



(<https://swayam.gov.in>)



([https://swayam.gov.in/nc\\_details/NPTEL](https://swayam.gov.in/nc_details/NPTEL))

teenakb1@gmail.com ↗

NPTEL (<https://swayam.gov.in/explorer?ncCode=NPTEL>) » Big Data Computing (course)



If already registered,  
click to check your  
payment status

## Course outline

[How does an  
NPTEL online  
course work? \(\)](#)

[Week-0 \(\)](#)

**Practice: Week 0:**  
[Assignment 0  
\(assessment?  
name=112\)](#)

[Week-1 \(\)](#)

[Week-2 \(\)](#)

[Week-3 \(\)](#)

[Week-4 \(\)](#)

[Week-5 \(\)](#)

[Week-6 \(\)](#)

[Week-7 \(\)](#)

[Week-8 \(\)](#)

[Text Transcripts \(\)](#)

[DOWNLOAD  
VIDEOS \(\)](#)

[Books \(\)](#)

# Week 0: Assignment 0

Your last recorded submission was on 2023-08-28, 20:41 IST

- 1) The maximum number of super keys for the relation schema R(E,F,G,H) with E as the key is 1 point

- 5
- 6
- 7
- 8

Yes, the answer is correct.  
Score: 1

Accepted Answers:  
8

- 2) Consider the following relational schemas for a library database: 1 point

Book (Title, Author, Catalog\_no, Publisher, Year, Price)

Collection (Title, Author, Catalog\_no)

with the following functional dependencies:

- I. Title Author --> Catalog\_no
- II. Catalog\_no --> Title, Author, Publisher, Year
- III. Publisher Title Year --> Price

Assume {Author, Title} is the key for both schemas. Which of the following statements is true?

- Both Book and Collection are in BCNF
- Both Book and Collection are in 3NF only
- Book is in 2NF and Collection is in 3NF
- Both Book and Collection are in 2NF only

Yes, the answer is correct.  
Score: 1

Accepted Answers:  
Book is in 2NF and Collection is in 3NF

- 3) Consider a B+-tree in which the maximum number of keys in a node is 5. What is the minimum number of keys in any non-root node ? 1 point

- 1
- 2
- 3
- 4

Yes, the answer is correct.  
Score: 1

Accepted Answers:  
2

- 4) Consider the join of a relation R , with a relation S . If R has m number of tuples and S



has n number of tuples then the maximum and minimum sizes of the join respectively are:

- m + n & 0
- mn & 0
- m + n & | m - n |
- mn & m + n

Yes, the answer is correct.

Score: 1

Accepted Answers:

*mn & 0*

5) Which one of the following is NOT a part of the ACID properties of database transactions ?

- Atomicity
- Consistency
- Isolation
- Deadlock-freedom

Yes, the answer is correct.

Score: 1

Accepted Answers:

*Deadlock-freedom*

6) In the IPv4 addressing format, the number of networks allowed under Class C addresses is

- $2^{14}$
- $2^7$
- $2^{21}$
- $2^{24}$

Yes, the answer is correct.

Score: 1

Accepted Answers:

*$2^{21}$*

7) One of the header fields in an IP datagram is the Time to Live (TTL) field. Which of the following statements best explains the need for this field ?

- It can be used to prioritize packets
- It can be used to reduce delays
- It can be used to optimize throughput
- It can be used to prevent packet looping

Yes, the answer is correct.

Score: 1

Accepted Answers:

*It can be used to prevent packet looping*

8) The address resolution protocol (ARP) is used for

- Finding the IP address from the DNS
- Finding the IP address of the default gateway
- Finding the IP address that corresponds to a MAC address
- Finding the MAC address that corresponds to an IP address

Yes, the answer is correct.

Score: 1

Accepted Answers:

*Finding the MAC address that corresponds to an IP address*

9) Consider different activities related to email:

m1: Send an email from a mail client to a mail server

m2: Download an email from mailbox server to a mail client

m3: Checking email in a web browser

Which is the application level protocol used in each activity?

- m1: HTTP m2: SMTP m3: POP
- m1: SMTP m2: FTP m3: HTTP
- m1: SMTP m2: POP m3: HTTP

**1 point**

**1 point**

**1 point**

**1 point**



m1: POP m2: SMTP m3: IMAP

Yes, the answer is correct.

Score: 1

Accepted Answers:

*m1: SMTP m2: POP m3: HTTP*

10) A process executes the code

`fork();`

`fork();`

`fork();`

The total number of child processes created is

3

4

7

8

**1 point**

Yes, the answer is correct.

Score: 1

Accepted Answers:

7

**Check Answers and Submit**

Your score is: 10/10



## **Quiz Assignment-I Solutions: Big Data Computing (Week-1)**

---

1. What are the three key characteristics of Big Data, often referred to as the 3V's, according to IBM?

- A) Viscosity, Velocity, Veracity
- B) Volume, Value, Variety
- C) Volume, Velocity, Variety
- D) Volumetric, Visceral, Vortex

**Solution:**

**C) Volume, Velocity, Variety**

**Explanation:**

**Volume:** Refers to the massive amount of data generated and collected from various sources. This includes both structured and unstructured data.

**Velocity:** Represents the speed at which data is generated, processed, and analyzed. It emphasizes the real-time nature of data and the need to handle and react to data quickly.

**Variety:** Encompasses the different types and formats of data, including structured, semi-structured, and unstructured data. This diversity challenges traditional data processing methods.

Option A is incorrect because "Viscosity" is not one of the 3V's, and "Veracity" relates to the accuracy and trustworthiness of data, not velocity.

Option B is incorrect because while "Volume" and "Variety" are correct, "Value" is not one of the 3V's.

Option D is incorrect because "Volumetric," "Visceral," and "Vortex" are not the terms used to describe the characteristics of Big Data according to IBM.

2. What is the primary purpose of the MapReduce programming model in processing and generating large data sets?

- A) To directly process and analyze data without any intermediate steps.
- B) To convert unstructured data into structured data.

- C) To specify a map function for generating intermediate key/value pairs and a reduce function for merging values associated with the same key.
- D) To create visualizations and graphs for large data sets.

Solution:

- C) To specify a map function for generating intermediate key/value pairs and a reduce function for merging values associated with the same key.

Explanation:

MapReduce is a programming model used for processing and generating large data sets. It involves two main steps: mapping and reducing. Users specify a map function that processes a key/value pair to generate a set of intermediate key/value pairs. The map function operates in parallel across the input data. The intermediate key/value pairs are then grouped by key and passed to a reduce function, which merges all intermediate values associated with the same intermediate key. This process allows for distributed and parallel processing of large datasets.

Option A is incorrect because MapReduce does involve intermediate steps (mapping and reducing) to process data.

Option B is incorrect because while MapReduce is used for processing unstructured data, its primary purpose is not to convert it into structured data.

Option D is incorrect because MapReduce is not primarily focused on creating visualizations and graphs; its main focus is on processing and generating large data sets using the map and reduce functions.

3. \_\_\_\_\_ is a distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of log data.

- A) Flume
- B) Apache Sqoop
- C) Pig
- D) Mahout

Solution:

- A) Flume

Explanation:

Flume is a distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of log data. It has a simple and very flexible architecture based on streaming data flows. It's quite robust and fail tolerant, and it's really tunable to enhance the reliability mechanisms, fail over, recovery, and all the other mechanisms that keep the cluster safe and reliable. It uses simple extensible data model that allows us to apply all kinds of online analytic applications.

4. What is the primary role of YARN (Yet Another Resource Manager) in the Apache Hadoop ecosystem?

- A) YARN is a data storage layer for managing and storing large datasets in Hadoop clusters.
- B) YARN is a programming model for processing and analyzing data in Hadoop clusters.
- C) YARN is responsible for allocating system resources and scheduling tasks for applications in a Hadoop cluster.
- D) YARN is a visualization tool for creating graphs and charts based on Hadoop data.

Solution:

C) YARN is responsible for allocating system resources and scheduling tasks for applications in a Hadoop cluster.

Explanation:

YARN, which stands for "Yet Another Resource Manager," is a key component of the Apache Hadoop ecosystem. Its primary role is resource management and job scheduling. YARN is responsible for efficiently allocating system resources, such as CPU and memory, to various applications running in a Hadoop cluster. It also handles the scheduling of tasks to be executed on different cluster nodes, ensuring optimal utilization of resources and improving overall cluster performance.

Option A is incorrect because YARN is not a data storage layer; it focuses on resource management and job scheduling.

Option B is incorrect because while YARN plays a role in supporting data processing and analysis, its main function is not to define a programming model.

Option D is incorrect because YARN is not a visualization tool; it is a resource management and scheduling technology.

5. Which of the following statements accurately describes the characteristics and functionality of HDFS (Hadoop Distributed File System)?
- A) HDFS is a centralized file system designed for storing small files and achieving high-speed data processing.
  - B) HDFS is a programming language used for writing MapReduce applications within the Hadoop ecosystem.
  - C) HDFS is a distributed, scalable, and portable file system designed for storing large files across multiple machines, achieving reliability through replication.
  - D) HDFS is a visualization tool that generates graphs and charts based on data stored in the Hadoop ecosystem.

**Solution:**

C) HDFS is a distributed, scalable, and portable file system designed for storing large files across multiple machines, achieving reliability through replication.

**Explanation:**

HDFS (Hadoop Distributed File System) is a fundamental component of the Hadoop framework. It is designed to store and manage large files across a distributed cluster of machines. The key features and functionality of HDFS include:

**Distributed and Scalable:** HDFS distributes data across multiple nodes in a cluster, allowing it to handle large datasets that range from gigabytes to terabytes, and even petabytes. It scales horizontally as more nodes are added to the cluster.

**Reliability Through Replication:** HDFS achieves reliability by replicating data blocks across multiple data nodes in the cluster. This replication ensures data availability even in the face of node failures.

**Single Name Node and Data Nodes:** Each Hadoop instance typically includes a single name node, which acts as the metadata manager for the file system, and a cluster of data nodes that store the actual data.

**Portability:** HDFS is written in Java and is designed to be portable across different platforms and operating systems.

**Option A is incorrect because HDFS is not centralized; it is distributed. It is also designed for storing large files rather than small files.**

**Option B is incorrect because HDFS is not a programming language; it is a file system.**

**Option D is incorrect because HDFS is not a visualization tool; it is a distributed file system for storing and managing data in the Hadoop ecosystem.**

6. Which statement accurately describes the role and design of HBase in the Hadoop stack?
- A) HBase is a programming language used for writing complex data processing algorithms in the Hadoop ecosystem.
  - B) HBase is a data warehousing solution designed for batch processing of large datasets in Hadoop clusters.
  - C) HBase is a key-value store that provides fast random access to substantial datasets, making it suitable for applications requiring such access patterns.
  - D) HBase is a visualization tool that generates charts and graphs based on data stored in Hadoop clusters.

**Solution:**

C) HBase is a key-value store that provides fast random access to substantial datasets, making it suitable for applications requiring such access patterns.

**Explanation:**

HBase is a NoSQL database that is a key component of the Hadoop ecosystem. Its design focuses on providing high-speed random access to large amounts of data. Key characteristics and roles of HBase include:

**Key-Value Store:** HBase stores data in a distributed, column-family-oriented fashion, similar to a key-value store. It allows you to look up data quickly using a key.

**Fast Random Access:** HBase is optimized for fast read and write operations, particularly random access patterns. This makes it suitable for applications that require quick retrieval of specific data points from massive datasets.

**Scalability:** HBase is designed to scale horizontally, allowing it to handle vast amounts of data by adding more nodes to the cluster.

Option A is incorrect because HBase is not a programming language; it's a database system.

Option B is incorrect because HBase is not a data warehousing solution; it's designed for real-time, random access to data rather than batch processing.

Option D is incorrect because HBase is not a visualization tool; it's a database system focused on high-speed data access.

7. \_\_\_\_\_ brings scalable parallel database technology to Hadoop and allows users to submit low latencies queries to the data that's stored within the HDFS or the Hbase without acquiring a ton of data movement and manipulation.

- A) Apache Sqoop
- B) Mahout
- C) Flume
- D) Impala

**Solution:**

- D) Impala

**Explanation:**

Cloudera, Impala was designed specifically at Cloudera, and it's a query engine that runs on top of the Apache Hadoop. The project was officially announced at the end of 2012, and became a publicly available, open source distribution. Impala brings scalable parallel database technology to Hadoop and allows users to submit low latencies queries to the data that's stored within the HDFS or the Hbase without acquiring a ton of data movement and manipulation.

## 8. What is the primary purpose of ZooKeeper in a distributed system?

- A) ZooKeeper is a data warehousing solution for storing and managing large datasets in a distributed cluster.
- B) ZooKeeper is a programming language for developing distributed applications in a cloud environment.
- C) ZooKeeper is a highly reliable distributed coordination kernel used for tasks such as distributed locking, configuration management, leadership election, and work queues.
- D) ZooKeeper is a visualization tool for creating graphs and charts based on data stored in distributed systems.

**Solution:**

- C) ZooKeeper is a highly reliable distributed coordination kernel used for tasks such as distributed locking, configuration management, leadership election, and work queues.

**Explanation:**

ZooKeeper is a distributed coordination service that provides a reliable and efficient way for coordinating various processes and components in a distributed system. It offers functionalities like distributed locking, configuration management, leader election, and work queues to ensure that distributed applications can work together effectively. ZooKeeper

acts as a central repository for managing metadata related to the coordination of these distributed tasks.

Option A is incorrect because ZooKeeper is not a data warehousing solution; its primary role is distributed coordination.

Option B is incorrect because ZooKeeper is not a programming language; it's a coordination service.

Option D is incorrect because ZooKeeper is not a visualization tool; it's focused on distributed coordination and management.

9. \_\_\_\_\_ is a distributed file system that stores data on a commodity machine. Providing very high aggregate bandwidth across the entire cluster.

- A) Hadoop Common
- B) Hadoop Distributed File System (HDFS)
- C) Hadoop YARN
- D) Hadoop MapReduce

Solution:

B) Hadoop Distributed File System (HDFS)

Explanation:

Hadoop Common: It contains libraries and utilities needed by other Hadoop modules.

Hadoop Distributed File System (HDFS): It is a distributed file system that stores data on a commodity machine. Providing very high aggregate bandwidth across the entire cluster.

Hadoop YARN: It is a resource management platform responsible for managing compute resources in the cluster and using them in order to schedule users and applications. YARN is responsible for allocating system resources to the various applications running in a Hadoop cluster and scheduling tasks to be executed on different cluster nodes

Hadoop MapReduce: It is a programming model that scales data across a lot of different processes.

10. Which statement accurately describes Spark MLLib?

- A) Spark MLLib is a visualization tool for creating charts and graphs based on data processed in Spark clusters.
- B) Spark MLLib is a programming language used for writing Spark applications in a distributed environment.
- C) Spark MLLib is a distributed machine learning framework built on top of Spark Core, providing scalable machine learning algorithms and utilities for tasks such as classification, regression, clustering, and collaborative filtering.
- D) Spark MLLib is a data warehousing solution for storing and querying large datasets in a Spark cluster.

Solution:

- C) Spark MLLib is a distributed machine learning framework built on top of Spark Core, providing scalable machine learning algorithms and utilities for tasks such as classification, regression, clustering, and collaborative filtering.

Explanation:

Spark MLLib (Machine Learning Library) is a component of the Apache Spark ecosystem. It offers a distributed machine learning framework that allows developers to leverage Spark's distributed computing capabilities for scalable and efficient machine learning tasks. Key features and roles of Spark MLLib include:

**Distributed Machine Learning:** MLLib provides a wide range of machine learning algorithms that are designed to work efficiently in a distributed environment. It enables the processing of large datasets across a cluster of machines.

**Common Learning Algorithms:** MLLib includes a variety of common machine learning algorithms, such as classification, regression, clustering, and collaborative filtering.

**Integration with Spark Core:** MLLib is built on top of Spark Core, which provides the underlying distributed processing framework. This integration allows seamless utilization of Spark's data processing capabilities for machine learning tasks.

Option A is incorrect because Spark MLLib is not a visualization tool; its focus is on distributed machine learning.

Option B is incorrect because Spark MLLib is not a programming language; it's a machine learning library.

Option D is incorrect because Spark MLLib is not a data warehousing solution; its primary purpose is machine learning on distributed data.

## Quiz Assignment-II Solutions: Big Data Computing (Week-2)

---

1. What is the primary purpose of the Map phase in the MapReduce framework?

- A) Combining and aggregating data.
- B) Storing intermediate results.
- C) Sorting and shuffling data.
- D) Applying a user-defined function to each input record.

**Solution:**

- D) Applying a user-defined function to each input record.

**Explanation:**

In the MapReduce framework, the Map phase is responsible for applying a user-defined function (often referred to as the "map function") to each input record individually. This function takes the input record as its input and produces a set of key-value pairs as intermediate outputs. The key-value pairs generated by the map function are then grouped by keys and passed on to the Shuffle and Sort phase, which is responsible for sorting and shuffling the data based on the keys. The Shuffle and Sort phase prepares the data for the subsequent Reduce phase.

Options A, B, and C are not accurate descriptions of the primary purpose of the Map phase:

- A) Combining and aggregating data is a task typically performed in the Reduce phase.
- B) Storing intermediate results is a function of the underlying MapReduce infrastructure and not the primary purpose of the Map phase itself.
- C) Sorting and shuffling data is performed by the Shuffle and Sort phase, which comes after the Map phase.

2. Which of the following statements about the components in the MapReduce framework is true?

Statement 1: The Job Tracker is hosted inside the master and it receives the job execution request from the client.

Statement 2: Task Tracker is the MapReduce component on the slave machine as there are multiple slave machines.

- A) Both statements are true.
- B) Only statement 1 is true.
- C) Only statement 2 is true.
- D) Both statements are false.

Solution:

- A) Both statements are true.

Explanation:

**Statement 1:** The Job Tracker is hosted inside the master and it receives the job execution request from the client.

This statement is true. In the Hadoop MapReduce framework, the Job Tracker is a master node component that manages and coordinates the execution of MapReduce jobs. It receives job execution requests from clients and is responsible for assigning tasks to Task Trackers.

**Statement 2:** Task Tracker is the MapReduce component on the slave machine as there are multiple slave machines.

This statement is also true. The Task Tracker is a slave node component in the Hadoop MapReduce framework. It runs on each slave machine and is responsible for executing individual tasks assigned to it by the Job Tracker. These tasks include both map tasks and reduce tasks.

In summary, both statements accurately describe the roles of the Job Tracker and Task Tracker components within the Hadoop MapReduce framework. Therefore, the correct answer is option A.

3. Which of the following is the slave/worker node and holds the user data in the form of Data Blocks?

- A. NameNode
- B. Data block
- C. Replication
- D. DataNode

Solution:

#### D. DataNode

Explanation:

In the Hadoop Distributed File System (HDFS), a DataNode is a slave/worker node responsible for storing the actual user data in the form of Data Blocks. DataNodes manage these data blocks and respond to requests for read and write operations. The NameNode, on the other hand, is the master node responsible for managing the metadata about the data blocks, such as their locations and replication status.

Option A (NameNode) is incorrect because the NameNode is the master node that maintains the metadata for the HDFS file system.

Option B (Data block) is not the correct answer as it refers to a unit of data storage within HDFS, but it is not a node itself.

Option C (Replication) is not the correct answer either. While replication involves creating multiple copies of data blocks for fault tolerance, it is not a node that holds user data.

Option D (DataNode) is the correct answer as it is the HDFS slave node responsible for storing user data in the form of data blocks.

4. The number of maps in MapReduce is usually driven by the total size of\_\_\_\_\_.

- A. Inputs
- B. Outputs
- C. Tasks
- D. None of the mentioned

Solution:

A. Inputs

Explanation:

In the MapReduce framework, the number of map tasks is usually determined by the total size of the inputs. The input data is divided into chunks, and each map task processes a portion of these chunks. The goal is to achieve parallelism by processing different portions of the input data simultaneously. Therefore, the larger the input data size, the more map tasks are typically used to process it efficiently.

Option B (Outputs) is incorrect because the number of map tasks is not driven by the outputs. Map tasks are responsible for processing input data, not generating outputs.

Option C (Tasks) is not the correct answer because while the number of tasks in a MapReduce job can include both map tasks and reduce tasks, the question specifically asks about the number of maps, which is driven by the input data size.

Option D (None of the mentioned) is not the correct answer either because the number of maps is indeed influenced by the size of the input data.

5. Identify the correct statement(s) in the context of YARN (Yet Another Resource Negotiator):

- A. YARN is highly scalable.
- B. YARN enhances a Hadoop compute cluster in many ways.
- C. YARN extends the power of Hadoop to incumbent and new technologies found within the data center.

Choose the correct option:

- A) Only statement A is correct.
- B) Statements A and B are correct.
- C) Statements B and C are correct.
- D) All statements A, B, and C are correct.

Solution:

- D) All statements A, B, and C are correct.

Explanation:

A. YARN is highly scalable.

This statement is correct. YARN (Yet Another Resource Negotiator) is designed to be highly scalable and is capable of managing resources efficiently in large clusters.

B. YARN enhances a Hadoop compute cluster in many ways.

This statement is correct. YARN enhances a Hadoop compute cluster by separating the resource management and job scheduling aspects from the original MapReduce framework. This separation allows for more diverse and efficient processing frameworks beyond MapReduce to coexist in the same Hadoop cluster.

C. YARN extends the power of Hadoop to incumbent and new technologies found within the data center.

This statement is correct. YARN's architecture allows it to support not only MapReduce but also other data processing frameworks like Spark, Tez, and more. This flexibility enables Hadoop clusters to accommodate a wide range of data processing technologies and workloads, making it a more versatile platform.

Therefore, all three statements A, B, and C are correct.

6. Which of the following statements accurately describe(s) the role and responsibilities of the Job Tracker in the context of Big Data computing?

- A. The Job Tracker is hosted inside the master and it receives the job execution request from the client.
- B. The Job Tracker breaks down big computations into smaller parts and allocates tasks to slave nodes.
- C. The Job Tracker stores all the intermediate results from task execution on the master node.
- D. The Job Tracker is responsible for managing the distributed file system in the cluster.

Choose the correct option:

- A) Only statement A is correct.
- B) Statements A and B are correct.
- C) Statements A, B, and C are correct.
- D) None of the statements are correct.

Solution:

B) Statements A and B are correct.

Explanation:

A. The Job Tracker is hosted inside the master and it receives the job execution request from the client.

This statement is correct. The Job Tracker is a component hosted on the master node in the Big Data computing framework. It is responsible for receiving job execution requests from clients and managing the execution of those jobs.

B. The Job Tracker breaks down big computations into smaller parts and allocates tasks to slave nodes.

This statement is correct. One of the main responsibilities of the Job Tracker is to divide large computations into smaller tasks and assign these tasks to slave nodes (Task Trackers) for execution in a distributed manner.

C. The Job Tracker stores all the intermediate results from task execution on the master node.

This statement is not correct. The Job Tracker is primarily responsible for coordinating job execution and managing task allocation, but it does not typically store intermediate results. Intermediate results are often stored on the slave nodes (Data Nodes) where the tasks are executed.

D. The Job Tracker is responsible for managing the distributed file system in the cluster.

This statement is not correct. The Job Tracker's main focus is on job execution and task allocation, rather than managing the distributed file system. The role of managing the file system is typically handled by other components like the NameNode in Hadoop's HDFS.

In summary, statements A and B accurately describe the role and responsibilities of the Job Tracker in the context of Big Data computing.

7. Consider the pseudo-code for MapReduce's WordCount example. Let's now assume that you want to determine the frequency of phrases consisting of 3 words each instead of determining the frequency of single words. Which part of the (pseudo-)code do you need to adapt?

- A) Only map()
- B) Only reduce()
- C) map() and reduce()
- D) None

Solution:

- A) Only map()

Explanation:

The map function takes a value and outputs key:value pairs.

For instance, if we define a map function that takes a string and outputs the length of the word as the key and the word itself as the value then

map(steve) would return 5:steve and

map(savannah) would return 8:savannah.

This allows us to run the map function against values in parallel.

So we have to only adapt the map() function of pseudo code.

8. How does the NameNode determine that a DataNode is active, using a mechanism known as:

- A. Heartbeats
- B. Datapulse
- C. h-signal
- D. Active-pulse

**Solution:**

A. Heartbeats

**Explanation:**

The NameNode in Hadoop's HDFS determines that a DataNode is active using a mechanism called "Heartbeats." DataNodes periodically send heartbeats to the NameNode to indicate that they are operational and reachable. This heartbeat mechanism allows the NameNode to keep track of the health and status of the DataNodes in the cluster. If the NameNode stops receiving heartbeats from a DataNode, it assumes that the DataNode is unavailable or has failed, and it starts the process of replicating the data stored on that DataNode to maintain data availability and fault tolerance.

Options B, C, and D are not accurate terms for describing the mechanism by which the NameNode determines the activity of DataNodes in HDFS. The correct term is "Heartbeats."

9. Which function processes a key/value pair to generate a set of intermediate key/value pairs?

- A. Map
- B. Reduce
- C. Both Map and Reduce
- D. None of the mentioned

**Solution:**

**A. Map**

**Explanation:**

In the context of the MapReduce programming model, the Map function processes input key/value pairs and generates a set of intermediate key/value pairs as output. The Map function applies a user-defined operation to each input record and produces intermediate key/value pairs based on the processing. These intermediate pairs are then grouped, shuffled, and sorted before being passed to the Reduce function for further processing.

The Reduce function, on the other hand, processes the intermediate key/value pairs generated by the Map function and performs operations to aggregate or combine the data based on the keys.

Option B (Reduce) is incorrect because the Reduce function processes intermediate data generated by the Map function, but it doesn't directly generate intermediate key/value pairs.

Option C (Both Map and Reduce) is incorrect because while both Map and Reduce functions are integral parts of the MapReduce process, the specific function responsible for generating intermediate key/value pairs is the Map function.

Option D (None of the mentioned) is not the correct answer as the Map function is indeed responsible for generating intermediate key/value pairs.

10. Which of the following options correctly identifies the three main components of the YARN Scheduler in Hadoop?

- A) Global Application Manager (GAM), Cluster Resource Tracker (CRT), Job Task Coordinator (JTC)
- B) Resource Monitor (RM), Cluster Supervisor (CS), Task Executor (TE)
- C) Global Resource Manager (RM), Per server Node Manager (NM), Per application (job) Application Master (AM)
- D) Central Resource Coordinator (CRC), Node Resource Manager (NRM), Application Controller (AC)

**Solution:**

- C) Global Resource Manager (RM), Per server Node Manager (NM), Per application (job) Application Master (AM)

**Explanation:**

**Global Resource Manager (RM):** Responsible for overall resource allocation and scheduling across the cluster.

**Per server Node Manager (NM):** Monitors resource usage on individual servers, reports it to the RM, and manages containers' lifecycle.

**Per application (job) Application Master (AM):** Manages the negotiation of resource containers with the RM and NM for its tasks. It also monitors task progress and handles task failures for a specific job.

Option A, B, and D are incorrect because they do not accurately identify the actual components of the YARN Scheduler in Hadoop.

### **Quiz Assignment-III Solutions: Big Data Computing (Week-3)**

---

1. What is the primary goal of Spark when working with distributed collections?

- A) Achieving real-time data processing
- B) Distributing data across multiple clusters
- C) Working with distributed collections as you would with local ones
- D) Optimizing data storage

**Solution:**

**C. Working with distributed collections as you would with local ones**

**Explanation:**

The primary goal of Spark is to allow users to work with distributed collections of data in a way that is similar to how they would work with local collections. This is achieved through the concept of Resilient Distributed Datasets (RDDs), which are immutable collections of objects spread across a cluster. RDDs can be built through parallel transformations like map, filter, etc., and they are automatically rebuilt on failure, ensuring fault tolerance. Users can also control the persistence of RDDs, such as caching them in RAM for faster access.

2. Which of the following is NOT a type of machine learning algorithm available in Spark's MLLib?

- A) Logistic regression
- B) Non-negative matrix factorization (NMF)
- C) Decision tree classification
- D) Convolutional neural network (CNN)

**Solution:**

**D) Convolutional neural network (CNN)**

**Explanation:**

**A) This is correct. MLLib includes logistic regression as a classification algorithm.**

- B) This is correct. MLlib includes non-negative matrix factorization (NMF) as a collaborative filtering algorithm.
- C) This is correct. MLlib includes decision tree classification as a classification algorithm.
- D) This is the correct answer. MLlib does not include convolutional neural networks (CNNs) in its library. CNNs are typically used for deep learning tasks and are not a part of Spark's MLlib.

3. Which of the following statements about Apache Cassandra is not accurate?

- A) Cassandra is a distributed key-value store.
- B) It was originally designed at Twitter.
- C) Cassandra is intended to run in a datacenter and across data centers.
- D) Netflix uses Cassandra to keep track of positions in videos.

Solution:

- B) It was originally designed at Twitter.

Explanation:

- A) This statement is correct. Apache Cassandra is a distributed NoSQL database that is often described as a distributed key-value store.
- B) This statement is not accurate. Cassandra was originally designed at Facebook, not Twitter.
- C) This statement is correct. Cassandra is designed to run in a datacenter and is capable of operating across multiple data centers for high availability and fault tolerance.
- D) This statement is correct. Netflix is known to use Cassandra to keep track of positions in videos, ensuring that users can resume playback where they left off.

4. Which of the following statements is true about Apache Spark?

Statement 1: Spark improves efficiency through in-memory computing primitives and general computation graphs.

Statement 2: Spark improves usability through high-level APIs in Java, Scala, Python, and also provides an interactive shell.

- A) Only Statement 1 is true.
- B) Only Statement 2 is true.
- C) Both Statement 1 and Statement 2 are true.
- D) Neither Statement 1 nor Statement 2 is true.

**Solution:**

C) Both Statement 1 and Statement 2 are true.

**Explanation:**

Statement 1: Spark improves efficiency through in-memory computing primitives and general computation graphs. This statement is true. Spark is designed to efficiently process data by utilizing in-memory computing, which reduces the need to read data from disk repeatedly. It also represents data processing workflows as general computation graphs, allowing for optimizations.

Statement 2: Spark improves usability through high-level APIs in Java, Scala, Python, and also provides an interactive shell. This statement is true. Spark provides high-level APIs in multiple programming languages (Java, Scala, Python) to make it accessible to a wide range of developers. Additionally, it offers an interactive shell (Spark Shell) for interactive data exploration and development, enhancing usability.

5. Which of the following statements is true about Resilient Distributed Datasets (RDDs) in Apache Spark?

- A) RDDs are not fault-tolerant but are mutable.
- B) RDDs are fault-tolerant and mutable.
- C) RDDs are fault-tolerant and immutable.
- D) RDDs are not fault-tolerant and not immutable.

**Solution:**

C) RDDs are fault-tolerant and immutable.

**Explanation:**

RDDs (Resilient Distributed Datasets) in Apache Spark are designed to be fault-tolerant and immutable:

Fault-tolerant: RDDs are resilient to node failures in a cluster. If a partition of an RDD is lost due to a node failure, Spark can recompute it using the lineage information, ensuring fault tolerance.

Immutable: Once an RDD is created, it cannot be modified. Any transformation applied to an RDD results in the creation of a new RDD, leaving the original RDD unchanged. This immutability property is crucial for fault tolerance and data consistency.

Options A and B are incorrect because RDDs are indeed fault-tolerant but are not mutable. Option D is incorrect because RDDs are fault-tolerant and immutable, as explained above.

6. Which of the following is not a NoSQL database?

- A) HBase
- B) Cassandra
- C) SQL Server
- D) None of the mentioned

**Solution:**

C) SQL Server

**Explanation:**

NoSQL, which stand for "not only SQL," is an alternative to traditional relational databases in which data is placed in tables and data schema is carefully designed before the database is built. NoSQL databases are especially useful for working with large sets of distributed data.

7. How does Apache Spark's performance compare to Hadoop MapReduce?

- A) Apache Spark is up to 10 times faster in memory and up to 100 times faster on disk.
- B) Apache Spark is up to 100 times faster in memory and up to 10 times faster on disk.
- C) Apache Spark is up to 10 times faster both in memory and on disk compared to Hadoop MapReduce.
- D) Apache Spark is up to 100 times faster both in memory and on disk compared to Hadoop MapReduce.

**Solution:**

The correct answer is B) Apache Spark is up to 100 times faster in memory and up to 10 times faster on disk compared to Hadoop MapReduce.

**Explanation:**

The biggest claim from Spark regarding speed is that it is able to "run programs up to 100x faster than Hadoop MapReduce in memory, or 10x faster on disk." Spark could make this claim because it does the processing in the main memory of the worker nodes and prevents the unnecessary I/O operations with the disks. The other advantage Spark offers is the ability to chain the tasks even at an application programming level without writing onto the disks at all or minimizing the number of writes to the disks.

8. \_\_\_\_\_ leverages Spark Core fast scheduling capability to perform streaming analytics.

- A) MLlib
- B) Spark Streaming
- C) GraphX
- D) RDDs

**Solution:**

B) Spark Streaming

**Explanation:**

Spark Streaming ingests data in mini-batches and performs RDD transformations on those mini-batches of data.

9. \_\_\_\_\_ is a distributed graph processing framework on top of Spark.

- A) MLlib
- B) Spark streaming
- C) GraphX
- D) All of the mentioned

Solution:

C) GraphX

Explanation:

GraphX is Apache Spark's API for graphs and graph-parallel computation. It is a distributed graph processing framework on top of Spark.

10. Which statement is incorrect in the context of Cassandra?

- A) It is a centralized key-value store.
- B) It is originally designed at Facebook.
- C) It is designed to handle large amounts of data across many commodity servers, providing high availability with no single point of failure.
- D) It uses a ring-based DHT (Distributed Hash Table) but without finger tables or routing.

Solution:

A. It is a centralized key-value store.

Explanation:

Cassandra is not a centralized key-value store; it is a distributed NoSQL database designed to handle large amounts of data across many commodity servers. It follows a decentralized architecture to provide high availability with no single point of failure.

## **Quiz Assignment-IV Solutions: Big Data Computing (Week-4)**

---

1. What is the primary role of Snitches in Apache Cassandra?

- A) Data encryption and decryption
- B) Load balancing and distribution of requests
- C) Data compression and decompression
- D) Managing database schema changes

**Solution:**

**B) Load balancing and distribution of requests**

**Explanation:**

Snitches in Apache Cassandra play a crucial role in determining the network topology and facilitating efficient routing of requests. They help Cassandra distribute replicas by grouping machines into data centers and racks. This information is essential for load balancing and ensuring that requests are routed to the appropriate nodes, which helps in achieving high availability, fault tolerance, and efficient data distribution across the cluster. While encryption, compression, and schema management are important aspects of a distributed database system like Cassandra, these functions are not the primary responsibility of Snitches.

2. Which statements about Cassandra and its Snitches are correct?

Statement 1: In Cassandra, during a write operation, when hinted handoff is enabled and if any replica is down, the coordinator writes to all other replicas and keeps the write locally until the down replica comes back up.

Statement 2: In Cassandra, Ec2Snitch is an important snitch for deployments, and it is a simple snitch for Amazon EC2 deployments where all nodes are in a single region. In Ec2Snitch, the region name refers to the data center, and the availability zone refers to the rack in a cluster.

- A) Only Statement 1 is correct.
- B) Only Statement 2 is correct.
- C) Both Statement 1 and Statement 2 are correct.

D) Neither Statement 1 nor Statement 2 is correct.

**Solution:**

C) Both Statement 1 and Statement 2 are correct.

**Explanation:**

Statement 1 is correct: In Cassandra, when hinted handoff is enabled and a replica is down during a write operation, the coordinator node writes to all other replicas that are available and keeps the write locally. It does this to ensure that the write is not lost and can be delivered to the down replica when it comes back up.

Statement 2 is correct: Ec2Snitch is indeed an important Snitch for deployments in Amazon EC2 environments. It is designed for use in Amazon EC2 deployments where all nodes are typically in a single region. In Ec2Snitch, the region name refers to the data center, and the availability zone refers to the rack within a cluster. This snitch helps Cassandra understand the network topology within Amazon EC2, facilitating efficient routing and data replication.

Therefore, both statements are correct.

3. ZooKeeper allows distributed processes to coordinate with each other through registers, known as \_\_\_\_\_.

- A) znodes
- B) hnones
- C) vnodes
- D) rnones

**Solution:**

A) znodes

**Explanation:**

ZooKeeper allows distributed processes to coordinate with each other through registers called "znodes." These znodes act as the basic building blocks in ZooKeeper, providing a distributed and hierarchical namespace where data can be stored and synchronized across a cluster of machines.

4. In Zookeeper, when a \_\_\_\_\_ is triggered the client receives a packet saying that the znode has changed.

- A) Event
- B) Row
- C) Watch
- D) Value

**Solution:**

C) Watch

**Explanation:**

ZooKeeper supports the concept of watches. Clients can set a watch on a znodes.

5. What does the CAP theorem, proposed by Eric Brewer and subsequently proved by Gilbert and Lynch, state about distributed systems?

- A) You can always achieve all three guarantees: Consistency, Availability, and Partition tolerance.
- B) In a distributed system, you can satisfy at most 3 out of the 3 guarantees.
- C) In a distributed system, you can satisfy at most 2 out of the 3 guarantees: Consistency, Availability, and Partition tolerance.
- D) The CAP theorem only applies to centralized systems, not distributed systems.

**Solution:**

C) In a distributed system, you can satisfy at most 2 out of the 3 guarantees: Consistency, Availability, and Partition tolerance.

**Explanation:**

The CAP theorem, proposed by Eric Brewer and subsequently proved by Gilbert and Lynch, states that in a distributed system, you can achieve at most two out of the three guarantees: Consistency, Availability, and Partition tolerance. This theorem highlights the trade-offs that need to be made when designing distributed systems. Depending on the system's requirements and the nature of network partitions, you may prioritize consistency and availability, consistency and partition tolerance, or availability and partition tolerance, but achieving all three simultaneously can be challenging or impossible in certain scenarios.

6. In Cassandra, what is the purpose of the "QUORUM" consistency level?

- A) It allows the client to read or write data on any server, regardless of whether it's a replica.
- B) It ensures strong consistency by requiring all replicas to respond, making it the slowest option.
- C) It provides the fastest response by caching writes and replying quickly to the client.
- D) It balances between consistency and availability by requiring a quorum of replicas across datacenters.

Solution:

- D) It balances between consistency and availability by requiring a quorum of replicas across datacenters.

Explanation:

In Cassandra, consistency levels allow clients to specify the level of consistency they require for read and write operations.

The "QUORUM" consistency level ensures a balance between consistency and availability. It requires a quorum of replicas to acknowledge the operation. A quorum is typically calculated as  $(N/2 + 1)$  replicas, where  $N$  is the total number of replicas.

This means that to achieve a successful read or write operation with "QUORUM," the client must receive acknowledgments from a majority of replicas, ensuring a level of data consistency while still allowing for reasonable availability and fault tolerance.

Options A, B, and C describe the characteristics of other consistency levels ("ANY," "ALL," and "ONE") in Cassandra, which have different trade-offs in terms of consistency, availability, and speed.

7. Which strong consistency model ensures that each operation by a client is visible instantaneously to all other clients, with real-time visibility?

- A) Sequential Consistency
- B) Linearizability
- C) Transaction ACID properties
- D) Transaction chains

**Solution:**

B) Linearizability

**Explanation:**

Linearizability (Option B) is a strong consistency model that guarantees that each operation by a client is visible instantaneously to all other clients, providing real-time visibility.

It ensures that operations appear to be executed instantaneously and in a total order, as if they were executed one after the other in real-time, even in a distributed system.

This strong consistency model is known for its strict and intuitive guarantees of visibility and consistency.

Options A, C, and D refer to other consistency models or properties, but they do not provide the same level of real-time visibility as linearizability.

8. What is the primary purpose of the gossip protocol in Cassandra?

- A) To facilitate data replication and consistency across the cluster.
- B) To synchronize the clocks of all nodes in the Cassandra cluster.
- C) To discover location and state information about other nodes in the cluster.
- D) To encrypt and secure communication between nodes in the cluster.

**Solution:**

C) To discover location and state information about other nodes in the cluster.

**Explanation:**

In Cassandra, the gossip protocol is primarily used for discovering location and state information about other nodes in the cluster.

It is a peer-to-peer communication protocol where nodes periodically exchange information about themselves and the nodes they are aware of in the cluster.

Through gossip, nodes learn about the status, health, and metadata of other nodes, helping Cassandra maintain an up-to-date and accurate view of the cluster's topology.

While data replication and consistency are important aspects of Cassandra's functionality, these are achieved through other mechanisms like partitioning and consistency levels, not directly through the gossip protocol.

9. What is the primary objective of the Network Topology Strategy in Cassandra?

- A) To determine the data types to be used in the database schema.
- B) To specify the consistency levels for read and write operations.
- C) To define data centers and the number of replicas to place within each data center.
- D) To optimize query performance by creating secondary indexes.

Solution:

- C) To define data centers and the number of replicas to place within each data center.

Explanation:

In Cassandra, the Network Topology Strategy is used to specify the organization of data centers and the number of replicas to place within each data center.

It is a crucial strategy for achieving fault tolerance, data availability, and disaster recovery by distributing data across multiple data centers and racks.

By defining the placement of replicas on distinct racks, the Network Topology Strategy helps ensure that data remains available even in the event of node or rack failures.

Options A, B, and D are not accurate descriptions of the Network Topology Strategy's primary objective. It focuses on data placement and replication, not data types, consistency levels, or query optimization.

10. How does ZooKeeper achieve high throughput values, including hundreds of thousands of operations per second for read-dominant workloads?

- A) By utilizing a lock-based approach to coordination.
- B) By employing eventual consistency for all operations.
- C) By exposing wait-free objects to clients.
- D) By using fast reads with watches and serving both from local replicas.

**Solution:**

- D) By using fast reads with watches and serving both from local replicas.

**Explanation:**

ZooKeeper achieves high throughput, including hundreds of thousands of operations per second for read-dominant workloads, through its use of fast reads with watches and serving both from local replicas.

Fast reads provide low-latency access to frequently read data, while watches allow clients to receive notifications of changes to data they are interested in, avoiding the need for frequent polling.

Serving both reads and watches from local replicas reduces the latency and load on the ZooKeeper ensemble, contributing to high throughput.

Options A, B, and C are not accurate descriptions of how ZooKeeper achieves high throughput. ZooKeeper's approach focuses on minimizing latency and optimizing read operations for distributed coordination.

1. Where are Bloom Filters generated and used in the context of HBase?
  - A) Bloom Filters are generated when an HFile is persisted and stored at the end of each HFile.
  - B) Bloom Filters are loaded into memory during HBase operations.
  - C) Bloom Filters allow checking on row and column levels within the HBase store.
  - D) Bloom Filters are useful when data is grouped and many misses are expected during reads.

**Solution:**

- A) Bloom Filters are generated when an HFile is persisted and stored at the end of each HFile.

**Explanation:**

In HBase, Bloom Filters are generated when an HFile is persisted, and they are stored at the end of each HFile. This allows HBase to use Bloom Filters to optimize read operations.

Option B is not entirely accurate because while Bloom Filters are used during HBase operations, they are generated and stored persistently in HFiles.

Option C is not a complete description of Bloom Filters; they are primarily used for filtering at the row level to determine whether data might exist for a particular row and column.

Option D mentions scenarios where Bloom Filters are useful but doesn't describe their generation and storage process.

2. What is the primary purpose of data streaming technologies?

- A) To transfer data in large, irregular chunks for batch processing.
- B) To ensure that data is transferred in a lossless manner over the internet.
- C) To process data as a continuous and steady stream.
- D) To reduce the growth of data on the internet.

**Solution:**

C) To process data as a continuous and steady stream.

**Explanation:**

Data streaming is a technique used to transfer data in a continuous and uninterrupted flow, allowing it to be processed as a steady and continuous stream.

Streaming technologies are essential for scenarios where data arrives continuously, such as real-time analytics, monitoring, and IoT applications.

Options A, B, and D do not accurately describe the primary purpose of data streaming technologies, which is to enable the processing of data as a continuous stream.

3. What is the primary purpose of column families in HBase?

- A) To group rows together for efficient storage.
- B) To define the data type of each column.
- C) To enable efficient grouping and storage of columns.
- D) To encrypt columns for enhanced security.

**Solution:**

C) To enable efficient grouping and storage of columns.

**Explanation:**

In HBase, column families are used to logically and physically group related columns together.

Columns within the same column family are stored separately from columns in other families, allowing for efficient storage and retrieval.

This grouping helps optimize storage, retrieval, and scans, as it allows HBase to read data efficiently at the column-family level.

Options A, B, and D do not accurately describe the primary purpose of column families in HBase, which is to enable efficient grouping and storage of columns.

4. Which of the following statements accurately describes HBase?

- A) HBase is a relational database management system (RDBMS) designed for structured data.
- B) HBase is a distributed Column-oriented database built on top of the Hadoop file system.
- C) HBase is a NoSQL database designed exclusively for document storage.
- D) HBase is a standalone database that does not require any distributed computing framework.

**Solution:**

B) HBase is a distributed Column-oriented database built on top of the Hadoop file system.

**Explanation:**

HBase is a distributed, scalable, and column-oriented database that is designed to run on top of the Hadoop Distributed File System (HDFS).

It is not a traditional relational database management system (RDBMS) like option A suggests.

While it is a NoSQL database, it is not designed exclusively for document storage, as option C implies.

Option D is incorrect because HBase is built for distributed computing and relies on Hadoop's infrastructure for distributed storage and processing.

5. What is a "Region" in the context of HBase?

- A) It refers to a single machine that holds an entire HBase table.
- B) It is a small chunk of data residing in one machine, part of a cluster of machines holding one HBase table.
- C) It represents the entire set of data in an HBase table.
- D) It is a backup copy of an HBase table stored on a remote server.

Solution:

B) It is a small chunk of data residing in one machine, part of a cluster of machines holding one HBase table.

Explanation:

In HBase, a "Region" is a small chunk of data that is part of a cluster of machines holding one HBase table.

HBase divides tables into regions to allow for horizontal scalability and efficient data distribution across the cluster.

Option A is not correct because a single machine typically does not hold an entire HBase table; rather, tables are divided into regions distributed across multiple machines.

Option C is not accurate; a table consists of multiple regions.

Option D is not a valid description of a region in HBase. Regions are not backup copies but rather partitions of data in the table.

6. In HBase, \_\_\_\_\_ is a combination of row, column family, column qualifier and contains a value and a timestamp.

- A) Cell
- B) Stores
- C) HMaster
- D) Region Server

Solution:

- A) Cell

Explanation:

Data is stored in HBASE tables Cells and Cell is a combination of row, column family, column qualifier and contains a value and a timestamp.

7. HBase architecture has 3 main components:

- A) Client, Column family, Region Server
- B) Cell, Rowkey, Stores
- C) HMaster, Region Server, Zookeeper
- D) HMaster, Stores, Region Server

**Solution:**

- C) HMaster, Region Server, Zookeeper

**Explanation:**

HBase architecture has 3 main components: HMaster, Region Server, Zookeeper.

- i. HMaster: The implementation of Master Server in HBase is HMaster. It is a process in which regions are assigned to region server as well as DDL (create, delete table) operations. It monitors all Region Server instances present in the cluster.
- ii. Region Server: HBase Tables are divided horizontally by row key range into Regions. Regions are the basic building elements of HBase cluster that consists of the distribution of tables and are comprised of Column families. Region Server runs on HDFS DataNode which is present in Hadoop cluster.
- iii. Zookeeper: It is like a coordinator in HBase. It provides services like maintaining configuration information, naming, providing distributed synchronization, server failure notification etc. Clients communicate with region servers via zookeeper.

8. What is the role of a Kafka broker in a Kafka cluster?

- A) A Kafka broker manages the replication of messages between topics.
- B) A Kafka broker allows consumers to fetch messages by topic, partition, and offset.
- C) A Kafka broker is responsible for maintaining metadata about the Kafka cluster.
- D) A Kafka broker is in charge of processing and transforming messages before they are consumed.

**Solution:**

- B) A Kafka broker allows consumers to fetch messages by topic, partition, and offset.

**Explanation:**

In Kafka, a Kafka broker is a server that stores and manages Kafka topics, allowing producers to write messages to topics and consumers to fetch messages from topics.

Option B is the correct choice because Kafka brokers play a central role in allowing consumers to fetch messages by specifying the topic, partition, and offset.

Kafka brokers do not manage message replication directly (Option A), maintain cluster metadata (Option C), or process/transform messages before consumption (Option D). Replication is managed by Kafka itself, metadata is managed collectively, and message processing typically occurs on the consumer side.

9. Which of the following statements accurately describes the characteristics of batch and stream processing?

Statement 1: Batch Processing provides the ability to process and analyze data at-rest (stored data).

Statement 2: Stream Processing provides the ability to ingest, process, and analyze data in-motion in real or near-real-time.

- A) Only Statement 1 is correct.
- B) Only Statement 2 is correct.
- C) Both Statement 1 and Statement 2 are correct.
- D) Neither Statement 1 nor Statement 2 is correct.

**Solution:**

- C) Both Statement 1 and Statement 2 are correct.

**Explanation:**

Statement 1 accurately describes Batch Processing, which involves processing and analyzing data that is already at rest or stored. Batch processing is typically used for tasks that do not require real-time processing, such as data warehousing, reporting, and ETL (Extract, Transform, Load) jobs.

Statement 2 accurately describes Stream Processing, which is designed for ingesting, processing, and analyzing data in-motion, often in real-time or near-real-time. Stream processing is used for applications like real-time analytics, monitoring, and processing of streaming data sources.

Both batch and stream processing have their unique use cases and characteristics, and both statements correctly capture these distinctions.

#### 10. What is Kafka Streams primarily used for?

- A) Kafka Streams is a Java library to process event streams live as they occur.
- B) Kafka Streams is a SQL database for querying event data.
- C) Kafka Streams is a message broker for pub-sub messaging.
- D) Kafka Streams is a distributed file storage system.

Solution:

- A) Kafka Streams is a Java library to process event streams live as they occur.

Explanation:

Kafka Streams is a Java library and framework for building real-time stream processing applications.

Option A accurately describes its primary purpose, which is to process event streams live as they occur, enabling developers to build applications that consume, process, and produce event data in real-time.

Kafka Streams is not a SQL database (Option B), message broker (Option C), or distributed file storage system (Option D). Its main focus is stream processing within the Kafka ecosystem.

## **Quiz Assignment-VI Solutions: Big Data Computing (Week-6)**

---

**1. What is Distributed K-Means Iterative Clustering?**

- A) A single-node clustering algorithm
- B) A clustering algorithm that uses distributed computing to improve scalability
- C) A supervised machine learning algorithm
- D) A dimensionality reduction technique

**Solution:**

**B) A clustering algorithm that uses distributed computing to improve scalability**

**Explanation:**

Distributed K-Means Iterative Clustering is a clustering algorithm that leverages distributed computing resources to improve the scalability of the traditional K-Means clustering algorithm. It allows the clustering process to be performed across multiple nodes or machines, making it suitable for handling large datasets and improving computational efficiency.

**2. What is the primary goal of the K-Means clustering algorithm?**

- A) Classification of data points
- B) Regression analysis
- C) Finding the nearest neighbor for each data point
- D) Partitioning data points into clusters based on similarity

**Solution:**

**D) Partitioning data points into clusters based on similarity**

**Explanation:**

The primary goal of the K-Means clustering algorithm is to partition a given dataset into clusters such that data points within the same cluster are more similar to each other than to

data points in other clusters. It does not involve classification, regression analysis, or finding nearest neighbors; its main purpose is unsupervised clustering based on similarity.

3. What is the main objective of using Parallel K-Means with MapReduce for Big Data Analytics?

- A) To reduce the dimensionality of the data
- B) To classify data points into predefined categories
- C) To efficiently cluster large datasets in a distributed manner
- D) To perform regression analysis on big data

Solution:

- C) To efficiently cluster large datasets in a distributed manner

Explanation:

The main objective of using Parallel K-Means with MapReduce for Big Data Analytics is to efficiently perform clustering on large datasets in a distributed and scalable manner, not to reduce dimensionality, classify data points, or perform regression analysis.

4. What is the primary goal of using Parallel K-Means with MapReduce in Big Data Analytics?

- A) To perform regression analysis
- B) To classify data points into predefined categories
- C) To efficiently handle large-scale clustering tasks
- D) To visualize data patterns

Solution:

- C) To efficiently handle large-scale clustering tasks

Explanation:

The primary goal of using Parallel K-Means with MapReduce in Big Data Analytics is to efficiently and effectively cluster large-scale datasets, not to perform regression analysis, classify data points, or visualize data patterns.

5. Which of the following tasks can be best solved using Clustering?

- A) Predicting the amount of rainfall based on various cues
- B) Training a robot to solve a maze
- C) Detecting fraudulent credit card transactions
- D) All of the mentioned

Solution:

C) Detecting fraudulent credit card transactions

Explanation:

Credit card transactions can be clustered into fraud transactions using unsupervised learning.

6. Identify the correct statement(s) in context of overfitting in decision trees:

Statement I: The idea of Pre-pruning is to stop tree induction before a fully grown tree is built, that perfectly fits the training data.

Statement II: The idea of Post-pruning is to grow a tree to its maximum size and then remove the nodes using a top-bottom approach.

- A) Only statement I is true
- B) Only statement II is true
- C) Both statements are true
- D) Both statements are false

**Solution:**

- A) Only statement I is true

**Explanation:**

With pre-pruning, the idea is to stop tree induction before a fully grown tree is built that perfectly fits the training data.

In post-pruning, the tree is grown to its maximum size, then the tree is pruned by removing nodes using a bottom up approach.

7. Identify the correct statement(s) in context of machine learning approaches:

Statement I: In supervised approaches, the target that the model is predicting is unknown or unavailable. This means that you have unlabeled data.

Statement II: In unsupervised approaches the target, which is what the model is predicting, is provided. This is referred to as having labeled data because the target is labeled for every sample that you have in your data set.

- A) Only Statement I is true
- B) Only Statement II is true
- C) Both Statements are false
- D) Both Statements are true

**Solution:**

- C) Both Statements are false

**Explanation:**

The correct statements are:

Statement I: In supervised approaches the target, which is what the model is predicting, is provided. This is referred to as having labeled data because the target is labeled for every sample that you have in your data set.

Statement II: In unsupervised approaches, the target that the model is predicting is unknown or unavailable. This means that you have unlabeled data.

8. What is the primary focus of Machine Learning?

- A) Accessing data from databases
- B) Extracting meaning from big data
- C) Learning from data
- D) Predicting future outcomes

Solution:

C) Learning from data

Explanation:

Machine Learning is a field of artificial intelligence that focuses on developing algorithms and models that enable computers to learn from data. It involves the creation of mathematical and statistical models that can identify patterns, make predictions, or make decisions without being explicitly programmed for each task. Here's a breakdown of the options and why "Learning from data" is the correct choice:

(A) Accessing data from databases: While accessing data is an important part of the ML process, it is not the primary focus. ML is more concerned with what you do with the data once you have it, such as training models and making predictions.

(B) Extracting meaning from big data: Extracting meaning from data is certainly a goal of ML, but it is a step within the broader process of learning from data. ML involves not only extracting meaning but also using that meaning to make predictions or decisions.

(D) Predicting future outcomes: Predicting future outcomes is one of the key applications of ML, but it is an outcome of the learning process. ML models learn from historical data to make predictions about future events. Therefore, predicting future outcomes is a result of learning from data.

9. Which of the following is an essential activity in the Machine Learning process?

- A) Writing code for specific tasks
- B) Designing graphical user interfaces
- C) Collecting and preprocessing data
- D) Creating beautiful data visualizations

**Solution:**

- C) Collecting and preprocessing data

**Explanation:**

Data is the fundamental building block of Machine Learning. Without high-quality, relevant data, it is impossible to train models that can make meaningful predictions or decisions.

Here's why data collection and preprocessing are crucial:

(A) Writing code for specific tasks: While writing code is important in the ML process, it is not the essential starting point. Before writing code for specific tasks, you need to have the right data to work with.

(B) Designing graphical user interfaces: Designing user interfaces may be necessary for deploying ML models in applications, but it is not a fundamental activity in the initial stages of the ML process.

(D) Creating beautiful data visualizations: Data visualization is a valuable skill for understanding data and presenting results, but it is typically a later-stage activity in the ML process, often used for communicating findings and insights.

Collecting data involves gathering relevant datasets that represent the problem you want to solve. Preprocessing data includes tasks like cleaning, transforming, and organizing the data to make it suitable for model training. These activities set the foundation for successful machine learning projects, as the quality of your data and how well it is prepared directly impact the performance and accuracy of the models you build.

10. Which distance measure calculates the distance along strictly horizontal and vertical paths, consisting of segments along the axes?

- A) Euclidean distance
- B) Manhattan distance
- C) Cosine similarity
- D) Minkowski distance

Solution:

B) Manhattan distance

Explanation:

Euclidean distance (option A) measures the straight-line distance (the shortest path) between two points in a two-dimensional or multi-dimensional space. It considers both horizontal and vertical movements as well as diagonal movements.

Manhattan distance (option B), on the other hand, calculates the distance by summing up the absolute differences in the x and y coordinates (or dimensions) between two points. It forces movement only along the axes (horizontal and vertical), resembling the grid-like layout of streets in Manhattan, hence the name.

Cosine similarity (option C) measures the cosine of the angle between two vectors in a multi-dimensional space. It doesn't measure distance in the same way as Manhattan or Euclidean distance but instead assesses the direction or orientation of vectors.

Minkowski distance (option D) is a generalized distance measure that includes both Manhattan (when the exponent is 1) and Euclidean (when the exponent is 2) distance as special cases. It allows you to control the level of emphasis on different dimensions based on the chosen exponent.

So, the correct answer is (B) Manhattan distance, which is explicitly defined to measure distances along strictly horizontal and vertical paths, making it ideal for scenarios where movement is constrained to grid-like patterns.

## **Quiz Assignment-VII Solutions: Big Data Computing (Week-7)**

---

1. What is a decision tree primarily used for in machine learning?

- A) Clustering
- B) Regression
- C) Classification
- D) Dimensionality reduction

**Solution:**

**C) Classification**

**Explanation:**

A decision tree is primarily used for classification in machine learning.

2. In the context of utilising a bagging-based algorithm like Random Forest for model development, which of the following statements is accurate?

- 1) Number of tree should be as large as possible
- 2) You will have interpretability after using Random Forest

- A) Only 1
- B) Only 2
- C) Both 1 and 2
- D) None of these

**Solution:**

**A) Only 1**

**Explanation:**

Since Random Forest aggregate the result of different weak learners, If It is possible we would want more number of trees in model building. Random Forest is a black box model you will lose interpretability after using it.

3. Which of the following statements about decision trees is accurate?
- A) Decision trees can automatically handle interactions of features and build complex functions involving multiple splitting criteria.
  - B) Decision trees are highly scalable for very large datasets with many features.
  - C) A single decision tree is typically a very good predictor.
  - D) Decision trees are not interpretable, and it is not possible to visualize the decision tree or analyze its criteria.

**Solution:**

- A. Decision trees can automatically handle interactions of features and build complex functions involving multiple splitting criteria.

**Explanation:**

A. Decision trees can automatically handle interactions of features and build complex functions involving multiple splitting criteria: This statement is generally true. Decision trees are capable of capturing interactions between features by recursively splitting the data based on different criteria at each node. They can create complex decision boundaries by combining multiple splitting criteria in a hierarchical manner. This is one of the strengths of decision trees in modeling complex relationships in the data.

B. Decision trees are highly scalable for very large datasets with many features: This statement is not entirely accurate. While decision trees are relatively fast to build, they may not be the most scalable option for very large datasets with many features. Random Forest and other ensemble methods are often preferred in such cases to improve predictive accuracy.

C. A single decision tree is typically a very good predictor: This statement is generally not true. While decision trees can model relationships in the data, a single decision tree is prone to overfitting, which can lead to poor predictive performance on new, unseen data. Ensemble methods like Random Forest, which combine multiple decision trees, are often used to enhance predictive power.

D. Decision trees are not interpretable, and it is not possible to visualize the decision tree or analyze its criteria: This statement is not accurate. Decision trees are known for their interpretability. They can be visualized to understand the splitting criteria at each node, the values in the leaves, and the decision-making process. This interpretability is one of the advantages of using decision trees in machine learning.

So, option A is the correct statement about decision trees.

4. Which statement accurately describes the role of bootstrapping in the context of the random forest algorithm?
- A) Bootstrapping generates replicas of the dataset without replacement, ensuring diversity in the random forest.
  - B) Bootstrapping is not used in the random forest algorithm, it is only employed in decision tree construction.
  - C) Bootstrapping produces replicas of the dataset by random sampling with replacement, which is essential for the random forest algorithm.
  - D) Bootstrapping creates additional features to augment the dataset for improved random forest performance.
- Solution:**
- C) Bootstrapping produces replicas of the dataset by random sampling with replacement, which is essential for the random forest algorithm.
- Explanation:**
- A. Bootstrapping generates replicas of the dataset without replacement, ensuring diversity in the random forest: This statement is not accurate. Bootstrapping involves random sampling with replacement, which means that some data points in the original dataset may appear multiple times in each replica. This process creates diversity in the data subsets used to train individual trees in a random forest.
  - B. Bootstrapping is not used in the random forest algorithm; it's only employed in decision tree construction: This statement is incorrect. Bootstrapping is a fundamental technique used in the random forest algorithm. It is used to create multiple random subsets of the dataset, each of which is used to train a separate decision tree in the forest.
  - C. Bootstrapping produces replicas of the dataset by random sampling with replacement, which is essential for the random forest algorithm: This statement is correct. Bootstrapping is a process in which random samples are drawn from the original dataset with replacement. These samples are used to create diverse training datasets for individual decision trees in the random forest. The combination of bootstrapping and ensemble averaging makes random forests robust and effective.
  - D. Bootstrapping creates additional features to augment the dataset for improved random forest performance: This statement is not accurate. Bootstrapping does not create additional features but rather generates multiple subsets of the original dataset, each with some degree of overlap with the others.

So, option C is the correct description of the role of bootstrapping in the random forest algorithm.

5. Which of the following statements about bagging is accurate in the context of understanding the random forest algorithm?

- A) Bagging is primarily used to average predictions of decision trees in the random forest algorithm.
- B) Bagging is a technique exclusively designed for reducing the bias in predictions made by decision trees.
- C) Bagging, short for Bootstrap Aggregation, is a general method for averaging predictions of various algorithms, not limited to decision trees, and it works by reducing the variance of predictions.
- D) Bagging is a method specifically tailored for improving the interpretability of decision trees in the random forest algorithm.

Solution:

C) Bagging, short for Bootstrap Aggregation, is a general method for averaging predictions of various algorithms, not limited to decision trees, and it works by reducing the variance of predictions.

Explanation:

A. Bagging is primarily used to average predictions of decision trees in the random forest algorithm: This statement is not accurate. While bagging is used in the random forest algorithm, it is not limited to decision trees. Bagging can be applied to a variety of algorithms to improve prediction accuracy by reducing variance.

B. Bagging is a technique exclusively designed for reducing the bias in predictions made by decision trees: This statement is not accurate. Bagging is primarily used to reduce the variance in predictions, not bias. It aims to improve the stability and generalization of predictive models, including decision trees.

C. Bagging, short for Bootstrap Aggregation, is a general method for averaging predictions of various algorithms, not limited to decision trees, and it works by reducing the variance of predictions: This statement is correct. Bagging involves creating multiple subsets of the dataset through bootstrapping (random sampling with replacement). These subsets are

used to train different models (which can be any algorithm), and their predictions are aggregated to reduce variance and enhance overall model performance.

D. Bagging is a method specifically tailored for improving the interpretability of decision trees in the random forest algorithm: This statement is not accurate. While bagging can be used in conjunction with decision trees in the random forest, its primary purpose is to enhance predictive accuracy and model robustness, not interpretability.

So, option C accurately describes bagging in the context of understanding the random forest algorithm.

6. In which of the following scenario a gain ratio is preferred over Information Gain?

- A) When a categorical variable has very small number of category
- B) Number of categories is the not the reason
- C) When a categorical variable has very large number of category
- D) None of the mentioned

Solution:

- C) When a categorical variable has very large number of category

Explanation:

When high cardinality problems, gain ratio is preferred over Information Gain technique.

7. Which of the following is/are true about Random Forest and Gradient Boosting ensemble methods?

- 1) Both methods can be used for classification task
- 2) Random Forest is use for classification whereas Gradient Boosting is use for regression task
- 3) Random Forest is use for regression whereas Gradient Boosting is use for Classification task

4) Both methods can be used for regression task

- A) 1 and 2
- B) 2 and 3
- C) 2 and 4
- D) 1 and 4

**Solution:**

- D) 1 and 4

**Explanation:**

Both algorithms are design for classification as well as regression task.

8. What is the primary benefit of bagging in machine learning?

- A) Bagging increases the complexity of individual models, leading to more accurate predictions.
- B) Bagging provides an averaging over a set of possible datasets, removing noisy and non-stable parts of models.
- C) Bagging generates diverse features to improve model robustness.
- D) Bagging ensures that models do not overfit the training data.

**Solution:**

- B) Bagging provides an averaging over a set of possible datasets, removing noisy and non-stable parts of models.

**Explanation:**

Bagging (Bootstrap Aggregation) is a technique in machine learning that involves creating multiple subsets of the original dataset through bootstrapping (random sampling with replacement). These subsets are used to train different models, and their predictions are combined or averaged to produce the final prediction. The primary benefit of bagging is that it reduces variance and increases model stability by averaging over these different datasets.

9. Hundreds of trees can be aggregated to form a Random forest model. Which of the following is true about any individual tree in Random Forest?

- 1) Individual tree is built on a subset of the features
  - 2) Individual tree is built on all the features
  - 3) Individual tree is built on a subset of observations
  - 4) Individual tree is built on full set of observations
- A) 1 and 3
- B) 1 and 4
- C) 2 and 3
- D) 2 and 4

**Solution:**

A) 1 and 3

**Explanation:**

Random forest is based on bagging concept, that consider fraction of sample and fraction of feature for building the individual trees.

10. Boosting any algorithm takes into consideration the weak learners. Which of the following is the main reason behind using weak learners?

Reason I: To prevent overfitting

Reason II: To prevent underfitting

- A) Reason I
- B) Reason II
- C) Both Reason I and Reason II
- D) None of the Reasons

**Solution:**

A) Reason I

**Explanation:**

To prevent overfitting, since the complexity of the overall learner increases at each step.  
Starting with weak learners implies the final classifier will be less likely to overfit.

---

1. Which of the following functionalities are provided by the Spark API for graph parallel computations?

- A) joinVertices
- B) subgraph
- C) aggregateMessages
- D) All of the above

**Solution:**

D) All of the above

**Explanation:**

**joinVertices:** This function is provided by the Spark API for graph parallel computations. It allows you to join the vertices of the graph with data from another DataFrame or RDD, updating the vertex attributes accordingly. This is useful for enriching the graph's vertex properties based on external data.

**subgraph:** The Spark API also provides the subgraph function, which allows you to create a subgraph of the original graph by specifying a subset of vertices and edges. This is useful when you want to work with a smaller portion of the graph for specific computations.

**aggregateMessages:** This is another functionality provided by the Spark API for graph parallel computations. It allows you to perform message aggregation across the edges of the graph and apply user-defined functions to update vertex attributes. It's a fundamental operation for graph algorithms such as PageRank and connected components.

So, all three options (i. joinVertices, ii. subgraph, iii. aggregateMessages) are provided by the Spark API for graph parallel computations, making option D the correct answer.

2. Which of the following statement(s) is/are true in the context of Apache Spark GraphX operators?

- A) Property operators modify the vertex or edge properties using a user-defined map function and produce a new graph.

B) Structural operators operate on the structure of an input graph and produce a new graph.

C) Join operators add data to graphs and produce new graphs.

D) All of the above

**Solution:**

D) All of the above

**Explanation:**

- A) Property operators in Apache Spark GraphX are used to modify the vertex or edge properties using a user-defined map function, and they produce a new graph with the updated properties. This allows you to transform the properties of vertices and edges in the graph.
- B) Structural operators work on the structure of the input graph, such as subgraph extraction, and they produce a new graph that represents the modified structure. These operators help you analyze and manipulate the structure of the graph.
- C) Join operators in GraphX are used to add data to graphs, typically by joining external datasets with vertices or edges. They produce new graphs with the additional data incorporated, allowing you to enrich the graph with external information.

So, all three statements (A, B, and C) are true, making option D the correct answer.

3. Why are graph-structured data increasingly common in data science contexts?

A) Because they are easy to model and analyze.

B) Because they are rare and unique in real-world applications.

C) Because they are commonly used to model communication between entities in various contexts.

D) Because they are limited to social networks only.

**Solution:**

C) Because they are commonly used to model communication between entities in various contexts.

**Explanation:**

Graph-structured data, represented as networks of nodes and edges, are becoming increasingly common in data science for several reasons:

**Communication Between Entities:** Graphs are an effective way to model and analyze communication and interactions between entities, such as people in social networks, computers in Internet communication, cities and countries in transportation networks, and corporations in financial transactions.

**Real-World Applications:** Graphs are not limited to a single domain. They are versatile and applicable to a wide range of real-world scenarios where entities and their relationships need to be understood and analyzed.

**Ubiquity:** Graphs are ubiquitous in the sense that they naturally represent relationships in various domains, making them a valuable tool for data scientists to gain insights and make predictions.

Option C accurately captures the reason for the increasing popularity of graph-structured data in data science contexts.

4. What is the purpose of indexing and bitmaps in the context of graph processing?

- A) To create colorful visual representations of graphs.
- B) To accelerate joins across graphs.
- C) To alphabetically organize graph data.
- D) To improve the aesthetics of graph presentations.

**Solution:**

- B) To accelerate joins across graphs.

**Explanation:**

Indexing and bitmaps play a crucial role in optimizing graph processing for various purposes:

- 1) **Accelerating Joins Across Graphs:** Indexing and bitmaps are used to speed up the process of joining data from different graphs or subgraphs. By creating efficient data structures and indices, it becomes quicker and more efficient to identify and combine related data elements from multiple sources, which is especially useful for complex analytical tasks.

- 2) Efficiently Constructing Subgraphs: Bitmaps, in particular, can be used to efficiently construct subgraphs by representing which vertices or edges belong to a specific subset. This allows for the rapid extraction of relevant portions of a graph without the need to traverse the entire graph, improving performance when dealing with large-scale graphs.

Options A, C, and D do not accurately describe the primary purpose of indexing and bitmaps in graph processing, making option B the correct answer.

5. What is the significance of substantial index and data reuse in graph processing?

- A) To create decorative elements within graphs.
- B) To increase the file size of graphs for better storage.
- C) To save memory and processing resources by reusing routing tables and edge adjacency information.
- D) To add redundancy to graphs for fault tolerance.

Solution:

- C) To save memory and processing resources by reusing routing tables and edge adjacency information.

Explanation:

Substantial index and data reuse in graph processing is essential for optimizing memory usage and processing efficiency. Here's why it's significant:

**Reuse Routing Tables Across Graphs and Subgraphs:** By reusing routing tables, you can avoid the unnecessary duplication of data structures, which saves memory and reduces the overhead of maintaining multiple copies of the same information. This is crucial when working with large graphs or subgraphs, as it helps conserve resources.

**Reuse Edge Adjacency Information and Indices:** Reusing edge adjacency information and indices means that you don't need to recreate these structures every time you perform graph operations. This leads to faster graph traversals and computations because you can build on existing data rather than starting from scratch.

Options A, B, and D do not accurately describe the significance of substantial index and data reuse in graph processing. Option C is the correct answer, as it highlights the benefits of saving memory and processing resources through reuse.

6. What is the purpose of the `outerJoinVertices()` operator in Apache Spark GraphX?
- A) It removes all vertices that are not present in the input RDD.
  - B) It returns a new graph with only the vertices from the input RDD.
  - C) It joins the input RDD data with vertices and includes all vertices, whether present in the input RDD or not.
  - D) It creates a subgraph from the input RDD and vertices.

**Solution:**

C) It joins the input RDD data with vertices and includes all vertices, whether present in the input RDD or not.

**Explanation:**

The `outerJoinVertices()` operator in Apache Spark GraphX serves the purpose of joining the input RDD data with the graph's vertices. However, it differs from a standard join in that it includes all vertices, whether they have matching data in the input RDD or not. The key points to note are:

It performs a join operation between the existing graph's vertices and the data in the input RDD.

It applies a user-defined `map()` function to all vertices, allowing you to update their properties.

Importantly, it includes vertices that may not have corresponding data in the input RDD, ensuring that all vertices are retained in the resulting graph.

Options A, B, and D do not accurately describe the purpose or behavior of the `outerJoinVertices()` operator, making option C the correct answer.

7. Which of the following statements regarding Apache Spark and its components are true?

S1: Apache Spark GraphX provides the following property operators - `mapVertices()`, `mapEdges()`, `mapTriplets()`.

S2: The RDDs in Spark depend on one or more other RDDs. The representation of dependencies between RDDs is known as the lineage graph. Lineage graph information is used to compute each RDD on demand, so that whenever a part of a persistent RDD is lost, the data that is lost can be recovered using the lineage graph information.

- A) Only S1 is true.
- B) Only S2 is true.
- C) Both S1 and S2 are true.
- D) Neither S1 nor S2 are true.

**Solution:**

C) Both S1 and S2 are true.

**Explanation:**

S1: Apache Spark GraphX does indeed provide property operators such as `mapVertices()`, `mapEdges()`, and `mapTriplets()`. These operators allow you to apply user-defined functions to vertices, edges, and combinations of both (triplets) to modify their properties. This statement is correct.

S2: In Spark, RDDs (Resilient Distributed Datasets) depend on one or more other RDDs, and the representation of dependencies between RDDs is known as the lineage graph. The lineage graph information is used to compute each RDD on demand. If a portion of a persistent RDD is lost due to a node failure, Spark can use the lineage graph information to recompute the lost data, ensuring fault tolerance and data recovery. This statement accurately describes the concept of lineage in Spark.

Therefore, both statements S1 and S2 are true, making option C the correct answer.

8. What does GraphX provide an API for expressing in terms of graph computation?

- A) Pregel algorithms.
- B) Genetic algorithms.
- C) Natural language processing.
- D) Image recognition.

**Solution:**

A) Pregel algorithms.

**Explanation:**

GraphX, a component of Apache Spark, provides an API for expressing graph computation that can model the "Pregel" abstraction. The Pregel abstraction, based on Google's Pregel paper, is a programming model for expressing distributed, iterative graph algorithms. It simplifies the development of large-scale graph algorithms by allowing developers to express algorithms as a series of iterations, where messages are sent along edges and aggregated at vertices. This approach is particularly well-suited for solving complex graph problems efficiently in distributed computing environments. Therefore, option A is the correct answer.

9. What are the characteristics of a Parameter Server in the context of distributed machine learning?

S1: Distributes a model over multiple machines.

S2: It offers two operations: (i) Pull for query parts of the model (ii) Push for update parts of the model.

A) Only S1 is true.

B) Only S2 is true.

C) Both S1 and S2 are true.

D) Neither S1 nor S2 are true.

**Solution:**

C) Both S1 and S2 are true.

**Explanation:**

S1: A Parameter Server is designed to distribute a machine learning model over multiple machines. This distribution allows for parallelism and distributed training, making it possible to work with large-scale machine learning models on clusters of machines.

S2: Parameter Servers typically offer two fundamental operations: "Pull" and "Push." These operations are essential for distributed machine learning. "Pull" is used to query parts of the model for inference or use, and "Push" is used to update parts of the model during training.

These operations facilitate the coordination and synchronization of model parameters across multiple machines in a distributed environment.

Both statements accurately describe the characteristics and functionality of a Parameter Server, making option C the correct answer.

10. What has driven the development of specialized graph computation engines capable of inferring complex recursive properties of graph structured data?

- A) Increasing demand for social media analytics
- B) Growing scale and importance of graph data
- C) Advances in machine learning algorithms
- D) Expansion of blockchain technology

Solution:

- B) Growing scale and importance of graph data

Explanation:

The correct answer is B) Growing scale and importance of graph data. The statement in the question highlights that the development of specialized graph computation engines has been driven by the increasing size and significance of graph data. Graph data structures are used to represent complex relationships and connections, such as social networks, recommendation systems, and knowledge graphs. As these datasets have grown in scale and importance, there has been a need for specialized engines capable of efficiently inferring complex recursive properties from them. These engines are essential for tasks like graph analytics, network analysis, and various applications in fields like social media, biology, and finance.

---