**Question # 1**

✏ Revisit

Why are big data applications liable to latency?

Choose the best option

○ Big data cannot use in-memory computing

○ Big data applications are still in the early stages of development.

○ The volume of big data is too large to be analyzed rapidly

○ Big data may reside in a different location from the application

**Question # 5**

🖊 Revisit

Choose the best option

Data Locality means:

○ Moving data to computation

○ Moving computation to data

○ Moving data to computation and Moving comp

○ None of the above

**Question # 4**                              ⏱ Revisit

Which of the following is fundamental datastructure of Spark?

**Choose the best option**

○ RDD

○ Dataframe

○ Dataset

○ None of the above

uestion # 3

⟳ Revisit

pically a Hbase Region server is collocated with

**Choose the best option**

○ HDFS namenode

○ HDFS datanode

○ As a client to HDFS server

○ Resource manager

**uestion # 2**

⟳ Revisit

mmand is used to know the status of all the daemons of hadoop.

**Choose the best option**

○ fsck

○ distcp

○ jps

○ None of the above

PG-DBDA_0921_230322 ⓘ

## Question # 8

Revisit

Which of the following is module for Structured data processing?

**Choose the best option**

○ GraphX

○ MLib

○ Spark SQL

○ Spark R

**Question # 7**                                    ✏ Revisit

What message is generated by a datanode to indicate its connectivity with namenode?

**Choose the best option**

○ Beap

○ Heartbeat

○ Analog Pulse

○ Map

**Question # 6**

⟳ Revisit

**Choose the best option**

_____ is not a component of spark.

- ○ SparkR
- ○ Mllib
- ○ GraphX
- ○ Sqoop

uestion # 17                                    ⟳ Revisit

ick the correct statement about hadoop when compared to RDBMS?

**Choose the best option**

○ Hadoop is suitable for read and write many times

○ Hadoop does ACID transactions

○ Hadoop works better on unstructured and semi-structured

○ Hadoop has higher data Integrity

**Question # 15**                                    ✏ Revisit          **Choose the best option**

Why do we use the SSH in Hadoop Cluster?

○  To perform the Passwordless authentication

○  To establish the communication between Master

○  Both of the above

○  None of The above

**Question # 21**

⟳ Revisit

Which of the following services is provided by YARN?

**Choose the best option**

○ Global resource management

○ Record reader

○ MapReduce engine

○ Data Mining

**Question # 19**

⟳ Revisit

Hive provides a SQL like language called

**Choose the best option**

○ SQL Hive

○ Hive QL

○ DB QL

○ Hive Data

**Question # 23**                                      ⟳ Revisit

Pick the false statement about hadoop:

**Choose the best option**

○ The main algorithm used in it is Map Reduce

○ It is a distributed framework

○ It runs with commodity hard ware

○ All are true

**Question # 29**

Which of the following is not true?

✏ Revisit

**Choose the best option**

○ In Pseudo distributed mode no daemons will be runni...

○ In fully distributed mode all the daemons will be runni...

○ In standalone no daemons will be running

○ In Pseudo distributed mode no daemons will be runni... daemons will be running on same machine

PG-DBDA_0921_230322 ⓘ

**Question # 27**                                    ⟳ Revisit

The main advantage of creating table partition is

**Choose the best option**

○ Effective storage memory utilization

○ faster query performance

○ Less RAM required by namenode

○ simpler query syntax

**Question # 28**

🖉 Revisit

**Choose the best option**

What can best be described as a programming model used to develop Hadoop-based applications that can process massive amounts of data?

○ Oozie

○ Mahout

○ MapReduce engine

○ All of the mentioned

**Question # 29**

Revisit

**Choose the best option**

Which of the following is not true?

- ○ In Pseudo distributed mode no daemons will be running

- ○ In fully distributed mode all the daemons will be running on same machine

- ○ In standalone no daemons will be running

- ○ In Pseudo distributed mode no daemons will be running and In fully distributed daemons will be running on same machine

**Question # 33**                              ⟳ Revisit          **Choose the best option**

Which property gets configured in mapred-site.xml?

○  Host and port where MapReduce job runs.

○  Java Environment variables.

○  Directory names to store hdfs files

○  Replication Factor

PG-DBDA_0921_230322 ⓘ

## Question # 35

Revisit

On dropping a managed/internal table

Choose the best option

○ The schema gets dropped without dropping t

○ The data gets dropped without dropping the s

○ An error is thrown

○ Both the schema and the data is dropped

## Question # 38

⟳ Revisit

**Choose the best option**

To retrieve all rows of a table in HBase, we use:

○ get

○ scan

○ put

○ select

**Question # 37**　　　　　　　　　　　　　　　✏ Revisit

What role does the map function play in a word count query?

**Choose the best option**

○ It sorts the words alphabetically and returns a list of the most freque

○ It creates a list with each word as a key and the number of occurre

○ It creates a list with each word as a key and every occurrence as v

○ It returns a list with each document as a key and the number of wo

PG-DBDA_0921_230322 ⓘ

Total 00:56:33
Section 00:56:33

Fini

Section 2 of 2

31   32   33   34   35   36   37   38   39   40   <   40 of 40   >   All   39

**Question # 40**

✎ Revisit

We have a 500 MB file and HDFS block size of 64 MB.

We are using hadoop fs -put command to write this file. Just after this command has finished writing 400 MB of this file, what would another user see when trying to access this file?

**Choose the best option**

○ They would see no content until the whole file written and closed.

○ They would see the current state of the file through the last completed block

○ They would see the current state of the file, up to the last bit written by the comman

○ They would see Hadoop throw a ConcurrentFileAccessException when they try to a this file.

**Question # 39**

⟲ Revisit

Which of the following deals with small files issue?

Choose the best option

○ Sequence files

○ Hadoop Archives

○ HBase

○ All of the above

## Question # 34

🖊 Revisit

**Choose the best option**

Apache HBase was modeled after Google's _____

○ Foundation DB

○ Big top

○ Big Tables

○ None of the above

**Question # 35**

🖉 Revisit

On dropping a managed/internal table

**Choose the best option**

○ The schema gets dropped without dropping the data

○ The data gets dropped without dropping the schema

○ An error is thrown

● Both the schema and the data is dropped

Clear Response

**Question # 36**                                      ⟳ Revisit          **Choose the best option**

If a Big data analyst were to analyze data from a database of call logs provided by a telecom          ○ Volume
service provider, which element of big data would be dealing with?
                                                                                                 ○ Variety

                                                                                                 ○ Velocity

                                                                                                 ○ Variable

**Question # 31**

Revisit

UDF stands for:

**Choose the best option**

- ◯ Universal Defined Function
- ◯ Unique Defined Function
- ◯ Universal Disk Format
- ◯ Unique Definition of Function

**Question # 32**

⟳ Revisit

Can you provide multiple input paths to map reduce jobs?

Choose the best option

○ Yes developers can add any number of input p

○ Yes, but limit is currently 10 input paths

○ No, hadoop works only on one input directory

○ None of the above

**Question # 26**                                    ⟳ Revisit        **Choose the best option**

The CREATE statement in Hive is related to:                          ○ DDL statements

                                                                     ○ DML statements

                                                                     ○ Session control statements

                                                                     ○ Embedded SQL statements

**Question # 22**                                    ⟳ Revisit

One of your HBase Region server (in a well configured in a sized HBase and Hadoop cluster) is reporting bad performance (slow response). What can be the possible reason?

**Choose the best option**

○ small rows and column names

○ uneven key space distribution

○ small column family name

○ None of the above

**Question # 25**

Hadoop framework is written in:

⟳ Revisit

**Choose the best option**

○ Java

○ Python

○ Scala

○ C++

**Question # 24**                              ⟳ Revisit                    **Choose the best option**

The parameter fs.default.name is set in _____ configuration file?

○  hadoop-env.sh

○  mapred-site.xml

○  core-site.xml

○  hdfs-site.xml

Question # 18                                          ⟳ Revisit

Which of the following statements are correct?

**Choose the best option**

○  Spark can run on the top of Hadoop

○  Spark can process data stored in HDFS

○  Spark can use Yarn as resource management layer

○  All of the above

**Question # 20**

⟳ Revisit

**Choose the best option**

In YARN, a container can run:

○ any application or job developed to run on YARN

○ only map and reduce tasks

○ only java applications

○ None of the above

PG-DBDA_0921_230322 🛈

Total 00:58:34
Section 00:58:34
Finish

Section 2 of 2

6   7   8   9   10   11   12   13   14   15   <   14 of 40   >   All   13   27

**Question # 14**

🔄 Revisit

What is default input format?

**Choose the best option**

○ The default input format is xml. Developers can specify other input formats as appro the xml is not correct input.

○ There is no default input format. The input format should always be specified.

○ The default input format is sequential input format. The data needs to be preproces before using the default input format.

○ The default input format is Text input format with byte offset as key and entire line

## Question # 16

⟳ Revisit

Which of the following are common feature of RDD and DataFrame?

**Choose the best option**

○ immutability

○ in-memory

○ resilient

○ All of the above

**Question # 11**

⟳ Revisit

Which of the following deals with feature "structured but not relational"?

Choose the best option

○ Structured data

○ Unstructured Data

○ Semi-structured data

○ None of the above

**Question # 10**

Revisit

**Choose the best option**

Which of the following characteristic does not belong to big data?

○ Volume

○ Variety

○ Velocity

○ Variable

**Question # 13**    ⟳ Revisit    **Choose the best option**

Hive uses _____ to store metadata:

○ Derby database

○ HiveQL

○ NoSQL

○ SQL

**Question # 12**                                    ↻ Revisit          **Choose the best option**

HBase is defined as _____                                          ○  Row oriented

                                                                     ○  Column oriented

                                                                     ○  Tuple oriented

                                                                     ○  None of the above

**Question # 9**    ✐ Revisit    **Choose the best option**

Which of the following is true regarding apache airflow?

○  open source

○  workflow management platform

○  data transformation pipeline ETL (Extract, Transfo

○  All of the above

रसी डैक
CDAC

Section 2 of 2

PG-DBDA_0921_230322 ⓘ

1   2   3   4   5   6   7   8   9   10   «   1 of 40   »   ⊕ All

Tot
Sect

**Question # 1**

⟳ Revisit

Why are big data applications liable to latency?

**Choose the best option**

○ Big data cannot use in-memory computing

○ Big data applications are still in the early stages of developm

○ The volume of big data is too large to be analyzed rapidly

○ Big data may reside in a different location from the applicatio