

Spark Interview Questions | 100 Toughest q&a 2021

The spark Interview Questions

100 Toughest spark Interview Questions 2021

Make yourself ready for your next interview with 100 Spark Questions and Answers for Job Interview

Spark Interview Questions (1-25)

1. Tell us something about Shark.

Answer: Shark is a beautiful program to work with many data users to understand just SQL for database management and are deficient in other programming languages. A shark is a tool that's been developed specifically for such men and women. This tool helps such database users quickly access Scala MLib capabilities through Hive like the SQL interface. A Shark is a tool that allows data users to operate Hive on Spark, all the while offering compatibility with Hive megastore, inquiries, and information.

2. Mention some instances in which Spark performed better than Hadoop in processing.

Response: There are several instances where Spark outperformed Hadoop:

- Sensor Data Processing –The exceptional feature of Apache Spark's In-memory computing functions best in this type of condition, as data must be recovered and

needs combination from different sources.

- For real-time querying of data, usually, Spark is favored over Hadoop.
- Stream Processing – Apache Spark is the optimal solution for transaction processes such as processing logs and discovering tips in live streams for alerts.

3. Do you know anything about Sparse Vector?

Response: A sparse vector has two parallel arrays –

- One for indices
- One for worth

We use these vectors for storing non-zero entrances to save space.

4. What do you know about RDD?

Answer: RDD is the abbreviation for Resilient Distributed Datasets. They're abstractions in Apache Spark representing the practice of information coming into the system in object format. The consumers use all RDDs for in-memory computations on large clusters, usually in a fault-tolerant way. All these databases have been read-only portioned, a set of records and have two categories–

- Immutable – It is not possible to alter RDD.
- Resilient – If a situation occurs when a node holding the partition fails, another node automatically chooses the information.

5. Do you know anything regarding transformations and activities connected to RDD?

Answer: Transformations are essential functions and executed on-demand to generate a brand new RDD. All these transformations are followed closely by user-defined actions. These transformations may include a filter, map, and reduceByKey.

Actions will be the outcomes of all sorts of RDD transformations and computations. After the consumer has done any steps, the data out of RDD returns to the local machine. Reduce, accumulate, and take are some examples of Actions.

6.Can you list the languages supported by Apache Spark for developing any huge data applications.

The Answer: The languages supported by Apache Spark for developing any big information applications are

1. Java
2. Python
3. Scala
4. R
5. Clojure

7.Are there some options for a user to use Spark to access and explore any external information saved in Cassandra databases?

Answer: Yes, it is possible to utilize Spark Cassandra Connector to test and access external information saved in Cassandra databases.

8.Is it possible for a user to use Apache Spark on Apache Mesos?

Response: We can implement Apache Spark on the hardware clusters handled by Apache Mesos. Additionally, this is among the attributes which make the Apache Spark quite popular.

9.Mention something about all the various cluster managers out there in Apache Spark.

Response: The three different clusters managers affirmed in Apache Spark are:

- Standalone deployments – These are ideal for new deployments that can only run and are extremely simple

to establish.

- Apache Mesos – This has rich resource scheduling capabilities designs, suitable for separate Spark and different applications. It's especially valuable when numerous users operate interactive shells, majorly since it climbs down the CPU allocation between orders.

10. Can Spark be connected to Apache Mesos with an individual?

Yes, an individual can connect Spark to Apache Mesos. To Be Able to connect Spark with Mesos, the user needs to follow the specified steps-

- The spark driver program has to be configured to connect to the Mesos. Any Spark binary package should be in a unique place accessible by Mesos.
- That is an alternative means to achieve the same. The user should set up Apache Spark at precisely the same place similar to the Apache Mesos and then configure the property `spark.Mesos.executor.home'` to point to the location where the installation will take place.

11. Can the consumer minimize data transfers when working with Spark? If so, how?

Response: Yes, any user gets the option to lessen the data transfers while working with Spark. Minimizing data transfers and escaping shuffling aids the user write Spark programs, which could be implementation needs quickly and trustworthy. The different ways to reduce data transfers if working with Apache Spark are:

Utilizing Broadcast Variable- Broadcast variables are designed to boost joins' efficacy between small and massive RDDs.

The most common means to minimize data transfers would be to steer precise operations ByKey, repartition, or any other similar function that triggers shuffles.

Utilizing Accumulators – Accumulators help the user upgrade the values of variables in parallel while simultaneously implementing it.



spark interview questions

12.Does an individual need broadcast variables while working with Apache Spark? If so, why?

Response: Broadcast factors are multiple-choice variables and are present within the memory card on each machine. When an individual is operating with Spark, he/she should use

broadcast factors to get rid of the necessity to send duplicates of a variable for every undertaking to process data faster. Broadcast variables also help store a lookup table within the memory to boost the recovery efficiency compared to an RDD lookup ().

13.Can a user implement Spark and Mesos, according to Hadoop?

Answer: Yes, it's possible to run Spark and Mesos using Hadoop by establishing all those respective services as a separate service on the device. The Apache Mesos functions as a unified schedule that assigns tasks to either Spark or Hadoop.

14.Do you understand anything about the lineage chart?

Answer: All the RDDs Offered in Spark solely depend on more than 1 RDD. The rendering of such dependencies between RDDs called the Lineage graph. The data offered by a Lineage chart calculate each RDD on demand to ensure that whenever a part of a continuous RDD is lost, the data lost can be retrieved with no fuss using the lineage graph info.

15.Can a user activates automatic clean-ups in Spark for handling accumulated metadata?

Answer: Yes, any user can trigger automatic clean-ups by setting the parameter ' Spark .cleaner.TTL'. Instead, the user can attain the same by dividing the long-running jobs into various batches and writing each of the intermediary results to the disk.

16.What do you know about the essential libraries which constitute the Spark Ecosystem?

Answer: Listed below are the essential libraries that make up a bulk of those Spark Ecosystem:

- 1.Spark Streaming – This library process real-time

streaming information.

- 2. Spark MLib- It is the Machine learning library at Spark and for studying algorithms like clustering.
- 3. Spark SQL – This library can help to execute SQL like queries on Spark data with regular visualization or BI tools.
- 4. Spark GraphX is the Spark API for parallel graph computations and basic operators such as joinVertices, subgraph, aggregate messages, etc.

17. Mention the advantages of using Spark following Apache Mesos.

Response: Spark, when combined with Apache Mesos, renders scalable partitioning of various Spark cases and dynamic partitioning between Spark and some other extensive information framework.

18. Mention the significance of Sliding Window operation.

Answer: The purpose called Sliding Window controls the transmission of data packets between a variety of computer networks. Spark Streaming library offers lots of windowed computations in which the transformations on RDDs facilitate a sliding window of information. Whenever the window slides, all of the RDDs that fall inside the particular window, combine. They manage to produce fresh RDDs of their windowed DStream.

19. What do you know about a DStream?

Answer: A Discretized Stream, commonly known as a DStream, is a succession of Resilient Distributed Databases (RDDs) representing a flow of information. Creation of DStreams from various sources such as Apache Kafka, HDFS, and Apache Flume.

All heads of DStreams have two surgeries –

- Output operations which write data to an outside system

- Transformations who Generate a new DStream

20. During running Spark applications, is it necessary for the consumer to set up Spark on all the available YARN cluster nodes?

Answer: It is one of the most striking features of Spark. It doesn't need installing when conducting a job under YARN or Mesos. That is because Spark can execute together with YARN or Mesos clusters without causing any change to your bunch.

21. What can you tell us about the Catalyst framework.

Response: A Catalyst framework is a new optimization framework within the Spark SQL. This special framework enables Spark to automatically transform SQL queries by adding further optimizations to construct a quicker processing system.

22. Can you mention that the companies which use Apache Spark within their respective production?

Response: Some of the companies that make use of this Apache Spark within their production are

1. Conviva
2. Pinterest
3. Shopify
4. Open-Table

23. What do you know about the Spark library that enables trusted file sharing in-memory speed across various cluster frameworks?

Response: The Tachyon is the Spark Library that lets trusted file sharing in-memory speed across different bunch frameworks.

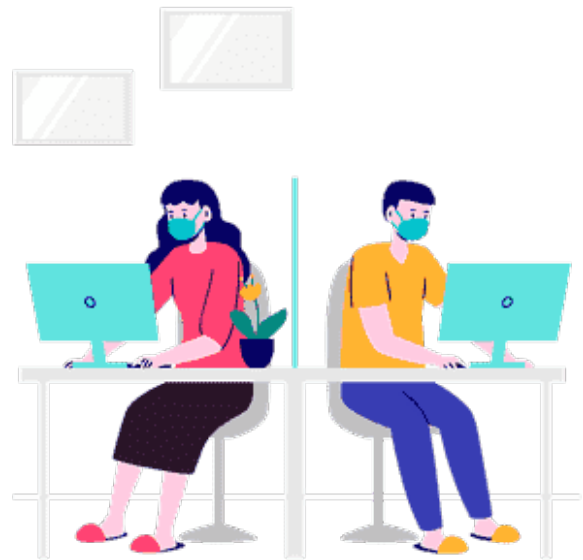
24. Do you know anything about BlinkDB? Why can it be used?

Response: BlinkDB is a query engine used for implementing interactive SQL queries on large quantities of data and condenses query outcomes marked with purposeful error bars. It is a beautiful instrument and helps users to balance query accuracy together with reaction time.

25. Can you differentiate between Hadoop and Spark concerning the simplicity of use.

Response: The user requires the Hadoop MapReduce for programming in Java that isn't easy. The developed Pig and Hive make programming using Java considerably easier. However, learning the syntax of Pig and Hive takes a whole lot of time. Spark has several interactive APIs for various languages such as Java, Python, or Scala and includes Spark SQL. It makes it relatively easier to utilize than Hadoop.

Spark Interview Questions (26-50)



Spark interview questions and answers

26. Mention the common mistakes that developers usually commit when conducting Spark software.

Answer: The most common mistakes that programmers usually commit when running Spark applications:

- Hitting the web service several times by employing multiple clusters.
- Run everything on a local node instead of distributing it.

27. Please mention the advantages of a Parquet file.

Answer: Parquet document is a columnar format file that helps the user to–

- Consumes less space
- Limit I/O operations
- Fetches only required columns

28. Mention the various information sources accessible SparkSQL.

Response: The various data sources accessible in SparkSQL are:

- Parquet document
- JSON Datasets

29. How can an individual execute Spark using Hadoop?

Answer: Spark was designed with its cluster management computation and utilized Hadoop for storage chiefly.

30. Mention the features of Apache Spark, which make it so popular.

Answer: The features of Apache Spark that make it so popular are:

- 1. Apache Spark provides innovative analytic options like graph algorithms, machine learning, streaming information, etc..
- 2. Apache Spark has good performance gains, as it helps to run an application in the Hadoop cluster ten times faster on disc and 100 times faster within the memory.
- 3. Apache Spark has built-in APIs in multiple languages like Java, Scala, Python, and R.

31. Tell us something about the Publish RDD.

Answer: We can do Particular operations on RDDs in Spark, working with the available key/value pairs. Pair RDDs allow countless users to get each key. Additionally, they have a `reduceByKey ()` method that gathers data based on each key and a `link ()` system that unites different RDDs collectively, based on the components that have the same key.

32. Between the Hadoop MapReduce and Apache Spark, which should we use for a job?

Answer: Choosing an application or development software depends on the given project scenario. Spark uses memory instead of network and disk I/O. But, Spark utilizes a large amount of RAM and requires a dedicated server to produce significant results. Therefore, the decision to use Hadoop or Spark varies sharply together with the organization's job and funding requirements.

Spark Interview Questions and answer Sample Video

33. Mention different types of transformations on DStreams.

Response: The different types of transformations on DStreams are:

- Stateful Transformations- performance of this batch depends on the intermediary results of the previous set.
- 0 Cases – Transformations that depend on sliding windows.
- Stateless Transformations- performance of the batch does not depend on the output of the previous set.
- 0 Cases – `map ()`, `reduceByKey ()`, `filter ()`.

34. Explain about the widespread usage cases of Apache Spark.

Answer: We use Apache Spark for:

- Interactive data analytics and processing.
- Sensor data processing

35.Can we use Apache Spark for Reinforcement Learning?

Response: No, a user Can't use Apache Spark for Reinforcement Learning. The Apache Spark works well for simple machine learning algorithms such as clustering, regression, and classification.

36.What do you understand about the Spark Core?

Answer: Spark Core is one of the features of Spark. It's all of Spark's basic functionalities, for example, interacting with storage systems, memory management, fault recovery, scheduling activities, etc.

37.Can the user remove the components with a key present in a different RDD?

Response: The user can remove the components with a key present in a different RDD using the subtract key () function.

38.Differentiate between persist() and cache() methods.

Response: The persistent () system allows the user to specify the storage level, whereas the method cache () uses the default storage amount.

39.Inform us about the numerous levels of persistence at Apache Spark.

Response: The Apache Spark automatically continues the intermediary data from several repeat surgeries. It depends on the users to call the persist() method on the RDD for reuse. The Spark has various persistence levels to keep several RDDs on the disc or inside the memory or combine both the disk and

the memory using different replication levels.

The Variety of storage/persistence levels in Spark are –

- OFF_HEAP
- MEMORY_ONLY
- MEMORY_AND_DISK
- MEMORY_ONLY_SER
- MEMORY_AND_DISK_SER, DISK_ONLY

40. How has the Spark been designed to handle monitoring and log in while in the Standalone mode?

Response: Spark has an online user interface for keeping an eye on the cluster in standalone mode that shows the audience and job statistics. The user log output for each job consists of the working directory of the slave nodes.

Is 41. Can the Apache Spark supply checkpointing to the user?

Answer: Lineage charts exist within Apache Spark to recuperate RDDs from a collapse. Nonetheless, this is time-consuming if the RDDs have lineage chains. Spark was provided with an API to get checkpointing, i.e., a REPLICATE flag to last. The choice on which the user decides data to the checkpoint. Checkpoints are helpful when the lineage graphs are extended and have complete dependencies.

42. How can the user start Spark tasks within Hadoop MapReduce?

Response: Using the SIMR (Spark in MapReduce), consumers can execute any Spark job inside MapReduce without using any admin rights.

43. How is Spark capable of using Akka?

Response: Spark uses Akka for scheduling. All the users

usually request for a job to learn after registering themselves. The master assigns the task. In this specific case, Spark uses Akka for messaging between the employees and the pros.

44. How is the consumer able to attain high availability in Apache Spark?

Answer: The user can Attain high availability in Apache Spark by applying the given methods:

- By implementing a single node recovery following the local file system.
- By using the StandBy Experts together with the Apache ZooKeeper.

45. How does Apache Spark achieve fault tolerance?

Answer: The data storage model at Apache Spark relies on RDDs. The RDDs help achieve fault tolerance through the lineage charts. The RDD always stores info about the best way best to build from other datasets. If there is a loss of any partition of the RDD due to the collapse, the lineage helps construct only that particular lost partition.

46. Explain the main components of a distributed Spark application

Response: The core components of any distributed Spark application are as follows:

- Executor – It comprises the worker processes that run the individual tasks of a Spark job.
- Driver- This contains the process that runs on the main() method of this app to make RDDs and perform transformations and activities on them.
- Cluster Manager- This is a pluggable component in Spark establish Executors and Drivers. The cluster manager allows the Spark to operate outside managers like Apache

Mesos or YARN from the backdrop.

47. Do you know anything about Lazy Evaluation?

Answer: When Spark runs on a particular dataset, it protects the directions and makes a note of it so that it does not forget. However, Spark does nothing about the instructions unless the user requests for the final result.

When a transformation such as the procedure `map()` runs an RDD, Spark doesn't execute the operation immediately. Assessment of All changes in Spark Tare happens after the consumer reacts. It helps to maximize the general data processing workflow.

48. What do you know about a worker node?

Answer: An employee node is a node that can run the Spark application code in a cluster. A worker node may have more than one procedure, easily configured by placing the `SPARK_WORKER_INSTANCES` property in the `spark-env. Sh` file. Just one employee node starts when the `SPARK_ WORKER_ INSTANCES` property particularizes.

49. Tell us something about SchemaRDD.

Response: A RDD that comprises row objects (wrappers around the whole string or integer arrays) with schema information about the sort of data in each column is called a SchemaRDD.

50. Mention the downsides of using Apache Spark over Hadoop MapReduce.

Answer: The Apache Spark doesn't perform very well for compute-intensive tasks and absorbs a high number of system tools. Apache Spark's in-memory capability causes a significant barrier to the cost-effective processing of big data. Spark includes its file management program and integrates with other cloud-based data platforms and Apache Hadoop.

51. Does the user need to set up Spark on all the YARN bunch nodes while running Apache Spark on YARN?

Answer: No, the consumer doesn't have to install Spark on all the nodes of a YARN bunch when running Apache Spark onto YARN because Apache Spark runs together with YARN.

52. Do you understand anything regarding the Executor Memory at a Spark application?

Response: Each Spark application has the same fixed heap size and a predetermined number of cores for a Spark executor. The heap size/ Spark executor memory controlled together with the Spark .executor. Memory property of this -executor-memory flag. Every Spark application has a single executor on every worker node. Executor memory is a measure of the magnitude of memory of the worker node that the application utilizes.

53. What's Spark Engine made to achieve?

Answer: The Spark engine accomplishes several tasks: creating schedules, distributing, and tracking all data applications across the Spark cluster.

54. Apache Spark is excellent at low-latency workloads like chart processing and machine learning. Elaborate on the reasons behind this.

Answer: The Apache Spark stores info in-memory for quicker model building and training. All Machine learning algorithms require numerous iterations to make an optimal model outcome. Likewise, graph algorithms browse through all of the nodes and edges. All these minimal latency workloads that require multiple iterations may result in more outstanding performance. Less disk controlled and access network traffic alters the entire equation when there's a lot of information to be processed.

55.Does the user needs to start Hadoop to conduct any Apache Spark Program?

Answer: No, beginning Hadoop is not required for the user to conduct any Spark application. Because there is no separate storage in Apache Spark, it uses the Hadoop HDFS. The information could be stored in a local file system and can be conveniently loaded from the local file system and processed.

56.Mention the default amount of parallelism in the Apache Spark.

Response: If the user does not explicitly state the parallelism level, then the number of partitions is considered the default degree of parallelism in Apache Spark.

57.Do you know anything about the frequent workflow of a Spark program?

Answer: Yes, the following is the typical workflow of the Spark Program:

- The first step at a Spark program is creating these input RDD's from external data.
- The various RDD transformations, including filter(), are used to develop new changed RDD's based on the company logic.
- The persist() method is employed for any intermediate RDD to be reused in the future.
- Eventually, the various RDD activities like first(), count() start the parallel computation. Later these are optimized and implemented by Spark.

58.In a Spark app, how does the user identify whether a specified operation is a Transformation or Action?

Answer: You can identify the operation based on the return type –

1.The operation is purely an Action if the return type is anything else than RDD.

2.The performance is Transformation when the return type is the same as the RDD.

59.What's a frequent mistake any Apache Spark developer usually makes while coping with Spark?

Response: Some of the common mistakes that all Apache Spark programmers make while coping with Spark are:

- Maintaining the necessary size of shuffle blocks.
- Trying to manage directed acyclic graphs (DAG's.)

60. It may sound tricky as a Spark interview question, but sometimes, the interviewer asks very basic ones. What is the primary difference between Spark SQL & HIVE?

Response: The following are the differences between Spark SQL and Hive:

- We can implement Any Hive question in Spark SQL; however, vice-versa is not valid.
- Spark SQL is faster than Hive.
- It isn't compulsory to make a metastore in Spark SQL, but it is mandatory to create a Hive metastore.
- Spark SQL deduces the schema while in Hive, the schema needs to be explicitly declared.

61.Mention the resources from where the Spark streaming component can process real-time information.

Answer: Normally, the users, apply Apache Flume, Apache Kafka, and Amazon Kinesis for Spark streaming to process real-time information.

62.Which are the companies that are presently utilizing Spark Streaming?

Response: Uber, Netflix, Pinterest are several companies that are presently using Spark Streaming.

63.What is the bottom layer of the abstraction in Spark Streaming API?

Response: DStream is the bottom layer of abstraction from the Spark Streaming API.

64.What do you know about receivers in Spark Streaming?

Answer: Receivers are particular entities in Spark Streaming that have data from various data sources and transfer them accordingly to Apache Spark. Receivers are often created by streaming contexts as long-running tasks on different executors and scheduled to function in a Round-Robin manner with each receiver taking a single core.

65.How is the user supposed to figure the number of executors necessary to do the real-time processing using Apache Spark? What factors need for deciding on the number of nodes for real-time processing?

Answer: Benchmarking the hardware, we calculate the number of nodes. While doing this, one must also consider numerous elements like optimum throughput (network rate), memory usage, the implementation frameworks used (YARN, Standalone, or Mesos), and contemplating the other jobs running inside these implementation frameworks together with Spark.

Scenario-Based Spark Interview Questions

66. Please elaborate change () work in Spark Streaming.

Response: The change () work in Spark Streaming enables the concerned developers to utilize Apache Spark transformations on the inherent RDD's for the flow.

The map() function in Hadoop is for element-to-element transform using the change () function. The real map() method works on the elements of Dream, while the change () method makes it possible for developers to operate with RDD's of this DStream. A map() method is a basic transformation, whereas the change () method is an RDD transformation.

67. What do you mean by Apache Spark?

Response: The Apache Spark is a fast, easy-to-use, and flexible data processing framework. It has an innovative implementation engine supporting cyclic data stream and in-memory computing. Spark can operate on Hadoop, run standalone or at the cloud, and obtain diverse data sources such as HDFS, HBase, Cassandra, etc.

68. Explain a few of the significant features of Spark.

Answer: Spark is now a favorite tool among developers because of the following features:

- Spark supports numerous analytic tools for interactive query analysis, real-time investigation, and graph processing
- The Spark has an interactive language casing as an independent Scala (Spark's language) interpreter.

69. Do you know anything about RDD?

It's a fault-tolerant collection of functional components that operate parallel. The partitioned information in RDD is immutable and dispersed. There are primarily two types of RDD:

1. Parallelized Dimensions: The present RDD's running parallel together.
2. Hadoop datasets: They perform a function on every record in HDFS or a different storage system.

70. Define Partitions.

Response: A partition is a smaller and logical branch of information similar to split' in the MapReduce procedure for programming. Partitioning is the procedure used to derive reasonable units of data from hastening the processing process. Everything in Spark is a partitioned RDD.

71. What are the operations backed by RDD?

Response: A RDD only supports the following two purposes:

72. What do you know about Transformations in Spark?

Answer: Transformations are acts applied to RDD, resulting in a different RDD. A transformation precisely doesn't execute until the action happens. The map() and also filter() methods are cases of the Transformation. A map() method former applies the function passed to it on every RDD section and outcomes into a different RDD. The filter() generates a new RDD by selecting elements to form the current RDD that pass purpose debate.

73. Elaborate on the concept of Actions.

Response: An act in Spark can help in restoring the information from RDD to the machine. The execution of any of the action is the result of all earlier created

transformations. The `reduce()` method is an activity that implements the role passed again and again till one value is finally remaining. The `take()` Action takes all values from an RDD to the local node.

74.Mention the functions of SparkCore.

Response: The SparkCore Functions as the base engine and performs several purposes, such as:

- Memory Administration
- Job scheduling
- Monitoring jobs
- Fault-tolerance

75.What do you know about RDD Lineage?

Response: Spark does not support data replication in memory. If any data is lost,t RDD lineage rebuilds it mechanically. An RDD lineage is a process that reconstructs lost data partitions and always recalls how to build from different datasets.

The spark Interview Questions (75-100)

Spark Interview Questions

76.Please explain Spark Driver.

Response: Spark Driver is the program that runs on the master node of the machine and can be used to announce transformations and Activity on information RDDs. The driver in Spark generates SparkContext, attached to a specified Spark Master. The driver also provides the RDD charts into the Spark Master, in which the standalone cluster manager runs.

77.Explain Hive on Spark.

Response: Hive contains significant aid for Apache Spark, but Hive implementation to Spark through the given piece of code:

```
Hive> place hive.execution.engine=spark;
```

78. Name a few of the commonly-used Spark Ecosystems.

1. Spark SQL (Shark)- for programmers.
2. SparkR to Market R in Spark engine.
3. GraphX for generating and computing charts.
4. Spark Streaming for processing live data streams.

79. Explain Spark Streaming

Answer: Spark supports flow processing. It is an extension of the Spark API and enables stream processing of live data streams.

We procure data from different sources such as Flume and HDFS.

Streamed and processed to document systems, databases, and live dashboards. It is close to batch processing as the input divides into channels such as batches.

80. What do you know about GraphX?

Response: Spark utilizes the application, GraphX, for chart processing to build and transform interactive graphs. The GraphX component enables the programmers to study structured data at scale.

81. Why is the MLlib required?

Response: The MLlib is an accessible machine learning library provided inside the Spark. It makes machine learning scalable and straightforward with shared learning algorithms and uses clustering, regression filtering, dimensional reduction, and alike.

82.What do you know about Spark SQL?

Answer: SQL Spark or Shark is a publication module introduced to utilize structured data and execute structured data processing. The crux of the Shark supports an altogether different RDD known as the SchemaRDD. The SchemaRDD comprises rows of items and schema objects specifying each column's data type in the row and is somewhat like a table in a relational database.

83.Next Spark Interview Question: Explain the Parquet file.

Response: A Parquet is a columnar format file supported by many different data processing methods. Spark SQL performs both the write and reads operations together with the Parquet file.

84.What are the file systems, do you think, supported by Spark?

- Neighborhood File system.
- S3

85.Elaborate on the Yarn.

Answer: The Yarn is one of the crucial features of Spark and is remarkably like Hadoop. It provides a central and source management platform to deliver accessible operations across a cluster efficiently. When the user uses Spark on Yarn, he must necessitate a binary distribution of info built on the Yarn support.

There are few things more fun than talking to seventh graders about their passions and vision for their future!
[#sparkinterview](#) [#PSDProud](#) [#icap](#) [#successinachangingworld](#)
[#connections](#) pic.twitter.com/SdIwgwvB2F

– Jesse Morrill (@Jessenm32) [September 18, 2018](#)

86. Kindly List the functions of Spark SQL.

Response: The Spark SQL is capable of accomplishing the following functions:

- Loading data from a variety of sources that are structured.
- Providing integration between the SQL and regular Python/Java/Scala code, together with the capacity to join RDDs and SQL tables, and expose customized SQL functions.
- Querying all the data using SQL statements, both inside a Spark app and outside tools connected to Spark SQL through standard database connectors (JDBC/ODBC).

87. Mention the advantages of using Spark as compared to MapReduce.

Answer: Spark has several benefits as compared to MapReduce:

- Spark implements the processing of about 10-100x quicker than Hadoop MapReduce due to the availability of in-memory processing. MapReduce uses persistent storage for any of the information processing tasks.
- Spark is capable of doing iterative computation while there's absolutely no pragmatic computing employed by Hadoop.

88. Is there some advantage of learning Hadoop MapReduce?

Answer: Yes, an individual should know the Hadoop MapReduce. It's a paradigm used by several large data tools, including Spark, and is extremely relevant when the data grows bigger and bigger. Tools such as Pig and Hive convert their queries into MapReduce phases to optimize them better.

89. What should you know about the Spark Executor?

Response: When a person connects the SparkContext into a

cluster manager, it acquires a Spark Executor on the cluster nodes. The Executors are Spark procedures that run computations and save the data on the worker node. The last task transfers to the Spark Executors because of their final performance.

90.Name the Kinds of Cluster Managers within Spark.

Response: Spark supports three major types of Cluster Managers:

- Standalone: Manager to establish a cluster.
- Apache Mesos: The generalized/commonly-used bunch manager also runs Hadoop MapReduce and other programs.
- Yarn: It is responsible for resource management in Hadoop.

91.What do you know about the worker node?

Response: The Worker node refers to any node that can run the application code in a bunch.

92.Do you know anything regarding the PageRank?

Response: The PageRank is the measure of every vertex from the graph and is one of the striking characteristics of this Graph in Spark.

93.Does the user install Spark on all Yarn cluster nodes while conducting Spark on Yarn?

Response: No, there is no compulsion regarding this since Spark runs along with Yarn.

The spark Interview Questions

94.Tell us about some of those demerits of using

Spark.

Response: Spark uses more storage space than Hadoop and MapReduce and hence, can cause specific problems. Developers must be mindful while conducting their applications in Spark. All Spark programmers must be sure that the work distributes evenly over multiple clusters.

95.How can an individual create RDD?

Response: Spark supplies two approaches to Make RDD:

From parallelizing a set on your Driver app.

```
val distIntellipaatData = sc.parallelize(IntellipaatData)
```

By loading the external dataset from external storage like HDFS and shared file system.

96.What are the key features of Apache Spark.

Answer: Apache Spark has the following key attributes:

- Spark can run along with a current Hadoop cluster utilizing YARN for resource scheduling.
- Speed: Spark runs around 100 times faster than Hadoop MapReduce for large-scale information processing. It can accomplish this rate through controlled partitioning and deftly manages information using walls that help parallelize dispersed data processing using minimal traffic.
- Real-Time Computation: Spark's computation in real-time. It has less latency due to its in-memory analysis. It is created for massive scalability and supports many computational models.
- Machine Learning: Spark's MLlib is your machine learning component used for big information processing. It eliminates the need to utilize numerous tools for processing and also for machine learning.
- Lazy Assessment: Apache Spark guarantees its test until

it's totally necessary and accelerates the whole process. For transformations, Spark adds to the DAG of computation, and when the driver requests some data, DAG gets executed.

97.How many languages supported by Apache Spark? What's the most popular language?

Response: Apache Spark supports the following four languages:

We can access the Scala shell. /bin/pyspark.

98.What is Resilient Distributed Dataset (RDD).

Answer: RDD stands for Resilient Distribution Datasets. It is a fault-tolerant collection of functional components that operate in parallel. There are two types of RDD:

- Parallelized Collections: The present RDDs are running parallel with one another.
- Hadoop Datasets: They play works on each file record in HDFS or alternative storage systems.

RDDs are elements of data. The memory [distributed](#) across several nodes save it. They're lazily evaluated in Spark, making Spark operate at a faster rate.

99.How do we create RDDs in Spark?

Answer: Spark provides two methods to Make RDD:

From parallelizing a collection on the Driver program. The goal of the method is to make use of SparkContext'parallelize' by applying the following part of code:

```
method value DataArray is equal to Array(2,4,6,8,10)
```

By loading the external dataset from external storage such as HDFS, HBase, shared file system.

100.Tell about the Executor Memory in any Spark application?

Response: Every Spark application has the same fixed heap size and a set number of cores for a Spark Executor. This heap size the Spark executor memory and can be controlled together with the Spark `.executor.memory` property of this `-executor-memory` flag. Every Spark application always has one Executor working on every worker node. An executor memory is a measure of the amount of memory that the worker node of any application will utilize.

Hope all these tips will help you to succeed in answering some tough spark Interview Questions.