



**BERLIN SCHOOL OF  
BUSINESS & INNOVATION**

**Essay / Assignment Title: Big Data Analytics with a Special Focus on  
Distributed File Systems**

**Programme title: MSc Data Analytics**

**Name: SMIT KEVADIYA KAMLESHBHAI**

**Year: 2023**

## CONTENTS



TABLE OF CONTENTS.....	02
INTRODUCTION.....	04
CHAPTER ONE.....	05
CHAPTER TWO.....	08
CHAPTER THREE.....	10
CONCLUSIONS.....	13
BIBLIOGRAPHY.....	14



**Statement of compliance with academic ethics and the avoidance of plagiarism**

I honestly declare that this dissertation is entirely my own work and none of its part has been copied from printed or electronic sources, translated from foreign sources and reproduced from essays of other researchers or students. Wherever I have been based on ideas or other people texts I clearly declare it through the good use of references following academic ethics.

(In the case that is proved that part of the essay does not constitute an original work, but a copy of an already published essay or from another source, the student will be expelled permanently from the postgraduate program).

Name and Surname (Capital letters):

SMIT KEVADIYA

.....

Date: 08/06/2023

## INTRODUCTION

In the digital era, big data has evolved as a crucial notion, referring to massive and complex datasets that surpass the capacity of standard data processing technologies. It includes structured, semi-structured, and unstructured data from many sources, and it is distinguished by its vast volume, fast rate of data production, and wide range of data kinds and formats. Big data brings with it new problems and opportunities for analysis and insight extraction.

Distributed File Systems (DFSs) play a significant role in managing and processing big data. DFSs, such as Hadoop Distributed File System (HDFS) and Google File System (GFS), are designed to store and manage data across multiple nodes in a distributed computing environment. They offer scalability, fault tolerance, and high throughput for handling large volumes of data.

DFSs consist of components like the NameNode, DataNodes, and client nodes, which work together to ensure efficient data storage and retrieval. DFSs provide critical support in the science of big data due to their ability to handle the unique challenges posed by large-scale datasets. They offer fault tolerance, ensuring data availability even in the event of node failures. DFSs enable parallel processing, allowing for efficient data analysis and processing across distributed nodes. They also provide data locality, minimizing data transfer costs and improving overall performance.

Hadoop is an open-source platform that uses HDFS as its distributed file system and the MapReduce programming language for distributed data processing. Hadoop's characteristics and qualities, such as scalability, fault tolerance, and distributed processing capabilities, make it a viable platform for big data learning and decision-making. Organizations may successfully manage and analyze large data by harnessing the capabilities of DFSs and Hadoop, gaining relevant insights and making educated choices. The combination of distributed file systems and the processing capacity of Hadoop enables effective management of large-scale datasets, supporting advances in big data research.

## CHAPTER ONE

Big data is fundamentally defined as large amounts of complicated and diverse information that cannot be properly handled or evaluated using typical data processing methods. It includes structured, semi-structured, and unstructured data from a variety of sources, including social media, sensors, online logs, and transactional systems. Big data is distinguished by its enormous volume, rapid velocity, and wide range of data kinds and formats.

Big data is defined structurally by its massive volume. Traditional data management solutions are incapable of handling the massive amounts of data created on a daily basis. Big data can be measured in terabytes ( $10^{12}$  bytes), petabytes ( $10^{15}$  bytes), and even exabytes ( $10^{18}$  bytes). This volume presents enormous storage, processing, and analytical issues. To handle the enormous volume of data involved, specialized infrastructure and technologies are necessary.

In addition to volume, big data is characterized by its high velocity. Data is generated and transmitted at an unprecedented rate in the digital age. Social media platforms, for example, generate vast amounts of data every second through user interactions, posts, and messages. This velocity necessitates real-time or near real-time processing and analysis to extract meaningful insights and respond to emerging trends and events. High-speed data processing systems and streaming technologies are essential to handle the continuous flow of data.

Big data is also defined by its diverse variety of data types and formats. It encompasses structured data, which fits neatly into traditional relational databases, as well as semi-structured and unstructured data. Semi-structured data includes formats such as XML and JSON, while unstructured data includes text, images, videos, and more. The variety of data sources and formats adds complexity to the data management and analysis process. Specialized techniques like natural language processing and computer vision are needed to effectively analyze unstructured data.

Big data may be mathematically represented using a variety of models and methodologies. One method is to think of huge data as high-dimensional spaces. Each data point represents an observation or item with a number of properties. Data can be represented in this context as vectors in an  $n$ -dimensional space, where  $n$  denotes the number of characteristics or features. To

analyze and identify patterns from data, mathematical processes such as clustering, classification, and regression can be used. To examine huge data, statistical approaches and algorithms are often used. Machine learning algorithms, for example, may find patterns in big and complicated information, generate predictions, and unearth insights. Deep learning algorithms, a type of machine learning, excel in analyzing unstructured data such as photos and text.

Furthermore, distributed computing frameworks such as Apache Hadoop and Apache Spark may be used to handle huge data. These frameworks enable distributed data processing across computer clusters, allowing for parallel execution and effective resource usage. This distributed processing capabilities is critical for dealing with the large-scale computations necessary for big data analysis.

Big data refers to large and complicated datasets that cannot be managed efficiently using typical data processing methods. It is distinguished by its vast volume, rapid velocity, and wide range of data kinds and formats. To meet the problems of big data, specialized infrastructure, technology, and mathematical models are required. Organizations may extract important insights and make data-driven choices from big data by utilizing distributed computing frameworks and employing statistical and machine learning approaches.

Large and complicated datasets that exceed the capability of typical data processing technologies are referred to as big data. It includes structured, semi-structured, and unstructured data from a variety of sources, including social media, sensors, and transactional systems. Big data is distinguished by its massive volume, rapid rate of data collection, and wide range of data kinds and formats.

Big data has a volume ranging from terabytes to petabytes and even exabytes, whereas standard data has a smaller dataset. Big data is created at a significantly faster rate, necessitating real-time or near real-time processing, whereas typical data is created at a slower rate. Furthermore, big data has a broader range of data kinds and formats than traditional data, which is often more organized.

The difference between big data and regular data is found in their magnitude and analytical methodologies. Big data goes beyond regular data, integrating large-scale datasets and necessitating specialized administration and analysis approaches. Organizations may use big data analytics to find patterns, trends, and insights that might otherwise be missed when dealing with smaller datasets.

While big data presents unique issues, the methodologies and technology used to analyze and manage big data may also be applied to other types of data. Big data approaches may be used by organizations to acquire complete insights from both forms of data. Big data analytics gives a larger view as well as the potential to identify hidden patterns, whereas traditional data analysis delivers concentrated and precise insights into specific parts of a business or topic.

In summary, big data represents large and complex datasets that exceed the capacities of traditional data processing methods. It encompasses vast volumes of data generated at a high velocity and exhibits a diverse variety of data types and formats. While big data and normal data have differences in scale and characteristics, they are interconnected, and techniques used for big data analysis can be beneficial in extracting insights from both types of data.

## CHAPTER TWO

Distributed File Systems (DFSs) are designed to store and manage data across several nodes or servers in a distributed computing environment. They provide a scalable and fault-tolerant architecture for effectively managing enormous amounts of data. In this review, we will extensively describe and examine DFSs, outlining their essential components and highlighting their benefits and challenges in comparison to other similar/relevant systems.

The components of DFSs typically include the following:

1. **NameNode:** The NameNode acts as the master node in the DFS architecture. It manages the file system namespace, maintains the metadata of files and directories, and coordinates access to data. The NameNode keeps track of the location of data blocks stored across DataNodes.
2. **DataNodes:** DataNodes serve as worker nodes in the DFS. They are responsible for storing and retrieving data blocks. DataNodes communicate with the NameNode to report their status and perform data-related operations.
3. **Client Nodes:** Client nodes interact with the DFS to access and manipulate data. They send requests to the NameNode for file operations and communicate directly with DataNodes for data read and write operations.

DFSs offer several advantages in handling big data:

1. **Scalability:** DFSs are designed to handle large-scale datasets by distributing data storage and processing across multiple nodes. This distributed architecture enables horizontal scalability, allowing organizations to seamlessly scale their storage and processing capabilities as data volumes grow. This scalability is crucial in managing and analyzing big data, where data sizes can reach terabytes, petabytes, or even exabytes.
2. **Fault Tolerance:** DFSs provide fault tolerance mechanisms to ensure data availability and reliability. Data replication across multiple DataNodes ensures that data can be retrieved



even if some nodes fail. The NameNode maintains information about data block locations, enabling recovery in case of node failures.

3. High Throughput: DFSs are built to handle high-volume data access and processing. DFSs can use parallel processing capabilities by spreading data over numerous nodes, enhancing overall system performance and minimizing data processing times.
4. Data Locality: DFSs optimize data access by storing data closer to where it is processed. This reduces data transfer costs and minimizes network congestion, leading to improved performance.

Despite their advantages, DFSs also face certain challenges:

1. Consistency and Coherency: Maintaining data integrity and coherency across dispersed nodes in DFSs can be difficult. To ensure that numerous clients accessing and updating the same data get consistent results, synchronization and coordination procedures are necessary.
2. Metadata Management: The centralized administration of information in the NameNode might create a scalability constraint as the number of files and directories grows. To overcome this difficulty, efficient information management approaches and distributed metadata systems are required.
3. Network Bandwidth and Latency: Data movement and communication between dispersed nodes rely on network capacity, which can be constrained by network congestion and delay. DFSs have increasing issues in optimizing data transmission and decreasing network overhead.

DFSs provide a more dispersed and scalable solution ideal for large data situations as compared to other similar/relevant systems such as Network Attached Storage (NAS) or Storage Area Networks (SAN). They are well-suited for handling large-scale datasets and facilitating parallel processing because to their fault tolerance, fast throughput, and data proximity.

The distributed nature of DFSs enables efficient data storage, retrieval, and processing over a network of nodes, allowing enterprises to manage and analyze huge data more effectively.

## CHAPTER THREE

Learning from big data in Hadoop involves leveraging its attributes and properties that are specifically designed to support the learning and decision-making processes. In this critical analysis, we will explore in detail how organizations can effectively learn from their big data in Hadoop and the key attributes and properties of Hadoop that enhance the learning and decision-making processes.

1. Scalability: Hadoop's scalability is a fundamental feature that allows it to manage massive amounts of data. Scalability is achieved through Hadoop's distributed design, which distributes data and compute over several nodes in a cluster. As the size of big data grows, Hadoop may scale horizontally by adding more nodes to the cluster. Because of its horizontal scalability, enterprises can store and handle huge volumes of data, meeting the ever-increasing demands of big data applications. Organizations can address the storage and processing requirements of their big data by successfully growing their Hadoop clusters, enabling extensive analysis and learning.
2. Distributed File System (HDFS): HDFS is Hadoop's distributed file system for storing and managing massive datasets across numerous nodes. Because it is designed for high throughput and fault tolerance, it is perfect for large data processing. HDFS splits files into data blocks and distributes them throughout the cluster's DataNodes. This distribution supports parallel processing by allowing each DataNode to process its local data blocks independently. Because HDFS is distributed, data may be retrieved and processed in a very efficient way, allowing for quicker analysis and learning from massive data. Additionally, HDFS provides fault tolerance through data replication. Data blocks are replicated across multiple DataNodes, ensuring data availability even in the event of node failures. This fault tolerance mechanism enhances the reliability of data access and processing, minimizing the risk of data loss and interruptions in the learning process.
3. MapReduce: MapReduce is a programming model and processing framework in Hadoop specifically designed for distributed data processing. It breaks down complex data

processing tasks into smaller, parallelizable units and distributes them across the cluster. The Map phase analyzes data across numerous nodes in parallel, with each node applying a specific function (map) to the input data and producing intermediate key-value pairs. To obtain the ultimate output, the Reduce step aggregates and merges the intermediate findings. MapReduce enhances data processing by exploiting Hadoop's distributed nature. It enables parallel calculations to be executed throughout the cluster, allowing for efficient processing of large-scale datasets. Organizations may use MapReduce to execute complex analytics, computations, and transformations on huge data, finding patterns, insights, and trends that can help with learning and decision-making.

4. Fault Tolerance: Fault tolerance is a critical attribute of Hadoop that ensures data reliability and availability. Hadoop achieves fault tolerance through various mechanisms.

- Data Replication: HDFS, Hadoop's distributed file system, replicates data across the cluster's DataNodes. Each data block is usually duplicated three times to provide redundancy and data availability even if a node fails. If a node fails, the system seamlessly switches data processing to other available replicas, eliminating data loss and disrupting the learning process.
- Job Tracker and Task Tracker: In Hadoop's MapReduce framework, the Job Tracker manages and tracks the execution of MapReduce jobs, while the Task Tracker oversees individual map and reduce tasks. If a Task Tracker fails, the Job Tracker automatically reassigns the failed tasks to other available Task Trackers, ensuring uninterrupted job execution. This fault tolerance mechanism prevents job failures and allows for continuous learning and analysis of big data.

5. Data-Locality:

Data locality is a key aspect of Hadoop that optimizes data processing performance by minimizing data transfer and network overhead. The principle of data locality ensures that data processing tasks are executed on the nodes where the data resides. This approach offers several advantages:

- Reduced Network Latency: By processing data locally, Hadoop minimizes the need to transfer large volumes of data across the network. This reduces network latency and accelerates data processing, leading to faster insights and decision-making.

- **Efficient Resource Utilization:** Data locality maximizes the utilization of cluster resources by leveraging the processing power available on each node. This distributed processing approach improves the overall performance and efficiency of big data analytics in Hadoop.

#### 6. Ecosystem of Tools:

Hadoop provides a rich ecosystem of tools and frameworks that complement its core components and enhance the learning and decision-making processes:

- **Apache Spark:** Spark is a fast and general-purpose cluster computing system that runs on top of Hadoop. It offers in-memory processing capabilities, allowing for faster data analytics and iterative computations. Spark provides high-level APIs in Java, Scala, Python, and R, making it easier for analysts and data scientists to leverage big data for learning and decision-making.
- **Hive and Pig:** Hive and Pig are high-level data query and scripting languages that make large data research easier. They provide SQL-like interfaces for interacting with Hadoop, allowing for data manipulations, querying, and summarization. Hive and Pig provide a user-friendly abstraction layer that enables businesses to extract insights and do ad-hoc analysis on massive datasets without requiring complicated programming.
- **HBase:** HBase is a distributed, scalable, NoSQL database that runs on Hadoop. It allows for random and real-time access to massive amounts of structured data, making it ideal for applications that demand quick data retrieval and storage. HBase enhances Hadoop's batch processing capabilities by allowing for interactive data exploration and real-time decision-making.

Hadoop's tool ecosystem extends its capabilities beyond basic data processing and storage, providing specialized capability for a variety of analytical applications. These technologies enable companies to draw meaningful insights from big data and assist the learning and decision-making processes by providing higher-level abstractions, sophisticated analytics, and data querying capabilities.

## CONCLUDING REMARKS

In conclusion, through the exploration of these tasks, we have gained valuable insights into the world of big data and its applications in Hadoop.

Firstly, we have learned that big data refers to extremely large and complex datasets that cannot be easily managed or processed using traditional data processing techniques. Big data is characterized by the three V's: volume, velocity, and variety. It requires specialized infrastructure, algorithms, and tools to extract meaningful insights and value from the data.

Secondly, we have analyzed Distributed File Systems (DFSs) and discovered that they provide a scalable and fault-tolerant solution for storing and processing big data. DFSs, such as Hadoop's HDFS, distribute data across multiple nodes in a cluster, enabling parallel processing and fault tolerance. This architecture is highly supportive in the science of big data as it facilitates efficient data storage, retrieval, and processing.

Lastly, we have critically examined how we can learn from our big data in Hadoop. Hadoop offers key attributes and properties that support the learning and decision-making processes. Its fault tolerance mechanisms ensure data reliability and availability, while data locality minimizes network overhead and improves processing performance. Additionally, Hadoop's ecosystem of tools provides advanced analytics capabilities, allowing organizations to derive valuable insights from big data.

Overall, the knowledge gained from these tasks highlights the importance of understanding and harnessing big data through technologies like Hadoop. By effectively managing and analyzing big data, organizations can unlock new opportunities, make informed decisions, and gain a competitive edge in today's data-driven world.

## BIBLIOGRAPHY

1. "Big Data: A Revolution That Will Transform How We Live, Work, and Think" by Viktor Mayer-Schönberger and Kenneth Cukier.
2. This book provides a comprehensive overview of big data, its fundamental concepts, and its impact on various aspects of society. "Hadoop: The Definitive Guide" by Tom White.
3. This book offers an in-depth exploration of Hadoop, covering its architecture, components, and practical usage for processing big data. "Distributed Systems: Principles and Paradigms" by Andrew S. Tanenbaum and Maarten van Steen.
4. This textbook provides a comprehensive introduction to distributed systems, including distributed file systems, and covers fundamental concepts and design principles. "Apache Hadoop Documentation" - Official Apache Hadoop website.
5. "Hadoop: The Definitive Guide" by Tom White (Online edition):  
Link: <https://www.oreilly.com/library/view/hadoop-the-definitive/9781491901687>
6. "Apache Hadoop Documentation" - Official Apache Hadoop website:  
Link: <https://hadoop.apache.org/docs>

## **APPENDIX (if necessary)**

