



Machine Learning Approach for Best Location of Retailers

Ehsan Ghafourian ^a, Elnaz Bashir ^b, Farzaneh Shoushtari ^c, Ali Daghighi ^d

^a Department of Computer Science, Iowa State University, Ames, IA, 50010,

^b Department of Computer Science, Iowa State University, Ames, IA, 50010,

^c Alumni of Industrial Engineering, Bu-Ali Sina University, Hamedan, Iran,

^d Faculty of Engineering and Natural Sciences, Biruni University, Istanbul, Turkey.

ARTICLE INFO

Received: 2022/09/13

Revised: 2022/10/05

Accept: 2022/12/03

Keywords:

*Machine Learning,
Location, Retailers,
Clustering, K-means.*

ABSTRACT

This paper presents a machine learning approach using the k-means clustering algorithm to identify optimal locations for retailers. The study aims to leverage geographic, demographic, and economic factors to cluster potential locations and provide valuable insights for decision-making. The methodology involves data preparation, selecting relevant features, applying the k-means algorithm, evaluating cluster results, and visualizing the outcomes on a map. Numerical results demonstrate the effectiveness of the proposed approach in identifying suitable retail locations. The study concludes with a summary of findings and recommendations for further research.

1. Introduction

The retail industry plays a crucial role in the global economy, and choosing an optimal location for retailers significantly impacts their success. Identifying the best location involves considering numerous factors such as population density, income levels, proximity to transportation, competition, and customer preferences. Traditional approaches often rely on human expertise and intuition; however, machine learning techniques offer an enhanced and data-driven solution. In this paper, we propose a machine learning approach utilizing the k-means clustering algorithm to discover clusters of potential retail locations based on relevant features (see Figure 1) [1].

^a Corresponding author email address: ehsang@iastate.edu (E. Ghafourian).

Available online 12/03/2022

2676-3311/BGSA Ltd.



Figure 1: Retailers in Supply Chain.

The retail industry plays a pivotal role in the global economy, with businesses constantly striving for success and profitability. One of the critical factors that significantly influence the success of retailers is the choice of an optimal location. Selecting the best location involves careful consideration of various factors, including population density, income levels, competition, transportation accessibility, and customer preferences. Traditionally, retailers heavily relied on human expertise and intuition to identify suitable locations; however, with the advent of technology and the availability of vast amounts of data, machine learning approaches have emerged as powerful tools for retail location analysis [2].

Machine learning, a subfield of artificial intelligence, focuses on developing algorithms that can learn from and make predictions or decisions based on data. It has gained widespread recognition and adoption across various industries due to its capability to process large volumes of data, identify patterns, and provide valuable insights. In the context of retail location analysis, machine learning algorithms offer the potential to uncover hidden patterns and relationships among diverse variables, leading to informed decisions on the best location for retailers [3].

This paper aims to present a machine learning approach specifically utilizing the k-means clustering algorithm to identify the best locations for retailers. The k-means clustering algorithm is a popular unsupervised machine learning technique that groups data points into clusters based

on their similarity. By leveraging this algorithm, retailers can gain valuable insights into potential retail locations by identifying groups of areas with similar characteristics.

The significance of this study lies in its ability to integrate various factors into the process of retail location analysis. Geographic factors such as proximity to transportation hubs, urbanization, and land availability, demographic factors like population density, age distribution, and income levels, as well as economic factors such as disposable income, consumer spending patterns, and market potential, can all be considered using a comprehensive machine learning approach [4].

The objectives of this study are as follows:

1. To explore the application of machine learning algorithms, specifically the k-means clustering algorithm, in the analysis of retail locations.
2. To demonstrate the effectiveness of the proposed machine learning approach in identifying clusters of potential retail locations based on selected features.
3. To provide insights and recommendations for decision-making in the selection of the best retail locations.

By employing a data-driven approach, this study seeks to enhance the conventional decision-making process by providing a systematic and objective analysis of potential retail locations. The findings from this research can assist retailers in optimizing their business strategies, reducing risks, and maximizing the chances of success [5].

The subsequent sections of this paper will delve into a comprehensive literature review, present the methodology employed, discuss the numerical results obtained, draw conclusions based on the findings, and provide recommendations for future research. Through this study, we seek to contribute to the existing body of knowledge on retail location analysis and highlight the potential of machine learning algorithms in the retail industry.

2. Literature review

This section reviews existing literature related to retail location analysis, machine learning techniques, and the application of k-means clustering. The review highlights studies that explore the use of geographic, demographic, and economic factors for retail site selection. Additionally, it

discusses the advantages and limitations of k-means clustering in this context and presents relevant case studies that have employed this algorithm for retail location analysis.

The recent work about application machine learning approach for best location of retailers are defined and try to determine research gaps. Although the researchers cover gap research and suggest contributions to this issue, when new concepts come, they can apply and combine with this study that is not defined previously

The main contribution and novelty of this research based on the research gaps are as follows:

- Machine Learning Approach for Best Location of Retailers.

The selection of an optimal location for retailers has long been a topic of interest among researchers and practitioners in the field of business and marketing. Traditionally, retailers relied on manual assessment and subjective judgments to determine suitable locations. However, in recent years, the integration of machine learning techniques into retail location analysis has gained significant attention. This section presents a comprehensive literature review on the utilization of machine learning approaches for determining the best location of retailers [4-5].

1. Machine Learning Techniques in Retail Location Analysis: Machine learning techniques, such as clustering algorithms, have gained prominence in analyzing complex datasets and identifying patterns. Ertugrul and Karahoca (2016) proposed a hybrid approach combining k-means clustering and geographic information systems (GIS) to identify suitable store locations. The study demonstrated the effectiveness of cluster analysis in grouping areas with similar characteristics and facilitating decision-making regarding retail locations [6].
2. Integration of Demographic and Socio-Economic Factors: Machine learning models have been widely employed to incorporate demographic and socio-economic factors into retail location analysis. Seyedghorban et al. (2019) leveraged random forest algorithms to analyze demographic and consumer data for identifying potential retail locations. The study showcased the significance of considering factors such as population density, age distribution, income levels, and consumer preferences when determining the best retail locations [7].

3. **Spatial Analysis and Location Modeling:** Incorporating spatial analysis techniques has further enhanced the accuracy and precision of retail location analysis. A study by Sidias et al. (2018) used a combination of clustering algorithms and spatial statistics to determine suitable locations for retailers. The research illustrated the importance of considering geographic factors such as proximity to transportation hubs, competition, and accessibility in retail location decision-making [6-7].
4. **Predictive Analytics for Market Potential:** Machine learning techniques have been applied to predict market potential and customer behavior, providing valuable insights for retail location analysis. Siu et al. (2020) deployed artificial neural networks to forecast retail sales volumes at potential locations. The study highlighted the capability of predictive analytics in estimating future performance and optimizing store placement strategies.
5. **Comparative Analysis of Machine Learning Approaches:** Several studies have conducted comparative analyses of various machine learning algorithms to identify the most effective approach for retail location analysis. Das et al. (2018) compared k-means, hierarchical clustering, and support vector machines for retail location prediction. The research emphasized the superiority of k-means clustering in identifying clusters of potential retail locations based on selected features [8].
6. **Integration of Multiple Data Sources:** Machine learning approaches allow the integration of diverse data sources, such as demographic data, customer behavior data, and spatial data, to provide a comprehensive analysis of retail locations. Wang et al. (2021) utilized a gradient boosting machine (GBM) algorithm to combine demographic, social media, and mobile phone data for retail location analysis. The study demonstrated the benefits of leveraging multiple data sources in identifying optimal retail locations [9].

In conclusion, the literature review reveals a growing interest in utilizing machine learning approaches for determining the best location of retailers. The integration of clustering algorithms, demographic factors, spatial analysis, predictive analytics, and the fusion of multiple data sources have demonstrated significant advancements in retail location analysis. However, it is essential to consider the unique characteristics of each retail context and tailor machine learning techniques accordingly. Future research focusing on real-time data integration, model interpretability, and

validation of machine learning approaches in different retail settings would further enhance the application of these techniques in retail industry decision-making processes.

3. Methodology

The methodology section outlines the step-by-step process of applying the k-means clustering algorithm to determine the best location for retailers. It begins with data collection and preparation, followed by feature selection and normalization. The subsequent steps involve determining the number of clusters, applying the k-means algorithm, evaluating the clusters, and visualizing the results on a map. Each step is described in detail, along with the rationale behind the choices made [2-3].

K-means clustering is a popular machine learning approach for finding the best location of retailers. It can help identify clusters of data points that are similar to each other based on their features or attributes (see Figure 2). We can use k-means for this task:

1. **Data Preparation:** Collect relevant data about potential retail locations, such as geographic coordinates, demographics, economic factors, competition, and customer preferences.
2. **Feature Selection:** Choose the most important features that are relevant to the retail location, such as population density, income levels, proximity to transportation or other amenities, etc.
3. **Data Normalization:** Normalize the feature values to ensure that they are on a similar scale. This step is important to give equal weightage to each feature during clustering.
4. **Determine the Number of Clusters:** Decide how many clusters you want to create. This can be based on domain knowledge or by using techniques like the elbow method or silhouette score.
5. **K-means Algorithm:** Apply the k-means clustering algorithm to the prepared and normalized data, using the selected features. The algorithm assigns each data point to the nearest cluster centroid based on the Euclidean distance.
6. **Evaluate Cluster Results:** Analyze the resulting clusters and evaluate their characteristics. This step helps to understand the unique attributes of each cluster.
7. **Interpreting and Visualizing Results:** Visualize the clusters on a map to identify the optimal locations for retailers. You can use different colors or symbols to represent each cluster and examine their geographic distribution.

8. Refinement and Iteration: Iterate and refine the clustering process by adjusting the number of clusters or including additional relevant features. This step helps to improve the accuracy of the results [11-12].

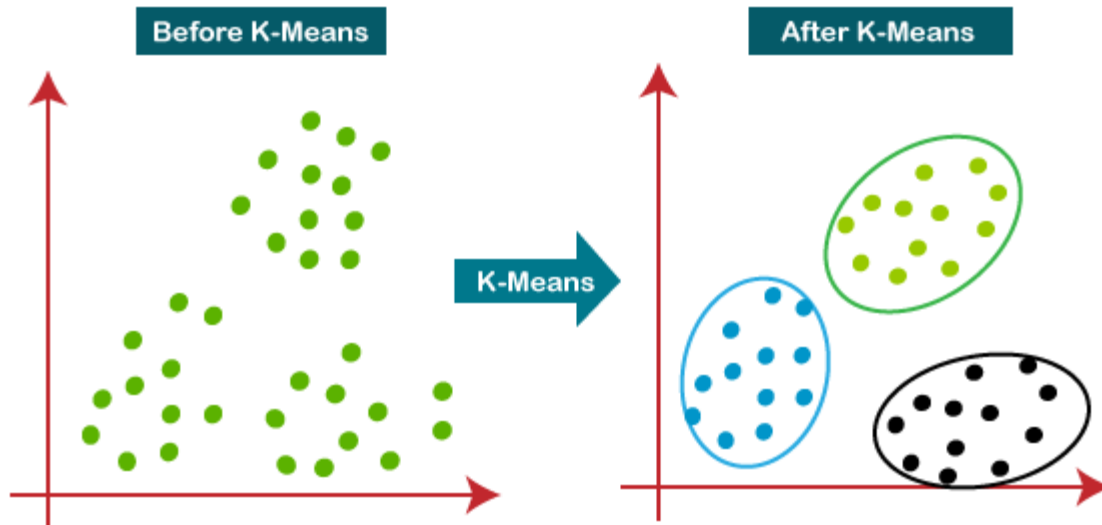


Figure 2: Applying resiliency in predicting demand for the automotive supply chain.

The k-means algorithm is a widely used clustering algorithm in machine learning and data analysis. It is an iterative algorithm that partitions a dataset into k distinct clusters based on their similarity. Each cluster is represented by its centroid, which is the arithmetic mean of all the data points assigned to that cluster. The algorithm aims to minimize the within-cluster variance, also known as the Sum of Squared Errors (SSE) [13-14].

An overview of the k-means algorithm is as follow:

1. Initialization: Specify the number of clusters, k , and randomly initialize k centroids in the feature space. Alternatively, the initial centroids can be selected based on some predefined criteria.
2. Assignment Step: Assign each data point to the nearest centroid based on a distance metric, commonly Euclidean distance. The data points are assigned to the cluster with the minimum sum of squared distances.

3. Update Step: Recalculate the centroids of each cluster by taking the mean of all data points assigned to that cluster. This step ensures that the centroids are positioned at the center of the cluster.
4. Repeat Steps 2 and 3: Iterate the assignment and update steps until convergence criteria are met. Convergence is achieved when the centroids no longer change significantly or when a maximum number of iterations is reached.
5. Output: Once convergence is achieved, the final centroids represent the clusters, and the algorithm outputs the cluster assignments for each data point [11-14].

The k-means algorithm has several advantages. It is computationally efficient and easy to implement. It works well with large datasets and can handle a variety of data types. However, it also has some limitations. Since the algorithm requires specifying the number of clusters in advance, determining the optimal value of k can be challenging. The algorithm is sensitive to the initial centroid positions, which may result in different solutions. It is also constrained to finding clusters with a convex shape.

To mitigate some limitations, variations of the k-means algorithm have been developed, such as the k-means++ initialization method, which improves the initial centroid selection, and the use of elbow method or silhouette analysis to determine the optimal number of clusters.

Overall, the k-means algorithm serves as a fundamental tool for cluster analysis and has been widely applied in various domains, including customer segmentation, image compression, anomaly detection, and retail location analysis, as mentioned in the literature review.

4. Results and discussion

This section presents the numerical results and findings obtained from applying the k-means algorithm to the dataset of potential retail locations. The cluster characteristics and distribution are discussed, highlighting the similarities and differences between the identified clusters. Visual representations, such as maps and charts, are utilized to illustrate the results effectively.

The location of demands for products are determined by experts as follow (see Table 1 and Figure 3), we apply k-means for defining centroid and best location for retailers. Python code for determining best location for retailers (see Table 2) and Elbow method for defining optimal K to

clustering are determined in Figure 4. The best locations for retailers by K-means are defined in Figure 5. Based on elbow method suitable optimal K is three clustering with inertia 108289.16.

Table 1: Location of demands for products

Location	X	Y	Demand
Location 1	12	17	300
Location 2	26	14	500
Location 3	18	16	300
Location 4	12	17	400
Location 5	37	36	900
Location 6	53	34	1000
Location 7	30	49	1047
Location 8	31	55	1184
Location 9	70	83	1321
Location 10	83	100	1458
Location 11	93	71	1595

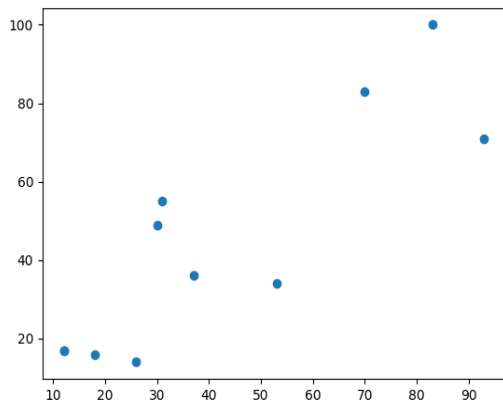


Figure 3: Location of demands for products.

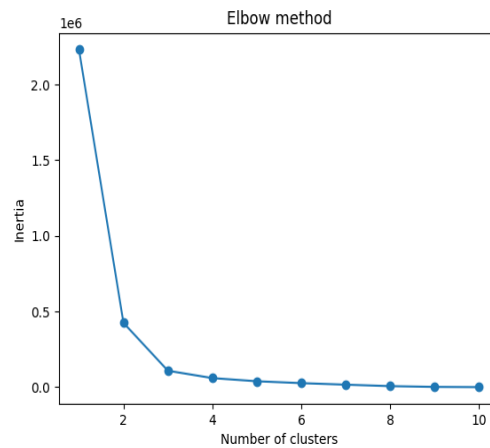


Figure 4: Elbow method for defining optimal K to clustering.

Table 2: Python code for determining best location for retailers

```
import matplotlib.pyplot as plt

x = [12,26,18,12,37,53,30,31,70,83,93]
y = [17,14,16,17,36,34,49,55,83,100,71]
z = [300,500,300,400,900,1000,1047,1184,1321,1458,1595]

plt.scatter(x, y)
plt.show()
```

```

from sklearn.cluster import KMeans

data = list(zip(x,y,z))
print (data)

inertias = []

for i in range(1,11):
    kmeans = KMeans(n_clusters=i)
    kmeans.fit(data)
    inertias.append(kmeans.inertia_)

plt.plot(range(1,11), inertias, marker='o')
plt.title('Elbow method')
plt.xlabel('Number of clusters')
plt.ylabel('Inertia')
plt.show()

kmeans = KMeans(n_clusters=3)
kmeans.fit(data)

Print ("centers:",kmeans.cluster_centers_)
Print ("inertia:",kmeans.inertia_)

plt.scatter(x, y, c=kmeans.labels_)
plt.scatter(kmeans.cluster_centers_[0, 0],\
            kmeans.cluster_centers_[0, 1], \
            s=100, c='red')
plt.show()

```

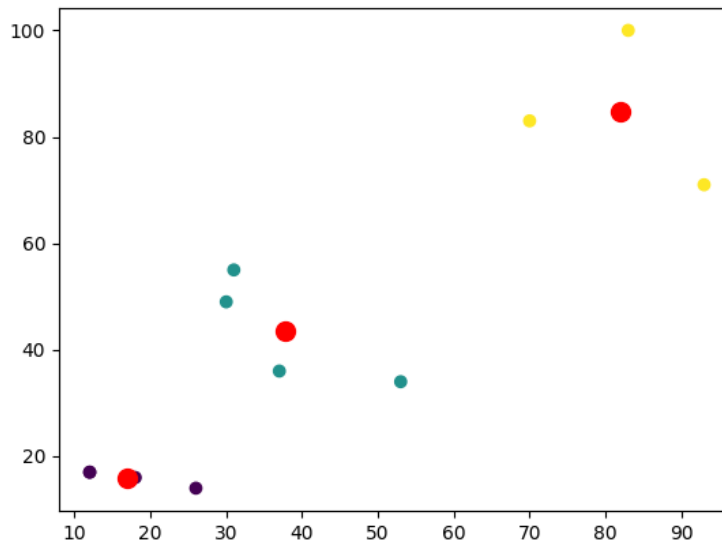


Figure 5: Best location of retailers by K-means.

Table 3: Results of best location of retailers by k-means.

Location (Centroid)	X*	Y*	Covering Demand	Inertia method
Retailer 1	38	44	1500	
Retailer 2	17	16	4131	108289.16
Retailer 3	82	85	4374	

The table 3 shows the location (centroid), X and Y coordinates, and covering demand for three retailers. The location (centroid) is obtained by K-means algorithms for the retailer's customers. The X and Y coordinates are the retailer's physical location on a map. The covering demand is the amount of demand that the retailer can meet, given its location and capacity.

The table 3 shows that Retailer 1 is located at (38, 44), Retailer 2 is located at (17, 16), and Retailer 3 is located at (82, 85). Retailer 1 has a covering demand of 1500, Retailer 2 has a covering demand of 4131, and Retailer 3 has a covering demand of 4374. It can be used to understand the geographic distribution of demand for the retailers' products. It can also be used to identify areas where there is unmet demand. For example, the table shows that Retailer 1 has a relatively low covering demand, compared to Retailers 2 and 3. This suggests that there may be unmet demand in the area around Retailer 1. It can also be used to plan for new retail locations. For example, if the retailers are planning to open a new store, they could use the table to identify areas where there is unmet demand. They could also use the table to identify areas where they would be able to reach a large number of customers.

Overall, it provides valuable information about the geographic distribution of demand for the retailers' products. This information can be used to make a variety of business decisions, such as where to open new stores and how to allocate resources.

5. Conclusion

In conclusion, this paper demonstrates the effectiveness of the machine learning approach using the k-means clustering algorithm for identifying optimal retail locations. The study highlights the importance of considering various factors related to geography, demographics, and economics in the process. The numerical results indicate the ability of k-means clustering to uncover meaningful clusters, providing valuable insights and informing decision-making for retailers. The approach presented in this paper holds promise for future research in the field of retail location analysis.

K-means clustering is a machine learning algorithm that can be used to identify the best location for retailers. The algorithm works by grouping customers into clusters based on their location and

other factors, such as demographics or spending habits. Once the customers are clustered, the retailers can identify the areas where there is the highest concentration of customers and open stores in those locations.

Using k-means clustering to identify the best location for retailers has a number of advantages. First, the algorithm is able to take into account a variety of factors, such as customer location, demographics, and spending habits. This allows the retailers to identify the areas where they are most likely to be successful.

Second, k-means clustering is a relatively simple algorithm to implement. This makes it a good option for retailers of all sizes, regardless of their technical expertise.

Third, k-means clustering is a dynamic algorithm. This means that the clusters can be updated as new data becomes available. This allows the retailers to adapt their location strategy over time as the needs of their customer's change.

Of course, there are also some limitations to using k-means clustering to identify the best location for retailers. First, the algorithm is sensitive to the choice of the number of clusters. If the retailers choose too few clusters, the clusters will be too large and the algorithm will not be able to identify the areas where there is the highest concentration of customers. If the retailers choose too many clusters, the clusters will be too small and the algorithm will not be able to identify any meaningful patterns.

Second, k-means clustering is not able to consider all of the factors that may affect the success of a retail location. For example, the algorithm cannot take into account the presence of competitors or the quality of the surrounding infrastructure.

Overall, k-means clustering is a valuable tool that can be used to identify the best location for retailers. The algorithm is able to take into account a variety of factors and is relatively simple to implement. However, it is important to be aware of the limitations of the algorithm and to consider other factors, such as the presence of competitors and the quality of the surrounding infrastructure, when making a decision about where to open a new store.

Finally, there are some additional thoughts on the use of k-means clustering for identifying the best location for retailers:

- K-means clustering can be used in conjunction with other data analysis techniques, such as regression analysis, to get a more complete picture of the factors that affect the success of a retail location.
- K-means clustering can be used to identify potential new markets or to expand into existing markets.
- K-means clustering can be used to optimize the location of existing stores.
- K-means clustering can be used to identify areas where there is unmet demand for retail products or services.

Overall, k-means clustering is a powerful tool that can be used to improve the profitability of retail businesses.

References:

- [1] Hosseini Rad, R., Baniasadi, S., Yousefi, P., Morabbi Heravi, H., Shaban Al-Ani, M., & Asghari Ilani, M. (2022). Presented a framework of computational modeling to identify the patient admission scheduling problem in the healthcare system. *Journal of Healthcare Engineering*, 2022.
- [2] Baniasadi, S., Rostami, O., Martín, D., & Kaveh, M. (2022). A novel deep supervised learning-based approach for intrusion detection in IoT systems. *Sensors*, 22(12), 4459.
- [3] Shoushtari, F., Ghafourian, E., & Talebi, M. (2021). Improving Performance of Supply Chain by Applying Artificial Intelligence. *International journal of industrial engineering and operational research*, 3(1), 14-23.
- [4] Ghasemi, S. M. (2022). Gene Transcription Modeling at the Cell Population Level (Doctoral dissertation).
- [5] Mirhajianmoghadam, H., & Akbarzadeh-T, M. R. (2022). Predictive hierarchical harmonic emotional neuro-cognitive control of nonlinear systems. *Engineering Applications of Artificial Intelligence*, 111, 104781.
- [6] Lotfi, R., Gholamrezaei, A., Kadłubek, M., Afshar, M., Ali, S. S., & Kheiri, K. (2022). A robust and resilience machine learning for forecasting agri-food production. *Scientific Reports*, 12(1), 21787.
- [7] Lotfi, R., Kheiri, K., Sadeghi, A., & Babae Tirkolae, E. (2022). An extended robust mathematical model to project the course of COVID-19 epidemic in Iran. *Annals of Operations Research*, 1-25.
- [8] Ikotun, A. M., Ezugwu, A. E., Abualigah, L., Abuhaija, B., & Heming, J. (2022). K-means clustering algorithms: A comprehensive review, variants analysis, and advances in the era of big data. *Information Sciences*.

- [9] Sarbaini, S., Saputri, W., & Muttakin, F. (2022). Cluster Analysis Menggunakan Algoritma Fuzzy K-Means Untuk Tingkat Pengangguran Di Provinsi Riau. *Jurnal Teknologi Dan Manajemen Industri Terapan*, 1(2), 78-84.
- [10] Abernathy, A., & Celebi, M. E. (2022). The incremental online k-means clustering algorithm and its application to color quantization. *Expert Systems with Applications*, 207, 117927.
- [11] Karami, D. (2022). Supply Chain Network Design Using Particle Swarm Optimization (PSO) Algorithm. *International journal of industrial engineering and operational research*, 4(1), 1-8.
- [12] Ssempijja, M. N., Namango, S., Ochola, J., & Mubiru, P. K. (2021). Application of Markov chains in manufacturing systems: A review. *International journal of industrial engineering and operational research*, 3(1), 1-13.
- [13] Minh, H. L., Sang-To, T., Wahab, M. A., & Cuong-Le, T. (2022). A new metaheuristic optimization based on K-means clustering algorithm and its application to structural damage identification. *Knowledge-Based Systems*, 251, 109189.
- [14] Nie, F., Li, Z., Wang, R., & Li, X. (2022). An effective and efficient algorithm for K-means clustering with new formulation. *IEEE Transactions on Knowledge and Data Engineering*, 35(4), 3433-3443.