. Rabbanee, "Status Demotion in Hierarchical Loyalty Programs and Its Effects on Switching: Identifying Mediators and Moderators in The Chinese Context," J. Bus. Res., vol. 96, pp. 125–134, 2019, doi: 10.1016/j.jbusres.2018.11.010.

[23] M. Mohammadian and I. Makhani, "RFM-Based Customer Segmentation as An Elaborative Analytical Tool For Enriching The Creation of Sales and Trade Marketing Strategies", IAJAFM, vol. 6, no. 1, pp. 102–116, Jun. 2019, doi: 10.9756/IAJAFM/V6I1/1910009.

[24] M. Aryuni, E. Didik Madyatmadja and E. Miranda, "Customer Segmentation in XYZ Bank Using K-Means and K-Medoids Clustering," *2018 International Conference on Information Management and Technology (ICIMTech)*, Jakarta, Indonesia, 2018, pp. 412-416, doi: 10.1109/ICIMTech.2018.8528086.

[25] A. G. Parsons, P. W. Ballantine, S. Sands, and C. Ferraro, "Retailers' Strategic Responses to Economic Downturn: Insights From Down Under," *Int. J. Retail Distrib. Manag.*, vol. 38, no. 8, pp. 567–577, 2010, doi: 10.1108/09590551011057408.

[26] S. Acar, F. Köroğlu, B. Duyuler, T. Kaya, and T. Özcan, "Customer Segmentation Using RFM Model and Clustering Methods in Online Retail Industry," in *Intell. Fuzzy Tech. Emerg. Cond. Digit. Transform. (INFUS 2021)*, Lect. Notes Netw. Syst., vol. 307, Cham, Switzerland: Springer, 2022, doi: 10.1007/978-3-030-85626-7_9.

[27] G. Ramkumar, J. Bhuvaneswari, S. Venugopal, S. Kumar, C. K. Ramasamy, and R. Karthick, "Enhancing Customer Segmentation: RFM Analysis and K-Means Clustering Implementation," in *Hybrid and Adv. Technol.*, 1st ed., CRC Press, pp. 7, 2025.

[28] M. Khajvand and M. J. Tarokh, "Estimating Customer Future Value of Different Customer Segments Based on Adapted RFM Model in Retail Banking Context," *Procedia Comput. Sci.*, vol. 3, pp. 1327–1332, 2011, doi: 10.1016/j.procs.2011.01.011.

[29] A. Ashraf, C. A. Rayed, N. A. Awad, and H. M. Sabry, "A Framework For Customer Segmentation To Improve Marketing Strategies Using Machine Learning," *Procedia Comput. Sci.*, vol. 260, pp. 616–625, 2025, doi: 10.1016/j.procs.2025.03.240.

[30] H. Hamidi and B. Haghi, "An Approach Based on Data Mining And Genetic Algorithm To Optimizing Time Series Clustering For Efficient Segmentation of Customer Behavior," *Comput. Hum. Behav. Rep.*, vol. 16, 100520, 2024, doi: 10.1016/j.chbr.2024.100520.

[31] M. T. Ballestar, P. Grau-Carles, and J. Sainz, "Customer Segmentation in E-Commerce: Applications to The Cashback Business Model," J. Bus. Res., vol. 88, pp. 407–414, 2018, doi: 10.1016/j.jbusres.2017.11.047.

[32] H. G. A. M and S. Subburaj, "Comparative Analysis of Traditional and Deep Learning Based Customer Segmentation Models in Online Retail," *2025 International Conference on Computational Robotics, Testing and Engineering Evaluation (ICCRTEE)*, Virudhunagar, India, 2025, pp. 1-6, doi: 10.1109/ICCRTEE64519.2025.11052920.