



CS-704 (Big Data)

Practical File

Submitted To :

Ms. Aruna Bajpai

Asst. Professor

Computer Science Department

Submitted By:

**Shivani
kushwah**

0905CS201162

Index

S.No.	Practical List	Date	Signature
1.	Case study on Big data and its challenges.		
2.	Case study on hadoop and MapReduce.		
3.	Case study on Hive and Pig.		
4.	Case study to perform analytics on No SQL and Spark technology.		
5.	Installation of hadoop in pseudo Distributed mode.		
6.	Word count MapReduce program to understand Mapreduce Paradigm.		
7.	Case study for solving big data analytics 1) AWS 2) Microsoft Azure		

PRACTICAL – 1

Case study on Big data and its challenges.

Title: Big Data: Challenges and Solutions

Introduction:

Big data has become an integral part of modern business and technology landscapes, revolutionizing the way organizations handle, analyze, and derive insights from vast amounts of data. This case study delves into the world of big data, highlighting its significance, key challenges, and potential solutions.

Background:

In recent years, the proliferation of digital technologies, the internet, and IoT devices has resulted in an exponential increase in the volume, velocity, and variety of data generated globally. Big data, characterized by its enormous size and complexity, poses both opportunities and challenges for organizations across various sectors.

Key Challenges:

1. Volume:

Big data is characterized by its sheer volume, often exceeding the processing capabilities of traditional databases. Organizations struggle to store, manage, and analyze petabytes of data efficiently.

2. Velocity:

The speed at which data is generated, processed, and analyzed is a significant challenge. Real-time analytics and decision-making require advanced technologies to keep pace with the constant influx of data.

3. Variety:

Big data encompasses diverse data types, including structured, semi-structured, and unstructured data. Integrating and analyzing data from different sources can be complex, requiring sophisticated tools and technologies.

4. Veracity:

The reliability and quality of data can be questionable, especially when dealing with unstructured data sources. Ensuring data accuracy and consistency is a crucial challenge for organizations aiming to derive meaningful insights.

5. Value:

Extracting valuable insights from big data is a persistent challenge. Many organizations face difficulties in translating massive datasets into actionable information that contributes to strategic decision-making.

Solutions:

1. Advanced Analytics:

Implementing advanced analytics tools, such as machine learning and artificial intelligence, enables organizations to derive meaningful insights from large datasets, predict trends, and make data-driven decisions.

2. Scalable Infrastructure:

Investing in scalable infrastructure, such as cloud computing, allows organizations to expand their data storage and processing capabilities based on demand. Cloud platforms offer cost-effective solutions for managing large volumes of data.

3. Data Integration Platforms:

Leveraging robust data integration platforms helps organizations handle a variety of data formats. These platforms facilitate seamless integration, transformation, and analysis of data from different sources.

4. Data Quality Management:

Establishing robust data quality management practices ensures the accuracy and reliability of big data. Regular data cleansing, validation, and quality checks are essential components of effective data governance.

5. Cross-functional Collaboration:

Promoting collaboration between IT, data science, and business teams is crucial for unlocking the value of big data. Cross-functional teams can work together to identify business-specific challenges and develop tailored solutions.

Conclusion:

While big data presents unparalleled opportunities for innovation and growth, organizations must navigate various challenges to harness its full potential. By implementing advanced technologies, scalable infrastructure, and effective data management practices, businesses can transform big data challenges into strategic advantages, gaining valuable insights and maintaining a competitive edge in the digital era.

PRACTICAL – 2

Case study on hadoop and MapReduce

Title: A Case Study - Hadoop and MapReduce

Introduction:

In the rapidly evolving landscape of big data, organizations are constantly seeking innovative solutions to handle and analyze massive datasets efficiently. This case study explores the implementation of Hadoop and MapReduce, two powerful technologies that have revolutionized data processing, storage, and analysis.

Background:

Traditionally, processing large datasets posed significant challenges, as conventional databases struggled to cope with the volume, velocity, and variety of data. Hadoop, an open-source framework, emerged as a game-changer by offering a scalable, distributed storage and processing solution. Complementing Hadoop, MapReduce, a programming model, enables parallel processing of data across a distributed cluster.

The Challenge:

Our organization, a leading e-commerce company, faced escalating challenges in handling vast amounts of customer data. The traditional relational database system was proving inefficient in managing the ever-growing datasets generated by user transactions, website interactions, and customer feedback.

Implementation:

1. Hadoop Framework:

The first step involved deploying Hadoop to create a distributed file system capable of storing and managing petabytes of data. This scalable infrastructure allowed seamless expansion based on the organization's evolving data requirements.

2. MapReduce Programming Model:

Leveraging the MapReduce programming model, our data engineering team developed algorithms to process and analyze data in parallel across the Hadoop cluster. This approach facilitated the efficient distribution of computing tasks, significantly reducing processing times.

Results:

1. Scalability:

Hadoop's distributed architecture allowed the organization to scale its storage and processing capabilities effortlessly. The system accommodated the growing volume of data without compromising performance.

2. Speed and Efficiency:

MapReduce improved data processing speed by breaking down complex tasks into smaller, parallelizable components. This resulted in a substantial reduction in the time required to process and derive insights from large datasets.

3. Cost-effectiveness:

The switch to Hadoop and MapReduce brought about cost savings by eliminating the need for expensive, proprietary hardware. The organization adopted a cloud-based Hadoop solution, optimizing costs based on usage.

4. Enhanced Analytics:

With the new infrastructure in place, the analytics team gained the ability to run complex queries and analyses on vast datasets. The insights derived from this enhanced analytics capability empowered decision-makers to make informed, data-driven choices.

Challenges Overcome:

1. Data Variety:

Hadoop's ability to handle various data formats, including structured, semi-structured, and unstructured data, addressed the challenge of data variety faced by the organization.

2. Processing Time:

The parallel processing capabilities of MapReduce significantly reduced the time required for complex data processing tasks, overcoming the challenge of processing large volumes of data within tight timelines.

Conclusion:

The successful implementation of Hadoop and MapReduce transformed our organization's data processing capabilities, providing scalable, efficient, and cost-effective solutions to handle big data challenges. By adopting these technologies, the organization has positioned itself at the forefront of data-driven decision-making, setting the stage for continued innovation and growth in the dynamic landscape of e-commerce.

Future Considerations:

As our organization continues to evolve, we are exploring advancements in the Hadoop ecosystem, such as Apache Spark, to further enhance real-time processing capabilities and accommodate the ever-expanding data needs of our business. The journey with Hadoop and MapReduce has laid a robust foundation for our data infrastructure, ensuring our readiness to tackle future challenges in the big data domain.

Practical – 3

Case study on Hive and Pig

Title: A Case Study : Hive and Pig

Introduction:

In the realm of big data analytics, organizations face the challenge of efficiently processing and analyzing large volumes of diverse data. This case study explores the implementation of Apache Hive and Apache Pig, two powerful tools that simplify and optimize data processing on Hadoop, providing scalable and user-friendly solutions for our organization.

Background:

Our organization, a global financial services provider, manages extensive datasets from various sources, including customer transactions, market data, and regulatory information. The need to process and derive insights from this vast and diverse data prompted the exploration of advanced tools for big data processing.

The Challenge:

Traditional data processing tools struggled to cope with the scale and complexity of our datasets. The organization needed a solution that could handle structured and semi-structured data efficiently, enabling data analysts to perform complex queries and analytics without extensive programming skills.

Implementation:

1. Apache Hive:

- **Data Warehousing** : Hive was deployed to create a data warehousing solution that provides a familiar SQL-like interface for data analysts. This allowed them to query and analyze structured data using Hive Query Language (HQL).
- **Metadata Management** : Hive's metadata management capabilities simplified data organization and facilitated seamless integration with existing business intelligence tools. This improved data governance and accessibility for the analytics team.

2. Apache Pig:

- **Data Flow Scripting** : Pig was adopted for its scripting language, Pig Latin, which enabled data analysts to create data processing pipelines without extensive programming knowledge. This streamlined the development of complex data processing workflows.
- **Flexibility and Extensibility** : Pig's flexibility allowed the organization to process diverse data formats, including unstructured and semi-structured data. Additionally, Pig's extensibility through User Defined Functions (UDFs) enabled custom processing for specific business requirements.

Results:

1. Improved Query Performance:

With Hive, data analysts gained the ability to run complex SQL-like queries on structured data, significantly improving query performance and reducing the time required for data retrieval.

2. Enhanced Data Processing Workflows:

Pig's data flow scripting simplified the creation of data processing workflows. Analysts could easily define and execute multi-stage data transformations, reducing development time and enhancing workflow efficiency.

3. Scalability and Adaptability:

Both Hive and Pig demonstrated scalability, allowing the organization to seamlessly adapt to growing data volumes. The tools proved effective in processing large datasets, ensuring consistent performance even as data demands increased.

4. User Empowerment:

The user-friendly interfaces of Hive and Pig empowered data analysts to perform complex data processing tasks without deep programming expertise. This democratization of data processing capabilities enhanced collaboration and knowledge sharing within the analytics team.

Challenges Overcome:

1. Data Variety:

Hive and Pig effectively addressed the challenge of handling diverse data types, including structured, semi-structured, and unstructured data, providing a unified platform for data processing.

2. User Skill Levels:

The intuitive interfaces of Hive and Pig mitigated the challenge of varying programming skill levels within the analytics team, allowing analysts with SQL or scripting experience to contribute effectively to big data processing tasks.

Conclusion:

The implementation of Hive and Pig has revolutionized our organization's approach to big data processing. By providing user-friendly interfaces, scalable solutions, and efficient processing capabilities, these tools have empowered our analytics team to derive valuable insights from diverse and extensive datasets. The success of this implementation positions our organization for continued innovation and growth in the dynamic landscape of financial data analytics.

Future Considerations:

As our organization continues to evolve, we are exploring advanced features and updates within the Hive and Pig ecosystems, ensuring that we stay at the forefront of big data processing capabilities. The success of this case study sets the stage for future enhancements and innovations in our data analytics journey.

PRACTICAL - 4

Case study to perform analytics on NoSQL and Spark technology.

Title : A Case Study on NoSQL and Spark Technology for Advanced Analytics

Introduction:

In the era of big data, organizations are continually seeking robust solutions for efficient data storage, processing, and analytics. This case study explores the integration of NoSQL databases and Apache Spark, two cutting-edge technologies, to perform advanced analytics on large and diverse datasets.

Background:

Our organization, a multinational e-commerce platform, faced challenges in handling vast amounts of data generated by user interactions, transactions, and product analytics. Traditional relational databases were struggling to scale with the exponential growth of data, necessitating the exploration of innovative technologies.

The Challenge:

The primary challenges included the need for a flexible and scalable database system capable of handling diverse data types, coupled with the requirement for a high-performance analytics engine that could process large volumes of data in real-time.

Implementation:

1. NoSQL Database:

A NoSQL database, specifically MongoDB, was chosen for its flexibility in handling unstructured and semi-structured data. MongoDB's document-oriented model accommodated the dynamic nature of e-commerce data, allowing for easy storage and retrieval.

2. Apache Spark:

Apache Spark was integrated to provide a powerful analytics engine capable of processing data in-memory and supporting real-time analytics. The Spark ecosystem, including Spark SQL and MLlib, was leveraged for querying, machine learning, and data processing tasks.

Results:

1. Flexible Data Modeling:

MongoDB's document-oriented structure allowed the organization to model and store complex data structures without the need for a predefined schema. This flexibility proved invaluable in adapting to the ever-changing nature of e-commerce data.

2. Real-time Analytics:

Spark's in-memory processing capabilities facilitated real-time analytics on streaming data. This empowered the analytics team to derive insights promptly, enabling quick responses to changing market dynamics and user behaviors.

3. Scalability:

Both MongoDB and Spark demonstrated scalability as data volumes increased. Horizontal scaling capabilities of MongoDB ensured seamless expansion, while Spark's distributed computing model allowed for the parallel processing of large datasets.

4. Predictive Analytics:

Spark's MLlib was utilized for predictive analytics, enabling the creation and deployment of machine learning models. The organization leveraged these models to predict user preferences, optimize product recommendations, and enhance the overall customer experience.

Challenges Overcome:

1. Data Variety:

NoSQL databases, with their schema-less design, effectively addressed the challenge of handling diverse data types, including user reviews, clickstream data, and product images.

2. Real-time Processing:

Spark's ability to process data in-memory and support real-time analytics overcame the challenge of delayed insights. This was particularly crucial in the dynamic e-commerce environment where timely decision-making is paramount.

Conclusion:

The integration of NoSQL and Spark technologies has revolutionized our organization's approach to data management and analytics. By combining the flexibility of NoSQL databases with the processing power of Spark, we have created a dynamic ecosystem that not only addresses our current challenges but also positions us for future growth and innovation in the rapidly evolving landscape of e-commerce analytics.

Future Considerations:

As our organization continues to evolve, we are exploring further advancements within the Spark ecosystem, including Spark Streaming for enhanced real-time processing and Delta Lake for improved data reliability. The success of this case study lays the foundation for continuous exploration and adoption of cutting-edge technologies to stay at the forefront of data analytics capabilities.

PRACTICAL – 5

Installation of Hadoop in pseudo Distributed mode.

Setting up Hadoop in pseudo-distributed mode allows you to run Hadoop on a single machine, simulating a multi-node cluster environment. This is useful for development, testing, and learning purposes. Below are the steps to install Hadoop in pseudo-distributed mode. These instructions assume you are using Linux.

Prerequisites:

1. Java Development Kit (JDK) installed. Hadoop requires Java.
2. SSH configured to allow passwordless login.

Steps:

1. Download Hadoop:

- Visit the official Apache Hadoop website: Hadoop Releases.
- Download the latest stable release. As of my knowledge cutoff in January 2022, you can use a command like:

```
wget https://downloads.apache.org/hadoop/common/hadoop-X.X.X/hadoop-X.X.X.tar.gz
```

Replace X.X.X with the version number you downloaded.

2. Extract the Tarball:

```
tar -xvzf hadoop-X.X.X.tar.gz
```

This will create a directory named hadoop-X.X.X.

3. Set Environment Variables:

- Edit your shell profile file (e.g., ~/.bashrc or ~/.zshrc) to add the following lines:

```
export HADOOP_HOME=/path/to/hadoop-X.X.X
export PATH=$PATH:$HADOOP_HOME/bin
```

- Update the environment by running:

```
source ~/.bashrc # or source ~/.zshrc
```

4.

Configure Hadoop:

- Navigate to the etc/hadoop/ directory within your Hadoop installation directory.
- Edit hadoop-env.sh to set the JAVA_HOME variable:

```
export JAVA_HOME=/path/to/your/java
```

- Copy mapred-site.xml.template to mapred-site.xml:

```
cp mapred-site.xml.template mapred-site.xml
```

- Open mapred-site.xml and add the following configuration:

```
<configuration>
  <property>
    <name>mapreduce.framework.name</name>
    <value>yarn</value>
  </property>
</configuration>
```

- Copy yarn-site.xml.template to yarn-site.xml:

```
cp yarn-site.xml.template yarn-site.xml
```

- Open yarn-site.xml and add the following configuration:

```
<configuration>
  <property>
    <name>yarn.nodemanager.aux-services</name>
    <value>mapreduce_shuffle</value>
  </property>
</configuration>
```

5. Format HDFS :

```
hdfs namenode -format
```

6. Start Hadoop:

Start the Hadoop services:

```
start-dfs.sh
start-yarn.sh
```

7. Verify Installation:

Open your web browser and visit the Hadoop web interface at <http://localhost:50070/>. This should show the HDFS status page.

Additionally, you can check the ResourceManager web interface at <http://localhost:8088/>.

8. Run a Sample MapReduce Job:

Hadoop comes with some sample MapReduce jobs. Run one to verify your

setup:

```
hadoop jar $HADOOP_HOME/share/hadoop/mapreduce/hadoop-mapreduce-examples-X.X.X.jar pi
```

Replace X.X.X with your Hadoop version.

Congratulations! You've successfully set up Hadoop in pseudo-distributed mode. This setup allows you to experiment with Hadoop and understand its basic functionality. For a production environment, a fully distributed setup with a cluster of machines is recommended.

PRACTICAL – 6

Word count MapReduce program to understand Mapreduce Paradigm.

Certainly! The classic Word Count program is a great way to understand the MapReduce paradigm. This program counts the frequency of each word in a given set of documents. Here's a simple example using Hadoop MapReduce in Java:

Mapper Class (WordCountMapper.java):

```
1  import java.io.IOException;
2  import java.util.StringTokenizer;
3
4  import org.apache.hadoop.io.IntWritable;
5  import org.apache.hadoop.io.LongWritable;
6  import org.apache.hadoop.io.Text;
7  import org.apache.hadoop.mapreduce.Mapper;
8
9  public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
10
11     private final static IntWritable one = new IntWritable(1);
12     private Text word = new Text();
13
14     public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
15         String line = value.toString();
16         StringTokenizer tokenizer = new StringTokenizer(line);
17
18         while (tokenizer.hasMoreTokens()) {
19             word.set(tokenizer.nextToken());
20             context.write(word, one);
21         }
22     }
23 }
```

Reducer Class (WordCountReducer.java):

```
1  import java.io.IOException;
2
3  import org.apache.hadoop.io.IntWritable;
4  import org.apache.hadoop.io.Text;
5  import org.apache.hadoop.mapreduce.Reducer;
6
7  public class WordCountReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
8
9      private IntWritable result = new IntWritable();
10
11     public void reduce(Text key, Iterable<IntWritable> values, Context context)
12         throws IOException, InterruptedException {
13         int sum = 0;
14         for (IntWritable value : values) {
15             sum += value.get();
16         }
17         result.set(sum);
18         context.write(key, result);
19     }
20 }
```

Driver Class (WordCountDriver.java):

```
1  import org.apache.hadoop.conf.Configuration;
2  import org.apache.hadoop.fs.Path;
3  import org.apache.hadoop.io.IntWritable;
4  import org.apache.hadoop.io.Text;
5  import org.apache.hadoop.mapreduce.Job;
6  import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
7  import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
8
9  public class WordCountDriver {
10
11      public static void main(String[] args) throws Exception {
12          Configuration conf = new Configuration();
13          Job job = Job.getInstance(conf, "word count");
14          job.setJarByClass(WordCountDriver.class);
15          job.setMapperClass(WordCountMapper.class);
16          job.setCombinerClass(WordCountReducer.class);
17          job.setReducerClass(WordCountReducer.class);
18          job.setOutputKeyClass(Text.class);
19          job.setOutputValueClass(IntWritable.class);
20          FileInputFormat.addInputPath(job, new Path(args[0]));
21          FileOutputFormat.setOutputPath(job, new Path(args[1]));
22          System.exit(job.waitForCompletion(true) ? 0 : 1);
23      }
24  }
```

How to Run:

1. Compile the Java files and create a JAR file.
2. Create a directory on HDFS and upload some text files to it.
3. Run the WordCount program with Hadoop:

```
hadoop jar WordCount.jar WordCountDriver input_directory output_directory
```

Replace WordCount.jar with the name of your JAR file, input_directory with the directory containing your input text files, and output_directory with the desired output directory.

This Word Count example demonstrates the key concepts of the MapReduce paradigm: the Mapper class processes input data and emits intermediate key-value pairs, the Combiner (optional) performs a local aggregation of the intermediate data, and the Reducer class processes the intermediate data and produces the final output.

PRACTICAL – 7

Case study for solving big data analytics

1) AWS

2) Microsoft Azure

Case Study: Big Data Analytics with AWS and Microsoft Azure

Challenges:

1. Data Integration and Variety:

The company grappled with integrating diverse financial data sources, including transaction records, market feeds, and customer interactions. The challenge was to harmonize this data for comprehensive analytics.

2. Scalability and Performance:

The existing on-premises infrastructure struggled to cope with the increasing volume of financial data. The company sought a scalable solution to handle peak loads during market events and high transaction volumes.

3. Real-time Analytics:

Real-time insights were crucial for risk management and decision-making in financial markets. The company needed a platform that could process and analyze streaming data in real-time.

4. Cost Optimization:

Cost-effectiveness was a key consideration. The company aimed to optimize infrastructure costs while ensuring a robust and high-performance analytics environment.

AWS Solution:

1. Data Integration and Variety:

- Utilized AWS Glue for automated ETL (Extract, Transform, Load) operations, ensuring seamless data integration from various sources.
- Implemented Amazon Athena for ad-hoc querying of data directly from Amazon S3, simplifying analytics on diverse datasets.

2. Scalability and Performance:

- Configured Amazon Redshift, a fully-managed data warehouse, for high-performance analytics and scalability during peak loads.
- Leveraged Amazon EMR (Elastic MapReduce) for distributed processing, allowing the company to scale resources based on demand.

3. Real-time Analytics:

- Employed Amazon Kinesis Data Streams to ingest and process streaming financial data in real-time.
- Integrated Kinesis Analytics for immediate insights, enabling timely risk assessments and decision-making.

4. Cost Optimization:

- Leveraged AWS Lambda for serverless computing, optimizing costs by aligning resources with actual usage.
- Utilized AWS Cost Explorer and AWS Budgets for continuous monitoring and optimization of infrastructure costs.

Microsoft Azure Solution:

1. Data Integration and Variety:

- Utilized Azure Data Factory for orchestrating and automating ETL workflows, ensuring efficient data integration.
- Implemented Azure Synapse Analytics (formerly SQL Data Warehouse) for data warehousing and analytics on diverse datasets.

2. Scalability and Performance:

- Configured Azure HDInsight for scalable big data analytics, leveraging Apache Spark and Apache Hive for processing large datasets.
- Implemented Azure Virtual Machines to scale computing resources based on demand.

3. Real-time Analytics:

- Deployed Azure Stream Analytics for real-time processing of streaming financial data, providing immediate insights.
- Utilized Azure Databricks for real-time analytics and collaborative data science.

4. Cost Optimization:

- Implemented Azure Functions for serverless computing, optimizing costs by executing functions in response to events.
- Utilized Azure Cost Management and Billing for continuous monitoring and optimization of infrastructure costs.

Results:

AWS:

- Achieved seamless data integration and variety handling with AWS Glue and Athena.
- Ensured high-performance analytics and scalability with Amazon Redshift and EMR.
- Enabled real-time analytics and decision-making with Amazon Kinesis.

Azure:

- Streamlined data integration with Azure Data Factory and Synapse Analytics.
- Achieved scalability and performance with HDInsight and Azure Virtual Machines.
- Enabled real-time analytics with Azure Stream Analytics and Databricks.

Conclusion:

By harnessing the capabilities of both AWS and Microsoft Azure, the Financial Services Conglomerate successfully addressed its challenges in big data analytics. The combined strengths of AWS and Azure provided a holistic solution, ensuring seamless data integration, scalability, real-time analytics, and cost optimization. This multi-cloud approach allowed the company to leverage the best features of each cloud platform, driving innovation and efficiency in their financial analytics processes.

Future Considerations:

The Financial Services Conglomerate plans to explore advanced analytics and machine learning services within both AWS and Azure to further enhance predictive modeling, fraud detection, and customer insights. The multi-cloud strategy ensures flexibility and adaptability for future advancements in the dynamic financial services landscape.

