

Review of Machine Learning Method for Resolving Issues of Big Data Analytics

M. Sharma^{1*}, I.S. Sohal², R.M.Singh³, A. Wadhwa⁴, D. Garg⁵

^{1,2,3,4,5}Dept. of Computer Science, Chitkara University, Punjab, India

Corresponding Author: mudita171492.cse@chitkara.edu.in

DOI: <https://doi.org/10.26438/ijcse/v7i4.559563> | Available online at: www.ijcseonline.org

Accepted: 11/Apr/2019, Published: 30/Apr/2019

Abstract—In this technology bound era, data analytics is a decisive way to deal with this enormous amount of data that is getting collected from various sources such as social media, banking, healthcare etc. With the growing volume of this data, it has been getting more and more difficult to analyze the same with the existing techniques. This is where the concept of Machine Learning (ML) has turned out to be an indispensable way for giving this data an intelligent structure i.e. by sorting the clusters of data into data sets and drawing associations from previous information. However, the traditional machine learning methods are not helpful in manipulating the data in a way that we require as we are advancing in these various fields involving big data. In our research we have reviewed the various ML algorithms and learning paradigms for handling the big data problems by associating them with the challenges of the 6 big data dimensions- Volume, Veracity, Velocity, Variety, Visualization and Value. We have studied the similar approach of research given by Alexander *et al.* and Gandomi and Haider. Adding on to their findings and methods we have considered two more V's – Visualization and Value and associated their characteristic challenges with the ML methods. We have mentioned the use of ML in preserving the privacy and security of the data as securing the data being generated is also a significant problem that needs to be addressed.

Keywords---Big data analytics, Machine Learning, V's of big data, algorithms, learning paradigms

I. INTRODUCTION

With the augmentation of technology in the contemporary society, where social media and the internet have made it possible to get a plethora of information at one's fingertips, there has been a marked rise in number of people using the internet and other electrical gadgets. The data thus obtained is in huge numbers and is referred as Big Data. This data is gathered from various sources such as mobile phone applications, emails, videos, click streams [17], social media, health sectors [2], education departments [3]. Analyzation of this data to find the possible patterns and correlations is termed as big data analysis. The characteristics of big data can be divided into categories which are known as the Vs of big data. Most commonly discussed ones are volume, velocity, variety, veracity and volume [11]. These dimensions make it easy to narrow down the challenges that big data faces. To understand the characteristics of the data, some statistical and geometric patterns need to be extracted [11]. Machine Learning is a method to tackle this problem by training the data and giving it an intelligent structure. As the technology has advanced, the conventional methods of machine learning are no longer ideal. Many researches have been presented to overcome the challenges using various approaches. In 2014 Sukumar [5] in his research highlighted

three important requirements for working with large datasets. They include scheming adaptive and highly scalable planning, developing ability and understanding analytical data elements before applying algorithms. In 2015, Najafabadi *et al.* [6] showed how deep learning could solve the general problems of machine learning for big data i.e. unstructured data, noisy data, high scalability of algorithms, streaming data, unlabeled data etc. presented their research focused on signal processing in machine learning. They focused on the five critical problems that include learning of large scale data, types of data, streaming data, uncertain data, low value data and related these to the characteristics of big data. Although these research papers are really informative, the lack of correlation of the challenges with their specific solution makes it difficult to draw conclusions. In 2105 Al-Jarrah *et al.* [4] surveyed machine learning for Big Data and their focus was to improve the effectiveness of wide-ranging systems and also the new algorithmic approaches for reducing memory being used. Al- Jarrah *et al.* mentioned the statistical aspect but the methods for minimizing computational complexity were not evaluated [9]. The approach used by Gandomi and Haider [10] to categorize the challenges according to the big data characteristics makes it easy to understand loopholes in big data analytics. However, they have not discussed the machine learning methods for the

problems. Alexander et al. [9] have used this approach in their research, of categorizing the characteristics of the big data V's, and how different machine learning methods can be used to address the challenges. We have used this same approach in our research work Alexander et al. have discussed about the 4 V's of big data- volume, variety, velocity, veracity, but there are two other important dimensions also that are visualization [33] and value [7]. We have studied their characteristic challenges also and associated them with the machine learning methods. Also security and privacy still remain an issue in Big data analytics. We have surveyed the possible machine learning algorithms resolving this issue which therefore also opens doors for future research work.

II. METHODOLOGY

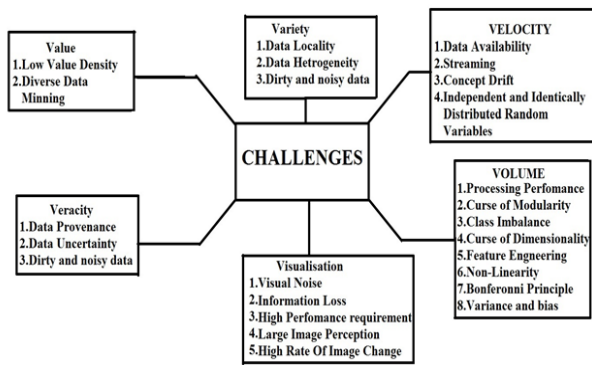


Fig. 1. Big Data characteristics along with associated challenges [9]

Fig. 1

A. Analytical Manipulation of Data

For data analytics using manipulation there are three methods.

These methods are shown in the fig2:

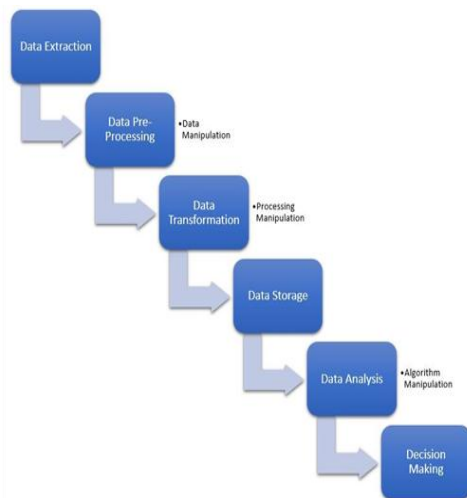


Fig. 2. Big Data Analysis Sequential Diagram [9]

Fig. 2

1. Data manipulations

- The process of organizing the data in order to make it easier to read, differentiate and locate is data manipulation. In order to do that. Two major aspects affecting the machine learning conventional methods in big data are high dimensionality of datasets and huge amount of sample datasets [9]. Therefore, three types of data manipulations – dimensionality reduction, instance selection and data clearing. [9].
- *Dimensionality reduction* - To map high measurement space to lower measurement space is called dimensionality reduction [9]. Techniques used for doing this are Hessian LLE, Isomap, PCA, locally linear embedding (LLE), and Laplacian Eigenmaps. [26]. Another method is Autoencoders for encoding datasets [25]. Dimensionality reduction aims to tackle the problems of large image perception, processing challenges and the dimensionality problems.
- *Instance selection*- Instance selection reduces size of the dataset and acts as a solution to imbalance in class [9]. Its approaches include progressive sampling, using domain knowledge, random selection, cluster sampling and genetic algorithm-based selection [12]. It tackles the problems of processing performance and modularity cure.
- *Data cleaning*- It refers to the processing of wrong or inaccurate information in data. The methods used for this are smoothing filters and wavelet transforms [13]. Even Autoencoders can be used in order to remove the noise from a corrupted input data [9]. It challenges the problems of dirty and noisy data.

2. Processing manipulation

This method centers on how data is stored and worked on [9]. The 3 phases in which the processing happens (is shown in fig 2) are data transformation, data storage, and data estimation. In order to capture the origin of the data so that it can be modified various techniques are used for e.g. Reduce and Map Provenance (RAMP) [14] but these methods have certain computational challenges. Parallelization is a good solution for the processing manipulations. Algorithms such as brute-force search and genetic algorithms, are parallel in their approach and may show a lot of performance improvement. Based on this parallelization, researchers have developed two methods- horizontal and vertical scaling paradigms.

- *Vertical scaling*- In this, graphic processing units (GPUs) are used as a machine learning approach in big data. GPUs were firstly structured for picture display, image processing, matrix operations, and vector operations are suited for such systems [9]. If they can be parallelised other machine learning methods can be implemented on them. However, other algorithms can

TABLE 1. MACHINE LEARNING APPROACHES IN THE FORM OF LEARNING PARADIGMS ALONG WITH THE CHALLENGES THEY ADDRESS [9]

APPROACHES		CHALLENGES																									
		VOLUME			VARIETY	VELOCITY	VERACITY	VISUALISATION		VALUE																	
		Processing Performance	Curse of Dimensionality	Class Imbalance	Curse of Dimensionality	Feature Engineering	Non-linearity	Bertrams' Principle	Variance and Bias	Data locality	Data Heterogeneity	Dirty and noisy Data	Data availability	Real time Processing/Streaming	Concept drift	E.t.d	Data Provenance	Data Uncertainty	Dirty and Noisy Data	Visual noise	Information Loss	High performance requirement	Large Image Perception	High storage usage	Low value density	Reverse Data Mining	
Deep Learning																											
Online Learning	+	+	.					+			.	+	.	+			.					+					
Local Learning	+	+	+					+	+													+					
Transfer Learning			+								+								
Lifelong Learning	+	+									+	+	+	+			.	.	.			+					
Ensemble Learning	+	+												+								+					

Fig. 4

IV. CONCLUSION

To conclude, we made an attempt to survey the existing challenges in the region of big data estimation. Nowadays, machine learning is gaining much attention of the researchers for being an important tool for data scientists to overcome the data related problems. So, we reviewed the latest and improved ML procedures that may be used to solve the problems in Big data analytics. In this age of technology, where data runs in massive numbers ranging between terabytes and zettabytes [11], proper handling of this data can be very beneficial not only for the scientific society but also help in solving real world challenges. This field has a lot of scope for research as even more challenges arise from the vast amounts of data. One such problem is of concept drift. Our research also gives a direction to the future research work that can be done to the fill the gaps by presenting new methods or even by the help of these existing methods.

ACKNOWLEDGMENT

We thank Dr Anshu Singla, Nature Master Classes(NMC) and all the volunteers for their involvement in this research project. We are highly indebted to Chitkara University for their instructions and their continuous management as well as for providing their support in the completion of the project.

REFERENCES

[1] L. Sweeney, "k-anonymity: A model for protecting privacy," International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, vol. 10, no. 05, 2002

[2] A.Vinothini and Dr.S.Baghavathi Priya "Survey of Machine Learning Methods for Big Data Applications" International Conference on Computational Intelligence in Data Science, 2017.

[3] Xi Fang and Juanjuan Wang, "The Application of Big Data Technology and Method in Moral Education in Colleges and a Universities", International conference on Intelligent Transportation, Big data and Smart City (ICITBS)

[4] S. R. Sukumar, "Machine Learning in the Big Data Era: Are We There Yet?" in Proceedings of the 20th ACM SIGKDD Conference on Knowledge Discovery and Data Mining: Workshop on Data Science for Social Good (KDD 2014), 2014.

[5] M. M. Najafabadi, F. Villanustre, T. M. Khoshgoftaar, N. Seliya, R. Wald, and E. Muharemagic, "Deep Learning Applications and Challenges in Big Data Analytics," Journal of Big Data, vol. 2, Feb. 2015.

[6] J. Qiu, Q. Wu, G. Ding, Y. Xu, and S. Feng, "A Survey of Machine Learning for Big Data Processing," EURASIP Journal on Advances in Signal Processing, vol. 67, 2016.

[7] J.Qiu et al., "A survey of machine learning for big data processing," Signal Processing for Big Data, 28 May,2016

[8] O. Y. Al-Jarrah, P. D. Yoo, S. Muhaidat, G. K. Karagiannidis, and K. Taha, "Efficient Machine Learning for Big Data: A Review," Big Data Research, vol. 2, no. 3, Apr. 2015.

[9] A. L'Heureux, K. Grolinger, H. F. Elyamany, and M. A. M. Capretz, "Machine Learning with Big Data: Challenges and Approaches," IEEE Access, vol. 5.

[10] A. Gandomi and M. Haider, "Beyond the Hype: Big data Concepts, Methods, and Analytics," International Journal of Information Management, vol. 35, Apr. 2015.

[11] S.Athmaja, M.Hanumanthapa and V.Kavitha, "A Survey Of Machine Learning Algorithms For Big Data Analytics," International Confrence on Innovations in Information, Embedded Communication Systems(ICIECS),2017

[12] H. Liu and H. Motoda, Instance Selection and Construction for Data Mining, vol. 608. Springer Science & Business Media, 2013.

[13] A. Buades, B. Coll, and J. Morel, "A Review of Image Denoising Algorithms, with a New One," Multiscale Modeling & Simulation, vol. 4, no. 2, 2005.

[14] H. Park, R. Ikeda, and J. Widom, "RAMP: A System for Capturing and Tracing Provenance in MapReduce Workflows.," Proceedings of the VLDB Endowment, vol. 4, no. 12, 2011.

[15] R. Agrawal and R. Srikant, "Privacy-preserving data mining," in ACM Sigmod Record, vol. 29, no. 2. ACM, 2000.

[16] K. Liu, H. Kargupta, and J. Ryan, "Random projection-based multiplicative data perturbation for privacy preserving distributed data mining," IEEE Transactions on knowledge and Data Engineering, vol. 18, no. 1, 2006.

[17] Seref Sagiroglu and Duygu Sinanc, "Big Data: A Review", IEEE journal,2013

[18] A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkitasubramaniam, "l-diversity: Privacy beyond k-anonymity," ACM Transactions on Knowledge Discovery from Data (TKDD), vol. 1, 2007.

[19] R. Salakhutdinov and G. E. Hinton, "Deep Boltzmann Machines," in Proceedings of the International Conference on Artificial Intelligence and Statistics, 2009.

[20] G. Hinton, "Deep Belief Nets," in Encyclopedia of Machine Learning, Springer, 2010, pp. 267-269.

[21] CW Tsai, CF Lai, MC Chiang, LT Yang, "Data mining for the internet of things: a survey". IEEE Commun Surv Tut 16(1), 77-97 (2014)

[22] X Wu, X Zhu, G Wu, W Ding, "Data mining with big data". IEEE Trans Knowl Data Eng 26(1), 97-107 (2014)

[23] U Fayyad, G Piatetsky-Shapiro, P Smyth, "From data mining to knowledge discovery in databases". AI Mag 17(3), 37-54 (1996)

[24] Mohammed Z. Omer and Hui Gao, "Privacy Preserving in Distributed SVM Data Mining on Vertical Partitioned Data", 3rd International Conference on Soft Computing & Machine Intelligence, 2016.

[25] M. M. Najafabadi, F. Villanustre, T. M. Khoshgoftaar, N. Seliya, R. Wald, and E. Muharemagic, "Deep Learning Applications and Challenges in Big Data Analytics," Journal of Big Data, Feb. 2015.

- [26] C. L. Philip Chen and C. Y. Zhang, “*Data-Intensive Applications, Challenges, Techniques and Technologies: a Survey on Big Data*,” Information Sciences, vol. **275**, **2014**.
- [27] S. G. Teo, S. Han, and V. C. Lee, “*Privacy preserving support vector machine using non-linear kernels on hadoop mahout*,” in Computational Science and Engineering (CSE), 2013 IEEE 16th International Conference on. IEEE, **2013**, pp. **941–948**.
- [28] Y. Rahulamathavan, S.Veluru, R. C.-W. Phan, J. A. Chambers, and M. Rajarajan, “*Privacy-preserving clinical decision support system using gaussian kernel-based classification*,” IEEE journal of biomedical and health informatics, vol. **18**, no. **1**, pp. **56–66**, **2014**.
- [29] Long Xu and Yihua Yan, “*Machine Learning for Astronomical Big Data Processing*”, IEEE,**2017**
- [30] A.Vinothini, S.Baghavati, “*Survey of Machine Learning for Big Data Applications*,” International Confrence on Computational Intelligence in Data Science(ICCIDS),**2017**
- [31] C.Augenstein, N.Spangenberg and Bogdan Franczyk, “*Applying Machine Learning to Big Data Streams*,”4th IEEE International Conference on Soft Computing and Machine Intelligence,**2017**