

# The Impact of Machine Learning Algorithms on Big Data Analysis

Gyanendra K. Gautam<sup>1</sup>, Devendra K. Mishra<sup>2</sup>

<sup>1</sup>Research Scholar, Computer Science & Engineering, Amity School of Engineering and Technology,

Amity University Madhya Pradesh Gwalior, MP, India

E-mail: [gyanendragautam26@gmail.com](mailto:gyanendragautam26@gmail.com)

<sup>2</sup>Associate Professor, Computer Science & Engineering, Amity School of Engineering and Technology, Amity University

Madhya Pradesh Gwalior, MP, India

[dkmishra@gwa.amity.edu](mailto:dkmishra@gwa.amity.edu)

## ABSTRACT

This study extends on the intricate relationship between Big Data and Machine Learning, focusing on how both novel technologies are working in synergy to provide new opportunities in various fields. Due to the ever-increasing volume, velocity, and variety of Big Data, its elements include rapidly turning into something that will need advanced tools and technologies for its processing, analysis et al and storage. Such complexities of Big Data are solved by machine Learning where powerful algorithms and learning techniques are employed to find patterns, enhance decisions and foster innovations. In this research, it is explored the background, problems, and characteristics of big data (5Vs), including the scope of machine learning, and its types like supervised, unsupervised, and reinforcement learning for assisting big data in various aspects. It follows up these technologies with real usages in health care, banking, retail, manufacturing and student telecommunication industries. The study also reviews trends that will impact the direction of Big Data and Machine Learning convergence including challenge of ethical AI, immediate overviews, and protests concerning data privacy and security, as well as expansion concerns. It also discusses some concerns related to interpretability, scalability, and diversity of data. It also provides some perspective of further investigations that may include: real time analytics, edge computing, and ethics in artificial intelligence. Finally, a tabular condemnation of the various Machine Learning techniques utilized in the setting of Big Data is provided.

**Keywords:** Machine Learning, Data Analytics, Big Data, Unsupervised Learning, Supervised Learning, Reinforcement Learning, Predictive Analytics, Scalability, Privacy, Edge Computing.

## 1. INTRODUCTION

A growing number of informed individuals in the IT and analytics communities are becoming watchful while using the term "big data." The size and complexity of big data make it hard for conventional systems and data-warehousing technologies to handle and process. Big data is produced by both humans and technology, as well as by the natural world. Large amounts of data, which might be structured, semi-structured, or unstructured from many sources are created as a result of the development of services and technologies[1].

### 1.1 Types of Big Data: - Big Data has been divided into following categories-

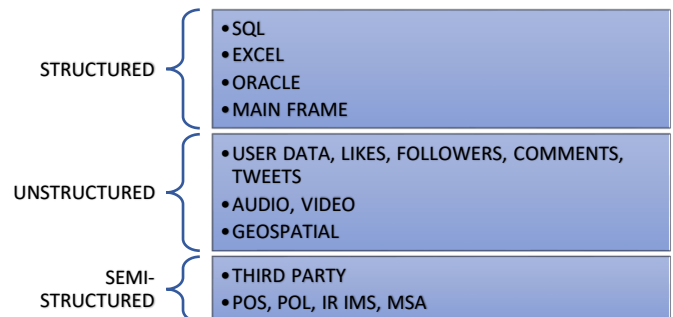


Fig. 1: Types of Big Data

**1.2 Big Data Characteristics (5Vs)** - There are 5Vs of Big Data that are:

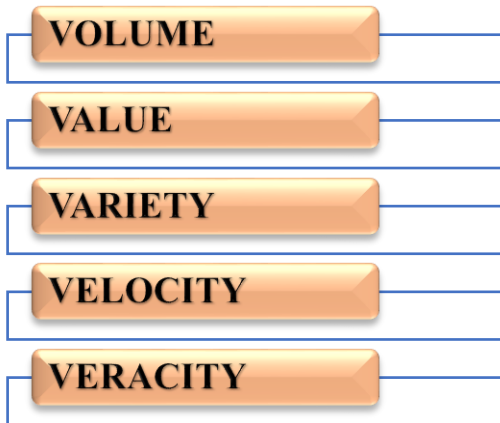


Fig. 2: Characteristics of Big Data

**1.3 Origin of Big Data:** Big Data can take place form different sources -

- (a) **Data Explosion:** With the advent of digitalization around the end of the 20th century, more and more data began accumulating at a rapid pace. This caused the generation of enormous tsunami of data that the rather rudimentary conventional data handling techniques could not cope including data from transactions, sensors, social media, and other sources.
- (b) **Internet Growth:** In the 1990s, the internet spiral began, which also accelerated the processes of data creation and data dissemination. The world wide web transformed into macromedia of both unstructured (text, multimedia) and structured (databases), which raised new challenges and opportunities for data management and analysis.
- (c) **Technological Advancements:** Organizations gained capacity to collect,

save, and work on more and more large amount of data sets due to processing power and storage capacity advancements. This also created challenges that led to the publishing of Big data with its volume, variety and velocity, distributed computing framework and scalable databases.

- (d) **Big Data Technologies:** Highlighted big Data problems also called for solutions such as Apache Hadoop, Apache Spark, and all.

**1.4 Applications of Big Data:** Some application areas of big data are as follows-

Table 1: The key applications of Big Data

Sector	Applications
Healthcare	<ul style="list-style-type: none"> <li>• Personalized Medicine</li> <li>• Medical Imaging Analysis</li> <li>• Drug Discovery</li> <li>• Healthcare Management</li> </ul>
Finance	<ul style="list-style-type: none"> <li>• Fraud Detection</li> <li>• Algorithmic Trading</li> <li>• Risk Management</li> <li>• Market Analysis.</li> </ul>
Retail	<ul style="list-style-type: none"> <li>• Customer Segmentation</li> <li>• Inventory Optimization</li> <li>• Dynamic Pricing</li> <li>• Recommendation Engines</li> </ul>
Manufacturing	<ul style="list-style-type: none"> <li>• Predictive Maintenance</li> <li>• Quality Control</li> <li>• Supply Chain Optimization</li> <li>• Energy Efficiency</li> </ul>
Telecommunications	<ul style="list-style-type: none"> <li>• Network Optimization</li> <li>• Service Personalization</li> <li>• Real-time Analytics</li> <li>• IoT Data Management</li> </ul>

**1.5 Challenges of Big Data:** Big data volume, velocity, variety, and authenticity provide a number of issues. The following are some major obstacles related to big data shows with block diagram.

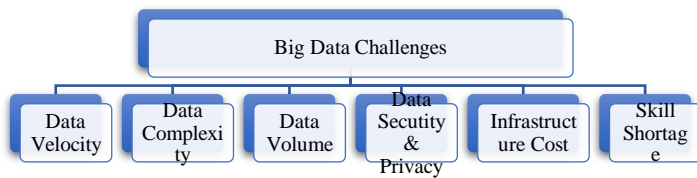


Fig. 3: Big Data Challenges Block Diagram

**1.6 Role of Big Data in Future:** Big Data will continue to shape the future in many different fields. The following crucial roles are anticipated of it:

- a) **Artificial Intelligence and Machine Learning:** Big Data will keep advancing the development of AI and ML. As the volume and variety of datasets increases, AI/ML algorithms will grow consequently and allow for better analysis, forecasting and systemization of the processes in the industries such as healthcare, finance, manufacturing and others.
- b) **Data Privacy and Security:** The more data there is the more the need to secure it and its contents. Advances in BD technologies will assist in increasing the encryption levels, addressing possible security threats, and ensuring that legal requirements (like GDPR, CCPA) are adhered to in order to safely extract valuable information.
- c) **Smart Cities and Urban Planning:** BD will support smart cities’ strategies by providing data integration techniques from other sources (including but not limited to sensors, social media, traffic flow). BD can provide a range of services such as resources

management, traffic management, safety services, and urban planning.

- d) **Healthcare and Precision Medicine:** BD will enhance healthcare systems as a result of development in genomics and personalized health care, and management of health at the population level. Through evaluating large scale medical databases, forecasted and resultant illnesses will also be treated, profiled and epidemiological tendencies will be determined at levels above customary measures.

## 2. TECHNIQUES OF MACHINE LEARNING

It is a One of the disciplines of AI that focuses on how to modify a computer to accomplish specific functions without any instructions by methods of algorithms and models is called machine learning. While Machine Learning (ML) systems learn from ‘education’ as stated above and this fact is important in their performance of such tasks as data processing, especially big and complex data.

### 2.1 Machine Learning Techniques in Big Data Analysis

- (a) **Supervised Learning:** It is the approach that provides a clear methodology for training algorithms on how to recognize patterns and forecast outcomes based on labelled datasets. These algorithms include the following:
  - (i) **Linear Regression:** It is a classical supervised approach used to establish the relationship that may exist between Two or more variables using Regression equation Models.

- (ii) **Support Vector Machines (SVM):** SVM is a common technique applied in machine learning which is defined as multi-class classifier for remote sensing image Classification and enhancing the classification accuracy of multitemporal satellite images. For a given sample, it is capable of regression and classification to determine the fitting hyperplane which best cuts the data space. The SVM has been demonstrated to effectively classify both the uniform and the diverse features in morphology in within remote [2].
- (iii) **Neural Networks:** Neural Networking is an effective computing modelling artificial brain, which applies neural interconnections for stimuli processing. Thanks to improved hardware and back propagation technologies, neural networks are becoming widely utilized for big data applications such as image and voice recognition[3].
- (b) **Unsupervised Learning:** This learning approach employs a computer to learn on its own without human involvement. The purposes of unsupervised learning are the restructuring of the input data into new attributes or a set of objects with the recognizable patterns. Typical algorithms are:
  - (i) **K-Means Clustering:** It teaches the algorithm to be able to perform autonomously on unlabelled, unclassified data and to permit the

computer to do so. That being said, the task of making a machine in this case is to categorize disordered information based on similarities, structures and discrepancies only.

- (ii) **Hierarchical Clustering:** In the field of data mining, cluster analysis is a set of techniques for the construction of nested orthogonal partitions. During the process every datum point is initialized as a single cluster and clusters which are in close proximity get merged stepwise until a satisfaction criterion is achieved.
- (iii) **Principal Component Analysis:** This approach was first developed in 1901 by the mathematician Karl Pearson. It doesn't allow too much circumstance that although the data from space of high axon number has been projected to lower space the lower space data must have its variances to be as high as possible.
- (c) **Reinforcement Learning:** Reinforcement learning would be one of the numerous challenges that involve an agent acting on its own, determining its environment and there is a feedback obtained in the form of reward which can either be good or bad. As an engagement goes on, an agent's goal is to increase this cumulative reward signal [4].

## 2.2 Applications of Machine Learning into Big Data

- (a) **Healthcare:** Making accountable predictions on the possible outcome of the patient, predicting illness and having customized therapy plans for the patients.

- (b) **Finance:** Segmentation and targeting of customers, trading using algorithms, management of risks, and fraud detection.
- (c) **Retail:** Control of the stock, analysis of customers' behaviour, and suggestion of systems.
- (d) **Transportation:** It is applied in the prediction of traffic jam, the most efficient routes and also self-driven automobiles.

### 2.3 Challenges and Future Directions

- (a) **Challenges:** Such application of technology in big data has implementation barriers and they are the present challenges which need to be tackled in order to avoid technology being a source of failure and frustration. Big data technologies are associated with problems like endless data that is acquired at rapid rates which needs to be stored and analysed.
  - (i) **Scalability:** Big data initiatives can grow fast. Cloud computing is the result's scalability of Big Data problem. It presents some drawback, such as managing and carrying out different tasks to meet each object of workload objective in an economical manner.
  - (ii) **Data Quality and Variety:** The huge amount of data generated can also have an effect on the correctness of data and quality. Making sure the data is appropriate, current, and full can be challenging given the volume of data being generated. This might lead to

assumptions or fallacy in the data analysis, generating inaccurate or insufficient discovery.

- (iii) **Interpretability:** Explanation of the big data interpretability: To communicate its conclusions properly, one has to understand complex models that are used. It is quite challenging to remain transparent without compromising privacy and proprietary information. Some of the significant challenges include a balance between accuracy and interpretability. Ensuring to retain interpretability while obtaining significant insights still remains an essential objective in making successful use of large data.
- (iv) **Privacy and Security:** Big data gives privacy and security challenges that including protect sensitive data, complying with laws such as GDPR, and reducing the likelihood of illegal access or data breaches. These tasks call for strong encryption, anonymization, and access control mechanisms.
- (b) **Future Instructions for use:** The growth of machine learning and artificial intelligence for deeper insights, edge computing for real-time analytics, and ethical concerns about data governance and privacy will probably be the main aim of big data in the coming times.
  - **Integration of Domain Knowledge:** with a view to enhance model performance and interpretability, domain-specific knowledge will be

progressively integrated into future machine learning systems.

- Automated Machine Learning:** The aim of the promptly growing field of automated machine learning research is to make machine learning obtainable to non-experts in the field by automating the processes required to build high-performance machine-learning pipelines for specific use cases[5].
- Ethical AI and Fairness:** with a view to ensure fairness in machine learning models and address biases in big data, further study and development are necessary.

- Advancements in Deep Learning:** As deep learning architectures and techniques continue to evolve, more complicated big data modelling will be possible, increasing scalability and accuracy.
- Edge Computing:** A promising approach is to directly install machine learning models on edge (Internet of Things) devices to process data locally instead of depending on centralized servers.
- Human-Machine Collaboration:** Improving human-in-the-loop systems, in which machine learning supports human judgment calls rather than taking over completely.

Table 2: Various paper of solving of big data problem using machine learning

Author (s)	Year	Title	ML Techniques Used	Application Area	Key Findings/Contributions
Meena et al. [6]	2020	Traffic Prediction for Intelligent Transportation System using Machine Learning	Deep Learning, Genetic Algorithms	Autonomous vehicles and traffic management	Developed a traffic prediction tool to support autonomous vehicle integration and improve traffic flow accuracy.
Shingate et al. [7]	2020	Adaptive Traffic Control System using Reinforcement Learning	Reinforcement Learning	Traffic light management	Optimized real-time traffic signal timing using deep neural networks and reinforcement learning.

Author (s)	Year	Title	ML Techniques Used	Application Area	Key Findings/Contributions
Singh et al. [8]	2021	Machine learning based distributed big data analysis framework for next generation web in IOT	Extreme Learning Machines (ELM), K-means, PCA	Internet of Things (IoT)	Achieved high accuracy in IoT data classification using scalable ML techniques.
Nassar & Kamal [9]	2021	ML & Big Data Analytics in Cybersecurity Threat Detection	Big Data Analytics, Machine Learning	Cybersecurity	Demonstrated real-time threat detection using Big Data and ML, and addressed privacy and ethical issues.
Punia et al. [10]	2021	Performance Analysis of ML Algorithms for Big	K-NN, Naïve Bayes, Decision Trees	Social media analytics	Evaluated different ML classification algorithms on social media data.

Author (s)	Year	Title	ML Techniques Used	Application Area	Key Findings/Contributions
		Data Classification			
Rahul et al. [11]	2021	Machine Learning Algorithms for Big Data Analytics	Supervised Learning, Reinforcement Learning	Industrial Applications	Addressed the effectiveness of ML in sectors such as banking, healthcare, and manufacturing.
Ang & Seng[12]	2021	ML and Big Data With Hyperspectral Information in Agriculture	ML, Deep Learning, Parallel Discriminant Analysis	Agriculture	Improved agricultural productivity by applying ML to hyperspectral data for crop management and soil monitoring.
Kumar et al. [13]	2022	Past, present, and future of sustainable finance: insights from big data analytics through machine learning of scholarly research	ML algorithms, Blockchain	Sustainable Finance	Proposed new research avenues for applying ML to green finance and carbon financing.
Manley et al.[14]	2022	A review of machine learning and big data applications in addressing ecosystem service research gaps	ML for uncertainty reduction	Cybersecurity	Addressed the importance of combining ML and Big Data for mapping ecosystem services and cyber security applications.
Deekshetha et al. [15]	2022	Traffic Prediction Using ML	Regression Models, TensorFlow	Intelligent Transportation Systems	Developed real-time traffic prediction system using ML libraries, improving traffic flow predictions.

Author (s)	Year	Title	ML Techniques Used	Application Area	Key Findings/Contributions
Khoshaba et al. [16]	2022	Implementation of machine learning techniques with big data and IoT to create effective prediction models for health informatics	Apache Spark, Apache Mahout	Big Data processing	Compared Apache Spark and Mahout for ML in Big Data environments, highlighting their impact on large dynamic datasets.
Zamani et al. [17]	2024	Multi-disease Prediction in Health Informatics	SSA-EL, DOA-EL, FPA-EL	Healthcare Informatics	Improved multi-disease prediction accuracy using heuristic ML models.
Tayseer et al. [18]	2024	IoT Integration for Machine Learning System using Big Data Processing	ML, Big Data Analytics	E-learning	Explored how Big Data and ML can enhance personalized e-learning experiences and security.
Adewusi et al. [19]	2024	Business Intelligence in the Era of Big Data	Predictive Analytics, ML	Business Intelligence	Highlighted the role of ML in improving business decision-making and gaining competitive advantages through data analytics.
Judjant et al. [20]	2024	Big Data Technology for Predicting Disease Spread	Integrated Data Sources, Big Data	Healthcare, Infectious Disease Control	Demonstrated improved predictions of disease spread using integrated Big Data platforms for quick decision-making in health crises.
Jahin et al. [21]	2024	Big Data and Supply	ML Algorithms, KPIs	Supply Chain Management	Introduced a framework for optimizing

Author (s)	Year	Title	ML Techniques Used	Application Area	Key Findings/Contributions
		Chain Management		ent	supply chain performance using ML and KPIs for better operational efficiency.

### 3. CONCLUSION

The combination of big data and machine learning has revolutionized many fields by promoting innovation and unprecedented insights. The former offers tools that would be able to handle and process the immense production of data every day, whereas the latter allows better quality and utility in these models. The blending of these elements paves the way, for progress in fields like automated technologies in cities and healthcare tailored to individual needs. There are still challenges to address in these areas though. Such, as making complex models easier to understand and ensuring scalability and privacy of data. Future potential for Big Data and Machine Learning is probably going to be concentrated on ethical challenges, real-time analytics, and insights into domain-specific breakthroughs. And as technology continues to advance, so will the benefits that Big Data and Machine Learning bring along for companies and society at large.

### REFERENCE

[1] Ishwarappa and J. Anuradha, "A brief introduction on big data 5Vs characteristics and hadoop technology," *Procedia Comput Sci*, vol. 48, no. C, pp. 319–324, 2015, doi: 10.1016/j.procs.2015.04.188.

[2] M. S. Chowdhury, "Comparison of accuracy and reliability of random forest, support vector machine, artificial neural network and maximum likelihood method in land use/cover classification of urban setting," *Environmental Challenges*, vol. 14, no. October 2023, p. 100800, 2024, doi: 10.1016/j.envc.2023.100800.

[3] Y. Lecun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015, doi: 10.1038/nature14539.

[4] Z. Ding, Y. Huang, H. Yuan, and H. Dong, "Introduction to reinforcement learning," *Deep Reinforcement Learning: Fundamentals, Research and Applications*, pp. 47–123, 2020, doi: 10.1007/978-981-15-4095-0\_2.

[5] M. Baratchi *et al.*, *Automated machine learning: past, present and future*, vol. 57, no. 5. Springer Netherlands, 2024. doi: 10.1007/s10462-024-10726-1.

[6] G. Meena, D. Sharma, and M. Mahrishi, "Traffic Prediction for Intelligent Transportation System using Machine Learning," *Proceedings of 3rd International Conference on Emerging Technologies in Computer Engineering: Machine Learning and Internet of Things, ICETCE 2020*, no. February 2020, pp. 145–148, 2020, doi: 10.1109/ICETCE48199.2020.9091758.

[7] Kranti Shingate, Komal Jagdale, and Yohann Dias, "Adaptive Traffic Control System using Reinforcement Learning," *International Journal of Engineering Research and*, vol. V9, no. 02, pp. 443–447, 2020, doi: 10.17577/ijertv9is020159.

[8] S. K. Singh, J. Cha, T. W. Kim, and J. H. Park, "Machine learning based distributed big data analysis framework for next generation web in iot," *Computer Science and Information Systems*, vol. 18, no. 2, pp. 597–618, 2021, doi: 10.2298/CSIS200330012S.

[9] A. Nassar and M. Kamal, "Traditional rule-based security systems, while effective to some extent, are insufficient to combat the dynamic and evolving strategies employed by cybercriminals," pp. 51–62, 2021.

[10] S. K. Punia, M. Kumar, T. Stephan, G. G. Deverajan, and R. Patan, "Performance analysis of machine learning algorithms for big data classification: MI and ai-based algorithms for big data analysis," *International Journal of E-Health and Medical Communications*, vol. 12, no. 4, pp. 60–75, 2021, doi: 10.4018/IJEHMC.20210701.0a4.

[11] K. Rahul, R. K. Banyal, P. Goswami, and V. Kumar, *Machine learning algorithms for big data analytics*, vol. 1227, no. January. Springer Singapore, 2021. doi: 10.1007/978-981-15-6876-3\_27.

[12] K. L. M. Ang and J. K. P. Seng, "Big data and machine learning with hyperspectral information in agriculture," *IEEE Access*, vol. 9, pp. 36699–36718, 2021, doi: 10.1109/ACCESS.2021.3051196.

[13] S. Kumar, D. Sharma, S. Rao, W. M. Lim, and S. K. Mangla, "Past, present, and future of sustainable finance: insights from big data analytics through machine learning of scholarly research," *Ann Oper Res*, 2022, doi: 10.1007/s10479-021-04410-8.

[14] K. Manley, C. Nyelele, and B. N. Egho, "A review of machine learning and big data applications in addressing ecosystem service research gaps," *Ecosyst Serv*, vol. 57, no. September, p. 101478, 2022, doi: 10.1016/j.ecoser.2022.101478.

[15] H. R. Deekshetha, A. V. Shreyas Madhav, and A. K. Tyagi, "Traffic Prediction Using Machine Learning," *Lecture Notes on Data Engineering and Communications Technologies*, vol. 116, pp. 969–983, 2022, doi: 10.1007/978-981-16-9605-3\_68.

[16] F. Khoshaba, S. Kareem, H. Awla, and C. Mohammed, "Machine learning algorithms in Bigdata Analysis and its applications: A review," *HORA 2022 - 4th International Congress on Human-Computer Interaction, Optimization and Robotic Applications, Proceedings*, no. July, pp. 1–8, 2022, doi: 10.1109/HORA55278.2022.9799848.

[17] A. S. Zamani, A. H. A. Hashim, A. S. A. Shatat, M. M. Akhtar, M. Rizwanullah, and S. S. I. Mohamed, "Implementation of machine learning techniques with big data and IoT to create effective prediction models for health informatics," *Biomed Signal Process Control*, vol. 94, no. April, 2024, doi: 10.1016/j.bspc.2024.106247.

[18] F. Tayseer *et al.*, "International Journal of INTELLIGENT SYSTEMS AND APPLICATIONS IN ENGINEERING IoT Integration for Machine

Learning System using Big Data Processing,” *Original Research Paper International Journal of Intelligent Systems and Applications in Engineering IJISAE*, vol. 2024, no. 14s, pp. 591–599, 2024, [Online]. Available: [www.ijisae.org](http://www.ijisae.org)

[19] Adebunmi Okechukwu Adewusi, Ugochukwu Ikechukwu Okoli, Ejuma Adaga, Temidayo Olorunsogo, Onyeka Franca Asuzu, and Donald Obinna Daraojimba, “Business Intelligence in the Era of Big Data: a Review of Analytical Tools and Competitive Advantage,” *Computer Science & IT Research Journal*, vol. 5, no. 2, pp. 415–431, 2024, doi: 10.51594/csitj.v5i2.791.

[20] Loso Judijanto, Hermansyah, K. P. Ningsih, D. Anurogo, and M. Firdaus, “The Role of Big Data Technology in Predicting and Managing the Spread of Infectious Diseases,” *J. of World Future Medicine, Health and Nursing*, vol. 2, no. 2, pp. 219–230, 2024.

[21] M. A. Jahin and I. A. Ridoy, “Big Data - Supply Chain Management Framework for Forecasting: Data Preprocessing and Machine Learning Techniques,” *SSRN Electronic Journal*, pp. 1–47, 2022, doi: 10.2139/ssrn.4076759.