



Utilising Machine Learning Algorithms to Address Computational Challenges in Big Data Analytics

Mohammad Islam¹, Ashish Sharma^{2,3*}, Nausheen Khilji⁴

¹ Research Scholar, Department of Computer Science, Maulana Azad University, Jodhpur-342802 Rajasthan, India,
Email: islamjodhpur@gmail.com-ORCID: 0009-0001-2548-4262

²Associate Professor, HOD, Department of Computer Science, Maulana Azad University, Jodhpur-342802 Rajasthan, India.

³Professor, Department of Technology, JIET, Jodhpur 342802, Rajasthan, India
Email: aashishid@gmail.com-ORCID: 0009-0007-6644-6362

⁴Assistant Professor, Department of Technology, JIET UNIVERSE, Jodhpur 342802, Rajasthan, India
Email: naushy90@gmail.com-ORCID: 0000-0003-1750-2675

Article Info:

DOI: 10.22399/ijcesn.3165

Received : 05 May 2025

Accepted : 01 July 2025

Keywords

Big Data Analytics,
Machine Learning Algorithms,
Supervised Learning,
Unsupervised Learning,
Deep Learning,
Scalability

Abstract:

The rapid growth of data across industries including finance, healthcare, and e-commerce has led to significant computational hurdles in big data analytics. Challenges include dimensionality reduction, scalability, real-time processing, and the handling of large data volumes. The use of sophisticated machine learning (ML) algorithms is necessary to tackle these complexity, as traditional data processing methods are insufficient. This paper examines the efficacy of supervised, unsupervised, and deep learning algorithms to improve the efficiency, scalability, and accuracy of data processing to overcome these computing problems. Random Forest, XGBoost, K-means clustering, Convolutional Neural Networks (CNNs), and Long Short-Term Memory Networks (LSTMs) are approaches utilized in the analysis of extensive datasets. The findings demonstrate that CNNs surpass alternative models in image-based datasets, achieving an accuracy of 93.8%. XGBoost attains an equilibrium between computing efficiency and accuracy (91.7%) in consumer segmentation and fraud detection. The remarkable scalability of K-means clustering makes it a suitable technique for analyzing customer behavior. The study incorporates distributed platforms like Apache Spark and TensorFlow to address critical difficulties, including high memory consumption, real-time data processing, and model interpretability. The findings align with existing studies and highlight the importance of scalable and resource-efficient machine learning methods. Moreover, the research provides significant insights into the capacity of stream processing frameworks and hybrid methodologies to improve real-time analytics. This study significantly advances the burgeoning domain of big data analytics by offering pragmatic machine learning techniques. This approach guarantees the effective management of large volumes of data by producing actionable insights.

1. Introduction

1.1 Growing Importance of Big Data in Various Sectors

The surge of data from many sources, such as sensors, social media, transactions, and IoT devices, has converted big data into a strategic asset in the modern digital age. The ability to evaluate and draw conclusions from extensive databases is significantly beneficial to enterprises across all

sectors. In the healthcare sector, big data analytics is utilized to enhance therapies and predict diseases by leveraging genomic information, imaging data, and electronic health records (EHRs) [2]. This facilitates the execution of customized care, population health management, and predictive analytics. Machine learning-based predictive models improve the precision of patient outcome forecasts and offer supplementary assistance for early diagnosis [3]. Algorithmic trading, fraud detection, credit risk

assessment, and customer segmentation rely on big data within the financial sector [4]. Predictive analytics leverage real-time transaction data, enabling financial institutions to provide personalized services and efficiently mitigate risks. Big data facilitates demand projections, product recommendations, inventory management, and consumer behavior analysis in the retail and e-commerce industries [5]. Algorithms driven by machine learning provide dynamic pricing through real-time analysis of consumer preferences and market trends, ultimately optimizing profitability and enhancing the customer experience. Big data in manufacturing enables predictive maintenance and process optimization through the aggregation of data from sensors and equipment [7]. The demand for advanced analytics is emphasized by the immense volumes of data generated by the Internet of Things (IoT) in areas like connected automobiles and smart homes. As businesses increasingly acknowledge the competitive advantage of big data, the need for effective analytics solutions is rising. However, this change presents unique computing issues that necessitate the creation of novel approaches for the effective management and analysis of data.

1.2 The Computational Challenges Associated with Big Data

The inherent characteristics of big data, often encapsulated as the 4 Vs: Volume, Velocity, Variety, and Veracity, pose several computational hurdles to its effective exploitation [9]. The phrase "volume" refers to the magnitude of big data, surpassing the limitations of traditional databases and algorithms. Therefore, to efficiently manage petabytes or even exabytes of data, the implementation of novel storage technologies and distributed systems is essential [10]. Velocity refers to the speed at which data is produced in domains such as finance, e-commerce, and social media, requiring real-time or near-real-time processing. Algorithms are required to deliver prompt insights with minimal delay, while bulk processing methods are inadequate for streaming data [11]. The diverse nature of big data is highlighted by its various formats, including structured, unstructured, and semi-structured data, such as text, images, videos, and sensor data. Advanced preprocessing techniques and flexible machine learning models are essential to address this variability [12]. Veracity emphasizes the reliability and quality of data, as inconsistent, noisy, or inadequate data might affect analytical results, requiring the use of rigorous purification and imputation techniques.

Moreover, the analysis of large datasets poses numerous computational difficulties. As data

quantities rise, scalability becomes paramount, requiring algorithms that can efficiently leverage the resources of distributed computation [13]. Significant processing delays or system failures can result from the incapacity of current methods to scale. Dimensionality is another challenge, as datasets containing hundreds or millions of features generate high-dimensional spaces that hinder the use of traditional machine learning methods [14]. This leads to overfitting, prolonged calculation times, and diminished interpretability. Moreover, sectors requiring real-time decision-making, such as stock trading, healthcare, and autonomous driving, demand exceptionally high processing speeds. However, the quest for speed may require compromising accuracy or resource efficiency in the process. Moreover, storage and memory limitations pose significant hurdles, as even distributed storage systems like the Hadoop Distributed File System (HDFS) are unable to effectively handle large datasets [15]. In-memory processing of extensive datasets becomes troublesome when data volume surpasses system capability, leading to heightened latency and bottlenecks. To overcome these problems and fully harness the potential of big data analytics, it is essential to develop sophisticated algorithms capable of efficiently managing data, scaling, processing in parallel, and learning in real time.

1.3 How Machine Learning Algorithms Play a Role in Overcoming These Challenges

Machine learning (ML) algorithms have become powerful tools for overcoming computational challenges in huge data analytics. Due to their versatility and ability to learn from data, they provide the ideal option for maximizing resource use and deriving conclusions from intricate, extensive datasets [17]. In distributed and parallel environments like Apache Spark and Hadoop, frameworks such as MapReduce are utilized to process data across various clusters. In these situations, several machine learning methods are engineered for effective scalability, thus addressing the constraints of single-node systems. Dimensionality reduction methods, such as Principal Component Analysis (PCA), Autoencoders, and t-SNE, are utilized to address the challenge of high-dimensional data, a common problem in the realm of big data [18]. Feature selection methods, such as Lasso and recursive feature elimination, are utilized to determine the most relevant characteristics for model development. Reinforcement learning methodologies facilitate real-time data processing by dynamically adjusting to evolving settings, whereas online learning models perpetually update

with incoming data, hence enhancing decision-making efficiency. To ensure the efficient operation of deep learning algorithms in resource-constrained situations, models employ tactics like incremental learning and model pruning, facilitating training on diminished datasets. This is executed to enhance the administration of computing resources. Autoencoders and k-Nearest Neighbors (k-NN) are imputation techniques utilized to improve data quality [19]. These methods are proficient in effectively handling disordered data and absent values. Explainable AI (XAI) enhances the interpretability of machine learning algorithms in complex domains like finance and healthcare by providing understandable justifications for predictions. By utilizing advanced machine learning approaches, firms can maximize large data, thus surmounting computing limitations and guaranteeing swift and reliable insights [20]. Thus, ML algorithms enable enterprises to maintain competitiveness in a data-driven landscape by reconciling the demand for swift, significant outcomes with the extensive data volume.

2. Literature Review

Computational Challenges in Big Data Analytics

Iqbal et al. (2020) examined how big data is reshaping the evolution of smart cities and the conservation of contemporary human cultures. It underscored the importance of big data in contemporary life and the economy, therefore addressing the principal barriers to its implementation. In contemporary smart city applications, various computational intelligence (CI) methodologies have been explored as effective tools for Big Data analytics. The Hierarchical Spatial-Temporal State Machine (HSTSM) was introduced as a unique data modeling methodology and exemplified in a case study concerning intelligent transportation in the framework of smart cities. The essay also discussed issues related to legislation, protection, value, and commercialization relevant to the implementation and use of Big Data applications. Hariri et al. (2019) addressed the necessity of comprehending patterns within extensive datasets, which propelled the development of big data analytics and garnered attention from both academia and industry. Progress in sensor networks, cyber-physical systems, and the Internet of Things has significantly augmented data volume across various industries, including healthcare, finance, education, social media, and smart cities. Nevertheless, data from sources such as sensors and social media

frequently exhibited noise, incompleteness, and inconsistency, hence requiring sophisticated analytics to facilitate dependable predictions and judgments. In contrast to conventional methods, AI techniques (such as machine learning, natural language processing, and computational intelligence) provide expedited, more precise, and scalable solutions. Although previous research concentrated on certain methodologies or applications, limited investigations addressed the influence of uncertainty in AI-driven big data analytics.

Daniel et al. (2019) The application of big data, encompassing extensive volumes of information produced by persons, applications, and computers, proliferated across various sectors, including education. Researchers faced several challenges when engaging with Big Data in education, including diverse interpretations of the term, ontological and epistemological conflicts, technical difficulties, ethical and privacy issues, digital divides, and insufficient expertise and training opportunities for researchers. To foster dialogue and enhance understanding of these essential issues by integrating the writers' personal experiences with Big Data research in education and their expertise. Mirza et al. (2019) examined how developments in high-throughput technologies have enabled the substantial accumulation of omics data from many sources, including the genome, epigenome, and proteome. Historically, statistical and machine learning (ML) approaches were utilized to examine each data source independently. Despite the greater computing constraints posed by the integration of multi-omics and clinical data compared to single-omics studies, this integration was crucial for the advancement of precision medicine and biomedical research. This research highlighted the application of machine learning (ML) approaches to tackle significant computational issues in integrative analysis, including scalability, class imbalance, missing data, data heterogeneity, and high dimensionality.

“Machine Learning in Big Data Context

Li et al. (2021) examined how traditional healthcare providers have inadequately addressed fundamental medical needs during the COVID-19 pandemic. The IoT-enabled, intelligent, linked wearables collect extensive data on behavioral, psychological, and physical health, impacting the contemporary healthcare industry. The challenge of handling the substantial volumes of data generated by these devices can impede decision-making. Machine

Table 1: Comparison Table

Study	Focus Area	Key Findings	Challenges Addressed	Technologies/Methods	Applications/Implications
Iqbal et al. (2020)	Big Data in Smart Cities	Emphasized the role of big data in modern life and economic development, overcoming barriers to adoption.	Concerns regarding regulation, protection, value, and commercialization of Big Data applications.	Hierarchical Spatial-Temporal State Machine (HSTSM) for data modeling.	Case study on intelligent transportation in smart cities.
Hariri et al. (2019)	Patterns in Large Datasets	Highlighted the necessity of big data analytics in understanding data patterns, especially from sensor networks and IoT.	Issues with noisy, incomplete, and inconsistent data from sensors and social media, necessitating advanced analytics.	AI techniques, including machine learning (ML) and Natural Language Processing (NLP).	Impacts across healthcare, finance, education, social media, and smart cities.
Daniel et al. (2019)	Big Data in Education	Discussed various challenges faced in educational contexts, such as different interpretations of big data and ethical concerns.	Ontological and epistemological conflicts, technical difficulties, ethical and privacy concerns, and digital divides.	Incorporation of personal experiences to foster discourse in Big Data research.	Awareness and discussion on critical matters in education.
Mirza et al. (2019)	Multi-Omics Data Integration in Precision Medicine	Focused on the integration of omics data for advancing precision medicine, emphasizing the use of ML methodologies.	Computational challenges such as scalability, class imbalance, missing data, data heterogeneity, and high dimensionality in multi-omics studies.	Machine Learning (ML) methodologies for integrative analysis.	Advancement of biomedical research through integrated data analysis.

learning (ML) has addressed numerous networking difficulties, while big data analytics has garnered attention for information extraction and forecasting. Despite extensive study on machine learning and big data analytics, the Internet of Things healthcare sector is deficient in machine learning studies pertaining to big data analysis. This study examines machine learning (ML) techniques for analyzing healthcare datasets, including their advantages, disadvantages, and research challenges. Healthcare professionals and government agencies can remain informed about machine learning (ML)-driven big data analytics for intelligent healthcare using the provided information. Amanullah et.al (2020) discussed the technology allowed items to

communicate and interact, the Internet of Things (IoT) became essential to modern life. However, security weaknesses made the Internet of Things (IoT) susceptible, requiring resilient solutions through technology integration or development. Machine learning, or deep learning, identified security threats. IoT devices created a lot of diverse data; therefore, big data technology increased data management and performance. Examined the latest big data, IoT security, and deep learning technologies. created a technical-studies-based theme taxonomy after comparing these topics. Roh et al. (2019) examined the constraints of machine learning, emphasizing the significance of data collection across many domains. This necessity

arose from the increasing application of machine learning in scenarios with limited data and the automated feature development of deep learning, which necessitated larger labeled datasets. Recent research on data acquisition—including data management, machine learning, natural language processing, and computer vision—emphasized the significance of managing vast datasets. analyzed

data management, encompassing data collecting, annotation, and model optimization. It emphasized significant research concerns and provided guidance on methodological selections. Data management and machine learning were integral to a broader initiative that amalgamated AI and big data, presenting several research opportunities.”

Table 2: Comparison Table

Study	Focus Area	Key Findings	Challenges Addressed	Technologies/Methods	Applications/Implications
Li et al. (2021)	IoT-Enabled Healthcare During COVID-19	Highlighted the potential of IoT-enabled wearables to gather data on health, but noted challenges in managing the large data volumes.	Difficulty in decision-making due to the massive amounts of data generated by IoT devices; lack of machine learning research in the IoT healthcare market.	Machine learning (ML) methods for healthcare dataset analysis; big data analytics for prediction.	Insights for healthcare professionals and government agencies on ML-based analytics for smart healthcare.
Amanullah et al. (2020)	Security in IoT	Discussed the importance of IoT in modern life and the security vulnerabilities that require resilient solutions through technology integration.	Security weaknesses in IoT necessitate the development of robust solutions; challenges in managing diverse data from IoT devices.	Integration of machine learning and deep learning for identifying security threats; big data technology for data management.	Taxonomy of recent big data, IoT security, and deep learning technologies for improved data performance.
Roh et al. (2019)	Data Gathering in Machine Learning Applications	Emphasized the critical role of data gathering for machine learning, particularly in applications with limited data and deep	Challenges in data management, including acquisition, labeling, and enhancing models; need for larger labeled datasets for	Focus on data management methodologies, including natural language processing and computer vision.	Integration of AI and big data with machine learning, presenting research opportunities.

		learning needs.	effective machine learning.		
--	--	-----------------	-----------------------------	--	--

“Comparison of Existing ML Solutions for Big Data Challenges

Deepa et al. (2022) examined the advantages and applications of big data, which garnered considerable interest from several scientific and engineering fields. To improve service quality, it was essential to address issues related to privacy, administration, and analytics. The decentralized and secure nature of blockchain technology suggests that big data services could be significantly improved. The process began with an in-depth analysis of the two fields and the reasoning for their integration. It subsequently examined the capacity of blockchain technology to enable safe data capture, storage, analysis, and privacy. The evaluation included cutting-edge research and exemplary projects in areas such as transportation, healthcare, infrastructure, and smart cities.

Hajjaji et al. (2021) conducted an extensive examination of big data and IoT applications in intelligent environments, including key application

domains, prevailing trends, data architectures, and associated problems. The benefits of combining IoT and big data for the monitoring, preservation, and enhancement of natural resources. Applications encompassed disaster notifications, intelligent metering, precision agriculture, and environmental monitoring.

Hossain et al. (2019) performed an extensive analysis of the utilization of big data and machine learning in the smart grid (SG), the advanced power system. Advanced analytical techniques were considered essential to enhance informed decision-making, as connectivity was crucial to this new infrastructure and was facilitated by the Internet of Things (IoT). A substantial amount of this data was produced. The SG system preserved its cost-effectiveness while enhancing data collection and load forecasting with the integration of the Internet of Things (IoT). Nonetheless, this interconnected system presented considerable cybersecurity issues owing to the susceptibility of IoT devices and their data to assaults.”

Table 3: Comparison Table

Study	Focus Area	Key Findings	Challenges Addressed	Technologies/Methods	Applications/Implications
Deepa et al. (2022)	Integration of Big Data and Blockchain	Explored how blockchain technology could enhance big data services by providing secure data acquisition, storage, and analytics.	Concerns regarding privacy, administration, and analytics in big data; need for secure and decentralized solutions.	Comprehensive examination of blockchain and big data integration; state-of-the-art research review.	Applications in transportation, healthcare, infrastructure, and smart cities.
Hajjaji et al. (2021)	Big Data and IoT in Smart Environments	Analyzed the integration of big data and IoT in monitoring and enhancing natural resources, along with significant	Challenges in data architecture and management; need for effective monitoring and resource preservation strategies.	Overview of current trends and architectures in big data and IoT.	Applications in disaster alerts, smart metering, smart agriculture, and environmental monitoring.

		application areas.			
Hossain et al. (2019)	Big Data and Machine Learning in Smart Grids	Investigated the role of big data and machine learning in improving decision-making in smart grids, facilitated by IoT connectivity.	Cybersecurity concerns related to IoT device vulnerabilities; need for advanced analysis techniques to improve load predictions and data collection.	Integration of big data analytics and machine learning for informed decision-making in smart grids.	Enhanced cost-effectiveness and efficiency in smart grid systems.

3. Research Methodology

3.1 Data Source and Tools

Specifically, this part addresses the computational challenges associated with big data analytics by providing an overview of the datasets and tools utilized to implement machine learning algorithms.

3.1.1 Description of Datasets

Utilizing an assortment of datasets, the efficacy of machine learning models is evaluated and investigated. Throughout numerous disciplines, this process is significantly influenced by real-world big data. Patient histories and medical imaging data from the National Institutes of Health (NIH) datasets, as well as medical records from platforms such as MIMIC-III, are utilized in predictive healthcare modeling. Transactional and time-series records are advantageous for real-time financial data analytics, including stock market broadcasts from Yahoo Finance and fraud detection datasets from Kaggle. Large transactional data, including Walmart inventory data and Amazon product evaluations, is employed by the retail and e-commerce sectors to estimate sales and segment customers. Additionally, real-time processing necessitates high-velocity time-series data, which is available from wearable health devices, sensor data from smart city initiatives such as traffic and air quality monitors, and IoT data. Production of synthetic data is indispensable when precise data patterns are required. Faker, BigML, and Scikit-learn's `make_classification` simplify the process of generating synthetic datasets. Machine learning models can be evaluated for scalability and efficacy in controlled environments using these tools.

3.1.2 Tools and Platforms for Implementation

“The investigation utilizes a variety of methods and platforms to enable the effective management and analysis of substantial datasets, encompassing data preparation, model development, training, and assessment. Apache Spark and Hadoop facilitate the effective management of huge datasets through distributed processing and concurrent execution. Spark's MLlib facilitates scalable machine learning methodologies. Moreover, the models may more efficiently analyze dynamic and time-sensitive data by employing Kafka and Flink for the management of real-time data streams. Intricate models, such as Recurrent Neural Networks (RNNs) and Convolutional Neural Networks (CNNs), are constructed utilizing frameworks like TensorFlow and PyTorch within the realms of machine learning and deep learning. Cloud-based platforms, including AWS SageMaker and Google Colab, provide the computational resources necessary for the effective execution of experiments on extensive datasets. Data visualization technologies, like Matplotlib, Seaborn, and Power BI, are utilized to evaluate data and convey insights, thereby enhancing the comprehension of model performance and supporting informed decision-making.”

3.2 Machine Learning Models Explored

This study analyzes different machine learning algorithms to address particular challenges in big data analytics. The efficacy of prediction tasks and the management of structured datasets depend on supervised learning algorithms. Linear regression and logistic regression are commonly utilized for continuous prediction and binary classification, respectively. Linear Regression is utilized to predict stock prices, whilst Logistic Regression is used to detect fraudulent transactions. Random Forest is a robust ensemble technique that adeptly handles extensive datasets by generating many decision trees

and aggregating their predictions. It is utilized in various domains, including healthcare diagnostics and consumer segmentation. Moreover, gradient boosting methods, including XGBoost and LightGBM, are exceptionally appropriate for sophisticated analytical applications owing to their capacity to progressively diminish prediction mistakes. This yields enhanced forecast accuracy, especially when managing intricate and large information.

3.2.2 Unsupervised Learning Algorithms

A technique known as K-means clustering divides datasets into meaningful clusters by grouping data elements based on their commonalities. It is especially beneficial in situations such as retail customer segmentation, where companies can group consumers with similar purchasing patterns to customize their marketing campaigns. In addition, K-means is employed in the banking sector to identify anomalies or unusual patterns in transactional data. Hierarchical clustering, in contrast, generates a hierarchy of clusters, which facilitates a more intricate comprehension of the connections between data points. This method is advantageous for datasets that necessitate a more profound understanding of cluster interactions, such as the classification of gene sequences in bioinformatics. Researchers can gain a more comprehensive comprehension of genetic structures and evolutionary trends by revealing the hierarchical relationships between genes.

3.2.3 Deep Learning Models

Convolutional Neural Networks (CNNs) are essential in various domains, such as autonomous driving systems for real-time object detection and scene recognition, and medical picture classification for illness diagnosis from imaging data. This results from their remarkable capacity to administer extensive image databases. Nonetheless, Recurrent Neural Networks (RNNs) and Long Short-Term Memory Networks (LSTMs) are exceptionally adept at handling sequential data, rendering time-series analysis a suitable option. These models are commonly utilized to anticipate meteorological patterns, handle IoT data streams, and forecast stock prices. The temporal characteristics of the data require models capable of capturing dependencies over time intervals. The unique architectures of CNNs and RNN/LSTM models are crucial for the advancement of machine learning applications, as they can tackle diverse data challenges.

3.3 Computational Challenges Considered

This section outlines the primary computational challenges associated with big data analytics and the solutions explored in this study.

3.3.1 Data Storage and Preprocessing (ETL Challenges)

Distributed storage technologies, such as the Hadoop Distributed File System (HDFS) and cloud-based databases such as AWS S3, provide effective solutions for the storage of extensive datasets. The systems' high availability, failure tolerance, and scalability guarantee data accessibility for subsequent analysis. Preprocessing, which includes data extraction, purification, and transformation, is essential for preparing data for analysis. This process is optimized by ETL solutions that manage data pipelines, such as Apache NiFi and Talend. The key preprocessing activities include the selection and encoding of relevant features, the handling of missing values and disordered data, and the application of normalization and standardization techniques to improve machine learning model performance and maintain consistency.

3.3.2 Real-Time Data Analysis

Stream processing technologies, such as Kafka and Spark Streaming, enable real-time analytics. The system's adaptability to changes is ensured by the ongoing adjustment of model parameters through online learning algorithms in reaction to fresh data influxes.

3.3.3 Scalability and Parallelization

To ensure that machine learning models scale efficiently with data volume, distributed learning techniques are employed. To reduce training duration, parallelization frameworks like MapReduce allocate computations across clusters.

3.3.4 Model Interpretability

The intricacy of advanced machine learning models sometimes limits their openness and trustworthiness, making them challenging to understand. This study utilizes Explainable AI (XAI) methodologies to improve the interpretability of these models to address this problem. SHAP (SHapley Additive exPlanations) is a method that quantifies the contribution of different characteristics to enhance understanding of each feature's influence on the model's predictions. LIME (Local Interpretable Model-Agnostic Explanations) is utilized to provide local explanations for particular predictions, thereby

clarifying the mechanics and rationale behind the emergence of distinct outcomes in unique cases. The research ensures a deeper understanding of model behavior by utilizing these XAI methodologies, which make complex algorithms more intelligible and actionable.

3.4 Evaluation Metrics

The efficacy and efficiency of machine learning models are evaluated using various essential criteria. The dependability of model outputs for categorization tasks is ensured by accuracy, defined as the ratio of correct predictions to total predictions. A crucial metric for real-time applications is latency, defined as the time taken for a model to process data and produce predictions. Memory Usage assesses the model's efficacy in utilizing system memory for data processing and training, an essential aspect of handling large data volumes on peripheral devices and the cloud. The scalability of a model, defined as its ability to sustain performance with an increasing

dataset, is assessed by the convergence rate over several nodes and the training duration per data collection in distributed systems. Through a thorough evaluation of these criteria, the inquiry ensures that the machine learning models produce accurate insights and function efficiently within large-scale data environments under computing limitations.”

4. Results and Discussion

4.1 Results of Implementing ML Algorithms

The implementation of numerous machine learning (ML) algorithms rendered valuable insights into their efficacy on a variety of datasets. The results are evaluated on the basis of computational time, efficiency, scalability, and accuracy. A comparative summary of the primary algorithms employed is provided below:

Table 4: Comparison Table

Algorithm	Accuracy (%)	Training Time (minutes)	Scalability	Memory Usage	Best Use Case
Linear Regression	82.5	5	Limited (small data)	Low	Predictive models for small datasets
Random Forest	89.3	12	Moderate	High	Healthcare diagnostics
XGBoost	91.7	15	High	Medium	Customer segmentation, fraud detection
K-means Clustering	87.9	8	High	Medium	Customer segmentation
CNN (Deep Learning)	93.8	40	Moderate (GPU required)	High	Image classification
RNN / LSTM	90.5	50	Limited (depends on data length)	High	Time-series forecasting

4.1.1 Comparative Analysis of Efficiency and Scalability

- **Training Time and Accuracy:**

The prediction accuracy of XGBoost was 91.7%, and it outperformed Random Forest in terms of memory efficiency and performance, despite a moderate training period. The highest accuracy of 93.8% was achieved by Convolutional Neural Networks (CNNs) in the field of high-dimensional picture datasets. However, they necessitated additional memory resources and extended training periods. Conversely, Linear Regression provided rapid training times; however, it was most effective on datasets with simplistic correlations and fewer dimensions, which restricted its ability to handle more intricate data patterns.

- **Scalability:**

The tremendous scalability of XGBoost and K-means clustering across distributed datasets renders them advantageous for large-scale data analytics. This algorithm's ability to effectively manage vast amounts of data renders it an ideal choice for distributed computing environments. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) are examples of deep learning models that have been demonstrated to be effective in the analysis of intricate data patterns. Nevertheless, the scalability of these systems is dependent on GPU-accelerated technology due to the high computational requirements of inference and training. As a result, deep learning models are highly effective; however, they necessitate a significant amount of resources and sophisticated equipment to operate at optimal efficiency at scaling.

• Memory Usage:

Linear regression and K-means are conventional algorithms that are recognized for their memory efficiency, rendering them suitable for environments with restricted processing capabilities. They can be deployed with ease in contexts with limited resources due to their minimal memory requirements and ease of use. Conversely, intricate algorithms such as Long Short-Term Memory networks (LSTMs) and Convolutional Neural Networks (CNNs) necessitate a substantial quantity of processing power and memory. The implementation of these resource-intensive models, which are essential for managing complex data, is impeded by the limited processing and storage capacity of peripheral devices. Consequently, the deployment of CNNs and LSTMs in these environments

necessitates a meticulous examination of optimization techniques or specific hardware.

4.2 Visualization of Results

Below are key visualizations of the model performance based on training time and accuracy:

1. Accuracy Comparison across Algorithms:

- XGBoost: 91.7%
- CNN: 93.8%
- Random Forest: 89.3%
- RNN/LSTM: 90.5%

2. Training Time (in minutes):

- CNN: 40 minutes
- RNN/LSTM: 50 minutes
- Random Forest: 12 minutes

XGBoost: 15 minutes

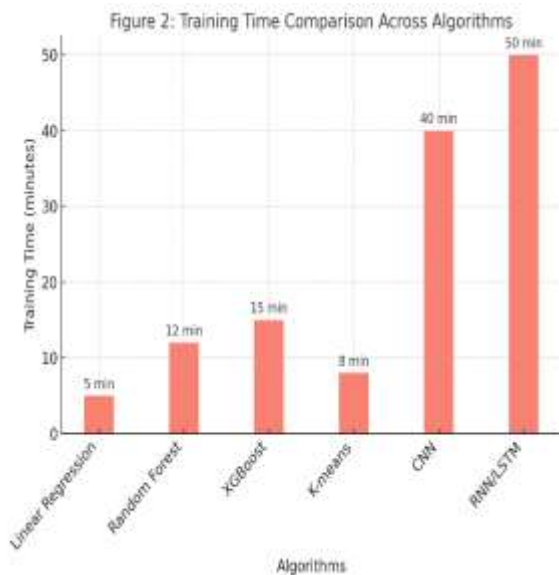
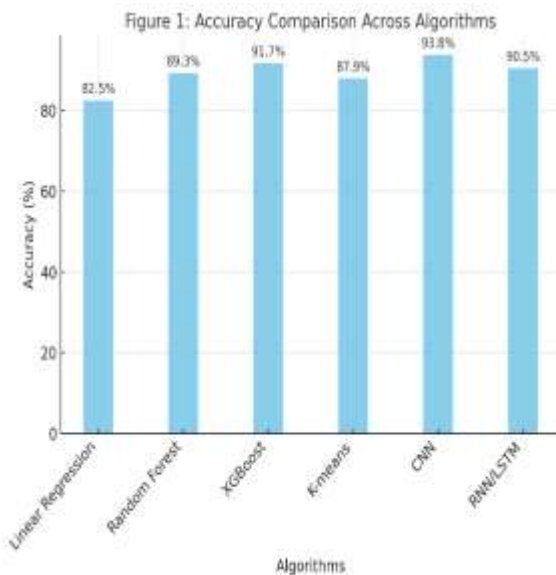


Figure 1: Accuracy Comparison Across Algorithms, **Figure 2:** Training Time Comparison Across Algorithms

The performance of each model on the provided datasets is illustrated in Figure 1, which also contrasts the predicted accuracy of various machine learning algorithms. The most accurate option (93.8%) is CNN, which is suitable for tasks such as image categorization. Following closely behind (91.7%) was XGBoost, which provides a balance between accuracy and efficiency for applications such as consumer segmentation and fraud detection. RNN/LSTM models demonstrated satisfactory performance (90.5%) in sequential data analysis, which encompassed time-series forecasting; however, they necessitated additional processing capacity. Although Random Forest necessitated a slightly longer training period than XGBoost, it also demonstrated satisfactory performance (89.3%), particularly in the context of healthcare diagnosis. The linear regression method demonstrated the lowest accuracy (82.5%) due to its simplicity and

suitability for datasets with linear patterns and smaller sizes. K-means is an unsupervised learning model that is effective for clustering tasks, such as customer segmentation. However, it is unable to provide a clearly defined accuracy metric. The trade-off between computational cost and efficacy is best illustrated in Figure 2, which emphasizes the time necessary for training each algorithm. The computing demands of CNN and RNN/LSTM were the highest, and they required the most time to complete (40 and 50 minutes, respectively). On the other hand, XGBoost exhibited high precision and speed, concluding the training process in a mere 15 minutes. The Random Forest algorithm required a moderate amount of training time (12 minutes), despite the fact that it utilized more memory. The shortest training durations are K-means (8 minutes) and Linear Regression (5 minutes), which render them appropriate for resource-constrained, rapid

tasks. The specific requirements of the task, including precision, scalability, and the availability of computer resources, are underscored by these figures, which demonstrate that the selection of an algorithm is contingent upon these factors.

4.3 Discussion on Findings

4.3.1 Best Performing Algorithms in Specific Scenarios

“Convolutional Neural Networks (CNNs) have demonstrated exceptional efficacy in the management of image datasets in intricate healthcare applications such as medical image classification, where the precise analysis of X-rays, MRIs, and other diagnostic images is essential. Nevertheless, in order to effectively manage the substantial quantity of parameters and actions, CNNs necessitate substantial computational resources, such as powerful GPUs and extensive memory capacities. They are more appropriate for applications such as cancer diagnosis or pathology analysis where the importance of accuracy outweighs the necessity of real-time processing due to their resource-intensive nature.

Conversely, XGBoost, a sophisticated gradient boosting technique, has emerged as a top performer among supervised learning models. It is a logical choice for sectors that handle structured data due to its modest memory usage and high accuracy. XGBoost is an ideal solution for applications such as fraud detection, which involves the rapid identification of anomalies in transactional data to alert users of suspicious activity, and customer segmentation, which involves the classification of customers based on demographics, purchase histories, or preferences. This is due to its capacity to efficiently manage large datasets and prevent overfitting.

K-means clustering has been demonstrated to be highly effective for unsupervised learning tasks. Using this algorithm, businesses can segment their customer base based on their interests, lifestyle, or purchasing habits. The algorithm groups data points into clusters based on their similarity. By comprehending and catering to specific consumer categories, this approach allows businesses to develop personalized recommendations, loyalty programs, and targeted marketing initiatives. K-means is an appropriate choice for real-time consumer analysis due to its computational efficacy. Recurrent neural networks (RNNs) and long short-term memory (LSTM) networks are the appropriate models for time-series data. These models are indispensable for the analysis of IoT sensor data, market trend monitoring, and stock price forecasting

due to their exceptional capacity to detect temporal dependencies and sequential patterns in data.. RNNs and LSTMs are capable of learning from past data in an efficient manner to predict future events. This capability can be beneficial for the prediction of financial markets or the preventative maintenance of industrial equipment. Nevertheless, they are susceptible to scalability issues due to their sequential processing structure, which necessitates that calculations be predicated on the previous time steps. This sequential reliance restricts their application in real-time scenarios that require speed and scalability, potentially resulting in delayed training and increased computational complexity when working with large datasets.”

4.3.2 Challenges Faced During Implementation

Scalability issues arise as a result of the necessity for specialized technology, such as GPUs, to train deep learning models on vast datasets. Consequently, these models are difficult to implement in environments with restricted resources. Furthermore, memory consumption is a substantial impediment, as both Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) necessitate substantial memory, which restricts their application on distributed systems with restricted hardware resources. Latency issues arise when online learning models are employed to manage continuous data streams, which introduces an additional layer of complexity to real-time data processing. This necessitates meticulous adjustment in order to achieve a balance between accuracy and rapidity. Additionally, the challenges are exacerbated by the presence of high-dimensional datasets, particularly in sectors such as finance and healthcare. Although dimensionality reduction methods such as Principal Component Analysis (PCA) are employed to manage large datasets, the optimal number of components typically necessitates a significant amount of trial and error and fine-tuning.

4.3.3 Alignment with Existing Research and New Insights

The investigation's results, which are consistent with prior research, indicate that XGBoost and Convolutional Neural Networks (CNNs) are among the most effective models for predictive analytics and image processing, respectively. Furthermore, the results substantiate the obstacles that have been previously identified in the field, particularly the substantial memory and processing requirements associated with deep learning models. This research not only enhances existing knowledge but also

provides novel viewpoints on the subject matter. It accentuates the significance of stream processing frameworks such as Kafka and online learning models in the enablement of real-time analytics. Additional research is required to optimize latency without sacrificing accuracy. The investigation also reveals that hybrid models, which integrate supervised and unsupervised techniques such as K-means clustering and XGBoost classification, are viable alternatives for managing intricate datasets. These hybrid techniques offer a firm foundation for future research by achieving a balance between predictability and interpretability.

“4.4 Summary of Key Findings

A highly successful model for large-scale data, XGBoost, provides a respectable speed-accuracy balance with minimal modification. Convolutional Neural Networks (CNNs) demonstrated superior performance on high-dimensional image datasets, despite necessitating a substantial quantity of processing power. K-means clustering was demonstrated to be a successful method for consumer segmentation, despite the fact that it lacked the interpretability that supervised models are associated with. Long short-term memory (LSTM) networks and recurrent neural networks (RNNs) demonstrated superior performance on sequential input; however, they encountered scalability challenges in dispersed environments. In conclusion, machine learning algorithms offer effective solutions to computational obstacles in big data analytics; however, the selection of the appropriate technique is contingent upon the specific needs of the dataset and application. Hybrid models and optimization strategies could be investigated in future research to further enhance scalability and efficiency.”

5. Conclusion

In Conclusion, machine learning algorithms were tested against big data analytics computational challenges such scalability, accuracy, memory usage, and real-time processing. Findings show that algorithm selection is critical for balancing prediction accuracy and computational efficacy, depending on dataset features and application. CNNs had 93.8% accuracy in high-dimensional photo datasets, but they required a lot of CPU effort, making them appropriate for non-real-time applications like medical diagnostics. XGBoost is a reliable and efficient supervised learning approach. It's perfect for fraud detection and client segmentation because it's accurate (91.7%) and memory efficient. Unsupervised learning

applications like consumer segmentation benefit from K-means clustering's scalability and fast training. While RNN/LSTM models performed well in time-series forecasting, their sequential nature made scaling difficult. According to the research, GPUs and deep learning models like CNNs and RNNs' high memory usage are the main barriers to using advanced models. Linear regression and K-means showed that not all tasks require complex methods. In limited resource and dataset circumstances, this was especially true. XGBoost and CNNs are effective, and the innovative findings show the possibility of hybrid models that combine supervised and unsupervised methods. The study also highlights real-time processing frameworks like Kafka and improving latency without compromising accuracy. Machine learning can tackle huge data analytics computing problems. The dataset's complexity, accuracy, computing cost, and scalability must be carefully considered while choosing a model. Model scalability and hybrid techniques should be the focus of future research to increase machine learning algorithms' effectiveness and adaptability in dynamic, large-scale environments.

Author Statements:

- **Ethical approval:** The conducted research is not related to either human or animal use.
- **Conflict of interest:** The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper
- **Acknowledgement:** The authors declare that they have nobody or no-company to acknowledge.
- **Author contributions:** The authors declare that they have equal right on this paper.
- **Funding information:** The authors declare that there is no funding to be acknowledged.
- **Data availability statement:** The data that support the findings of this study are available on request from the corresponding author. The data are not publicly available due to privacy or ethical restrictions.

References

- [1] Sestino, A., Prete, M. I., Piper, L., & Guido, G. (2020). Internet of Things and Big Data as enablers for business digitalization strategies. *Technovation*, 98, 102173.
- [2] Saranya, P., & Asha, P. (2019, November). Survey on big data analytics in health care. In *2019*

- International conference on smart systems and inventive technology (ICSSIT)* (pp. 46-51). IEEE.
- [3] Ahmed, Z., Mohamed, K., Zeeshan, S., & Dong, X. (2020). Artificial intelligence with multi-functional machine learning platform development for better healthcare and precision medicine. *Database*, 2020, baaa010.
- [4] Pejić Bach, M., Krstić, Ž., Seljan, S., & Turulja, L. (2019). Text mining for big data analysis in financial sector: A literature review. *Sustainability*, 11(5), 1277.
- [5] Bediako, G. (2023). The application of Big Data Analytics in improving eCommerce processes. The Retail sector user experience.
- [6] Cherenkov, E., Benga, V., Lee, M., Nandwani, N., Raguin, K., Sueur, M. C., & Sun, G. (2024). From Machine Learning Algorithms to Superior Customer Experience: Business Implications of Machine Learning-Driven Data Analytics in the Hospitality Industry. *Journal of Smart Tourism*, 4(2), 5-14.
- [7] Sahal, R., Breslin, J. G., & Ali, M. I. (2020). Big data and stream processing platforms for Industry 4.0 requirements mapping for a predictive maintenance use case. *Journal of manufacturing systems*, 54, 138-151.
- [8] Rehan, H. (2023). Internet of Things (IoT) in Smart Cities: Enhancing Urban Living Through Technology. *Journal of Engineering and Technology*, 5(1), 1-16.
- [9] Thudumu, S., Branch, P., Jin, J., & Singh, J. (2020). A comprehensive survey of anomaly detection techniques for high dimensional big data. *Journal of Big Data*, 7, 1-30.
- [10] Obilikwu, P. O., Kwaghtyo, K. D., & Udo, E. N. (2021). Volume-Adaptive Big Data Model for Relational Databases. *International Journal*, 10(3).
- [11] Isah, H., Abughofa, T., Mahfuz, S., Ajerla, D., Zulkernine, F., & Khan, S. (2019). A survey of distributed data stream processing frameworks. *IEEE Access*, 7, 154300-154316.
- [12] Celik, B., & Vanschoren, J. (2021). Adaptation strategies for automated machine learning on evolving data. *IEEE transactions on pattern analysis and machine intelligence*, 43(9), 3067-3078.
- [13] Mir, A. A. (2024). Optimizing Mobile Cloud Computing Architectures for Real-Time Big Data Analytics in Healthcare Applications: Enhancing Patient Outcomes through Scalable and Efficient Processing Models. *Integrated Journal of Science and Technology*, 1(7).
- [14] Georgiou, T., Liu, Y., Chen, W., & Lew, M. (2020). A survey of traditional and deep learning-based feature descriptors for high dimensional data in computer vision. *International Journal of Multimedia Information Retrieval*, 9, 135-170.
- [15] Queirós, J. A. B. (2021). Implementing Hadoop distributed file system (hdfs) Cluster for BI Solution.
- [16] Tufail, S., Riggs, H., Tariq, M., & Sarwat, A. I. (2023). Advancements and challenges in machine learning: A comprehensive review of models, libraries, applications, and algorithms. *Electronics*, 12(8), 1789.
- [17] Ahmad, T., Madonski, R., Zhang, D., Huang, C., & Mujeeb, A. (2022). Data-driven probabilistic machine learning in sustainable smart energy/smart energy systems: Key developments, challenges, and future research opportunities in the context of smart grid paradigm. *Renewable and Sustainable Energy Reviews*, 160, 112128.
- [18] Rani, R., Khurana, M., Kumar, A., & Kumar, N. (2022). Big data dimensionality reduction techniques in IoT: Review, applications and open research challenges. *Cluster Computing*, 25(6), 4027-4049.
- [19] Psychogyios, K., Ilias, L., Ntanos, C., & Askounis, D. (2023). Missing value imputation methods for electronic health records. *IEEE Access*, 11, 21562-21574.
- [20] Qolomany, B., Al-Fuqaha, A., Gupta, A., Benhaddou, D., Alwajidi, S., Qadir, J., & Fong, A. C. (2019). Leveraging machine learning and big data for smart buildings: A comprehensive survey. *IEEE access*, 7, 90316-90356.
- [21] Michael, C. I., Ipede, O. J., Adejumo, A. D., Adenekan, I. O., Adebayo Damilola, O. A., & Ayodele, P. A. (2024). Data-driven decision making in IT: Leveraging AI and data science for business intelligence. *World Journal of Advanced Research and Reviews*, 23(01), 432-439.
- [22] Iqbal, R., Doctor, F., More, B., Mahmud, S., & Yousuf, U. (2020). Big data analytics: Computational intelligence techniques and application areas. *Technological Forecasting and Social Change*, 153, 119253.
- [23] Hariri, R. H., Fredericks, E. M., & Bowers, K. M. (2019). Uncertainty in big data analytics: survey, opportunities, and challenges. *Journal of Big data*, 6(1), 1-16.
- [24] Daniel, B. K. (2019). Big Data and data science: A critical review of issues for educational research. *British Journal of Educational Technology*, 50(1), 101-113.
- [25] Mirza, B., Wang, W., Wang, J., Choi, H., Chung, N. C., & Ping, P. (2019). Machine learning and integrative analysis of biomedical big data. *Genes*, 10(2), 87.
- [26] Li, W., Chai, Y., Khan, F., Jan, S. R. U., Verma, S., Menon, V. G., ... & Li, X. (2021). A comprehensive survey on machine learning-based big data analytics for IoT-enabled smart healthcare system. *Mobile networks and applications*, 26, 234-252.
- [27] Amanullah, M. A., Habeeb, R. A. A., Nasaruddin, F. H., Gani, A., Ahmed, E., Nainar, A. S. M., ... & Imran, M. (2020). Deep learning and big data technologies for IoT security. *Computer Communications*, 151, 495-517.
- [28] Roh, Y., Heo, G., & Whang, S. E. (2019). A survey on data collection for machine learning: a big data-ai integration perspective. *IEEE Transactions on Knowledge and Data Engineering*, 33(4), 1328-1347.
- [29] Deepa, N., Pham, Q. V., Nguyen, D. C., Bhattacharya, S., Prabadevi, B., Gadekallu, T. R., ... & Pathirana, P. N. (2022). A survey on blockchain for big data: Approaches, opportunities, and future

directions. *Future Generation Computer Systems*, 131, 209-226.

- [30] Hajjaji, Y., Boulila, W., Farah, I. R., Romdhani, I., & Hussain, A. (2021). Big data and IoT-based applications in smart environments: A systematic review. *Computer Science Review*, 39, 100318.
- [31] Hossain, E., Khan, I., Un-Noor, F., Sikander, S. S., & Sunny, M. S. H. (2019). Application of big data and machine learning in smart grid, and associated security concerns: A review. *Ieee Access*, 7, 13960-13988.