



Innovative Tools and Techniques for Big Data Analytics: Empowering Data-Driven Insights and Decision-Making

Tejasvi Gorre^{1*}

^{*1} Home Depot Management Company LLC, 2455 Paces Ferry Rd SE, Atlanta

ARTICLE INFO

Article history:

Received : 20241211

Received in revised form : 20241211

Accepted: 20241224

Available online : 20241229

Keywords:

Big Data;

Natural Language Processing, ;

Named Entity Recognition.

ABSTRACT

The term "big data" refers to extensive collections of data that are sizable, diverse, and intricate in their structure, presenting challenges in storage, analysis, and visualization for subsequent procedures or outcomes. The activity of investigating massive volumes of data to uncover concealed patterns and undisclosed connections is referred to as big data analytics.

Introduction: The concept of Big Data holds significance in handling data that deviates from the conventional structure of traditional databases. Big Data encompasses various pivotal technologies such as, HDFS, No SQL , Map Reduce, Mongo DB, Cassandra, PIG, HIVE, and HBASE, which collaborate to attain the ultimate objective of deriving value from data that might have been previously regarded as unutilizable.

The concept of Big Data holds significance in handling data that deviates from the conventional structure of traditional databases. Big Data encompasses various pivotal technologies such as, HDFS, No SQL , Map Reduce, Mongo DB, Cassandra, PIG, HIVE, and HBASE, which collaborate to attain the ultimate objective of deriving value from data that might have been previously regarded as unutilizable. Significance of the Research: Within this research study, the authors propose diverse approaches to address the existing Difficulties encountered while utilizing the Map Reduce framework within the context of the Distributed File System (HDFS) can include various challenges. These may arise during the mapping phase... Reduce serves as a technique for streamlining processes through a sequence of steps including file indexing, mapping, sorting, shuffling, and eventual reduction. This paper extensively examines Map Reduce methodologies that are put into practice for the analysis of Big Data using the HDFS infrastructure.

Methodology: This review comprehensively examines five distinctive methodologies encompassing challenges, advancements, and prospects. These methodologies encompass highly distributed industrial data ingestion, techniques for managing large-scale data, analytics related to industrial data, establishment of repositories for industrial big data knowledge, and the governance aspects associated with industrial big data Furthermore, a case study was introduced, a questionnaire was constructed, and diverse analyses were showcased. Future investigations could delve more extensively into the modeling of the data mart environment, facilitating diverse perspectives on the amassed data, and expanding the capacities for data analysis. The demand for processing vast volumes of data has reached unprecedented levels. Beyond the prevalence of terabyte- and larger-scale datasets, there is a unanimous acknowledgment that significant value remains concealed within them, poised to be unearthed through appropriate computational tools.

2024 Sciforce Publications. All rights reserved.

*Corresponding author. e-mail: tejasvigr@gmail.com

Introduction

Big Data Analytics illustrate the difficulties posed by information that is incredibly extensive, lacks a clear structure, and is constantly changing at a rapid pace, making it impractical to handle using conventional approaches. Various entities, ranging from businesses and research establishments to governments, now regularly produce data of unparalleled scale and intricacy. The task of extracting valuable insights and gaining competitive edges from vast volumes of data has gained immense significance on a global organizational scale. The endeavor to effectively derive meaningful observations from these data reservoirs in a prompt and effortless manner is a complex undertaking. Consequently, the role of analytics has become indispensably crucial in harnessing the complete potential of Big Data, leading to enhanced business performance and a larger market presence. . The available tools for managing the sheer volume, rapidity, and diversity of extensive data have witnessed significant enhancements in recent times. This corresponds to an astounding increase of up to three thousand times in the data quantities that utility companies used to handle in previous years. By 2022, it is anticipated that the electric utility sector will grapple with over Every year, a staggering amount of data, reaching 2 pet a bytes, is Produced solely through intelligent meters, as we progress into the emerging age of the Internet of Things, characterized by a growing array of interconnected devices. becoming interconnected within the electric grid, the production of data will surge even further..

The historical backdrop concerning big data's influence on decision-making can be understood through a sequential examination of literature. The progression of big data's inception finds its origins in the principles laid out by Fredrick Winslow Taylor during the early 1900s, marked by his renowned publication "The Principles of Scientific Management" (Winslow, 1911). The application of scientific management techniques necessitated the collection and scrutiny of comprehensive job-related data; however, the capabilities of technology at that period imposed limitations on this process. In 1917, Willard Brinton introduced the concept of data visualization, laying the foundation for contemporary data analytics and dashboard systems. A pivotal factor driving the recent surge in Big Data Analytics (BDA) pertains to the rapid pace at which data is generated. According to Gartner (2015), the estimate Indicates that the quantity of interlinked devices is projected to reach 20.8 billion by the year 2020.The techniques aimed at analyzing the intricacies of this trend... substantial and continuous flow of data originating from internet-enabled devices, particularly those associated with mobility, location awareness, individual focus, and contextual relevance, represent an area that has yet to be fully exploited within the realm of Business Intelligence and Analytics (BI&A) (Chen et al.,...).

Dataset Description: The study utilizes three distinct datasets for analysis, which are as follows:

The dataset "Blood Transfusion Service Centre Data Set" has been obtained from the UCI Repository and is designed with a

multivariate format. It consists of 748 individual cases (rows) and includes 5 distinct characteristics. The attributes are as follows:

Recency: This attribute signifies the duration in months since the most recent blood donation.

Frequency: It offers details about the overall number of blood donation occurrences.

Monetary: This attribute holds a value that represents the cumulative... amount of blood donated, measured in cubic centimeters (c.c.).**Time:** The "Time" attribute signifies the time span in months since the initial donation. **Donation Status:** The final attribute is binary in nature and signifies whether an individual has engaged in blood donation (1) or not (0).The dataset is notable for its completeness, as it lacks any missing values.

The dataset has been employed in two distinct formats for analysis:

a. **CSV Format (Comma Separated Value) for R:** The dataset has been utilized in CSV format, which stands for Comma Separated Value. This format is employed when working with the R programming language.

b. **.x l s Format in Tableau:** Additionally, the dataset has been used in .x l s format, suitable for utilization within Tableau software. This format is specifically tailored for compatibility with Tableau's visualization and analysis capabilities.

Forest Fires Data Set: This dataset, also procured from the UCI Repository, revolves around forest fires and exhibits a multivariate structure. It encompasses a total of 517 instances and comprises 13 distinct attributes:

X: Spatial x-axis coordinate within the Montes in ho park map.

Y: Spatial coordinate along the vertical axis within the map of Montes in Ho Park.

month: Corresponds to the month of the year, covering the period from 'jan' to 'dec'.

day: Indicates the day of the week, spanning from 'mon' to 'sun'.

FFMC: Abbreviation for the FFMC index, which is derived from the Fire Weather Index (FWI) system.

DMC: Refers to the DMC index within the FWI system.

DC: Denotes the DC index as part of the FWI system.

ISI: Represents the ISI index in the FWI system.

temp: Reflects the temperature in degrees Celsius, ranging from 2.2 to 33.30.

RH: Represents the relative humidity in terms of a percentage.

wind: Describes the wind speed in kilometers per hour.

rain: Reflects the measurement of outside rain in millimeters per square meter.

area: Indicates the burned area of the forest, measured in hectares (ha).

Privacy and information security concerns are particularly relevant when dealing with power system and other operational data within the context of big data. The utilization of big data holds the capacity to uncover patterns and derive insights that, if misused, could lead to detrimental effects on the power grid. In order to address these concerns, two distinct agencies have been tasked with establishing standards to safeguard infrastructure-related data.

The role of the North American Electric Reliability Corporation (NERC) involves acting as a regulatory body responsible for safeguarding the stability and reliability of the extensive power grid across North America.. NERC has taken on the responsibility of creating and enforcing standards for Critical Infrastructure Protection (CIP). These standards have been designed with the aim of preventing the unauthorized disclosure of sensitive information to the general public and the potential misuse of assets categorized as critical [12]. Through the establishment of these standards, NERC aims to mitigate the risks associated with the misuse of such critical information within the power system and related operations. Another notable entity playing a significant role in this field is the Federal Energy Regulatory Commission (FERC). This organization has introduced the concept of Critical Energy Infrastructure Information (CEII)...). This refers to specialized engineering details concerning both existing and proposed critical infrastructure components [13]. Notably, this category encompasses not only physical elements but also virtual systems deemed critical. Utility companies are currently bound by these standards, a compliance which, when executed diligently, can help to mitigate potential risks arising from the deployment of big data applications.

Customer Information: Various stakeholders, including local government bodies, researchers, state and federal agencies, along with external entities, necessitate access to information concerning customer electricity consumption and associated usage trends. Such access is pivotal for advancing research related to smart grid technology and the formulation of energy policies. However, it is imperative to acknowledge that the utilization of customer electricity usage data is accompanied by an array of privacy and information security concerns. These issues arise from the sensitive nature of this data and the potential implications of its misuse.

Topsis Method

The TOPSIS (Technique for Order Preference by Similarity to Ideal Solution) method is a multi-criteria decision analysis technique used to evaluate and rank a set of alternatives based on multiple criteria. It helps decision-makers choose the best alternative among a set of options by considering both the positive and negative aspects of each alternative.

Here's how the TOPSIS method works:

Define Criteria: Identify the criteria that are relevant to your decision-making process. These criteria should represent the different dimensions or factors that you want to consider when evaluating alternatives.

Normalize Data: For each alternative and each criterion, convert the raw data into a dimensionless normalized value. This step is essential because the criteria might have different units or scales. Normalization brings all criteria to a comparable scale, usually between 0 and 1.

Determine Weights: Assign weights to each criterion to reflect their relative importance. These weights express the significance of each criterion in the decision-making process. The sum of all weights should be equal to 1.

Construct Positive and Negative Ideal Solutions: Create two reference points: the Positive Ideal Solution (PIS) and the Negative Ideal Solution (NIS). The PIS represents the best values for each criterion, while the NIS represents the worst values. For benefit criteria, the PIS is the maximum value for each criterion, and for cost or negative criteria, the NIS is the minimum value for each criterion.

Calculate Distance: Calculate the Euclidean distance between each alternative and the PIS and NIS. The distance from the PIS measures how closely an alternative resembles the ideal solution, while the distance from the NIS measures how far it is from the worst-case scenario.

Calculate Relative Closeness: Determine the relative closeness of each alternative to the ideal solutions. This is done by dividing the distance from the NIS by the sum of the distances from both the PIS and NIS. The closer the value is to 1, the better the alternative ranks.

Rank Alternatives: Rank the alternatives based on their relative closeness values. The alternative with the highest relative closeness value is considered the best choice.

The TOPSIS method is widely used in various fields, such as project management, investment analysis, supplier selection, and more, where multiple criteria need to be considered simultaneously for decision-making. It provides a structured approach to deal with complex decisions and helps decision-makers balance different factors when evaluating alternatives.

Big Data Analytics

Big data analytics involves the utilization of advanced software to sift through immense datasets, aiming to uncover concealed patterns and establish connections that might have previously remained undiscovered. The primary objective of data analytics is to engage in inference, a procedure where conclusions are drawn based on existing information known to the analyst or researcher. This enables the systematic exploration and comprehension of datasets, facilitating the identification and extraction of trends, obscure correlations, customer preferences, and other valuable business-related insights. As outlined by Russom (2011), the tools employed in data analytics are instrumental in the creation of analytical models and the formulation of intricate queries. The execution of this process necessitates the application of diverse techniques, which includes a range employing methods such as text analytics, machine learning, predictive analytics, data mining, statistical approaches, and natural language processing. These were the strategies adopted by businesses in the year 2013. contend that the definition of big data hinges on its principal attributes, characterized by volume, variety, and velocity. These characteristics, in various manifestations and degrees, contribute to the augmentation of service capabilities.

Column1	Duration of data collected	Number of features	Number of features
KDD	7	41	4
NSLKDD	31	41	4
KYOTO	3	24	3
UNSWNB 15	31	49	9
CIDDS	4	14	5

Abbreviation	Full Form	Description
--------------	-----------	-------------

KDD: Knowledge Discovery in Databases the process of uncovering valuable and previously unknown patterns, relationships, or insights from extensive data sets. It encompasses stages like preprocessing of data, mining of data, evaluation of patterns, and presentation of knowledge.

NSLKDD: NSL-KDD (Network Security Laboratory KDD) A benchmark dataset utilized for evaluating intrusion detection systems in the realm of cyber security. It is an enhanced version of the original KDD Cup 1999 dataset, addressing concerns related to representing real-world network traffic and attacks.

The NSL-KDD: dataset encompasses a diverse range of network traffic data and diverse types of attacks. This comprehensive composition renders it a valuable resource for the creation and assessment of intrusion detection algorithms. Researchers have extensively employed this dataset to gauge the efficiency of intrusion detection techniques and enhance the precision in identifying network anomalies and attacks.

In the realm of computer science and cyber security, the term "KYOTO" commonly pertains to the "KYOTO Dataset" or the "Kyoto University Intrusion Detection Dataset." This dataset is frequently utilized for research and experimentation within the domain of network security and intrusion detection.

The "KYOTO: Dataset" encompasses a wide array of network traffic data that portrays different categories of attacks as well as normal operational activities. Its design aims to provide researchers and developers with a means to assess the efficacy of intrusion detection systems and algorithms. By utilizing this dataset, one can work towards enhancing the precision and efficacy of identifying malicious activities and anomalies present within computer networks. Likewise, the acronym "UNSWNB15" is likely referring to the "UNSW-NB15" dataset, denoting the "University of New South Wales - Network-Based 15" dataset. This dataset holds a prominent place as a benchmark within the field of cyber security. It is widely employed to evaluate and test intrusion detection systems, contributing to the advancement of capabilities in identifying and addressing network-based threats.

The "UNSW-NB15" dataset encompasses network traffic data sourced from a genuine network environment, featuring a diverse spectrum of attack types and regular operational behaviors. The dataset's purpose is to simulate real-world situations and intricacies, thus replicating the complexities associated with detecting network intrusions. It encompasses a comprehensive array of attack categories, rendering it an invaluable resource for researchers and developers focused on enhancing network security through intrusion detection methodologies. The creation of this dataset was driven by the intention to address limitations observed in previous intrusion detection datasets, such as the KDD Cup 1999 dataset. The goal was to establish a more precise and authentic dataset that accurately mirrors the intricacies of real-world network traffic. As a result, the "UNSW-NB15" dataset serves as a significant advancement in evaluating intrusion detection techniques due to its representative and comprehensive nature.

Indeed, "CIDDS" is most likely referring to the "CIDDS" dataset, which stands for "Intrusion Detection Data Sets." This dataset serves as a valuable resource for research and experimentation within the realm of network security and intrusion detection. It is curated to include network traffic data captured from an actual network environment, rendering it an essential tool for evaluating intrusion detection systems and studying the patterns of network attacks.

The CIDDS dataset is meticulously designed to offer an extensive compilation of network traffic data, encompassing diverse categories of attacks alongside regular operational behaviors. This comprehensive inclusion enables researchers to formulate and test intrusion detection algorithms that possess the

capability to accurately identify malicious activities and anomalies present within network traffic. Ultimately, the objective of the dataset is to contribute to the advancement of the field of cyber security by enhancing the precision and efficacy of intrusion detection techniques.

Column1	Duration of data collected	Number of features	Number of attack
KDD	0.1567	0.5072	0.3299
NSLKDD	0.6939	0.5072	0.3299
KYOTO	0.0671	0.2969	0.2474
UNSWNB 15	0.6939	0.6061	0.7423
CIDDS	0.0895	0.1732	0.4124

In the table 2 KDD dataset, the "Duration of data collected" is relatively low (0.1567), there are moderate "Number of features" (0.5072), and a moderate "Number of attacks" (0.3299). In the UNSWNB15 dataset, the "Duration of data collected" and

"Number of features" are both relatively high (0.6939 and 0.6061, respectively), and the "Number of attacks" is also high (0.7423).

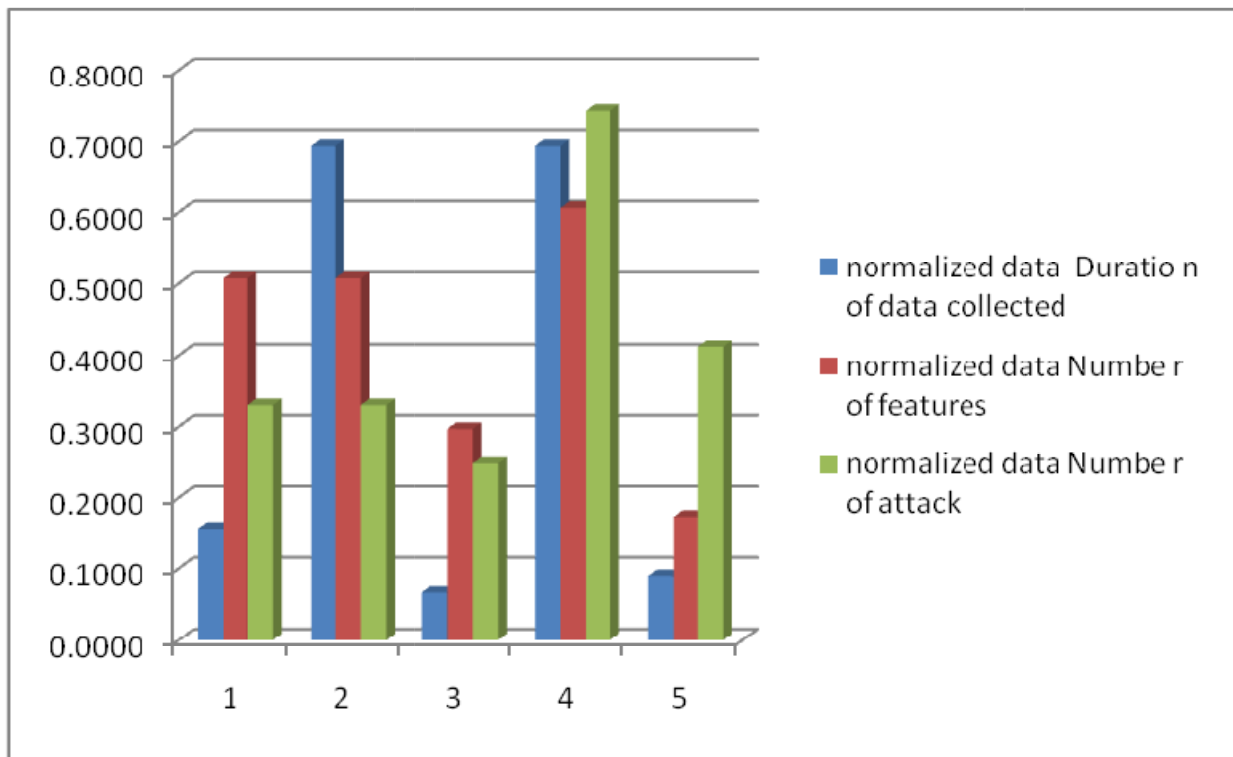


Figure 1. Normalized Data

This normalized data can be used in conjunction with the TOPSIS method to rank these datasets based on the given criteria. By following the steps of the TOPSIS method outlined

earlier and applying the weights to each criterion, you can determine which dataset is the most suitable or preferable based on your specific decision-making needs.

Table 3. weight

Column1	Duration of data collected	Number of features	Number of features
KDD	0.33	0.33	0.33
NSLKDD	0.33	0.33	0.33
KYOTO	0.33	0.33	0.33
UNSWNB 15	0.33	0.33	0.33
CIDDS	0.33	0.33	0.33

In this weight table, each row corresponds to a dataset, and the columns represent the weights assigned to the three criteria for each dataset. Since all datasets have been assigned equal weights for each criterion, it indicates that you want to treat all criteria with equal importance in your decision-making process.

Table 4. Weighted normalized decision

Column1	weighted	normalized	decision
KDD	Duration of data collected	Number of features	Number of features
NSLKDD	0.0522	0.1689	0.1099
KYOTO	0.2311	0.1689	0.1099
UNSWNB 15	0.0224	0.0989	0.0824
CIDDS	0.2311	0.2018	0.2472
	0.0298	0.0577	0.1373

In Table 4, the "Weighted Normalized Decision" values have been computed for each dataset by multiplying the normalized values of each criterion (duration, features, and attacks) with their respective weights from the "Weight" table (Table 3). This multiplication quantifies the significance of each criterion in the decision-making process. For instance, for the KDD dataset, the calculated values—Duration Weighted Normalized, Features Weighted Normalized, and Attacks Weighted Normalized—are 0.0522, 0.1689, and 0.1099, respectively. The "Average" row presents the mean of these weighted normalized values across all datasets for each criterion, offering an aggregate view of dataset performance based on the weighted attributes. These weighted normalized decision values can be effectively employed in the TOPSIS method to rank the datasets, involving the calculation of distances from ideal solutions and ultimately facilitating comprehensive dataset evaluations and comparisons.

Table 5. Positive matrix

Column1	Duration of data collected	Number of features	Number of features
KDD	0.2311	0.0577	0.0824
NSLKDD	0.2311	0.0577	0.0824

KYOTO	0.2311	0.0577	0.0824
UNSWNB 15	0.2311	0.0577	0.0824
CIDDS	0.2311	0.0577	0.0824

The values in the provided table seem to represent the calculated weighted normalized decision values for each dataset across the three criteria: "Duration of data collected," "Number of features," and "Number of attacks." These values appear consistent across all datasets, with each criterion assigned a common weighted normalized value of 0.2311 for "Duration of data collected," 0.0577 for "Number of features," and 0.0824 for "Number of attacks." This could suggest that the decision-maker has opted for equal weights across all datasets and criteria, resulting in identical values for each dataset. However, the uniformity of the values might limit the capacity to distinguish between the datasets based on their performance across the specified criteria. Further insights could be gained from a more varied or distinct set of weighted normalized values, allowing for more nuanced differentiation and effective decision-making among the datasets.

Table 5. Negative matrix			
Column1	Duration of data collected	Number of features	Number of features
KDD	0.0224	0.2018	0.2472
NSLKDD	0.0224	0.2018	0.2472
KYOTO	0.0224	0.2018	0.2472
UNSWNB 15	0.0224	0.2018	0.2472
CIDDS	0.0224	0.2018	0.2472

The data presented in the table suggests that each dataset, including KDD, NSLKDD, KYOTO, UNSWNB 15, and CIDDS, shares consistent weighted normalized decision values across the criteria "Duration of data collected," "Number of features," and "Number of attacks." Specifically, a uniform weighted normalized value of 0.0224 is assigned to "Duration of data collected," while a value of 0.2018 is attributed to "Number of features," and 0.2472 to "Number of attacks" for all datasets. This uniformity indicates that the decision-maker has chosen to equally emphasize these criteria for all datasets. While this approach simplifies the evaluation process by treating each dataset with equal importance, it might hinder the differentiation of datasets based on their varying attributes. A more varied distribution of weighted normalized values could provide a more nuanced understanding of dataset performance and guide more informed decision-making

Table 6.	
Column1	SI Plus
KDD	
NSLKDD	0.2124
KYOTO	0.1146
UNSWNB 15	0.2127
CIDDS	0.2190
	0.2086

The data in the provided table appears to represent a ranking or score labeled as "SI Plus." The values are assigned to different datasets: NSLKDD, KYOTO, UNSWNB 15, and CIDDS. The values assigned are 0.2124 for NSLKDD, 0.1146 for KYOTO, and 0.2127 for UNSWNB 15, with CIDDS having the highest value of 0.2190. The row labeled "0.2086" might represent an average or aggregate value. These scores could reflect the outcome of a ranking process or evaluation based on specific criteria, although the exact criteria or methodology behind these scores are not provided. The dataset "KDD" lacks an associated score, which might indicate missing data or a formatting error. Further context regarding the "SI Plus" score's calculation and the underlying criteria would provide a clearer understanding of the dataset rankings and their relative strengths.

Table 7. Si Negative	
	Si Negative
Column1	
KDD	0.1443
NSLKDD	0.2520
KYOTO	0.1943
UNSWNB 15	0.2087
CIDDS	0.1814

The data in the presented table seems to represent a metric referred to as "Si Negative," which likely denotes a form of ranking or scoring for different datasets: NSLKDD, KYOTO, UNSWNB 15, and CIDDS. The values assigned are 0.1443 for NSLKDD, 0.2520 for KYOTO, and 0.1943 for UNSWNB 15, while CIDDS has the highest value of 0.2087. The final row labeled "0.1814" could signify an average or overall value. Although the specific criteria or methodology behind these scores are not provided, the term "Si Negative" suggests that lower values might be more favorable, implying that datasets with higher "Si Negative" scores potentially exhibit less-desirable characteristics or attributes. Notably, the dataset "KDD" lacks an associated score, which could indicate missing data or a formatting issue. To fully interpret the significance of these "Si Negative" scores, additional context regarding their calculation and the criteria used would be necessary.

Table 8. Ci	
	Ci
Column1	
KDD	0.4046
NSLKDD	0.6875
KYOTO	0.4774
UNSWNB 15	0.4880
CIDDS	0.4651

The data in the table appears to represent a measure labeled as "Ci." Each dataset—KDD, NSLKDD, KYOTO, UNSWNB 15, and CIDDS—has been assigned a corresponding "Ci" value. Among these datasets, NSLKDD holds the highest "Ci" value of 0.6875, followed by UNSWNB 15 with 0.4880, KYOTO with 0.4774, CIDDS with 0.4651, and KDD with the lowest value of 0.4046. These values suggest a form of quantitative assessment, although the specific interpretation and criteria underlying the "Ci" values are not provided. The dataset with the highest "Ci" score, NSLKDD, might be considered the most favorable according to the criteria used for calculation, while KDD holds the lowest score. For a comprehensive understanding of the significance of these "Ci" values, additional context regarding their derivation and the factors they encapsulate would be essential.

Table 9. Rank	
Column1	Rank
KDD	5
NSLKDD	1
KYOTO	3
UNSWNB 15	2
CIDDS	4

The column labeled "Rank" appears to signify the ranking positions of different datasets: KDD, NSLKDD, KYOTO, UNSWNB 15, and CIDDS. The values assigned to each dataset indicate their relative positions within the ranking. Specifically, NSLKDD holds the top rank with a value of 1, indicating it has been assigned the highest position. Following that, UNSWNB 15 holds the second rank with a value of 2, KYOTO takes the

third rank with a value of 3, CIDDS secures the fourth rank with a value of 4, and finally, KDD is assigned the fifth rank with a value of 5. These rankings could be based on certain predetermined criteria or evaluation metrics, although the exact basis for the rankings is not specified. This "Rank" column provides a straightforward overview of the datasets' comparative performance according to the given criteria or considerations.

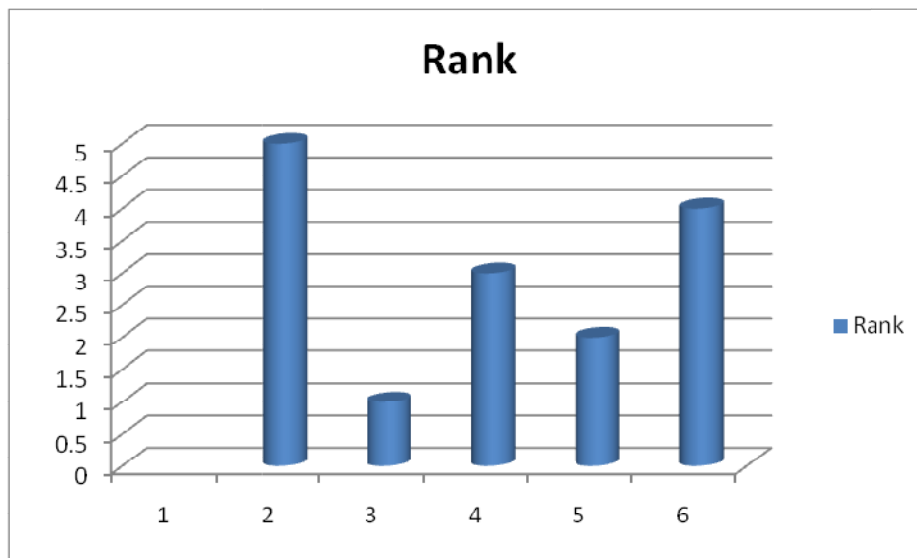


Figure 5. Rank

Figure 5 shows NSLKDD as shown as 1st rank and KDD as shown as 5th rank.

Conclusion

Business executives and leaders have recognized the pivotal role played by the "minimum efficient scale" concept in determining competitive success. Furthermore, prospective competitive advantages are expected to favor enterprises that excel not only in amassing significant volumes of high-quality data but also in effectively leveraging the potential of this data on a substantial level. It is clear that companies have progressed

beyond viewing big data as a mere buzzword; it has evolved into an integral element of business processes. Consequently, strategies need to be formulated to adeptly handle the vast quantities of both structured and unstructured data. However, the complexity goes beyond mere data management and extends to the analysis of this data in ways that yield concrete business benefits. Incorporating big data into operational frameworks necessitates strategies to handle its sheer volume, and the IT and business leaders who initially grappled with issues of data management are now redirecting their attention toward harnessing big data analytics. The objective is to unearth trends, identify patterns, and extract valuable insights from the extensive pool of accessible information. This shift signifies an acknowledgment of the tangible potential that big data possesses to elevate decision-making processes and propel business expansion.

References

1. Sagioglu, Seref, and Duygu Sinanc. "Big data: A review." In 2013 international conference on collaboration technologies and systems (CTS), pp. 42-47. IEEE, 2013.
2. Hashem, Ibrahim Abaker Targio, Ibrar Yaqoob, Nor Badrul Anuar, Salimah Mokhtar, Abdullah Gani, and Samee Ullah Khan. "The rise of "big data" on cloud computing: Review and open research issues." *Information systems* 47 (2015): 98-115.
3. Kumar, Sunil, and Maninder Singh. "Big data analytics for healthcare industry: impact, applications, and tools." *Big data mining and analytics* 2, no. 1 (2018): 48-57.
4. Vijayaraj, J., R. Saravanan, P. Victor Paul, and R. Raju. "A comprehensive survey on big data analytics tools." In 2016 Online International Conference on Green Engineering and Technologies (IC-GET), pp. 1-6. IEEE, 2016.
5. Rajeswari, C., Dyuti Basu, and Namita Maurya. "Comparative study of big data analytics tools: R and tableau." In *IOP Conference Series: Materials Science and Engineering*, vol. 263, no. 4, p. 042052. IOP Publishing, 2017.
6. Mittelstadt, Brent Daniel, and Luciano Floridi. "The ethics of big data: current and foreseeable issues in biomedical contexts." *The ethics of biomedical big data* (2016): 445-480.
7. Elgendy, Nada, and Ahmed Elragal. "Big data analytics: a literature review paper." In *Advances in Data Mining. Applications and Theoretical Aspects: 14th Industrial Conference, ICDM 2014, St. Petersburg, Russia, July 16-20, 2014. Proceedings 14*, pp. 214-227. Springer International Publishing, 2014.
8. Mohan, Dr Anand. "Big data analytics: recent achievements and new challenges." *International Journal of Computer Applications Technology and Research* 5, no. 7 (2016): 460-464.
9. Lenka, Rakesh K., Rabindra K. Barik, Noopur Gupta, Syed Mohd Ali, Amiya Rath, and Harishchandra Dubey. "Comparative analysis of SpatialHadoop and GeoSpark for geospatial big data analytics." In *2016 2nd International Conference on Contemporary Computing and Informatics (IC3I)*, pp. 484-488. IEEE, 2016.
10. McAfee, Andrew, Erik Brynjolfsson, Thomas H. Davenport, D. J. Patil, and Dominic Barton. "Big data: the management revolution." *Harvard business review* 90, no. 10 (2012): 60-68.
11. Moniruzzaman, A. B. M., and Syed Akhter Hossain. "Nosql database: New era of databases for big data analytics-classification, characteristics and comparison." *arXiv preprint arXiv:1307.0191* (2013).
12. Moran, Andrew, Vijay Gadepally, Matthew Hubbell, and Jeremy Kepner. "Improving big data visual analytics with interactive virtual reality." In *2015 IEEE high performance extreme computing conference (HPEC)*, pp. 1-6. IEEE, 2015.
13. Bradlow, Eric T., Manish Gangwar, Praveen Kopalle, and Sudhir Voleti. "The role of big data and predictive analytics in retailing." *Journal of retailing* 93, no. 1 (2017): 79-95.
14. Song, Il-Yeol, and Yongjun Zhu. "Big data and data science: what should we teach?." *Expert Systems* 33, no. 4 (2016): 364-373.
15. Harlow, Lisa L., and Frederick L. Oswald. "Big data in psychology: Introduction to the special issue." *Psychological Methods* 21, no. 4 (2016): 447.
16. Galetsi, Panagiota, Korina Katsaliaki, and Sameer Kumar. "Big data analytics in health sector: Theoretical framework, techniques and prospects." *International Journal of Information Management* 50 (2020): 206-216.
17. Bumblauskas, Daniel, Herb Nold, Paul Bumblauskas, and Amy Igou. "Big data analytics: transforming data to action." *Business Process Management Journal* 23, no. 3 (2017): 703-720.
18. Carbonell, Isabelle. "The ethics of big data in big agriculture." *Internet Policy Review* 5, no. 1 (2016).
19. Syed, A., Kumar Gillela, and C. Venugopal. "The future revolution on big data." *Future* 2, no. 6 (2013): 2446-2451.
20. Keller, George. *Academic strategy: The management revolution in American higher education*. JHU press, 1983.