NAME: SHAIK ABDUL KHADAR JILANI          ROLLNO: DXC262AB12038

BATCH: DXC-262-Analytics-B12-Azure          SUBMISSION: 7-6-2022

COMPANY: DXC TECHNOLOGY

DAY 7

**1) Explain what are various components of SPARK with block diagram?**

**explain functionality of every components?**

Components of Apache Spark

GraphX

| Spark Core | Spark SQL | Spark Streaming | MLlib | GraphX |
|---|---|---|---|---|

**SPARK CORE**

-It is the base engine for large-scale parallel and distributed data processing

-It is responsible for

1)Memory management

2)Fault recovery

3)scheduling distributing and monitoring jobs on a cluster

4)Interacting with storage system

**SPARK SQL**

-Spark SQL framework component is used for structured and semi-structured data
 processing

**Spark streaming**

-It is a lightweight API that allows developers to perform batch processing and
 real-time streaming of data with ease

-provides secure, reliable and fast processing of live data streams

**MLlib**

-IT is a low-level machine learning library that is simple to use,is scalable,and
 compatible with various programming languages.

-MLlib eases the deployment and development of scalable machine learning algorithms

**GraphX**

-GraphX is Spark's own Graph Computation Engine and data store

-Provides a uniform tool for ETL

-Exploratory data analysis

-Interactive graph computations

**2) Explain Spark core in details & how RDD is related to Spark core - explain with Spark program ?**

SPARK CORE

-It is the base engine for large-scale parallel and distributed data processing

-It is responsible for

1)Memory management

2)Fault recovery

3)scheduling distributing and monitoring jobs on a cluster

4)Interacting with storage system

RDD

Spark Core is embedded with RDDs(Resilient Distributed Dataset) an immutable

fault-tolerant, distributed collection of objects that can be operated on in parallel

RDD

**1)transformation**

These are operations (such as map, filter, join, union) that are performed on an RDD that

yields a new RDD containing the result

E.g.,

val x=sc.textfile  (RDD will be created)

val Y=x.map

val z=y.filter

**2)Action**

These are operations(such as reduce, first, count,) that return a value after running a computation on an RDD

Eg: z.count()  ,x.count()

**3) Explain various MLlib algorithms Spark is supporting ?**

MLlib

-IT is a low-level machine learning library that is simple to use, is scalable and

compatible with various programming languages.

-MLlib eases the deployment and development of scalable machine learning algorithms

-It contains machine learning libraries that have an implementation of various machine

learning algorithms

1)Clustering

Clustering is an unsupervised learning problem whereby we aim to group subsets of entities with one another based on some notion of similarity. Clustering is often used for exploratory analysis and/or as a component of a hierarchical supervised learning pipeline (in which distinct classifiers or regression models are trained for each cluster).

2)Classification

The spark.mllib package supports various methods for binary classification, multiclass classification and regression analysis. Some of the most popular algorithms in classification are Random Forest, Naive Bayes, Decision Tree, etc.

3)Collaborative Filtering

Collaborative filtering is commonly used for recommender systems. These techniques

aim to fill in the missing entries of a user-item association matrix. spark.mllib

currently supports model-based collaborative filtering, in which users and products

are described by a small set of latent factors that can be used to predict missing

entries.

**4) Explain benefits Spark SQL & how relational data will be inserted into SPARK?**

Spark SQL, it is a module of Apache Spark that analyses the structured data. It provides Scalability, it ensures high compatibility of the system. It has standard connectivity through JDBC or ODBC. Thus, it provides the most natural way to express the Structured Data.

1) Integrated

Apache Spark SQL mixes SQL queries with Spark programs. With the help of Spark SQL, we can query structured data as a distributed dataset (RDD). We can run SQL queries alongside complex analytic algorithms using tight integration property of Spark SQL.

2) Unified Data Access

Using Spark SQL, we can load and query data from different sources. The Schema-RDDs lets single interface to productively work structured data. For example, Apache Hive tables, parquet files, and JSON files.

3) High compatibility

In Apache Spark SQL, we can run unmodified Hive queries on existing warehouses. It allows full compatibility with existing Hive data, queries and UDFs, by using the Hive fronted and MetaStore.

4) Standard Connectivity

It can connect through JDBC or ODBC. It includes server mode with industry standard JDBC and ODBC connectivity.

5) Scalability

To support mid-query fault tolerance and large jobs, it takes advantage of RDD model. It uses the same engine for interactive and long queries.

6) Performance Optimization

The query optimization engine in Spark SQL converts each SQL query to a logical plan. Further, it converts to many physical execution plans. Among the entire plan, it selects the most optimal physical plan for execution. Read more about Apache Spark performance tuning techniques in detail.

7) For batch processing of Hive tables

We can make use of Spark SQL for fast batch processing of Hive tables.

Apache Spark has multiple ways to read data from different sources like files, databases etc. But when it comes to loading data into RDBMS(relational database management system), Spark supports only Append and Overlay of the data using dataframes. Spark dataframes do not support Updating of data into a database.

- Read data from a CSV file
- Create a database schema and table in MySQL db
- Load spark dataframe data into a database.
- Update database table records using Spark

**5) Explain Spark streaming in detail?**

-It is a lightweight API that allows developers to perform batch processing and

 real-time streaming of data with ease

-provides secure, reliable, and fast processing of live data streams

| Input data Stream | → | Streaming | → Batches of input data → | Engine | → Batches of processed data → |

**6)Explain SPARK architecture? what is Master - Slave architecture?**

SPARK ARCHITECTURE

Master Node

Driver Program

SparkContext

Cluster Manager

Worker Node

Executor

Cache

Task    Task

As the name suggests, the master-slave is a database architecture divided into a master database and slave databases. The slave database serves as the backup for the master database. The master database is the keeper of the data resources and also the place where all the writing requests are performed.

**7)Explain various cluster managers in SPARK?**

-A job is split into multiple tasks that are distributed over the worker node

-When an RDD is created in Spark context, it can be distributed across various nodes

-Worker nodes are slaves that run different task

Spark Cluster Managers

1)Standalone mode

By default applications submitted to the standalone mode cluster will run in FIFO order
and each application will try to use all available nodes

2)MESOS

Apache Mesos is an open-source project to manage computer clusters and can also run
Hadoop application

3)Hadoop Yarn

Apache YARN is the cluster resource manager of Hadoop 2.Spark can be run on YARN

4)kubernetes

Kubernetes is an open source system for automating deployment,scaling and management
of containerized applications.

## 8) Explain with screenshots & steps how to create Cosmos DB?

Go to portal.azure.com

Search for cosmosDB and select Azure Cosmos DB



Click on create

-it will show Core(SQL),Azure Cosmos DB API for MongoDB,Cassandra,Azure Table, Gremlin(Graph) options to create

In that select Core(SQL) to create



-give accountname dxc and create

Click next till you get this page keep it all default and click next we will get Your deployment is complete



## 9) Explain with screenshots & step how to insert data into Cosmos DB?

After our deployment is complete go to resource

-click on Data Explorer

And the click on New Container

-Fill the New Container details and click on OK



-After that go to sports->items and click New item

-Write the data with key and value like the below



And then save it the data is stored



**10)Explain with screenshots & step how to create Azure SQL Db & also explain how to insert data into Azure SQL D?**

-First search the sql database and then we will get the home page like the below one

-Click on Create option we will go to the basics where we will give the details

-give server name and location

-click Configure database give Max vCore 1  GB and Apply

Keep backup storage redundancy as Geo-redundant backup storage



After that review + create and Our Deployment is complete



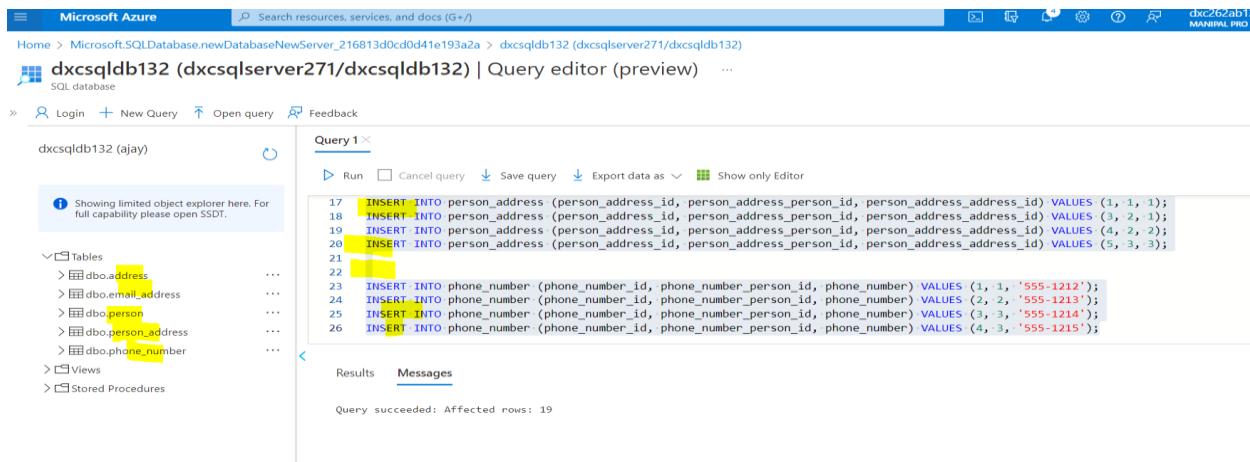To insert data in Azure Sql after deployment go to

-Query editor (preview)

Give your login id and password and click ok

Create the tables whatever we want and run the query tables will be created

After that insert the values we need in the table we get tables with inserted values



To see the inserted values in the table go to tables and write the query which is select * from table_name