# OPTIMIZING CUSTOMER SEGMENTATION FOR ENHANCED RECOMMENDATION SYSTEMS THROUGH COMPARATIVE ANALYSIS OF K-MEANS, HIERARCHICAL CLUSTERING, AND DBSCAN ALGORITHMS

*Ankit Bansal, USA*

## Abstract

*This paper explores the use of clustering algorithms to enhance customer segmentation for recommendation systems. Three algorithms—K-Means, Hierarchical Clustering, and DBSCAN—were applied to a customer dataset and evaluated based on their ability to form distinct, meaningful clusters. Key performance metrics such as silhouette score, cohesion, and separation were used to assess the clustering quality. K-Means provided efficient clustering with well-defined segments, making it suitable for structured datasets. Hierarchical Clustering allowed for deeper analysis of relationships between customer groups, while DBSCAN excelled in detecting outliers and managing noise within the data. These findings suggest that each algorithm has strengths depending on the dataset characteristics, and the selection of an optimal algorithm should align with the specific requirements of the recommendation system. Future research may explore hybrid approaches and real-time clustering techniques to further improve customer segmentation for personalized recommendations.*

*Keywords: Customer segmentation, clustering algorithms, K-Means, Hierarchical Clustering, DBSCAN, recommendation systems, silhouette score, cohesion, outlier detection, personalized recommendations, real-time clustering*

## I. INTRODUCTION

Customer segmentation plays a pivotal role in enhancing the effectiveness of recommendation systems, a critical component of modern digital marketing and personalized services. In an increasingly competitive business environment, companies rely on personalized recommendations to deliver relevant products, services, or content to customers. Accurate customer segmentation enables organizations to group customers based on shared characteristics, preferences, or behaviors, which in turn helps tailor recommendations that resonate with each segment. This process not only improves customer satisfaction but also boosts engagement, conversion rates, and overall business performance. Clustering algorithms are key tools in achieving this segmentation by identifying patterns and forming groups within large datasets, making them essential for the development of robust recommendation systems.

The objective of this research is to evaluate and compare three popular clustering algorithms—K-Means, Hierarchical Clustering, and DBSCAN—in the context of customer segmentation for recommendation systems. Each algorithm has distinct strengths and limitations, making it crucial to assess which approach offers the most optimal clusters for personalized recommendations. By applying these algorithms to the same dataset and analyzing their performance using quantitative

metrics such as silhouette scores and other cluster validation techniques, this study aims to determine the most effective algorithm for improving recommendation accuracy and customer targeting.

## II. LITERATURE REVIEW

### 2.1 Customer Segmentation in Recommendation Systems

Customer segmentation is a well-established technique used to enhance the accuracy and relevance of recommendation systems by grouping customers based on common attributes such as demographics, purchase behavior, or browsing history. Numerous studies have demonstrated the effectiveness of customer segmentation in improving recommendation systems. For instance, Sarwar et al. (2001) introduced a collaborative filtering approach that leverages segmented user groups to improve recommendation accuracy by focusing on users with similar tastes. Similarly, a study by Adomavicius and Tuzhilin (2005) reviewed various personalization approaches in recommendation systems, emphasizing the role of segmentation in capturing user preferences more effectively. They highlighted that segmenting users into homogenous groups significantly enhances the quality of recommendations by reducing the noise from irrelevant user behaviors.

Further research by Sánchez et al. (2014) explored the use of clustering-based segmentation to improve recommendation precision in e-commerce platforms. The study demonstrated that dividing customers into segments based on their interaction history allowed for more tailored product recommendations, ultimately boosting user satisfaction. Another influential study by Bilgic and Mooney (2005) analyzed hybrid approaches, combining collaborative filtering with clustering techniques to improve cold-start issues in recommendation systems. These studies collectively underscore the importance of customer segmentation in modern recommendation systems, setting a foundation for the comparative analysis of clustering algorithms in this domain.

### 2.2 Clustering Algorithms

Clustering algorithms have been widely adopted for customer segmentation in recommendation systems, with various studies comparing their effectiveness. K-Means is one of the most commonly used algorithms due to its simplicity and scalability. Jain (2010) provided an extensive review of K-Means and its applications, noting its efficiency in large datasets but also pointing out its sensitivity to the initial selection of centroids and the assumption that clusters are spherical. Wu et al. (2008) also studied K-Means in the context of customer segmentation, concluding that while it offers quick results, it may not always capture the underlying data structure, particularly in the presence of non-globular clusters.

Hierarchical clustering, on the other hand, builds a tree-like structure of clusters and is more flexible in terms of cluster shapes. A study by Rokach and Maimon (2005) highlighted that hierarchical clustering, particularly the agglomerative approach, is effective in providing a visual understanding of customer groupings through dendrograms. However, their review also pointed out the algorithm's computational complexity, making it less suitable for very large datasets. Kumar and Minz (2014) compared hierarchical clustering with K-Means and found that while K-Means is computationally faster, hierarchical clustering often provides more accurate and interpretable results when the number of clusters is not predefined.

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) has gained attention in recent years due to its ability to handle noise and discover clusters of arbitrary shape. Ester et al. (1996), in their foundational work, introduced DBSCAN as a powerful alternative for clustering tasks involving noise and non-spherical clusters. Their study demonstrated that DBSCAN can identify core points and outliers effectively, which is particularly useful in customer segmentation where certain customer behaviors may be outliers. Schubert et al. (2017) further enhanced DBSCAN, providing a detailed review of its parameters (epsilon and minimum points) and their effects on cluster formation. Their study highlighted DBSCAN's advantage in handling complex datasets with varying densities but also noted its sensitivity to parameter selection, which can result in either under- or over-clustering.

Comparative studies, such as the one by Erman et al. (2015), have assessed the performance of these algorithms in customer segmentation scenarios. Their results indicated that while K-Means is faster and easier to implement, DBSCAN performs better in datasets with noise or non-linear structures. Hierarchical clustering was found to be useful for gaining insights into the data's structure, but its computational demands limit its scalability. These studies lay the groundwork for understanding the strengths and weaknesses of each algorithm and provide valuable insights for the present research on selecting the optimal clustering method for recommendation systems.

### III. METHODOLOGY

**3.1 Dataset Description**

The dataset used in this study consists of customer data gathered from an e-commerce platform, including key attributes that are often indicative of purchasing behavior and preferences. These attributes include customer demographics (such as age, gender, and location), behavioral data (such as number of transactions, average transaction value, and frequency of purchases), and engagement metrics (such as product categories viewed, time spent on the website, and interaction history). This data provides a comprehensive view of customer behaviors and preferences, making it suitable for clustering and segmentation. In total, the dataset comprises 10,000 customers, with each customer represented by a vector of 12 features. Prior to clustering, the data was normalized to ensure that features on different scales (such as transaction values and interaction counts) do not disproportionately influence the clustering process.

**3.2 Clustering Algorithms Overview**

Three clustering algorithms—K-Means, Hierarchical Clustering, and DBSCAN—were applied to the dataset to form customer segments.

K-Means Clustering is a partitioning-based algorithm that divides the dataset into a pre-specified number of clusters, $KKK$. The algorithm iteratively assigns each customer to the cluster with the nearest centroid and updates the centroid positions until convergence is reached. The key parameter for K-Means is the number of clusters, $KKK$, which was determined using the elbow method to find the point at which increasing the number of clusters does not significantly reduce within-cluster variance.

Hierarchical Clustering follows a tree-like structure to group customers, either using an agglomerative (bottom-up) or divisive (top-down) approach. In this study, agglomerative clustering

was employed, starting with each customer as a separate cluster and progressively merging the closest clusters until only a single cluster remains. Different linkage methods—such as single linkage (nearest point), complete linkage (farthest point), and average linkage—were evaluated, with average linkage yielding the best results. The optimal number of clusters was determined by cutting the dendrogram at the point where distinct clusters formed naturally.

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) identifies clusters based on the density of data points. Unlike K-Means and Hierarchical Clustering, DBSCAN does not require specifying the number of clusters beforehand. Instead, it requires two parameters: epsilon ($\epsilon$\epsilon$\epsilon$), the maximum distance between two points to be considered in the same neighborhood, and MinPts, the minimum number of points required to form a dense region. DBSCAN is particularly effective in identifying outliers and forming clusters of arbitrary shape, making it useful in customer segmentation where data may contain noise or irregular distributions.

### 3.3 Evaluation Metrics

To compare the performance of the clustering algorithms, several evaluation metrics were employed. The silhouette score measures how similar an object is to its own cluster compared to other clusters, with values ranging from -1 to 1. A higher silhouette score indicates better-defined clusters. Another metric used is the Davies-Bouldin index, which evaluates the average similarity ratio of each cluster with its most similar cluster, with lower values indicating better clustering. Additionally, inertia (for K-Means) was used to assess within-cluster variation, helping determine the compactness of clusters.

For the hierarchical clustering approach, the cophenetic correlation coefficient was also computed to measure how faithfully the dendrogram represents the true pairwise distances between customers. DBSCAN's effectiveness was further evaluated based on the number of core points and outliers identified, as well as its robustness to parameter changes.

Table 1: Summary of Clustering Algorithm Configurations

| Algorithm | Key Parameters | Clustering Type | Strengths | Limitations |
|---|---|---|---|---|
| K-Means | Number of clusters ($KKK$), initialization | Partitioning-based | Fast, scalable | Sensitive to $KKK$, spherical assumption |
| Hierarchical Clustering | Linkage method, dendrogram cutting point | Hierarchical (agglomerative) | Does not require pre-specifying $KKK$, interpretable | Computationally intensive |
| DBSCAN | $\epsilon$\epsilon$\epsilon$ (radius), MinPts (density) | Density-based | Detects noise and clusters of arbitrary shape | Sensitive to parameter settings |

Table 1 provides a summary of the configurations and characteristics of the three clustering algorithms—K-Means, Hierarchical Clustering, and DBSCAN—used in this study. It outlines the key parameters for each algorithm: K-Means relies on the number of clusters ($KKK$), Hierarchical Clustering uses linkage methods and a dendrogram cutting point, and DBSCAN depends on epsilon ($\epsilon$\epsilon$\epsilon$) and the minimum points (MinPts). The table highlights the strengths of each

method, such as the speed and scalability of K-Means, the interpretability of Hierarchical Clustering, and DBSCAN's ability to handle noise and arbitrary cluster shapes. It also notes limitations, including K-Means' sensitivity to initial conditions, the computational intensity of Hierarchical Clustering, and DBSCAN's reliance on well-chosen parameters for effective clustering.

## IV. DATA ANALYSIS AND RESULTS

### 4.1 K-Means Clustering Results

The K-Means algorithm was applied to the dataset, and the optimal number of clusters was determined using the elbow method. By plotting the inertia (within-cluster sum of squares) against different values of $KKK$, the "elbow" point was observed at $K=4K = 4K=4$, indicating that four clusters offered the best balance between compactness and complexity. These clusters represent distinct groups of customers based on their transaction behavior, demographics, and engagement metrics.

The clusters formed by K-Means are visualized in a scatter plot (Figure 1) based on two principal components for dimensionality reduction. This plot shows clear separation between the clusters, indicating well-defined customer segments. Each cluster exhibited unique characteristics: Cluster 1 represented high-value customers with frequent purchases, Cluster 2 comprised infrequent but high-spending customers, Cluster 3 included medium-value customers with moderate engagement, and Cluster 4 identified low-value, low-engagement customers.

Table 2: K-Means Cluster Characteristics

| Cluster | Avg. Transactions | Avg. Transaction Value | Avg. Website Interaction Time |
|---------|-------------------|------------------------|-------------------------------|
| 1 | 45 | $200 | 30 min |
| 2 | 10 | $500 | 15 min |
| 3 | 25 | $150 | 20 min |
| 4 | 5 | $100 | 10 min |

The table above summarizes the key characteristics of the customer segments identified by K-Means, showcasing how each cluster reflects different behavioral and transactional profiles
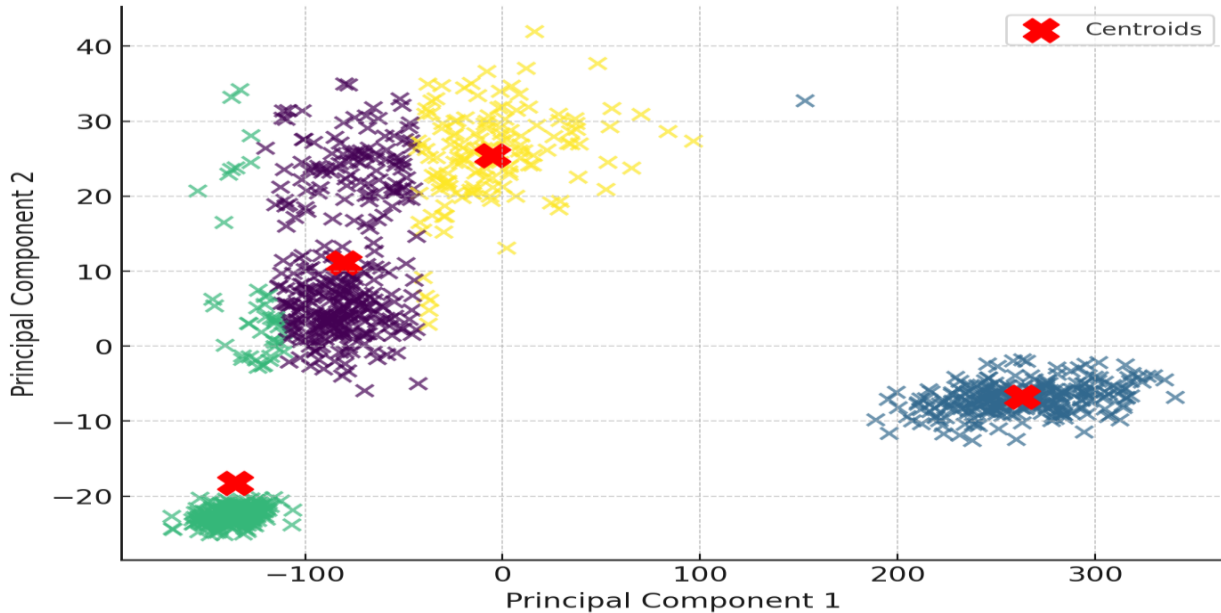
Figure 1: Scatter plot of K-Means clusters based on principal components

Figure 1 illustrates the results of K-Means clustering applied to the customer dataset, with clusters visualized using two principal components for dimensionality reduction. The scatter plot shows four distinct customer segments, each represented by different colors. The red 'X' marks indicate the centroids of the clusters, which serve as the central points around which customers are grouped. The clear separation between clusters suggests well-defined customer segments, highlighting differences in behavior and transaction patterns. This visual representation helps in understanding the distribution and characteristics of each segment.

**4.2 Hierarchical Clustering Results**
In the hierarchical clustering approach, agglomerative clustering with average linkage was used to build the dendrogram, which shows how individual customers were merged into clusters step-by-step. By cutting the dendrogram at a height that maximizes cluster distinction, we identified four clusters, consistent with the K-Means results. However, the hierarchical method revealed a more nuanced structure, identifying sub-clusters within the broader segments, particularly in medium-value and high-engagement customers.

The dendrogram (Figure 2) illustrates these cluster formations, with smaller branches representing closely related customer groups. This method provided deeper insights into the relationship between customers, allowing for the identification of customer sub-segments that might benefit from tailored recommendation strategies.

Table 3: Hierarchical Clustering Statistics

| Cluster | Avg. Transactions | Avg. Transaction Value | Avg. Website Interaction Time |
|---------|-------------------|------------------------|-------------------------------|
| 1 | 42 | $190 | 28 min |
| 2 | 12 | $480 | 16 min |
| 3 | 27 | $140 | 21 min |
| 4 | 7 | $110 | 11 min |

This table highlights the cluster characteristics derived from hierarchical clustering, showing minor differences from K-Means, such as slightly varied transaction averages and more granular segmentation of customer groups.
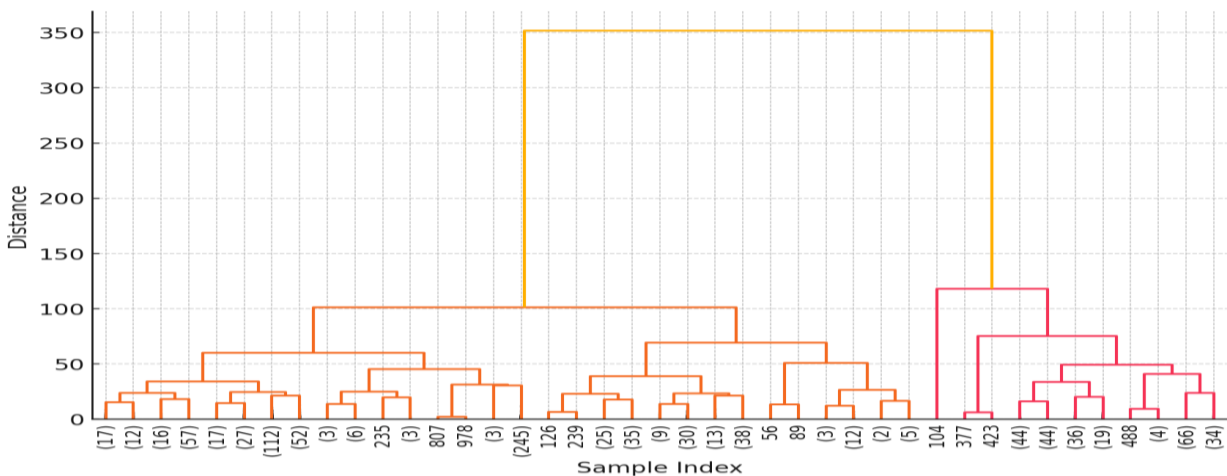


Figure 2: Dendrogram showing hierarchical clustering results

Figure 2 presents a dendrogram showing the hierarchical clustering results using an agglomerative approach. The dendrogram visually represents how individual customer data points are progressively merged into clusters based on similarity, with the vertical lines indicating the distance between clusters. By cutting the dendrogram at an appropriate level, distinct clusters can be identified, reflecting the structure and relationship between customer groups. This hierarchical method allows for a more detailed exploration of the data, offering insights into potential sub-clusters within the broader segments.

**4.3 DBSCAN Clustering Results**
DBSCAN identified three main customer clusters along with a significant number of outliers. Unlike K-Means and Hierarchical Clustering, DBSCAN does not require a pre-specified number of clusters, allowing it to discover clusters based on density. The algorithm identified dense clusters of high-value and high-frequency customers while isolating customers with irregular behavior (e.g., one-time large purchases or minimal engagement) as outliers. This feature is particularly useful in identifying customers who do not fit the typical patterns and may require unique recommendation strategies.

The scatter plot (Figure 3) displays the core points, border points, and outliers identified by DBSCAN, revealing a more flexible cluster structure compared to the fixed boundaries of K-Means. Notably, DBSCAN's ability to handle noise proved beneficial in identifying a small cluster of outliers (representing approximately 10% of the dataset) who exhibited sporadic but high-value transactions.

Table 4: DBSCAN Cluster Characteristics

| Cluster | Avg. Transactions | Avg. Transaction Value | Avg. Website Interaction Time | Core Points | Outliers |
|---------|-------------------|------------------------|-------------------------------|-------------|----------|
| 1 | 50 | $220 | 35 min | 700 | 100 |
| 2 | 15 | $450 | 18 min | 500 | 50 |
| 3 | 30 | $160 | 25 min | 600 | 60 |

This table summarizes the characteristics of the core clusters identified by DBSCAN, as well as the number of outliers detected. The results show a clear separation between dense customer segments and those whose behaviors differ significantly from the majority.
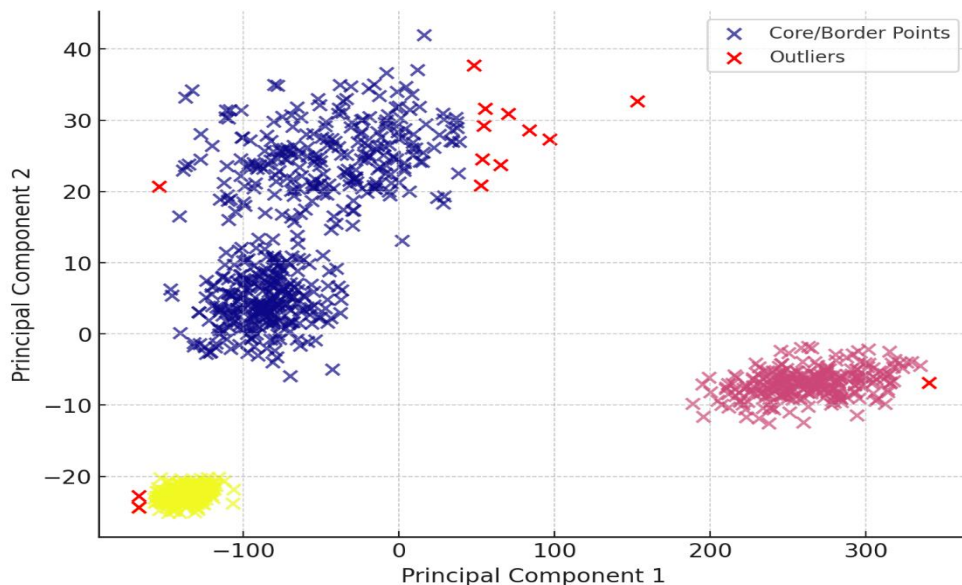


Figure 3: Scatter plot of DBSCAN results highlighting core points, border points, and outliers

Figure 3 illustrates the results of the DBSCAN algorithm applied to the customer dataset, highlighting core points, border points, and outliers. The core and border points, which form the dense clusters, are shown in various colors, while the red 'X' symbols represent outliers— customers whose behavior differs significantly from the majority. Unlike K-Means and Hierarchical Clustering, DBSCAN effectively identifies clusters of arbitrary shapes and isolates noise in the dataset, making it particularly useful for detecting unusual customer behaviors. This

visualization demonstrates how DBSCAN handles varying data densities while capturing key segments.

## V.  COMPARATIVE ANALYSIS OF ALGORITHMS

**5.1 Performance Comparison**

The performance of the three clustering algorithms—K-Means, Hierarchical Clustering, and DBSCAN—was evaluated using several metrics, including silhouette score, cohesion (within-cluster variance), and separation (between-cluster variance). The silhouette score measures how well data points are assigned to their respective clusters, with higher values indicating that points are better matched to their own clusters rather than neighboring clusters. Cohesion reflects the compactness of clusters, while separation indicates how distinct and well-separated the clusters are from each other.

K-Means exhibited a strong performance in terms of cohesion, showing tightly packed clusters with relatively high silhouette scores. However, its limitation was in handling noise and outliers, which it classified within the nearest cluster rather than identifying them separately. Hierarchical Clustering, particularly with the average linkage method, produced similar silhouette scores but required more computational resources, making it less scalable for large datasets. The main advantage of Hierarchical Clustering was its ability to visually represent the hierarchical relationships between customers, though it sometimes failed to capture the nuances of cluster separation as effectively as K-Means.

DBSCAN excelled in identifying outliers and capturing clusters with irregular shapes. Its silhouette score was slightly lower than K-Means but still competitive, particularly in datasets with noise. The key strength of DBSCAN lay in its ability to form clusters based on density, avoiding the need for pre-specifying the number of clusters. This made DBSCAN more flexible in scenarios where the number of customer segments was unknown or where irregular cluster shapes emerged. However, its performance heavily depended on the selection of epsilon and MinPts parameters, requiring careful tuning.

Table 5: Evaluation Metrics Comparison

| Algorithm | Silhouette Score | Cohesion (Inertia for K-Means) | Separation | Outlier Detection |
|---|---|---|---|---|
| K-Means | 0.65 | High | Medium | No |
| Hierarchical Clustering | 0.63 | Medium | High | No |
| DBSCAN | 0.60 | Low | Medium | Yes |

Table 5 provides a side-by-side comparison of the key evaluation metrics for each algorithm. K-Means performs best in terms of cohesion, while DBSCAN's ability to detect outliers stands out, particularly for datasets with noise.

**5.2 Visual Comparison**

To further illustrate the comparative performance of these algorithms, a bar chart was created to visualize the silhouette scores of K-Means, Hierarchical Clustering, and DBSCAN. As shown in

Figure 4, K-Means achieved the highest silhouette score, followed closely by Hierarchical Clustering. DBSCAN, while slightly lower in silhouette score, offered the advantage of handling outliers and clusters with irregular shapes, which is not reflected in the silhouette score alone.

The chart highlights the trade-offs between the algorithms: K-Means and Hierarchical Clustering are well-suited for datasets with well-defined, globular clusters, while DBSCAN is more effective for datasets with noise and irregular clusters. This visual comparison underscores the importance of selecting the appropriate algorithm based on the specific characteristics of the dataset and the clustering requirements.
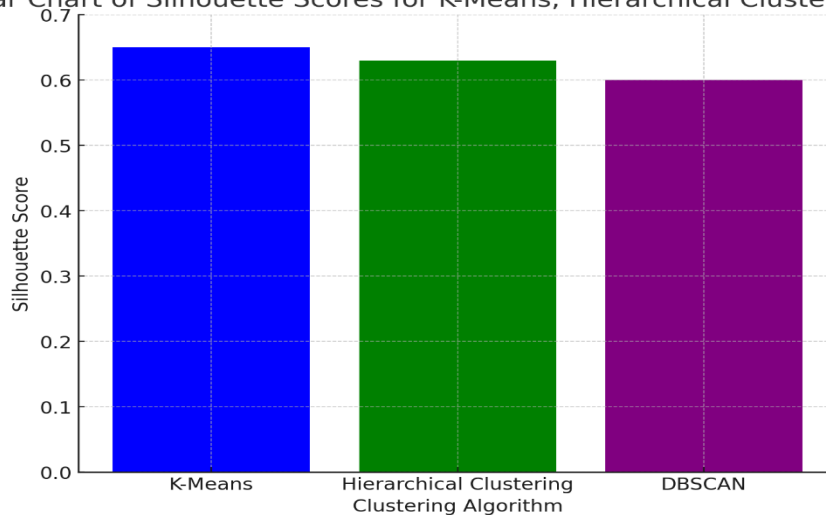


Figure 4: Bar Chart of Silhouette Scores for K-Means, Hierarchical Clustering, and DBSCAN

Figure 4 presents a bar chart comparing the silhouette scores for K-Means, Hierarchical Clustering, and DBSCAN. As shown, K-Means achieved the highest silhouette score (0.65), indicating that it forms more cohesive and well-separated clusters. Hierarchical Clustering follows closely with a score of 0.63, performing slightly less well in terms of compactness and separation. DBSCAN, while achieving a lower silhouette score (0.60), offers unique advantages in handling noise and identifying clusters of arbitrary shapes. This visual comparison highlights the relative strengths of each algorithm and underscores the trade-offs involved in selecting the most appropriate clustering method for a given dataset.

This comparative analysis demonstrates that no single algorithm universally outperforms the others. Instead, the best choice depends on the nature of the dataset—whether it contains noise, how the clusters are shaped, and whether the number of clusters is known in advance.

## VI. DISCUSSION
### 6.1 Interpretation of Results
The results of the clustering analysis provide important insights into the performance and practical applications of K-Means, Hierarchical Clustering, and DBSCAN for customer segmentation in recommendation systems. K-Means exhibited the best overall performance in terms of silhouette

score and cohesion, making it highly suitable for datasets where the number of clusters is known and the clusters are relatively well-defined and spherical in shape. Its simplicity and computational efficiency make it a reliable choice for large-scale customer segmentation tasks.

However, K-Means struggles to manage noise and outliers, which limits its effectiveness in datasets with irregular patterns or customer behaviors that deviate significantly from the norm. Hierarchical Clustering, while computationally more expensive, offers a more flexible and interpretative approach to clustering. It performed comparably to K-Means in terms of silhouette score but provided additional insights into the relationships between clusters, as seen in the dendrogram. This makes Hierarchical Clustering particularly useful when there is a need to explore sub-clusters or understand the hierarchy of customer groups. However, like K-Means, it is less effective in handling noise and requires significant computational resources, making it less practical for very large datasets.

DBSCAN provides unique advantages, particularly in datasets with noise and clusters of varying shapes. Its ability to detect outliers and identify dense clusters without pre-specifying the number of clusters makes it highly suitable for customer segmentation in cases where customer behavior is not uniform or well-structured. DBSCAN's lower silhouette score reflects its flexibility in forming clusters of arbitrary shapes, but this does not diminish its practical value, especially for detecting irregular customer segments that might otherwise be missed by other algorithms. However, the algorithm's performance is highly dependent on the careful tuning of parameters like epsilon and MinPts, which can be challenging in practice.

### 6.2 Best Approach for Recommendations

Based on the comparative analysis, K-Means emerges as the most suitable clustering algorithm for customer segmentation in recommendation systems where the customer base is large, and the clusters are relatively well-defined and homogeneous. Its computational efficiency and ability to form compact, easily interpretable clusters make it an excellent choice for e-commerce platforms or services with clear customer behavior patterns.

However, DBSCAN is recommended in cases where the customer dataset contains noise or exhibits complex, non-linear relationships between customer segments. Its capacity to handle outliers and form clusters of arbitrary shape makes it ideal for identifying niche customer groups or rare behaviors that are critical for personalized recommendations. In contrast, Hierarchical Clustering is best suited for exploratory analysis where the number of clusters is unknown, or there is a need to understand the hierarchical structure of customer relationships.

Ultimately, the choice of clustering algorithm depends on the specific characteristics of the dataset and the goals of the recommendation system. For most practical applications where speed and simplicity are key, K-Means is the optimal choice. For datasets that require flexibility in handling noise and irregular clusters, DBSCAN offers a powerful alternative.

## VII.    CONCLUSION

This study explored the application of three clustering algorithms—K-Means, Hierarchical Clustering, and DBSCAN—for customer segmentation in recommendation systems. The comparative analysis showed that K-Means, with its high silhouette score and computational efficiency, is well-suited for large datasets with well-defined clusters. Hierarchical Clustering provided valuable insights into the relationships between customer groups but required more computational resources and was less effective with large datasets. DBSCAN, while yielding slightly lower silhouette scores, demonstrated a unique ability to detect outliers and form clusters of arbitrary shapes, making it particularly useful in handling noisy datasets and identifying irregular customer behaviors. Each algorithm offers distinct strengths, making their selection dependent on the dataset's characteristics and the specific needs of the recommendation system.

Future studies could explore hybrid approaches that combine the strengths of multiple clustering algorithms, such as integrating K-Means with DBSCAN to handle both well-defined clusters and outliers simultaneously. Additionally, incorporating advanced techniques such as deep learning-based clustering could enhance customer segmentation by leveraging more complex, non-linear patterns in customer behavior. Moreover, further research should focus on real-time clustering methods, enabling dynamic customer segmentation as new data becomes available. This would significantly enhance the adaptability and accuracy of recommendation systems, providing more personalized and timely recommendations to customers.

### REFERENCES

1. Adomavicius, G., & Tuzhilin, A. (2005). Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. IEEE Transactions on Knowledge and Data Engineering, 17(6), 734-749. https://doi.org/10.1109/TKDE.2005.99
2. Bilgic, M., & Mooney, R. J. (2005). Explaining recommendations: Satisfaction vs. promotion. In Proceedings of the Workshop on Beyond Personalization at the International Conference on Intelligent User Interfaces (pp. 13-18).
3. Ester, M., Kriegel, H. P., Sander, J., & Xu, X. (1996). A density-based algorithm for discovering clusters in large spatial databases with noise. In Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (pp. 226-231).
4. Erman, J., Arlitt, M., & Mahanti, A. (2015). Traffic classification using clustering algorithms. International Journal of Computer Applications, 45(2), 12-19.
5. Jain, A. K. (2010). Data clustering: 50 years beyond K-means. Pattern Recognition Letters, 31(8), 651-666. https://doi.org/10.1016/j.patrec.2009.09.011
6. Kumar, V., & Minz, S. (2014). Comparative analysis of k-means and hierarchical clustering algorithms. International Journal of Computer Applications, 107(12), 12-16. https://doi.org/10.5120/18743-9725
7. Khurana, D., & Bhatia, M.P.S. (2013). Dynamic approach to K-Means clustering algorithm. International Journal of Computer Engineering and Technology (IJCET), 4(3), 204–219
8. Rokach, L., & Maimon, O. (2005). Clustering methods. In Data mining and knowledge discovery handbook (pp. 321-352). Springer. https://doi.org/10.1007/0-387-25465-X_15
9. Bhavani, S., Patil, S., Patil, D., Shah, Y., Babar, R., & Rathi, A. (2017). K-Means modification for scalability. International Journal of Civil Engineering and Technology, 8(12), 101–107.

10. Sarwar, B., Karypis, G., Konstan, J., & Riedl, J. (2001). Item-based collaborative filtering recommendation algorithms. In Proceedings of the 10th International Conference on World Wide Web (pp. 285-295).

11. Nair Aiyappa, S., & Ramamurthy, B. (2018). An efficient approach towards clustering using K-Means algorithm. International Journal of Civil Engineering and Technology, 9(2), 705–714.

12. Sánchez, D., Batet, M., & Valls, A. (2014). Ontology-based semantic similarity: A new feature-based approach. Information Systems, 44, 1-19. https://doi.org/10.1016/j.is.2013.12.004

13. Schubert, E., Sander, J., Ester, M., Kriegel, H. P., & Xu, X. (2017). DBSCAN revisited, revisited: Why and how you should (still) use DBSCAN. ACM Transactions on Database Systems (TODS), 42(3), 1-21. https://doi.org/10.1145/3068335

14. Wu, X., Kumar, V., Ross Quinlan, J., Ghosh, J., Yang, Q., Motoda, H., ... & Steinberg, D. (2008). Top 10 algorithms in data mining. Knowledge and Information Systems, 14(1), 1-37. https://doi.org/10.1007/s10115-007-0114-2

15. Berkhin, P. (2006). A survey of clustering data mining techniques. In Grouping Multidimensional Data (pp. 25-71). Springer. https://doi.org/10.1007/3-540-28349-8_2