



# The Important of Big Data in Machine Learning

Luu Tan Thanh <sup>a\*</sup>, Nguyen Quang Dat <sup>a</sup>, Vu Hoang <sup>a</sup>  
and To Hien Huy Hieu <sup>b</sup>

<sup>a</sup> Hanoi University of Science, Vietnam National University, Hanoi, Vietnam.  
<sup>b</sup> Le Quy Don Technical University, Vietnam.

## Authors' contributions

*This work was carried out in collaboration among all authors. All authors read and approved the final manuscript.*

## Article Information

DOI: <https://doi.org/10.56557/jobari/2024/v30i68952>

## Open Peer Review History:

This journal follows the Advanced Open Peer Review policy. Identity of the Reviewers, Editor(s) and additional Reviewers, peer review comments, different versions of the manuscript, comments of the editors, etc are available here: <https://prh.ikpress.org/review-history/12489>

**Short Communication**

**Received: 28/08/2024**  
**Accepted: 02/11/2024**  
**Published: 20/11/2024**

## ABSTRACT

The integration of Big Data and Machine Learning is revolutionizing various industries by enabling smarter decision-making, enhancing automation, and improving predictive analytics. Big Data plays a pivotal role in Machine Learning by providing large volumes of diverse, real-time data that fuel the learning process. The availability of vast amounts of data allows Machine Learning models to be trained on complex patterns, leading to better accuracy, improved generalization, and more reliable predictions. Moreover, Big Data facilitates the use of advanced techniques like deep learning, which require massive datasets for tasks such as image recognition, natural language processing, and recommendation systems.

However, the role of Big Data extends beyond mere volume; it also offers variety, providing diverse datasets essential for building more robust and adaptable Machine Learning models. The continuous stream of real-time data enables dynamic learning, while data diversity enhances model versatility. Despite these benefits, challenges such as data quality, processing scalability, and

\*Corresponding author: E-mail: [nguyenquangdat@hus.edu.vn](mailto:nguyenquangdat@hus.edu.vn);

privacy concerns must be addressed. In summary, Big Data significantly amplifies the capabilities of Machine Learning by enhancing model performance, driving innovations, and enabling applications across domains such as healthcare, finance, retail, and autonomous systems.

*Keywords: Big data; machine learning (ML); model performance; autonomous systems.*

## 1. INTRODUCTION

In today's rapidly evolving world of data-driven technologies, Big Data serves as a foundational element that is significantly advancing the field of Machine Learning (ML). As businesses and industries generate data at an unprecedented scale, the integration of Big Data with ML is reshaping various sectors by enabling a new level of precision in predictions, optimizing decision-making processes, and automating tasks that were once dependent on manual interventions. This convergence leverages the enormous volumes, variety, and velocity of data, transforming it into actionable insights that fuel innovations across domains such as healthcare, finance, retail, autonomous systems, and beyond (Mayer-Schönberger and Cukier 2013).

Big Data's impact on ML extends far beyond simply providing a large amount of information for training models. It facilitates the creation of complex, high-dimensional datasets that are essential for training sophisticated ML models like deep learning networks, which require extensive data to recognize patterns and improve generalization. The diversity in data types – ranging from structured datasets (e.g., relational databases) to unstructured content (e.g., social media posts, images, videos) – enhances the versatility of ML applications, making it possible for algorithms to tackle tasks such as natural language processing, image recognition, and predictive maintenance with greater accuracy.

The benefits of Big Data in Machine Learning are multi-faceted. On one hand, the availability of large-scale datasets helps in reducing issues like overfitting by providing more scenarios for models to learn from, thereby enhancing their performance on unseen data. On the other hand, it enables the use of advanced techniques like transfer learning, where pre-trained models on massive datasets can be fine-tuned for specific tasks, boosting efficiency and performance in specialized applications. Additionally, Big Data facilitates real-time learning and dynamic decision-making, making it possible to implement systems that adapt to changing conditions and evolving trends (Géron 2019).

However, the integration of Big Data and ML also presents unique challenges that must be addressed to maximize their potential. Issues related to data quality, such as the presence of noise and missing values, require robust data preprocessing techniques. The sheer volume of data necessitates scalable processing solutions, often involving distributed computing frameworks like Hadoop and Spark. Moreover, the rapid generation of data (velocity) calls for efficient data management strategies to handle the continuous influx of information.

The infrastructure supporting Big Data also plays a crucial role in ML. Technologies such as NoSQL databases, cloud computing platforms, and parallel processing frameworks allow for the storage and processing of vast datasets, making it feasible to conduct ML experiments at scale. Furthermore, the emergence of new computing paradigms like edge computing and federated learning is helping to distribute the data processing load, enabling ML models to learn from decentralized data sources while maintaining privacy and reducing latency (Dasgupta 2018).

As the synergy between Big Data and Machine Learning continues to evolve, it is driving transformative applications across industries. In healthcare, for example, predictive analytics using Big Data is improving patient outcomes by enabling personalized treatment plans based on individual health data. In finance, Big Data-powered ML models help detect fraudulent activities in real time, while in retail, it supports personalized marketing strategies that cater to individual consumer preferences. Autonomous systems, such as self-driving cars, rely on the vast amounts of sensory data processed by ML algorithms to navigate and make decisions (Minelli et al., 2013).

The future of data-driven innovation lies in the effective use of Big Data to not only train and improve Machine Learning models but also to derive meaningful insights that can guide strategic decision-making. As organizations continue to collect and analyze massive datasets, the role of Big Data in shaping the next

generation of intelligent systems will only become more pronounced, setting the stage for a more data-centric world where Machine Learning thrives on the power of information (ML Journey 2024).

## 2. MATERIALS AND DISCUSSION

Big data refers to extremely large and complex datasets that are difficult to process, store, or analyze using traditional data management tools and techniques. These datasets come from a wide range of sources, including social media, sensors, web traffic, transaction records, and more (Mayer-Schönberger and Cukier 2013, Géron 2019). The three main characteristics of big data, often referred to as the 3 Vs, are:

1. **Volume:** The sheer amount of data, often in petabytes or exabytes.
2. **Velocity:** The speed at which new data is generated, collected, and processed.
3. **Variety:** The wide range of data types, including structured data (like databases), semi-structured data (like JSON), and unstructured data (like images, videos, or text).

Sometimes, additional Vs such as **Veracity** (referring to the accuracy and trustworthiness of data) and **Value** (the usefulness of the data) are also mentioned.

Uses of Big Data (Mayer-Schönberger and Cukier 2013)

Big data is used in various industries for a wide range of purposes, including (Minelli et al., 2013):

- **Predictive analytics:** Predicting trends, customer behavior, or potential failures in machinery.
- **Business intelligence:** Gaining insights from customer data, improving products, and services.
- **Healthcare:** Analyzing patient data to improve diagnostics and treatment.
- **Marketing:** Tailoring ads and campaigns to specific audiences using data-driven insights.

### Technologies for Handling Big Data:

To manage and process big data, specialized tools and platforms have been developed, such as (Dasgupta 2018):

- **Hadoop:** A distributed storage and processing framework.
- **Spark:** A fast, in-memory data processing engine.
- **NoSQL Databases:** Like MongoDB or Cassandra, for handling large volumes of semi-structured or unstructured data.
- **Cloud Platforms:** Like AWS, Google Cloud, and Azure, which offer scalable storage and analytics tools.

Big data helps organizations make data-driven decisions and gain insights that were previously impossible due to data limitations.



Fig. 1. Salient feature of big data

**Big Data's role in machine learning:** Big Data plays a critical role in Machine Learning by providing the vast amounts of data required to train models effectively, leading to more sophisticated and accurate predictions (Géron 2019). The availability of large datasets allows machine learning models, especially those based on deep learning, to learn from a wider array of patterns and make better generalizations, which results in enhanced accuracy and performance. Furthermore, Big Data helps reduce the risk of overfitting that can occur when models are trained on smaller datasets, as it allows the model to encounter a greater variety of scenarios and noise, improving its ability to generalize to new, unseen examples.

Big Data also facilitates the handling of complex and high-dimensional data, benefiting from diverse data types such as structured, unstructured, and semi-structured formats including text, images, audio, video, and sensor data. This diversity expands the range of applications for machine learning, enabling models to tackle tasks such as natural language processing, image recognition, and recommendation systems with greater effectiveness. When dealing with tasks that involve a large number of features or dimensions, like image recognition or natural language processing, Big Data provides the necessary scale to train models to manage this complexity effectively.

The ability of Big Data to boost deep learning models is particularly significant, as deep learning architectures like Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) require extensive labeled datasets to achieve high accuracy. Additionally, even when datasets are not fully labeled, Big Data supports unsupervised and semi-supervised learning techniques, enabling models to derive insights and structure from vast amounts of unlabeled data.

Moreover, Big Data enhances personalization and predictive power in machine learning applications. The ability to analyze and process large volumes of data in real-time is essential for applications such as recommendation systems, fraud detection, and predictive maintenance. With more information about user behavior, models can provide personalized recommendations, advertisements, or experiences, thus delivering a more tailored and engaging user experience. The large-scale data

availability enables a nuanced understanding of users and facilitates more accurate and timely predictions.

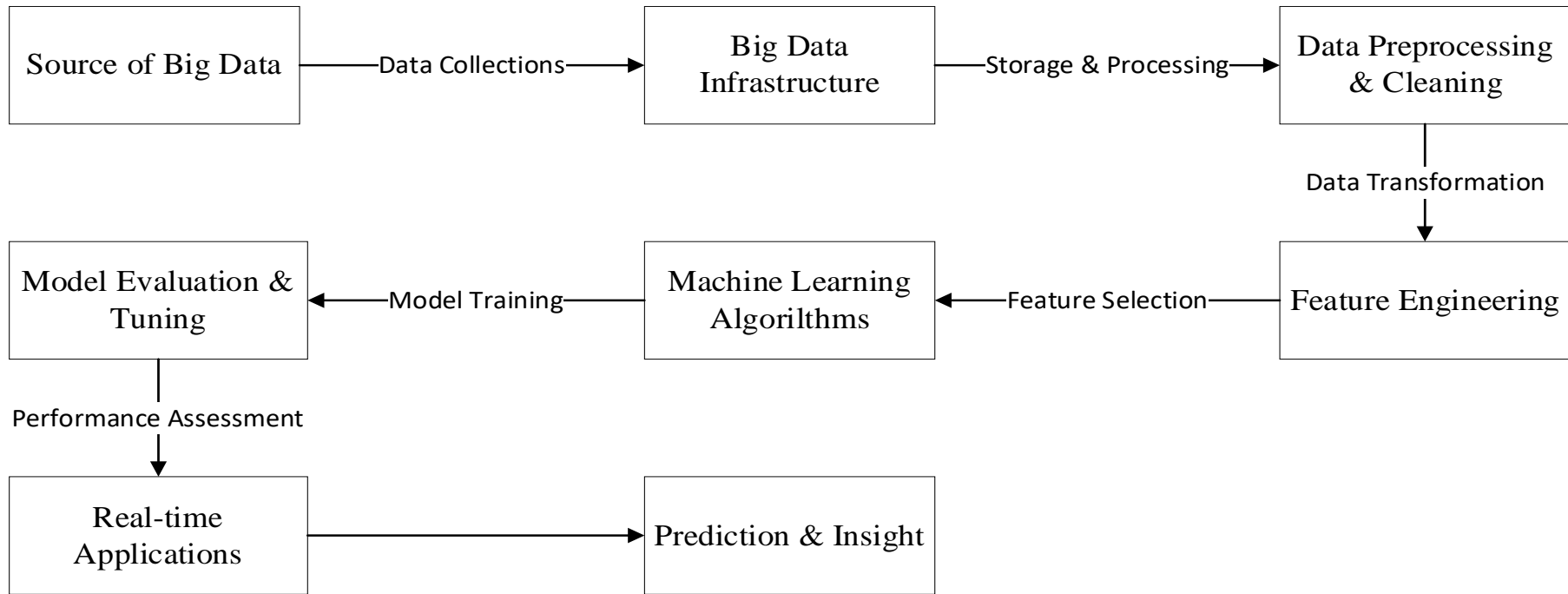
Another crucial contribution of Big Data to machine learning is in enabling transfer learning and model pretraining. Models can be pre-trained on large datasets, such as ImageNet for computer vision tasks or GPT for natural language processing, before being fine-tuned for more specific tasks with smaller, task-focused datasets. This approach allows for the development of generalized models that can be adapted to specific tasks with minimal additional training, thereby increasing efficiency and reducing resource requirements.

Big Data's role extends to supporting data-driven decision-making by empowering organizations to derive actionable insights from customer preferences, operational inefficiencies, or market trends. This integration not only informs better strategic decisions but also helps automate various business processes, including predictive analytics, customer service via chatbots, and anomaly detection, all powered by the large-scale data available.

In the realm of AI development, Big Data serves as the foundational fuel for machine learning algorithms that enable systems to "learn" and adapt over time. This is especially important for AI applications such as autonomous vehicles and smart home devices, which rely on extensive data to continuously improve their capabilities.

Big Data also contributes to handling noisy and incomplete data, as the sheer volume of information allows machine learning models to become more robust by learning from imperfect data. Even if individual data points are flawed, models can still identify useful patterns by averaging over the large datasets.

Finally, scalability and computational advancements in Big Data technologies such as Hadoop, Spark, and distributed machine learning frameworks like TensorFlow and PyTorch support the efficient processing of large datasets (Leskovec et al., 2014). These tools allow for parallel processing, which significantly reduces training time for machine learning models and facilitates rapid experimentation. The combination of Big Data and distributed computing enables machine learning algorithms to handle the massive scale of data required for advanced applications, pushing the boundaries of what these technologies can achieve.



**Fig. 2. Role of Big Data in Machine Learning**

To help you better visualize the role of Big Data in Machine Learning, here is a suggestion for a basic diagram structure:

**In there:**

- **Sources of Big Data:** Big data comes from many different sources such as IoT sensors, social networks, financial transactions, system logs and other data collection systems.
- **Big Data Infrastructure:** Platforms like Hadoop, Spark, NoSQL, and Cloud Storage help store and process large volumes of data efficiently.
- **Data Preprocessing & Cleaning (Data preprocessing and cleaning):** Raw data often contains many errors, noise, and omissions, so it needs to be cleaned and converted into a form that can be used in machine learning.
- **Feature Engineering:** Transform data into features that can be used to train machine learning models.
- **Machine Learning Algorithms:** Machine learning uses algorithms (supervised learning, unsupervised learning, deep learning) to learn from big data.
- **Model Evaluation & Tuning (Model evaluation and tuning):** After training, the model needs to be evaluated and fine-tuned to optimize performance (precision, recall, etc.).
- **Real-time Applications:** Once the model is optimized, it can be used for prediction or decision making in applications such as recommendation systems, fraud detection, and autonomous systems (Kumar et al., 2018).

### 3. CONCLUSION

Big Data and Machine Learning form a synergistic relationship that drives advancements across a wide range of industries, transforming how organizations operate and make decisions (Mayer-Schönberger and Cukier 2013). The massive amounts of data generated from various sources—such as IoT devices, social media, and business transactions—serve as the foundational raw material for machine learning algorithms. This wealth of data allows machine learning models to learn from diverse and complex datasets, leading to the development of more accurate, reliable, and insightful predictions. By leveraging Big Data, machine learning can

achieve a level of precision that is essential for tasks like detecting fraudulent activities, forecasting demand, or recommending personalized products and services.

Machine learning, in turn, is one of the most effective tools for unlocking the potential of Big Data. It not only identifies patterns and trends within the data but also enables organizations to automate decision-making processes and optimize operations in real time (Géron 2019). The ability to process and analyze large datasets empowers companies to derive actionable insights that can enhance customer experiences, reduce operational costs, and drive product innovation. This makes Big Data and machine learning indispensable for businesses aiming to stay competitive in today's data-driven world.

Moreover, the combined power of Big Data and Machine Learning extends beyond traditional business applications. In fields such as healthcare, it can lead to early detection of diseases and personalized treatment plans, while in the field of autonomous vehicles, it contributes to real-time navigation and safety. The integration of these technologies is also critical in tackling global challenges like climate change, where predictive models can help optimize resource usage and anticipate environmental risks (Minelli et al., 2013, Hastie et al., 2009).

In conclusion, the convergence of Big Data and Machine Learning is paving the way for a future where data-driven innovation is at the heart of progress. Their complementary nature not only drives technological advancements but also provides the tools necessary for solving some of the world's most complex problems, making them essential components of the modern technological landscape (Salkuti 2000). Together, they empower industries to not just react to change, but to predict and shape the future.

### DISCLAIMER (ARTIFICIAL INTELLIGENCE)

Author(s) hereby declare that NO generative AI technologies such as Large Language Models (ChatGPT, COPILOT, etc) and text-to-image generators have been used during writing or editing of this manuscript.

### COMPETING INTERESTS

Authors have declared that no competing interests exist.

## REFERENCES

- Mayer-Schönberger, V., & Cukier, K. (2013). *Big data: A revolution that will transform how we live, work, and think*. Houghton Mifflin Harcourt.
- Géron, A. (2019). *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build intelligent systems*. O'Reilly Media.
- Minelli, M., Chambers, M., & Dhiraj, A. (2013). *Big data, big analytics: Emerging business intelligence and analytic trends for today's businesses*. Wiley.
- Leskovec, J., Rajaraman, A., & Ullman, J. (2014). *Mining of massive datasets* (2nd ed.). Cambridge University Press.
- Dasgupta, N. (2018). *Practical big data analytics: Hands-on techniques to implement big data analytics using Hadoop, Spark, NoSQL, and data science*. Packt Publishing.
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: Data mining, inference, and prediction* (2nd ed.). Springer.
- ML Journey. (2024, November 4). *Big data vs machine learning: Differences and interplay*. ML Journey. <https://mljourney.com/big-data-vs-machine-learning-differences-and-interplay/>
- Kumar, A., Agarwal, S., & Gupta, A. (2018). A comparative study on big data analytics. *International Journal of Applied Research*, 4(10), 210–213. <https://www.allresearchjournal.com/archives/2018/vol4issue10/PartC/10-1-27-785.pdf>
- Salkuti, S. R. (2020). A novel approach for big data analytics in predictive systems. *International Journal of Electrical and Computer Engineering*, 10(5), 575–580. <https://ijece.iaescore.com/index.php/IJECE/article/view/19184/13525>

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of the publisher and/or the editor(s). This publisher and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.

© Copyright (2024): Author(s). The licensee is the journal publisher. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

*Peer-review history:*  
The peer review history for this paper can be accessed here:  
<https://prh.ikpress.org/review-history/12489>