

Optimizing Big Data Analysis with Machine Learning: Clustering, Visualization, and Insight Extraction

Rajesh Govind Talekar^{1,3}, Avinash Vasantrao Khambayat²

¹Research Scholar, Department of Mathematics, Sandip University, Nashik, 422213, Maharashtra, India

²Professor, Department of Mathematics, Sandip University, Nashik, 422213, Maharashtra, India

³Assistant Professor in Mathematics, Department of Applied Science and Humanities, Pimpri Chinchwad College of Engineering, Nigdi, Maharashtra, India

Email: ¹rajesh.talekar@pccoepune.org,

²avinash.khambayat@sandipuniversity.edu.in

ARTICLE INFO

ABSTRACT

Received: 28 Dec 2024

Revised: 19 Feb 2025

Accepted: 28 Feb 2025

Big Data analysis is understanding valuable insights within large datasets. The current analysis, which looks into how machine learning optimizes Big Data analysis concerning clustering, visualization, and insight extraction, explores this area. We show how clustering techniques work in several areas with the application on publicly available Mall Customer and Iris datasets. For this purpose, clustering methodologies such as K-Means, Hierarchical Clustering, and DBSCAN may then be used to cluster the data and recognize the patterns. Forms of reducing the feature space and enhancing visualization are Principal Component Analysis (PCA), t-SNE, and dendrograms, which improve interpretability and give clear representations of complex patterns. The assessment on the ideal number of clusters is based on the Silhouette score. In addition to these, some classification methods such as Support Vector Machine (SVM), Logistic Regression, and K-Nearest Neighbour (K-NN) are used for classifying the data into several classes. The results show how clustering creates advantages in decision-making based on the economic and biological spheres. By discussing the scalability challenge and optimization techniques in processing large datasets, our future work is taken up. This is to show for machine learning-based clustering in getting necessary forward inputs to make this possible across different sectors.

Keywords: Machine Learning, Clustering Techniques, Big Data, Mall Customer Dataset, Classification, Visualization.

1. INTRODUCTION

The amount of information generated from a variety of sources such as social media, online transactions, IoT devices, healthcare systems, financial markets, and scientific research is unprecedented in the present era of digitization. This vast and complex set of data is commonly referred to as "Big Data", and this type of data necessitates efficient techniques for processing, storing, and analyzing the same. The ability to draw inferences from Big Data is a must-have for decision-making in many areas such as scientific research, marketing, healthcare, and finance. However, the extreme volume, velocity, and variety present major challenges for interpretability, scalability, and computational efficiency. Often, the effectiveness of traditional data processing methods gets thwarted since they rely heavily on organized datasets, preset schemas, and limited processing capabilities to manage big data. With the number and complexity of data increasing, traditional statistical techniques and rule-based approaches have become practically useless in spotting trends and making data-driven decisions. This has, therefore, led to a surge in the demand for machine-learning (ML) techniques, which offer automated, scalable, and adaptive solutions to big-data challenges. Clustering lets us bring together similar data points based on their inherent characteristics. By coupling this clustering with activity visualization and dimensionality reduction techniques, one can gain further insight into massive datasets.

1.1 Overview of Big Data Analysis: Importance and Challenges in Handling Large Datasets

In order for enterprises to handle and interpret huge and complicated datasets and derive valuable information, big data analysis is an essential field. Big Data analysis is becoming more and more significant in a number of fields [4]:

- **Business and Marketing:** Companies analyze consumer behavior, preferences, and purchasing patterns to develop targeted marketing strategies and optimize customer engagement.
- **Healthcare:** Big Data analytics helps in predicting disease outbreaks, personalizing treatments, and improving patient care.
- **Finance:** Financial institutions use predictive analytics for risk assessment, fraud detection, and portfolio optimization.
- **Scientific Research:** Large-scale datasets in physics, climate science, and genomics necessitate sophisticated analytical methods for testing hypotheses and identifying patterns.

Big Data analysis has a number of drawbacks despite its advantages.

- **Volume:** A massive spurt of data resonates with requisite computing frameworks and storage arrangements. Such volumes of data are mainly way too big to be archived in regular databases.
- **Velocity:** Fraud detection and stock market forecasting require real-time data processing. The higher order of data streaming demands higher-order processing techniques to handle the same.
- **Variety:** Data might be semi-structured (JSON, XML), unstructured (text, photos, videos), or organised (databases). It can be challenging to integrate these different types.
- **Veracity:** It is important to ensure the quality, consistency, and correctness of the data because improper data such as noisy and incomplete data can produce misleading insight.
- **Scalability:** Analytical models must scale well without performance degradation when data volume increases. This challenge is met through distributed processing frameworks like Hadoop and Spark and parallel computation.
- **Interpretability:** As the raw numerical output is often devoid of intuitive meaning, the explanation capability and visualization techniques would be needed to derive meaningful insights from complex data.

Machine learning techniques are essential for optimizing Big Data analysis in order to overcome these issues.

1.2 Role of Machine Learning in Optimization: How ML Techniques Enhance Big Data Insights

Unlike the traditional methods of big data analysis where human intervention was essential, now machine learning is providing scalable and automated and completely adaptive solutions for pattern recognition and data processing. ML algorithms become instrumental in finding hidden patterns, trends, and correlations in massive datasets, thus contributing to strategic plans and informed decision-making [5].

Automating Data Processing and Analysis: With large datasets, the magic happens within a black box. Here, man has no interaction and no work involved at all for the purposes of training. That is how the entire process of training happens. While pattern recognition and adaptation by ML models occur over an automatic course, traditional mechanisms require rules to be static and statistical assumptions predetermined making it obviously more stable and applicable to complex and all kinds of changing datasets.

Clustering for Pattern Discovery: In order to segment and categorize Big Data, clustering—among the most successful ML approaches for unsupervised learning—groups similar data points according to their attributes. Several popular clustering algorithms are as follows:

- **K-Means Clustering:** A centroid-based approach that partitions data into k clusters, minimizing intra-cluster variance.
- **Hierarchical Clustering:** creates a dendrogram, or tree-like structure, to show nested clusters; this is helpful for figuring out hierarchical links in data.
- **DBSCAN (Density-Based Spatial Clustering of Applications with Noise):** effectively manages noise and finds clusters of any shape by identifying them based on density.

Dimensionality Reduction for Visualization and Interpretation: Big Data often possesses high-dimensional features that complicate interpretation and visualization. Dimensionality reduction techniques like Principal Component

Analysis (PCA) and t-SNE help reduce the data complexity while keeping the essential patterns. These techniques help map out the high-dimensional datasets in a way that eases human understanding and decision making.

Enhancing Insight Extraction through Advanced Visualization: It is said that complex digits can be translated into understandable figures through these visualization techniques. Visualization tools motivated by machine learning enable researchers and companies to:

- Detect trends and anomalies in data.
- Understand customer segmentation and preferences.
- Interpret biological and scientific data patterns.
- Communicate findings effectively to stakeholders.

Optimization for Scalability and Performance: Effective computational methods are necessary for big data analysis in order to manage massive datasets without consuming excessive amounts of memory or processing time. ML-based improvements consist of:

- **Parallel Computing:** Running ML algorithms on distributed systems like Apache Spark to accelerate processing.
- **GPU Acceleration:** Leveraging GPUs for deep learning-based clustering and high-performance computations.
- **Feature Selection and Engineering:** Reducing dimensionality by selecting the most relevant features for clustering and insight extraction.

The fusion of ML into big data research is taking the art of handling and evaluating very large datasets into a perfect system for extraction of real information. Although visualization techniques such as PCA or t-SNE provide intuitive interpretation of data, actual pattern detection and segmentation is represented with clustering techniques such as hierarchical clustering, K-means, and DBSCAN. While some of the challenges faced are due to scalability, quality of data, and computational efficiency, ML-based approaches still remain the major optimizers for Big Data analysis across various fields. The study will illustrate the efficacy of machine-learning-based clustering in combination with visual techniques using the Iris and Mall Customer datasets. It is anticipated that some significant insights for consumer behavior analysis and biological classification will be derived through these techniques. Thus, the findings would contribute to advancing the much more extensive field of big data analytics by providing better methodologies for handling humongous and complex datasets. The study will therefore select the Iris dataset (biological categorization dataset) and the Mall Customer dataset (consumer behavior dataset) to optimize Big Data analysis through clustering, visualization, and insight extraction.

1.3 Objectives of the Study

- **Clustering Techniques:** For testing effectiveness of clustering techniques like Hierarchical Clustering, K-Means and DBSCAN in grouping of data and recognizing patterns within Mall Customers and Iris datasets, disclosing its applicability in both economic and biological domains.
- **Optimal Cluster Selection:** Determining the optimal number of clusters using the Silhouette score and other metrics so that truly accurate segmentation of the datasets can lead to insight extraction.
- **Classification and Visualization:** To understand the data in a classification framework using SVM, Logistic Regression, K-NN, among classification techniques; conversely, with that understanding, data visualization becomes an interpretative tool, involving methods like PCA, t-SNE, and dendrograms.
- **Comparative Analysis and Identifying Drawbacks:** Comparativeness analysis will be done involving the clustering and classification techniques, including their merits and demerits, to find possible scalability problems and optimization strategies for the efficient handling of large datasets.

1.4 Structure of the Paper

The structure of the paper is as follows: A review of machine learning-based Big Data analyses is presented in Section 2, wherein various clustering strategies, visualization methods, and insight extraction approaches are discussed. Section 3 describes the proposed methodology, consisting of data pre-processing, clustering approaches, methods for selecting the number of clusters, and classification techniques, while Section 4 discusses the results, including data visualization, clustering, and classification performance. Conclusions and further research avenues are presented in Section 5, thereby bringing the study to an end.

2. LITERATURE REVIEW

With Big Data multiplying exponentially, there has been an increased adoption of machine learning techniques for drawing insights from complex datasets. These approaches are broadly divided into those classification approaches that try to classify data into predefined classes for predictive analysis and clustering approaches that group data on the basis of inherent patterns for exploratory analysis. This section reviews the existing literature on these methodologies with reflections on applications, advantages, and disadvantages in Big Data analysis while investigating new research gaps and opportunities.

Machine Learning Based Classification Approaches in Bigdata

[6] One of the papers called "Big Data Analysis Techniques" by Dhawas Pranali et al. makes a good analysis of various techniques with regard to preparation, mining, ML, and visualization. This research paper aims at introducing many fundamental techniques of preprocessing, which are very much useful for many noise, incomplete and high-dimensional data at the moment, since they are extremely important in performance enhancement. It also studies and investigates the applications of different algorithms of data mining and ML for pattern detection and predictive analytics. It addresses altogether the most sophisticated visualization tools that will make decision making much better and even more understandable. These problems are big in Big Data analytics, and it serves good insight into developing data-driven tactics across many disciplines.

Mallikarjuna Paramesha et al. [7] The improvement of Business Intelligence (BI) combining Big Data Analytics, AI, ML, IoT, and Blockchain was looked over in the report. The report speaks across industries how these technologies are supporting data-driven decision-making, process automation, and real-time analytics. It presents Blockchain for data security and transparency while classifying different Big Data methodologies, AI-based prediction models, and IoT data collection approaches. The research proposes optimizing BI while shedding light on data integration, scalability, and computational complexity challenges. This research systematically formulates a framework to adapt innovative technologies for business analytics.

Eid A et al. [8] focused on using Big Data, Cloud Computing, and AI for enhancing the decision-making and predictive analysis processes. The framework proposed in this manuscript relies on Hadoop and consisted of AI guided learning algorithms for effective data processing and predictive modeling. Strategic planning and performance optimization was attempted with the Sport AI Model (SAIM), and real gains concerning outcome predictions in sport became available through comprehensive simulations. According to this report, sports analytics would now be based on real-time analysis of data and precise decision making by means of AI and big data, which indicates a paradigm shift. These results illustrate the essence with which sports intelligence could be improved through such AI-enabled Big Data frameworks.

Sadat Lavasani et al. [9] provided deeper insights into the developments and applications of Big Data Analytics (BDA), highlighting its growing significance in industries. This study offers a systematic review of Big Data platforms, tools, and techniques, classifying them according to their functions. The important features and substance of Big Data are also examined, with a focus on its applicability to process engineering. Its application in chemical industries, which shows its ability to optimize engineering processes, is a noteworthy area of study. By advocating for the use of BDA, this study provides a workable vision for researchers and industries, stimulating more research and development in data-driven decision-making.

Wengang Zhang et al. [10] conducted a thorough analysis of the uses of optimization techniques, DL, and ML in geoscience and geoengineering. The work demonstrates how ML approaches, through the analysis of large and intricate geospatial datasets, enhance subsurface modeling, geological hazard prediction, and resource exploration. In order to improve computational accuracy and efficiency in geoscientific models, optimization procedures are incorporated. The study also highlights the necessity of hybrid models and interdisciplinary approaches in order to solve data uncertainty, scalability, and interpretability challenges. Future developments will focus on AI-driven geoscience applications, which will stimulate creativity in geological risk assessment and sustainable resource management.

Chinta Swetha et al. [11] concentrated on how to improve model performance and data processing efficiency by integrating machine learning approaches into Big Data contexts. The volume, velocity, and diversity of Big Data can overwhelm traditional analytical techniques, making optimized frameworks necessary for successful deployment.

Studies already conducted show that in order to optimize prediction accuracy and decision-making, organized methodologies are required for data pretreatment, algorithm selection, and performance evaluation. Nevertheless, there are still gaps in the standardization of these approaches for many applications. Case studies have illustrated the usefulness of ML in big data, highlighting its potential for better insights and developments across the industry.

According to Torre-Bastida et al. [12], bio-inspired computation has been investigated in terms of its role in Big Data fusion, storage, processing, learning, and visualization functions that enable it to handle sophisticated and large-scale datasets. This article reviews modern nature-inspired algorithms that help data processing efficiency and model flexibility, including genetic algorithms, swarm intelligence, and neural networks. Other issues tackled by the study include real-time learning, scalability, and data heterogeneity, with suggestions for bio-inspired solutions to improve these aspects of computability. The research provides insight into how biologically inspired methods in different sectors can maximize Big Data analytics while carving future paths.

Boppiniti Sai Teja and others [13] have analyzed the relationship between Big Data and ML and suggested techniques to process and analyze big data efficiently. They include optimization techniques, such as parallel computing, distributed systems, and feature engineering, and discuss a few major issues such as data scalability, computational efficiency, and real-time processing. They focus on the importance of applying appropriate algorithms of machine learning for various datasets in a manner that is accurate and efficient. This work proves to be useful to Big Data analytics academicians and practitioners when it comes to computational bottlenecks and hence promotes the development of scalable and very efficient solutions that are based on data.

Verbeeck N et al. [14] study unsupervised machine learning methods for exploratory data analysis in imaging mass spectrometry (IMS), including how it can be used for pattern recognition and feature extraction. Important clustering and dimensionality reduction methods such as K-Means, hierarchical clustering, and t-SNE help to extract important biological information from complex spectrum datasets in this study. This research will show how machine learning can enhance the interpretation of IMS data through considerations for data preprocessing, noise reduction, and computational efficiency. This improvement results in broadening the application of mass spectrometry in biomedical science and enhancing biomolecular analysis using unsupervised learning.

Zhou et al. [15] explored the prospects and challenges of applying machine learning to Big Data, focusing on the management of massive, high-dimensional datasets. The three major hurdles in Big Data environments: scalability, computational efficiency, and data heterogeneity; form the core contents of the paper. Maximizing model accuracy and data processing to be considered is established through advancements in deep learning, feature selection, and distributed computing frameworks. Moreover, discussions concerning ethical matters such as data bias and privacy have been included. This present paper propagates basic knowledge about the role of machine learning in analytics of Big Data while addressing the major challenges that could hinder its effectiveness in a variety of applications.

Praful Bharadiya [16] analyzed both the way Machine Learning (ML) ultimately changes Business Intelligence (BI)—essentially, how it influences predictive analytics, automation, and data-driven decision-making. It was stated that ML generally improves real-time data processing, anomaly detection, and consumer behavior analysis to help organizations better optimize their plans and operations. BI applications are described in the context of common ML approaches: supervised and unsupervised learning, DL, and reinforcement learning. Issues like data privacy, scalability, and computing efficiency are also considered throughout the study. The report describes continuous development in machine learning and its increasing importance to enhance the BI framework for competitiveness and company performance.

Clustering Approaches in Bigdata

Rehman A et al. [17] studied Big Data Analytics (BDA) usage in healthcare, focusing on emerging trends, challenges, and possible improvements in healthcare. For making data-driven decisions, the research explored BDA in improving prognostics, patient monitoring, and personalized treatment through machine learning and artificial intelligence. Data privacy, security, interoperability, and scalability are the major roadblocks towards seamless integration in systems of healthcare. Authors emphasize that these constraints can be resolved with advanced data processing methods and regulatory frameworks. This study prepares the foundation for subsequent advances of medical informatics by showing the revolutionary effects of BDA in health.

Akash Tayal et al. [18] have incorporated Big Data Analytics, Machine Learning, and Hybrid Meta-heuristics to optimize layouts for sustainable manufacturing. A four-stage methodology has been devised in this paper to deal with Sustainable Facility Layout Problem (SFLP) under uncertainty of demand. The article takes up the issue of the dimensionality curse by using K-Means clustering, Data Envelopment Analysis (DEA), and mathematical modeling for ranking and forecasting effective facility layouts. Energy consumption and CO₂ emissions are estimated in the study to ensure sustainability. A fictitious data-based case study showcases the relevance of the model and illustrates how it can be adjusted to various social, political, and economic conditions in facility construction.

Rahat Iqbal et al. [19] focuses on advancements and applications concerning the integration of Big Data Analytics (BDA) and Computational Intelligence (CI) within Cyber-Physical Systems (CPS). The study mainly analyzes how BDA can leverage predictive analytics, anomaly detection, and real-time processing in CPS environments. Furthermore, it cites several possible computational techniques-including deep learning, swarm intelligence, and machine learning-as enabling automation, security, and efficiency. Also taken into consideration are some problems, such as the heterogeneity of data, a lack of scalability, and the presence of cyber risks. The study further highlighted the possible radical benefits that BDA and CI could offer in optimizing CPS across sectors, such as healthcare, smart cities, and the manufacturing industry, as the major application showcase.

M. Feng et al. [20], realizing the importance of Big Data Analytics (BDA) in crime analysis and prediction, have applied data mining, deep learning, and exploratory data analysis on crime datasets from Philadelphia, Chicago, and San Francisco. Law enforcement decision-making should be aided in some way through statistical analysis and visualization for highlighting important crime trends. The study established that in the face of three years of optimal training, the Prophet model and Keras stateful LSTM perform better than ordinary neural networks. In helping police departments track crime, predict incidents, and allocate resources, these findings edify the idea of BDA as a tool to improve public safety strategies.

Table 1: Comparative Analysis of Literature Review

Author Name & Ref. No.	Methodology Used	Datasets Used	Advantages	Results
Dhawas Pranali et al. [6]	Data preprocessing, data mining, ML, visualization	Various market intelligence datasets	Comprehensive data processing & visualization techniques	Enhanced ML model accuracy & interpretability
Mallikarjuna Paramesha et al. [7]	AI, ML, IoT, Blockchain integration in BI	Business intelligence datasets	Enhanced BI through AI, ML & Blockchain	Improved real-time data processing in BI
Eid et al. [8]	Cloud computing & AI-based sports prediction model	Sports data	Improved sports outcome predictions	Higher accuracy in sports analytics
Sadat Lavasani et al. [9]	Big Data analytics in process engineering	Process engineering datasets	Optimized process engineering insights	Better process optimization & sustainability
Bharadiya Jasmin Praful [10]	ML techniques for BI transformation	Business intelligence datasets	Enhanced business analytics decision-making	More effective BI-driven decision-making
Rehman et al. [11]	Big Data analytics for healthcare enhancement	Healthcare datasets	Improved disease prediction & healthcare analytics	Better healthcare insights and patient monitoring
Wengang Zhang et al. [12]	ML, DL, and optimization in geoscience	Geoscience datasets	Better geological predictions & risk analysis	Enhanced modeling accuracy in geoenvironment

Chinta Swetha [13]	ML algorithms in Big Data analytics framework	Predictive insights datasets	Enhanced predictive modeling for Big Data	Improved predictive insights
Torre-Bastida et al. [14]	Bio-inspired computation for Big Data fusion	Bioinformatics & computational datasets	Efficient data fusion and real-time learning	Faster, scalable Big Data processing
Boppiniti Sai Teja [15]	Big Data & ML strategies for efficient processing	Large-scale datasets	Optimized data processing for large datasets	Optimized ML applications for Big Data
Akash Tayal et al. [16]	Big Data, ML, meta-heuristic, DEA for sustainable manufacturing	Manufacturing layout data	Sustainable and energy-efficient layout modeling	Efficient sustainable layout design
Verbeeck et al. [17]	Unsupervised ML for imaging mass spectrometry	Imaging mass spectrometry data	Advanced pattern discovery in mass spectrometry	More accurate biomolecular analysis
Rahat Iqbal et al. [18]	BDA & computational intelligence for cyber-physical systems	Cyber-physical system datasets	Optimized automation & anomaly detection	Enhanced real-time analytics & system efficiency
M. Feng et al. [19]	Big Data mining for crime trend forecasting	Crime data from multiple cities	Crime trend analysis and real-time forecasting	Increased crime trend prediction accuracy
Lina Zhou et al. [20]	ML applications & challenges in Big Data	General Big Data applications	Addressing scalability and efficiency in Big Data ML	Addressed scalability, bias, & efficiency issues

This segment lays down a review of the research on the detection of artificial videos and reviews indeed with machine learning models, deep learning models, significant natural language processing approaches, and generation processes for datasets.

3. METHODOLOGY

The study methodology is specifically designed for Big Data analysis through clustering and classification followed by visualization techniques in the most efficient manner. The study draws upon the publicly available data sets such as the Mall Customer data set and Iris data set from two different domains. Then the data is preprocessed to handle missing values, normalize the features, and ensure their suitability for analysis. Optimization techniques in particular for huge data sets would increase the computational scalability and efficiency of the processing. Clustering techniques would, therefore, help in gathering descriptions of data into meaningful clusters, indicating the inherent patterns and structures. Next, the classification techniques put the data into predetermined classes and allow a predictive analysis. Visualization methods help interpret and make clear representations out of the complex patterns. The results are finally evaluated using measures like accuracy, Silhouette score, and optimal number of clusters, thus providing a comparative assessment of the effectiveness of the techniques used. This methodology assures a comprehensive inquiry into Big Data analysis while also addressing issues of scaling and optimization.

Dataset Description

Iris Dataset [21]: The Iris dataset contains a unique household name in statistical classification or machine learning. It contains 150 iris flower samples, which can be classified into three categories: Iris setosa, Iris versicolor, and Iris virginica. Every sample can be described using four numerical characteristics, namely sepal length, sepal width, petal length, and petal width (in centimeters).

Mall Customer Dataset: The Mall consumer dataset is a very useful dataset for the fields of marketing and business analytics. After all, it is widely used for consumer segmentation and behavioral analysis. Such a dataset contains five variables per record containing the information on 200 mall patrons. The dataset is frequently used to group clients according to their spending patterns, income levels, and age demographics using clustering techniques. Table 2 shows the Mall Customer dataset parameters.

Table 2: Mall Customer Dataset Features and Description.

Features	Description
Customer ID	A unique identifier for each customer.
Gender	The gender of the customer (Male/Female).
Age	The customer's age (in years).
Annual Income	The customer's annual income in thousands of dollars
Spending Score (1-100)	a figure that is given to consumers according to their shopping habits and shows how much they spend at the mall.

3.1 Clustering in Bigdata

3.1.1 Load Dataset

The clustering methodology begins by providing the dataset as input to the model. Figure 1 is an example of the Mall Customer Dataset.

	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
0	Male	19	15	39
1	Male	21	15	81
2	Female	20	16	6
3	Female	23	16	77
4	Female	31	17	40

Fig 1: Mall Customer Dataset Sample

3.1.2 Feature Visualization:

After providing the dataset as input, the next step is to visualize the relationships between features. Feature visualization techniques, such as scatter plots, pair plots, and correlation matrices, are employed to explore the distribution, patterns, and interactions within the data. Figure 2 illustrated Scatter plot of Annual Income (in k\$) vs Spending Score (1-100) from the Mall Customer Dataset. It shows how customers are categorized according to their spending patterns and income. Customer segmentation analysis may benefit from the distribution's indication of discrete clusters that show diverse consumer purchasing patterns across income groups.



Fig. 2: Scatter Plot

3.1.3 Dimensionality Reduction:

Dimensionality reduction methods like PCA and t-SNE are also applied to project high-dimensional data into lower-dimensional spaces, enabling clearer visualization of clusters and trends. This step aids in understanding the underlying structure of the data, identifying potential outliers, and guiding the selection of appropriate clustering and classification techniques. Figure 3 displays the data distribution according to gender (male vs. female) using a 2D PCA projection of features. PCA helps identify patterns by reducing the dimensionality of the dataset while maintaining significant volatility. The points' natural group structures are suggested by their clustering, which facilitates customer segmentation and analysis.

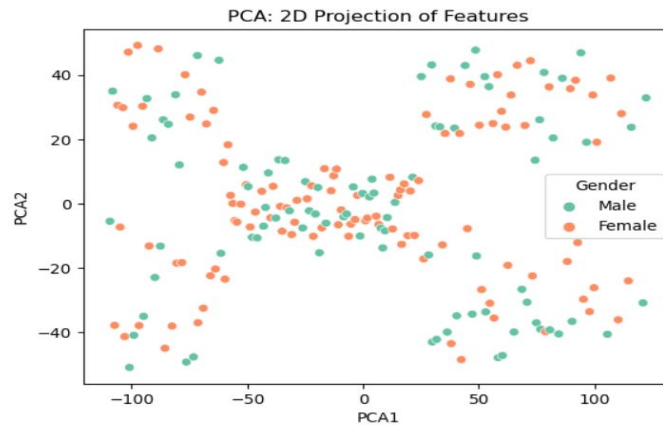


Fig. 3: Feature Projection Using PCA

In Figure 4, the 3D scatter plot shows Age, Annual Income (in k\$), and Spending Score (1-100) by Gender from the Mall Customer Dataset. Male and female consumers are represented by different colors, which makes it easier to see how much money people of different ages and income levels spend. The distribution points to possible customer segmentation clusters by indicating different spending tendencies. Understanding customer behavior for focused marketing tactics is made easier with the help of this visualization.

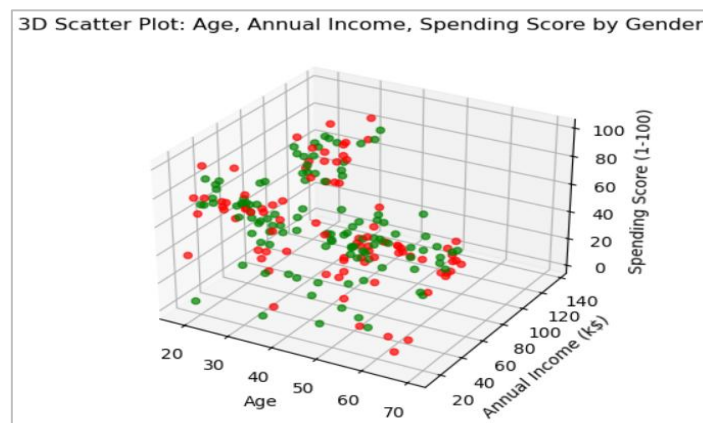


Fig.4: 3D Scatter Plot of MALL Dataset

3.1.4 Apply Clustering Technique:

The next step involves applying clustering techniques to group the data into meaningful clusters. Clustering techniques, such as Hierarchical Clustering, K-Means Clustering, and DBSCAN, reveal patterns in structures of the dataset. Each is specifically chosen depending on how well it fits the characteristics of the data and what the user wants to achieve with the implementation. To make interpretation easier, dendrograms (for hierarchical clustering), cluster plots, and 2D/3D scatter plots (using dimensionality reduction techniques such as PCA or t-SNE) are produced. The visualizations help to see the distribution of clusters, assess the quality of the results, and get insight into the underlying data patterns. This step is essential in unveiling hidden relationships and preparing the data for further analyses.

a. Hierarchical Clustering

Hierarchical clustering refers to the tree method in clustering, which forms a clustered structure in layers without necessarily specifying the number of clusters (k) needed. The two approaches of hierarchical clustering are agglomerative (bottom-up) and divisive (top-down). In agglomerative clustering, the clustering starts with each point as a cluster and combines them based on their similarity, while in the divisive approach, clustering starts with a single cluster and breaks it down in an iterative way. Hierarchical clustering presents the entire clustering process by a dendrogram that shows relationships between data points. The structure of the dendrogram allows for exploratory data analyses and is therefore especially useful for understanding hierarchical relationships among species since the K-Means necessitates a known fixed number of clusters. Versicolor and Virginica remain close, showing overlapping feature spaces, while Iris Setosa tends to form its branch due to its dissimilar attributes. For that matter, a primary benefit of hierarchical clustering is that it allows the visualization of the whole formation of cluster structures without making assumptions about the number of clusters. But it turns ineffective for huge datasets due to its limitation in computational complexity ($O(n^2)$). Furthermore, it is less flexible than K-Means in handling high-dimensional data and sensitive to noise. Nevertheless, it holds ample power when being applied to biological and exploratory data.

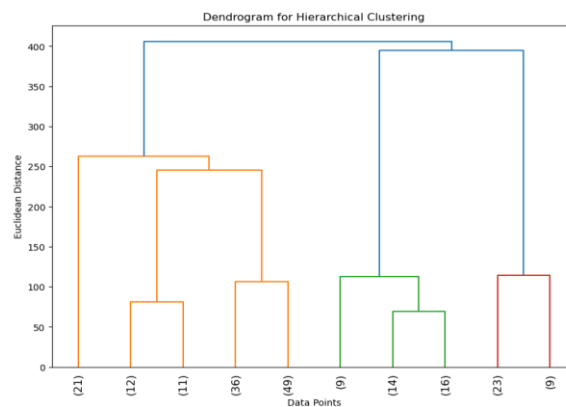


Fig. 5: Dendrogram for Hierarchical Clustering

Algorithm: Agglomerative Hierarchical Clustering

Input: Dataset, Linkage method (e.g., Single, Complete, Average)

Output: Hierarchical cluster tree (Dendrogram)

Steps:

- Start with each data point as an individual cluster.
- Compute pairwise distances between all clusters using the chosen linkage method:
 - Single Linkage: Minimum distance between two clusters.
 - Complete Linkage: Maximum distance between two clusters.
 - Average Linkage: Mean distance between all points in two clusters.
- Merge the two closest clusters based on the computed distances.
- Repeat until all points are merged into a single cluster.
- Construct a dendrogram to visualize hierarchical relationships.
- Cut the dendrogram at the desired level to determine the final clusters.

b. K-Means Clustering

K Means clustering is a common unsupervised method of machine learning which helps to divide a particular dataset in k number of groups based on the similarity of features within the groups. The process is an iterative procedure that uses the minimum distance between a data point and its closest centroid, followed by a move to update the centroid based on the mean of the associated data points. Mostly effective in identifying true structure when the dataset clusters are spherical and far apart from each other.

The study of K-Means found and important patterns among the data of Mall Customer and Iris datasets. Based on sepal and petal lengths and widths according to K-Means differentiation between the three flower species (Setosa,

Versicolor, and Virginica) sensorily interpreted. The algorithm shares three breeding hybrid plants into groups concerning a natural species category when $k=3$. K-Means will segregate the Mall Customer database into the customer's income-generating bracket of funds and spending behavior patterns. By this means of setting $k=5$, the algorithm now categorizes people into otherwise distinct groups such as average consumers, low-income shoppers, and big spenders for their focused marketing campaigns. This K-Means has few drawbacks, as well; it is not robust for non-spherical clusters, is sensitive to outliers, and requires a fixed number of clusters (k). But because of its ease of use and computational speed, it is a popular clustering method in business intelligence and data analytics.

Algorithm: K-Means Clustering

Input: Dataset X , Number of clusters k Maximum iterations

Output: Cluster assignments and centroids

Steps:

- Initialize k cluster centroids randomly from the dataset.
- Repeat until convergence or maximum iterations:
- Assign each data point to the nearest centroid using Euclidean distance:

$$d(x_i, c_j) = \sqrt{\sum (x_i - c_j)^2}$$

- Update each centroid as the mean of all assigned points:

$$c_j = \frac{1}{|C_j|} \sum_{x_i \in C_j} x_i$$

- Check for convergence (i.e., centroids do not change significantly).
- Return final cluster assignments and centroids.
- Calculate Time Complexity: $O(n \cdot k \cdot i)$

Where, n is the number of data points, k is the number of clusters, and i is the number of iterations.

i. Find Best Number of Cluster using Elbow Method

A well-liked method for figuring out the ideal number of clusters in clustering algorithms like K-Means is the Elbow Method. It assists in determining the point at which the performance of the model is not appreciably enhanced by the addition of more clusters. Figure 6 shows the Optimal Value of K using Elbow Method.

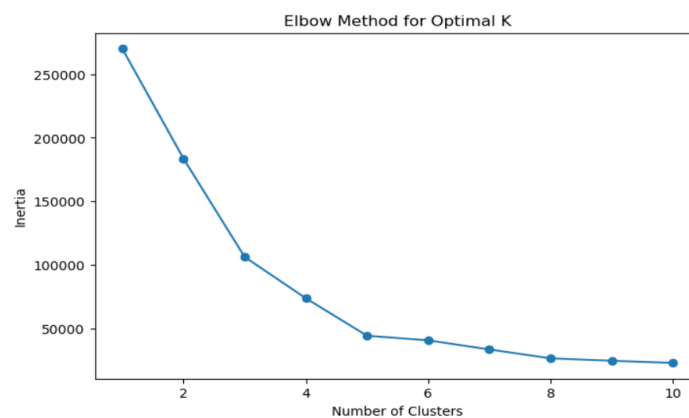


Fig.6: Optimal K Value using Elbow

ii. Find Best Number of Cluster using Silhouette Method

Another well-liked method for figuring out the ideal number of clusters in clustering algorithms like K-Means is the Silhouette Method. In contrast to the Elbow Method, which is dependent on Within-Cluster-Sum of Squares (WCSS), the Silhouette Method measures how well each data point has been assigned to its cluster with respect to other clusters and hence offers a more refined and valid way of calculating the optimal number of clusters.

Silhouette Score for K-Means: 0.47

Fig.7: Optimal K value using Silhouette Score

c. DBSCAN (Density-Based Clustering)

DBSCAN applies a very powerful unsupervised machine learning method for discovering any possible-shaped clusters and also for noise detection in data sets. It ignores the need for users to tell the number of clusters to create, but requires two parameters to make it work: minPts (minimum points required to form a cluster) and radius ϵ (epsilon - neighborhood radius) so as to classify points based on density connection compared with K-Means. Points that are within designated ϵ -radius of a dense zone contribute to building clusters while points that do not meet the criteria are denoted as noise.

This dataset is segmented in this study according to annual income and expenditure score thanks to the usage of DBSCAN on consumers' data. Its power to detect outliers in data and discover clusters of any shape gives it a notable advantage as compared to the K-Means algorithm bringing excellent results for customer segmentation. Customers may devote more or less of their earnings depending on the shopping situations. These different behaviors would result in non-circular clusters as assumed by K-Mean, where high-income individuals could have different spending habits but group together at different densities. DBSCAN filters out anomalies, including erratic purchasing patterns, and assists in identifying distinct client segments. DBSCAN has certain limitations even though it works well for detecting non-linear patterns and eliminating noise. Performance is greatly impacted by the choice of ϵ and minPts , and the method may have trouble with clusters with different densities. Additionally, distance-based clustering may not work as well with high-dimensional data. Notwithstanding these drawbacks, DBSCAN is a great method for consumer segmentation in real-world settings with noisy data and unevenly shaped clusters.

Algorithm:

Input: Dataset X , Neighborhood radius ϵ , Minimum points minPts

Output: Cluster labels with noise points identified

Steps:

- Mark all data points as unvisited.
- For each unvisited point:
 - If the point has at least minPts neighbors within radius ϵ , mark it as a core point and start a new cluster.
 - If the point is a border point (fewer than minPts neighbors but within an existing cluster), assign it to that cluster.
 - If the point is an outlier (noise), mark it as noise.
- Expand the cluster recursively by adding all density-reachable points.
- Repeat until all points are assigned.
- Return final cluster assignments.
- Time Complexity: $O(n \log n)$ (using spatial indexing like KD-Trees)

Table 2: Comparison of Clustering Algorithms

Algorithm	Type	Advantages	Limitations
K-Means	Centroid-based	Fast, efficient, and works well with spherical clusters	Requires predefined k, sensitive to outliers
Hierarchical	Tree-based	No need to define k, provides dendrogram visualization	Computationally expensive, not scalable for large datasets
DBSCAN	Density-based	Detects arbitrary-shaped clusters, identifies noise points	Struggles with varying densities and high-dimensional data

These algorithms help analyze and extract meaningful insights from the Iris Dataset and Mall Customer Dataset, each with unique advantages depending on the clustering goal.

3.2 CLASSIFICATION

3.2.1 Load Dataset

Measurements of sepal length, sepal width, petal length, and petal width for numerous iris flower species are included in the Iris Dataset example, which is displayed in Figure 8. In contrast to species like as Iris-versicolor and Iris-virginica, the species Iris-setosa in this sample is distinguished by its shorter petals and sepals. Machine learning classification and clustering tasks frequently use this dataset.

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa

Fig. 8: IRIS Dataset

3.2.2 Features Visualization

Figure 9 illustrates the count distribution of species in the Iris dataset, showing three classes: Iris-setosa, Iris-versicolor, and Iris-virginica. Each species has 50 samples, indicating a balanced dataset with equal representation of each class. This balanced distribution ensures fair model training and helps avoid class imbalance issues during classification tasks.

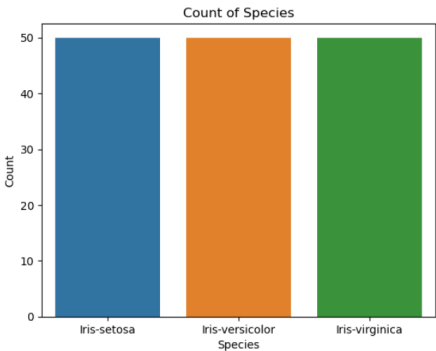


Fig. 9 Label Count of All Classes

Figure 10 shows a violin plot of petal length by species in the Iris dataset, comparing Iris-setosa, Iris-versicolor, and Iris-virginica. The plot reveals distinct differences in petal length distributions among the species. Iris-setosa has the smallest and least variable petal length, while Iris-virginica exhibits the largest and most variable petal lengths. This clear separation indicates petal length is a strong distinguishing feature for classification tasks.

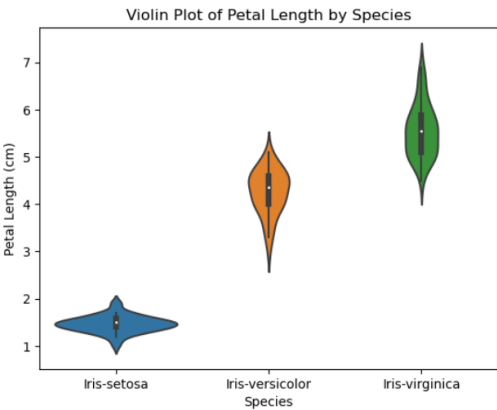


Fig. 10 Violin Plot by Species

In the figure 11, the species *Iris-setosa*, *Iris-versicolor*, and *Iris-virginica* are shown by colour-coded dots in a scatter plot of Sepal Length vs Petal Length for the Iris dataset. The species are clearly distinguished in the plot, especially *Iris-setosa*, which has noticeably shorter petal lengths. Regional overlap is seen between *Iris-versicolor* and *Iris-virginica*, however *Iris-virginica* typically has longer sepal and petal lengths. A positive correlation between sepal and petal lengths is shown by the upward trend, indicating that these features have strong predictive power for species categorization using ML models such as K-Means and Hierarchical Clustering.

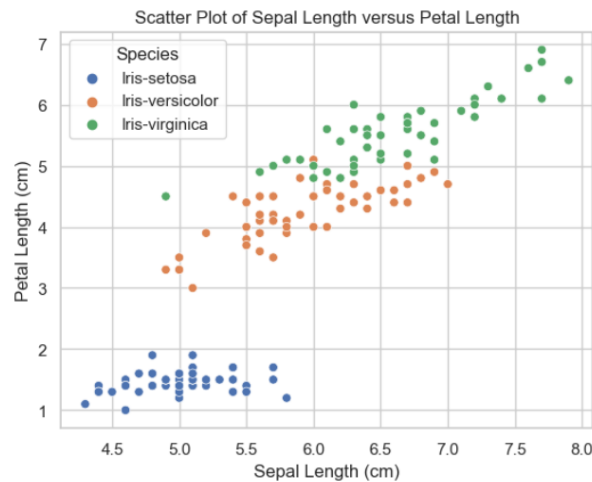


Fig. 11 Scatter Plot

3.2.3 Dimensionality Reduction for Visualization

After data preprocessing the next step is to reduce features, here PCA and tsne feature reduction algorithms are applied. In machine learning and data analysis, dimensionality reduction [23] is an essential preprocessing step, especially when working with high-dimensional datasets. It improves visualizations, enhances modeling efficiency, and reduces computational complexity. In particular, t-distributed stochastic neighbor embedding (t-SNE) and principal component analysis (PCA) are the two methods most commonly used in this study for dimensionality reduction and visualization.

Principal Component Analysis (PCA): Principal Component Analysis (PCA) was a popular linear technique for dimensionality reduction that kept the most variation in high-dimensional data while embedding it into a lower-dimensional space. The dimensionality reduction of PCA is done by finding a new coordinate system composed of a set of orthogonal axes, called principal components, that best account for the greatest variability in the dataset.

Algorithm Steps

- Standardize the dataset to have zero mean and unit variance.
- Compute the covariance matrix to understand feature relationships.
- Calculate the eigenvalues and eigenvectors of the covariance matrix.
- Select the top k principal components based on the highest eigenvalues.
- Transform the original dataset into a new subspace defined by these principal components.

t-SNE (t-distributed Stochastic Neighbor Embedding)

The non-linear dimensionality reduction method t-SNE is primarily intended for visualization. It differs from PCA in that t-SNE usually emphasizes maintaining local structure and distances among data points even more than it maintains variance. For this purpose, it uses probabilistic similarity metrics to map high-dimensional data into a fairly low-dimensional space, usually 2D or 3D.

Algorithm Steps

- Compute pairwise similarities in high-dimensional space using a probability distribution.
- Convert similarities into a low-dimensional probability distribution.
- Minimize the difference (KL-divergence) between the two distributions using gradient descent.
- Project the data onto a 2D or 3D plane while retaining local clusters.

3.2.3 Apply Machine Learning Models

After preprocessing and reduction in features comes the step of classification. Here, SVM, LR, and KNN [22] [24] [25] are examples of classification algorithms applied to create the prediction models. The following provides information on the details of these algorithms and their application.

SVM: With applications ranging from regression to classification and outlier detection, Support Vector Machine (SVM) is a strong supervised machine learning technique. An SVM works by finding the optimal hyperplane to distinguish data points of different classes; the support vectors are the closest data points from each class, and the hyperplane is selected to maximize the margin between them. In discriminating the classes, SVM uses kernel methods such as linear kernel, polynomial kernel, radial basis function (RBF) kernel, and sigmoid kernel to map input characteristics into higher-dimensional spaces. SVM handles linear and non-linear classification problems effectively. Because of its versatile nature, SVM is an excellent machine learning technique to apply to complex datasets where non-linear decision boundaries exist. In this research, SVM has been shown to classify clusters in datasets such as Iris and Mall Customer datasets with high accuracy and precision in pattern recognition.

Logistic Regression: Logistic regression is a statistical machine-learning technique predominantly applied for binary classification tasks. Interestingly, as the name suggests, logistic regression is a classification model and not really regression model. Essentially, the logistic (sigmoid) function is used to estimate the probability of occurrence of a binary event. Using the logistic curve, it then fits into data by evaluating how independent variables (also called features) influence the dependent variable (also called the class label). Logistic Regression is used here for predicting class labels for the Iris dataset and for customer segmentation in the Mall Customer dataset, where decent accuracy was achieved in all metrics of evaluation.

The logistic function is defined as:

$$P(y = 1|x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n)}}$$

The output is a probability value between 0 and 1, and predictions are typically made by applying a threshold (e.g., 0.5) to classify observations.

K- Neareast Naibhours: A straightforward yet effective instance-based learning approach for characterization and regression problems is KNN. The k nearest data points (neighbors) to a test instance are identified by the non-parametric, distance-based KNN model, which then assigns the majority class (for classification) or averages the values (for regression) to provide predictions. Although Euclidean distance is commonly applied to estimate the distance between points, depending on the type of dataset, Manhattan, Minkowski, or Hamming distances may also be employed. In this study, KNN is applied to both datasets to identify patterns and classify observations based on customer spending behavior and Iris species characteristics.

4 RESULTS ANALYSIS

4.1 Clustering

i. Hierarchical Clustering

Figure 12 illustrates Annual Income (k\$) versus Spending Score (1-100) to show the outcomes of Hierarchical Clustering applied to the Mall Customer Dataset. Different clusters are distinguished along the colour gradient bar, and the points are colored according to their cluster assignments. The algorithm highlights a variety of consumer subgroups, including high-income, high-spending customers and low-income, low-spending customers, by grouping clients with similar spending behavior and income levels. Businesses can find target groups for personalized marketing with the use of hierarchical clustering. Better insights into consumer habits and possible marketing strategies for various customer profiles are made possible by the cluster separation, which shows distinct behavioral patterns.

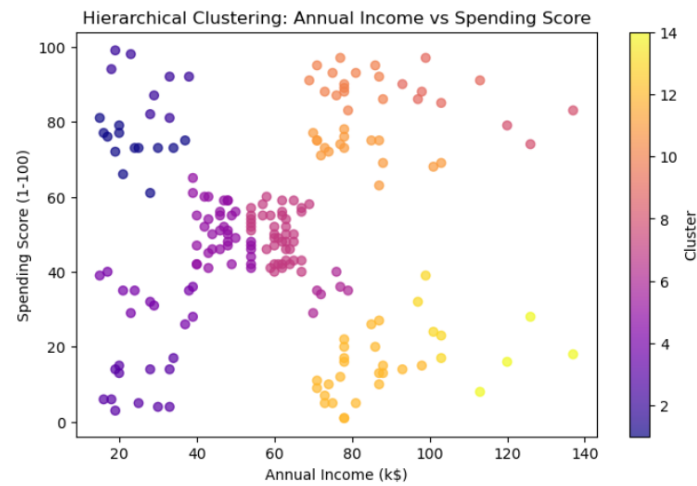


Fig.12: Hierarchical Clustering Results

ii. K- Means Cluster

Figure 13 shows the Annual Income (k\$) versus Spending Score (1-100) as a consequence of K-Means Clustering applied to the Mall Customer Dataset. Different colors are used to illustrate the distinct clusters that the data points are grouped into. Using similar spending habits and income levels, K-Means groups customers into high-income, high-spending customers, low-income, low-spending customers, and others with moderate income and spending behaviors. Businesses can use this segmentation to find target audiences for marketing strategies. For insightful market analysis and personalized marketing initiatives, the algorithm's effectiveness in classifying similar consumer behaviors is demonstrated by the clear distinction between clusters.



Fig. 13 K-Means Clustering Result

4.2 Comparative Analysis of Classification Models

In Figure 14, KNN, SVM, and LR are evaluated using Accuracy, Precision, Recall, and F1-Score in order to provide a Comparative Analysis of Classification Models. Perfect classification performance is demonstrated by the LR model, which beats the other models, attaining 100% across all metrics. Following closely behind, the SVM model exhibits great classification abilities with 99% accuracy, recall, and F1-score, and 100% precision. And lastly, KNN outperformed in terms of accuracy, recall, and F1 score which stood all at 97%, but had precision at 98%, suggesting a slightly reduced level of efficacy. With this consideration, the greatest model fit for the specified classification problem would be a Logistic Regression.

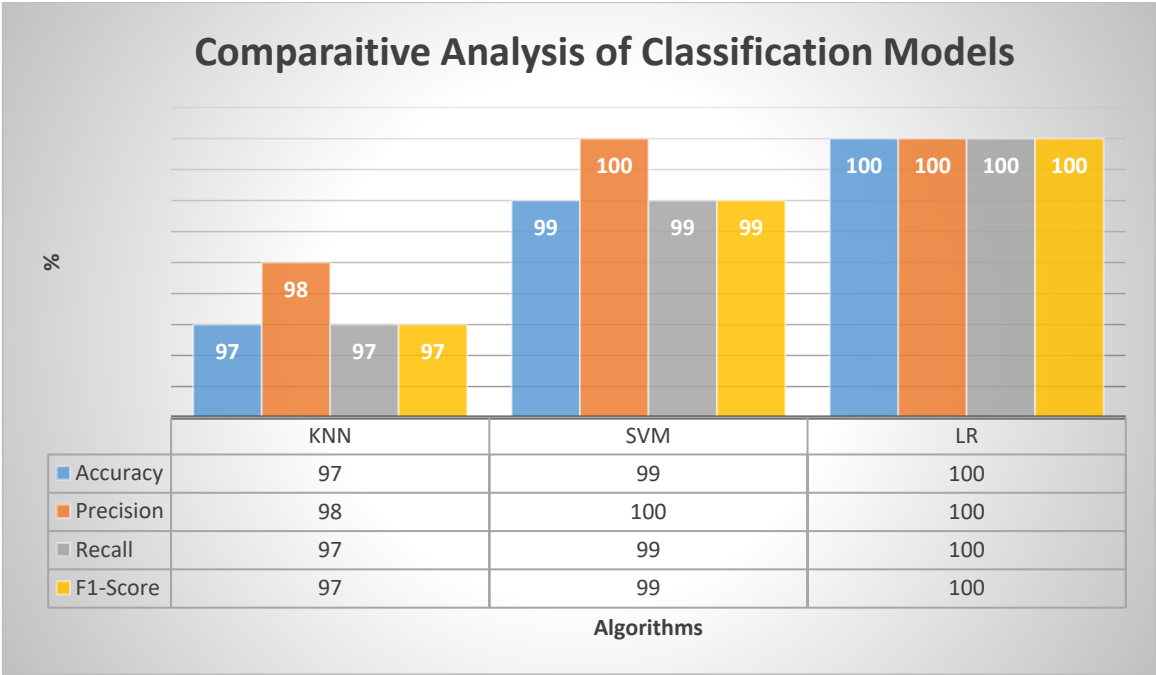


Figure 14. Comparative Analysis of Classification Models

5 CONCLUSION

This work reflects the converting of machine learning in the optimization of Big Data analysis through clustering, visualization, and insight extraction. By implementing clustering techniques such as K-Means, Hierarchical Clustering, and DBSCAN on datasets like Mall Customer and Iris, meaningful patterns and groups were identified. Visualization approaches such as PCA, t-SNE, and dendrogram contributed to interpretability, supported by the usage of the Silhouette score that ensures a strong discrimination for the number of clusters. Classification algorithms, which were used to classify the data, include Logistic Regression (LR), Support Vector Machine (SVM), and K-Nearest Neighbours (KNN). Among these, the LR turned out to be the most performing, with 100% accuracy, precision, recall, and F1-score. The results of this finding demonstrate the utility of clustering in supporting decision-making that is data-driven, especially in the economic and biological areas. The work also discusses some of the challenges in scalability and suggests ways for optimization with regard to efficiently working on large datasets. The theoretical implications of this research exhibit ways for enhancement of data-driven initiatives across various industries through machine learning-based clustering and classification for benefits in targeted marketing, personalization, and better market analysis.

REFERENCES

[1] Badawy, M., Ramadan, N. & Hefny, H.A. Big data analytics in healthcare: data sources, tools, challenges, and opportunities. *Journal of Electrical Systems and Inf Technol* 11, 63 (2024). <https://doi.org/10.1186/s43067-024-00190-w>

[2] Tosi, D., Kokaj, R. & Rocchetti, M. 15 years of Big Data: a systematic literature review. *J Big Data* 11, 73 (2024). <https://doi.org/10.1186/s40537-024-00914-9>

[3] Peng, Benji & Pan, Xuanhe & Wen, Yizhu & Bi, Ziqian & Chen, Keyu & Li, Ming & Liu, Ming & Niu, Qian & Liu, Junyu & Wang, Jinlang & Zhang, Sen & Xu, Jiawei & Feng, Pohsun. (2024). Deep Learning and Machine Learning, Advancing Big Data Analytics and Management: Handy Appetizer. 10.48550/arXiv.2409.17120.

[4] Domingue, J., Lasierra, N., Fensel, A., van Kasteren, T., Strohsbach, M., Thalhammer, A. (2016). Big Data Analysis. In: Cavanillas, J., Curry, E., Wahlster, W. (eds) *New Horizons for a Data-Driven Economy*. Springer, Cham. https://doi.org/10.1007/978-3-319-21569-3_5

[5] Gandomi, A.H.; Chen, F.; Abualigah, L. Machine Learning Technologies for Big Data Analytics. *Electronics* 2022, 11, 421. <https://doi.org/10.3390/electronics11030421>

- [6] Dhawas Pranali, et al. "Big Data Analysis Techniques: Data Preprocessing Techniques, Data Mining Techniques, Machine Learning Algorithm, Visualization." *Big Data Analytics Techniques for Market Intelligence*, edited by Dina Darwish, IGI Global, 2024, pp. 183-208. <https://doi.org/10.4018/979-8-3693-0413-6.ch007>
- [7] Mallikarjuna Paramesha, Nitin Liladhar Rane, & Jayesh Rane. (2024). *Big Data Analytics, Artificial Intelligence, Machine Learning, Internet of Things, and Blockchain for Enhanced Business Intelligence*. *Partners Universal Multidisciplinary Research Journal*, 1(2), 110–133. <https://doi.org/10.5281/zenodo.12827323>
- [8] Eid, A. , Miled, A. , Fatnassi, A. , Nawaz, M. , Mahmoud, A. , Abdalla, F. , Jabnoun, C. , Dhibi, A. , Allan, F. , Elhossiny, M. , Belhaj, S. and Mohamed, I. (2024) Sports Prediction Model through Cloud Computing and Big Data Based on Artificial Intelligence Method. *Journal of Intelligent Learning Systems and Applications*, 16, 53-79. <https://doi.org/10.4236/jilsa.2024.162005>
- [9] Sadat Lavasani, Mitra, Raeisi Ardali, Nahid, Sotudeh-Gharebagh, Rahmat, Zarghami, Reza, Abonyi, János and Mostoufi, Navid. "Big data analytics opportunities for applications in process engineering" *Reviews in Chemical Engineering*, vol. 39, no. 3, 2023, pp. 479-511. <https://doi.org/10.1515/revce-2020-0054>
- [10] Bharadiya, Jasmin Praful. "The role of machine learning in transforming business intelligence." *International Journal of Computing and Artificial Intelligence* 4.1 (2023): 16-24. <https://doi.org/10.33545/27076571.2023.v4.i1a.60>
- [11] Rehman, A., Naz, S. & Razzak, I. Leveraging big data analytics in healthcare enhancement: trends, challenges and opportunities. *Multimedia Systems* 28, 1339–1371 (2022). <https://doi.org/10.1007/s00530-020-00736-8>
- [12] Wengang Zhang, Xin Gu, Libin Tang, Yueping Yin, Dongsheng Liu, Yanmei Zhang, Application of machine learning, deep learning and optimization algorithms in geoenvironment and geoscience: Comprehensive review and future challenge, *Gondwana Research*, Volume 109, 2022, Pages 1-17, ISSN 1342-937X, <https://doi.org/10.1016/j.gr.2022.03.015>
- [13] Chinta Swetha, Integrating Machine Learning Algorithms in Big Data Analytics: A Framework for Enhancing Predictive Insights (September 01, 2021). *IJARESM*, Volume 9, Issue 10, October-2021, Pp. 2145-2161, Available at SSRN: <https://ssrn.com/abstract=5046555> or <http://dx.doi.org/10.2139/ssrn.5046555>
- [14] Torre-Bastida, A.I., Díaz-de-Arcaya, J., Osaba, E. et al. Bio-inspired computation for big data fusion, storage, processing, learning and visualization: state of the art and future directions. *Neural Comput & Applic* (2021). <https://doi.org/10.1007/s00521-021-06332-9>
- [15] Boppiniti, Sai Teja. "Big Data Meets Machine Learning: Strategies for Efficient Data Processing and Analysis in Large Datasets." *International Journal of Creative Research In Computer Technology and Design* 2.2 (2020).
- [16] Akash Tayal, Arun Solanki, Simar Preet Singh, Integrated frame work for identifying sustainable manufacturing layouts based on big data, machine learning, meta-heuristic and data envelopment analysis, *Sustainable Cities and Society*, Volume 62, 2020, 102383, ISSN 2210-6707, <https://doi.org/10.1016/j.scs.2020.102383>
- [17] Verbeeck, N., Caprioli, R.M. and Van de Plas, R. (2020), Unsupervised machine learning for exploratory data analysis in imaging mass spectrometry. *Mass Spec Rev*, 39: 245-291. <https://doi.org/10.1002/mas.21602>
- [18] Rahat Iqbal, Faiyaz Doctor, Brian More, Shahid Mahmud, Usman Yousuf, Big Data analytics and Computational Intelligence for Cyber–Physical Systems: Recent trends and state of the art applications, *Future Generation Computer Systems*, Volume 105, 2020, Pages 766-778, ISSN 0167-739X, <https://doi.org/10.1016/j.future.2017.10.021>
- [19] M. Feng et al., "Big Data Analytics and Mining for Effective Visualization and Trends Forecasting of Crime Data," in *IEEE Access*, vol. 7, pp. 106111-106123, 2019, <https://doi.org/10.1109/ACCESS.2019.2930410>
- [20] Lina Zhou, Shimei Pan, Jianwu Wang, Athanasios V. Vasilakos, Machine learning on big data: Opportunities and challenges, *Neurocomputing*, Volume 237, 2017, Pages 350-361, ISSN 0925-2312, <https://doi.org/10.1016/j.neucom.2017.01.026>
- [21] Omelina, Lubos & Goga, Jozef & Pavlovicova, Jarmila & Oravec, Miloš & Jansen, Bart. (2021). A survey of iris datasets. *Image and Vision Computing*. 108. 104109. 10.1016/j.imavis.2021.104109.
- [22] Mane, Deepak & Ashtagi, Rashmi & Suryawanshi, Ranjeetsingh & Kaulage, Anant & Hedao, Anushka & Kulkarni, Prathamesh & Gandhi, Yatin. (2024). Diabetic Retinopathy Recognition and Classification Using Transfer Learning Deep Neural Networks. *Traitement du Signal*. 41. 2683-2691. 10.18280/ts.410541.
- [23] Y. Dongre, N. Ranjan, P. Niranjane, M. Gulhane, Y. Gandhi, and P. Karmore, "Optimizing resource allocation in cloud-based information systems through machine learning algorithms," *J. Inf. Syst. Eng. Manag.*, vol. 10, no. 1s, 2025. [Online]. Available: <https://www.jisem-journal.com/>

- [24] Shanthi Kunchi, Vijaya N. Aher, Sharayu Ikhar, Kishor Pathak, Yatin Gandhi, and Kirti Wanjale. 2024. Risk Factor Prediction for Heart Disease Using Decision Trees. In Proceedings of the 5th International Conference on Information Management & Machine Intelligence (ICIMMI '23). Association for Computing Machinery, New York, NY, USA, Article 110, 1–6. <https://doi.org/10.1145/3647444.3647937>
- [25] Khetani, Vinit, et al. "Cross-domain analysis of ML and DL: evaluating their impact in diverse domains." International Journal of Intelligent Systems and Applications in Engineering 11.7s (2023): 253-262.