# A Survey Of Machine Learning Techniques In Big Data Analytics

Dr. S. Balamurugan* Dr.G.Anantharaj** & V. Ratnavali***

*Asst. Prof., Research Department of Computer Science, Kamaraj College (Autonomous), Thoothukudi – 628003, Tamilnadu, India (Affiliated to Manonmaniam Sundaranar University)
**Asst.Prof &Head, Department of Physical Education, Kamaraj College (Autonomous), Thoothukudi -628003, Tam Tamilnadu, India.
***Asst. Prof., Research Department of Computer Science, Kamaraj College (Autonomous), Thoothukudi – 628003, Tamilnadu, India (Affiliated to Manonmaniam Sundaranar University)

* balamurugan@kamarajcollege.ac.in, **ananthphs@gmail.com ***vprn93@gmail.com

**ABSTRACT**

Data is everywhere; everything is data. In traditional days we stored structured and smaller data using centralized database SQL, MySQL, and Oracle. In the modern era, we need to store structured, semi-structured, and unstructured data in a database to analyse real-time data. Big data analytics helps retrieve information from big and complex datasets. As the volume of data increases, conventional data processing becomes inefficient. We use machine learning algorithms to identify hidden patterns, predict outcomes, and make decisions. This paper reviews three algorithms—Decision Tree, Support Vector Machine (SVM), and K-Means clustering—used in big data analytics. We used a machine learning algorithm to analyse big data effectively. We analyse the results and discuss future enhancements of the algorithms for real-time big data processing, scalability improvement, and integration with deep learning techniques.

**Keywords:** Scalability, Big Data, Machine Learning, Deep Learning, Decision Tree, K-Means

## 1.Introduction

In the modern era, large volumes of complex data are generated and handled across domains such as social media, IoT devices, healthcare systems, and the financial sector [5]. We use big data analytics techniques to improve raw data effectively. To analyse patterns and decision-making, we use machine learning and data mining techniques. [7] Traditional databases are inefficient for analysing big data layers, so we are using big data to analyse large sets of complex data.

## 2. Limitation of Traditional Techniques

Large volumes of data like terabytes and petabytes are inefficient in traditional databases. When the amount of data increases, scalability is inefficient. It is inefficient in handling large-scale datasets. It requires more time for analysis, retrieval, and storage. It is unsuitable for real-time application. Large sets of data like images, text, audio, and video are difficult to handle in this traditional data processing. Its cost is very high for storing vast amounts of data. It fails to work in real-time data applications. Complex datasets and hidden insights are impossible with this type of technique.

## 3. Role of Machine Learning

Machine learning plays a significant role in big data analytics for analysing real-time data. It automatically learns from data [11]. Machine learning models learn from historical data, identify hidden patterns, and make decisions. By learning from data, they predict outcomes effectively. This process can be compared to how parents train toddlers to speak: children listen, learn, and gradually begin to use new words. Similarly, machine learning techniques learn from data and recognize patterns. Training machines with data enables accurate predictions.

"Children learn from past experiences and correct their mistakes. For example, when children touch a pooja lamp, they realize that it is hot and may burn their fingers; therefore, they avoid touching it in the future and can predict the outcome. Similarly, a machine learns from past data and generates predictions. This learning process enables effective solutions in big data analysis [8], [4]."

## 4.Why These Three Algorithms

In this survey paper, we analyse the decision tree, support vector machine, and K-means algorithms because of their significant roles in big data analytics. Decision trees are easy to interpret, support vector machines provide high accuracy, and K-means efficiently performs clustering on

large datasets. These algorithms are widely used and represent different machine learning approaches.

## 4.1 Decision Tree

The decision tree algorithm plays an important role in machine learning [6]. It makes decisions by asking a sequence of questions structured in the form of a tree. The tree is split based on conditions such as yes/no or true/false, and the final output is produced at the end of the decision path.

The algorithm learns from data. For example, in healthcare applications, it uses symptoms such as *"Fever?"*, *"High blood pressure?"*, and *"Sugar level?"* to predict diseases. By learning from past data, it identifies patterns and predicts outcomes effectively. Decision tree algorithms are widely used in applications such as banking, healthcare, and e-commerce.

A decision tree is constructed using a recursive splitting method, where
- The root node represents the first decision.
- An internal node represents a condition or feature.
- Leaf node represents the final output

### 4.1.1 Splitting criteria

Medical Example: Disease Test Prediction (Positive / Negative)

Dataset (Simple)

| Fever | Cough | Test Result |
| --- | --- | --- |
| Yes | Yes | Positive |
| Yes | Yes | Positive |
| Yes | No | Negative |
| No | No | Negative |

A medical diagnosis dataset is used to explain the splitting criteria of a decision tree algorithm. The dataset includes patient symptoms such as fever and cough, with the test outcome classified as positive or negative. Different attributes are evaluated to determine the optimal split. When the data is divided based on the fever attribute, the resulting subsets remain impure. However, splitting the dataset using the cough attribute produces pure subsets, where each child node contains a single class. Therefore, the cough attribute provides maximum information gain and is selected as the best splitting criterion. This example demonstrates how decision trees reduce uncertainty to improve classification accuracy in healthcare applications.

Several studies have surveyed decision tree algorithms and reported their effectiveness and limitations in machine learning and big data analytics. It handles both numerical and categorical data without requiring data normalization or scaling.

### 4.1.2 Drawback

The decision tree algorithm is overfitting [12], where the model learns noise from the training data and performs poorly on unseen data. They are also sensitive to noisy and imbalanced datasets, as small changes in data can significantly alter the tree structure. When applied to large and high-dimensional datasets, the tree may become complex and computationally expensive, reducing scalability. Furthermore, single decision tree models generally provide lower predictive accuracy compared to ensemble-based approaches.

## 5.SVM

Support Vector Machine (SVM) is an effective supervised learning algorithm commonly employed in big data analytics for handling complex and high-dimensional datasets. [13] The algorithm identifies an optimal separating hyperplane by maximizing the margin between different classes, which leads to better generalization on unseen data. [9] Through the application of kernel functions, SVM is capable of modelling nonlinear data distributions efficiently. Nevertheless, the high computational cost associated with training SVM models restricts their scalability in large-scale environments, making distributed and approximate learning strategies necessary.

## 5.1 Working Principle of Support Vector Machine

Support Vector Machine operates by identifying an optimal separating hyperplane that distinguishes data points belonging to different classes. The primary objective of the algorithm is to maximize the margin between classes, which enhances generalization performance on unseen data. The data points that lie closest to the decision boundary are known as support vectors, and they play a crucial role in defining the position of the hyperplane. By focusing on margin maximization, [10] SVM achieves robust classification even in high-dimensional feature spaces.

## 5.2 Kernel Functions in SVM

One of the key strengths of SVM is its ability to handle nonlinear data using kernel functions. The kernel trick enables SVM to implicitly transform data into a higher-dimensional space where linear separation becomes feasible. Kernel functions are linear, polynomial, radial basis function (RBF), and sigmoid kernels. [5] These kernels allow SVM to capture complex patterns in data without explicitly increasing computational dimensionality, making it suitable for diverse big data applications.

## 5.3 Advantages of SVM in Big Data Analytics

SVM offers high classification accuracy, particularly in high-dimensional datasets common in big data environments. [3] It is effective in handling complex and nonlinear relationships through kernel-based learning. Additionally, margin maximization reduces the risk of overfitting, making SVM a reliable choice for applications such as text classification, bioinformatics, and medical diagnosis.

### 5.4 Limitation
Despite its advantages, SVM faces scalability challenges when applied to very large datasets due to high computational and memory requirements. The selection of an appropriate kernel and tuning of hyperparameters can also be complex and time-consuming. Furthermore, SVM models lack interpretability compared to rule-based algorithms such as decision trees, which can limit their use in explainable AI applications.

## 6. K-Means Clustering in Big Data Analytics
### 6.1. Overview of K-Means Algorithm
K-means is an unsupervised machine learning algorithm widely used for clustering large datasets in big data analytics. [2] The algorithm partitions data points into $K$ distinct clusters by minimizing the distance between data points and their corresponding cluster centroids. Due to its simplicity and computational efficiency, K-Means is commonly applied in exploratory data analysis and large-scale data processing tasks.

### 6.2. Working Principle of K-Means
The K-Means algorithm begins by selecting $K$ initial centroids, either randomly or using heuristic methods. Data points are assigned to clusters based on their minimum distance to the centroids. After assignment, centroids are updated by computing the mean of all data points within each cluster. This iterative process continues until cluster assignments stabilize or convergence criteria are met.

### 6.3 K-Means in Big Data Analytics
K-Means is particularly suitable for big data analytics because of its scalability and ease of implementation. It can be efficiently parallelized using distributed computing frameworks such as Hadoop and Apache Spark, enabling clustering of large-scale datasets. K-means is widely used in applications including customer segmentation, image compression, and market basket analysis.

### 6.4 Advantages of K-Means
- Simple and easy to implement
- Computationally efficient for large datasets
- Scales well with distributed processing frameworks
- Effective for discovering hidden patterns in data

### 6.5 Limitations of K-Means
- Requires predefined number of clusters ($K$)
- Sensitive to initial centroid selection
- Performs poorly with non-spherical or uneven cluster sizes

Not suitable for categorical data without preprocessing prediction, and personalized learning systems.

### • Table: Comparative Analysis of Machine Learning Algorithms

| Algorithm | Learning Type | Key Strength | Major Limitation |
|---|---|---|---|
| Decision Tree | Supervised | High interpretability | Overfitting |
| Support Vector Machine (SVM) | Supervised | High classification accuracy | Scalability issues |
| K-Means | Unsupervised | Efficient data clustering | Requires predefined $K$ |

### 6.6 Applications of Machine Learning in Big Data Analytics
Machine learning techniques are widely applied in healthcare for disease diagnosis and patient data analysis. In banking, these algorithms support fraud detection, risk assessment, and customer behaviour analysis. Social media platforms utilize machine learning to analyse user interactions, identify trends, and provide personalized recommendations. In education, machine learning aids in performance evaluation, learning outcome

### 6.7 Challenges and Open Issues
Despite their effectiveness, machine learning algorithms face several challenges in big data analytics. Handling large-scale datasets remains difficult due to high computational and storage requirements. Real-time data processing demands fast and scalable learning models. Parameter tuning significantly influences model performance but is often complex and time-consuming.

Additionally, a trade-off exists between interpretability and accuracy, as highly accurate models are often less transparent.

### 7. Conclusion and Future Enhancements
This survey reviewed the role of machine learning techniques in big data analytics, with a focus on Decision Tree, Support Vector Machine, and K-Means algorithms. The study highlighted their working principles, advantages, limitations, and application domains. While Decision Trees provide interpretability, SVM achieves high accuracy for complex and high-dimensional data, and K-Means offers efficient clustering for large datasets. However, challenges related to scalability, real-time processing, parameter optimization, and the trade-off between interpretability and accuracy remain unresolved. Future enhancements should emphasize the development of hybrid models, scalable and distributed learning frameworks, and explainable machine learning

approaches to improve performance and reliability in large-scale big data applications.

## References

[1]   X. Liu and L. Song, "Hybrid SVM-K-means approach for big data pattern recognition," IEEE Access, vol. 12, pp. 5100–5115, 2024.

[2] R. Singh and R. K. Singh, "A novel hybrid clustering approach for big data analytics," IEEE Access, vol. 10, pp. 12345–12356, 2022.

[3] S. Kumar and N. Gupta, "Explainable machine learning in big data environments," IEEE Access, vol. 10, pp. 10252–10269, 2022.

[4] Z. K. Gao, "Editorial: Big data learning and discovery," IEEE Access, vol. 9, pp. 1–4, 2021.

[5] A. Kumar and S. Rao, "Machine learning techniques for big data analytics," IEEE Access, vol. 8, pp. 102345–102356, 2020.

[6] R. Sharma, "A survey on decision tree algorithms,"in Proc. IEEE Int. Conf. Data Science, Chennai, India, 2019, pp. 45–50.

[7] S. Tao, Z. Sun, and Z. Sun, "An improved intrusion detection algorithm based on GA and SVM," IEEE Access, vol. 6, pp. 13624–13631, 2018.

[8] Z. Sun, K. Hu, T. Hu, J. Liu, and K. Zhu, "Fast multi-label low-rank linearized SVM classification based on approximate extreme points," IEEE Access, vol. 6, pp. 42319–42326, 2018.

[9] A. Gandomi and M. Haider, "Beyond the hype: Big data concepts, methods, and analytics," Int. J. Information Management, vol. 35, no. 2, pp. 137–144, 2015.

[10] S. Chen, H. Mao, and Y. Liu, "Big data: A survey," Mobile Networks and Applications, vol. 19, no. 2, pp. 171–209, 2014.

[11] J. Han, M. Kamber, and J. Pei,Data Mining: Concepts and Techniques,3rd ed., San Francisco, CA, USA: Morgan Kaufmann, 2012.

[12] J. R. Quinlan, "Induction of decision trees,"Machine Learning, vol. 1, no. 1, pp. 81–106, 1986.

[13] V. Vapnik,The Nature of Statistical Learning Theory,New York, NY, USA: Springer, 1995.