# TRAINING ASSIGNMENT

# (7th JUNE 2022)

NAME – ROHIT ARORA            ROLLNO – DXC-262AB-1209

BATCH – DXC-262-ANALYTICS-B12-AZURE            COMPANY – DXC TECHNOLOGY

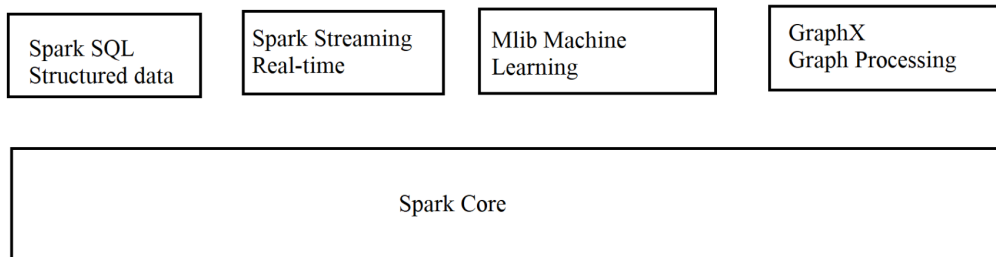EMPLOYEE DOMAIN – AZURE ANALYTICS            TRAINER NAME – MR. AJAY KUMAR

DATE OF SUBMISSION – 7th JUNE 2022            NO.OF QUESTIONS: 10

1. Explain what are various components of SPARK with block diagrams?
explain the functionality of every component??

Components of Apache Spark:



Spark Core - It is the base engine for large scale distributed and parallel data
    processing in Spark. It is responsible for memory management , fault
    recovery , job scheduling to different clusters and interacting with
    storage systems.

Spark SQL - Spark SQL is embedded with RDD , resilient distributed database with
    fault tolerance and is immutable with parallel processing features.

Spark Streaming - Spark streaming provides real time data processing because of a
    lightning fast clustering algorithm.

Spark Mlib - It provides and executes most of all the machine learning algorithms like
    clustering , classification and regression.

Spark GraphX - It provides all the data analysis features like ETL , EDA , reporting
    using graphs.

2.  Explain Spark core in details & how RDD is related to Spark core.Explain with Spark program ?

All the functionalities being provided by Apache Spark are built on the highest of the Spark Core. It delivers speed by providing in-memory computation capability. Spark Core is the foundation of parallel and distributed processing of giant dataset. It is the main backbone of the essential I/O functionalities and significant in programming and observing the role of the spark cluster. It holds all the components related to scheduling, distributing and monitoring jobs on a cluster, Task dispatching, Fault recovery. The functionalities of this component are:

·        It contains the basic functionality of spark. (Task scheduling, memory management, fault recovery, interacting with storage systems).
·        Home to API that defines RDDs.

Resilient Distributed Datasets (RDD) is a fundamental data structure of Spark. It is an immutable distributed collection of objects. Each dataset in RDD is divided into logical partitions, which may be computed on different nodes of the cluster.

·        TRANFORMATION – operations like map, filter performed on RDD yielding new RDD. Eg- val(x)=sc.textfile -àval y=x.map()

·        ACTION – Operations such as reduce,first,count that returns a value after running a computation RDD. Eg- z.count()

3.  Explain various Mlib algorithms Spark is supporting ?

Spark MLlib is used to perform machine learning in Apache Spark. MLlib consists of popular algorithms and utilities. MLlib in Spark is a scalable Machine learning library that discusses both high-quality algorithm and high speed. The machine learning algorithms like regression, classification, clustering, pattern mining, and collaborative filtering. Lower level machine learning primitives like generic gradient descent optimization algorithm are also present in MLlib.

        The popular algorithms and utilities in Spark MLlib are:

·        Basic Statistics.
·        Regression.
·        Classification.
·        Recommendation System.
·        Clustering.
·        Dimensionality Reduction.
·        Feature Extraction.
·        Optimization.

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

In reference to 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.  It provides a programming abstraction called DataFrames and can also act as a distributed SQL query engine. It enables unmodified Hadoop Hive queries to run up to 100x faster on existing deployments and data.

There are multiple methods to insert relational data in spark. For that first we need to open Google Collabs. Import PySpark. Then type the following commands and press Shift+Enter

5. Explain Spark streaming in detail ?

Apache Spark Streaming is a scalable fault-tolerant streaming processing system that natively supports both batch and streaming workloads. Spark Streaming is an extension of the core Spark API that allows data engineers and data scientists to process real-time data from various sources including (but not limited to) Kafka, Flume, and Amazon Kinesis. This processed data can be pushed out to file systems, databases, and live dashboards. Its key abstraction is a Discretized Stream or, in short, a DStream, which represents a stream of data divided into small batches. DStreams are built on RDDs, Spark's core data abstraction. This allows Spark Streaming to seamlessly integrate with any other Spark components like MLlib and Spark SQL. Spark Streaming is different from other systems that either have a processing engine designed only for streaming.
INPUT DATA STREAM àSPARK STREAMING àBATCH OF INPUT DATA àSPARK ENGINE à BATCHES OF PROCESSED DATA .

6. Explain SPARK architecture? what is Master - Slave architecure ?

The master slave architecture is there is a master node and in that master node all the worker nodes work

- 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 tasks

7. Explain various cluster managers in SPARK?

Apache Standalone - Application runs on FIFO mode when run in standalone and consume all resources available.
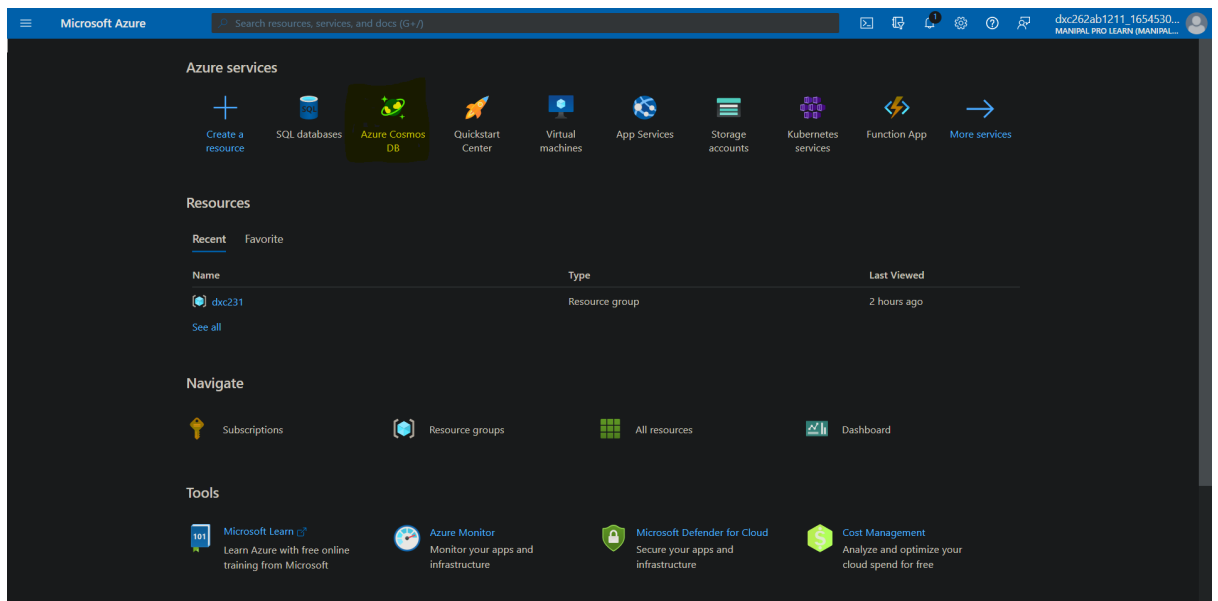
Apache Mesos - is an open source project to manage computer clusters
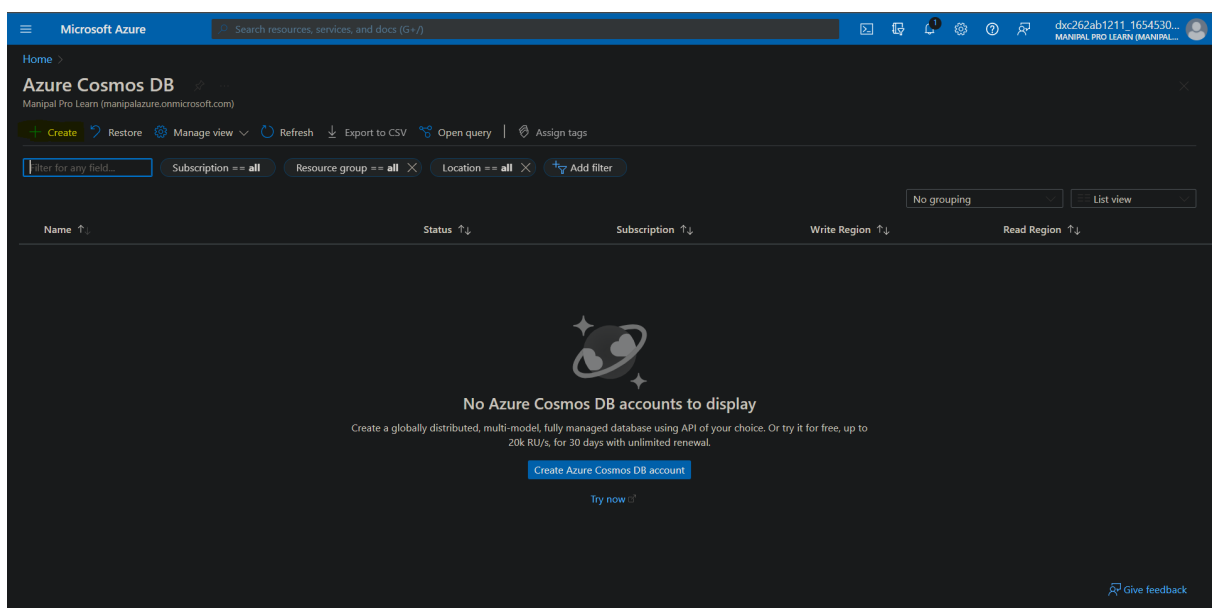
Apache Yarn - It is a cluster manager in hadoop 2 , spark can run on yarn.

Kubernetes - It is used for automating deployment , scaling and managing the containerized applications.

8. Explain with sceenshots & steps how to create Cosmos DB ?

To create Cosmos DB firstly we need to login into the azure cloud platform after that we need to select Azure Cosmos DB from resources

After selecting that we need to click on +Create in order to create a new CDB database



After we will click on the kind of database you want to create here we will create SQL supported Document database

After that we need to enter the details regarding the database



After that validate and review the details and click of create after validating

Through this we have created a Cosmos DB database with name as dxc101



After this our Database will be ready

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

To insert data into the database of cosmos DB , open the azure cloud platform and navigate towards the Azure cosmos DB and click on the database you want to work with

After choosing the database click on data explorer



After clicking on that click on the new container and enter the database you want to create and the tables in that

After the tables are created you need to navigate to the table and click on + new item



Add the data you want to insert and click on Save

Now through this you have successfully insert a data into Cosmos DB

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

To create a SQL database and insert data into it we need to first navigate to the SQL database on Azure Cloud Platform



After you are into the SQL database you need to click  + create

After clicking on create enter the details regarding the resource group , server and database id , name  and click on create
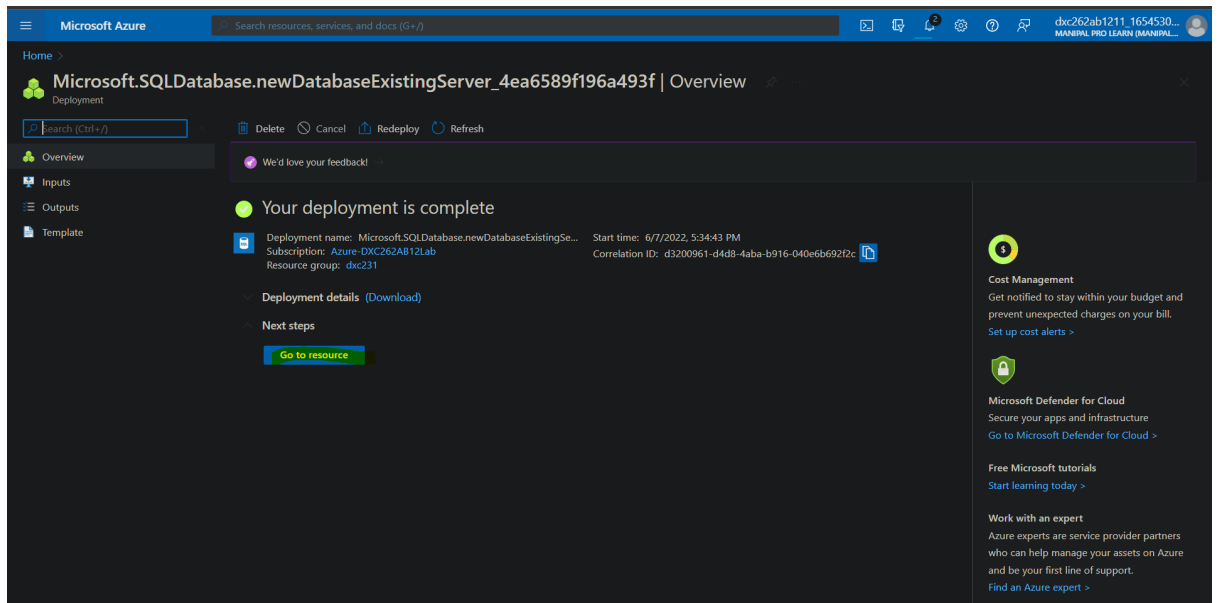
After that you need the validate the data that you have entered proceed and click create after validating
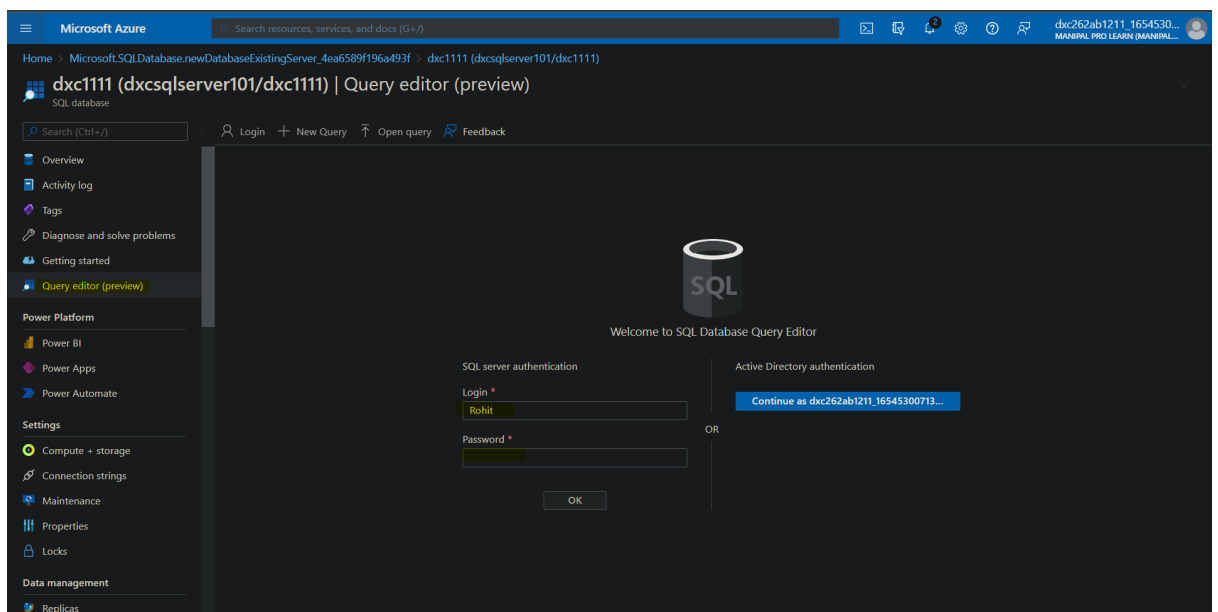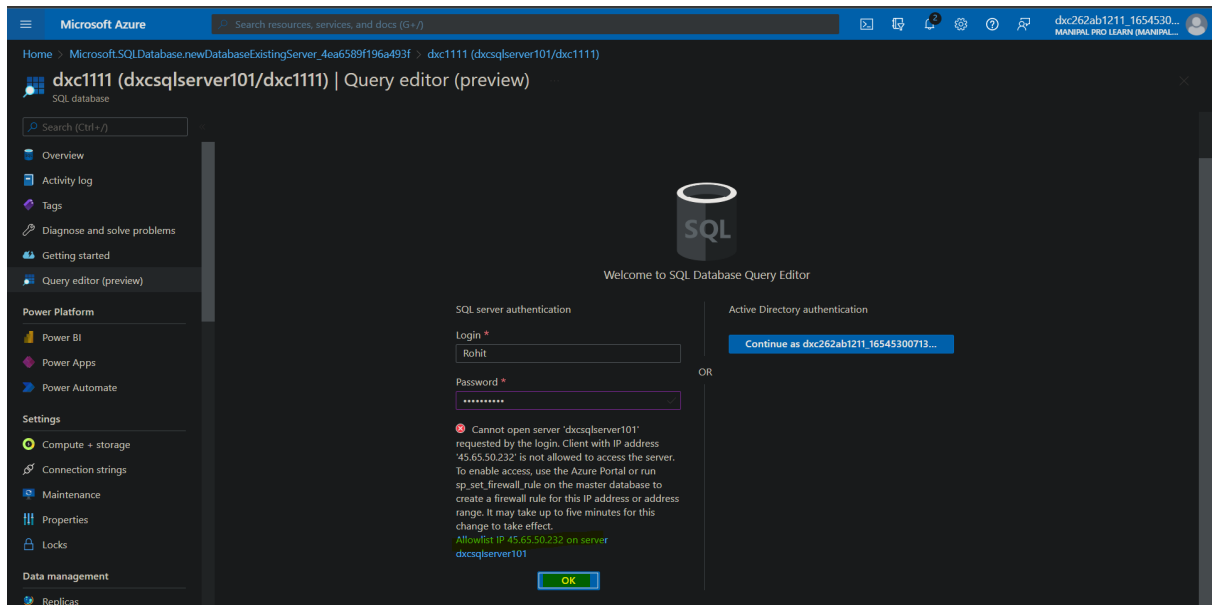


It wil; take some time to deploy the database



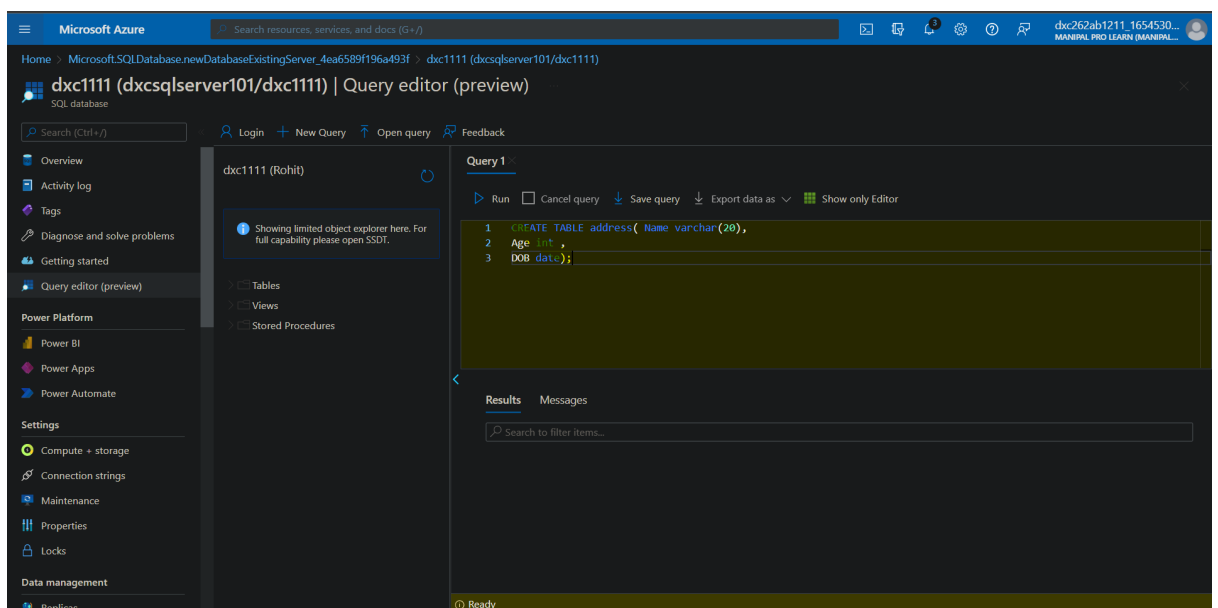After the deployment is done click on Go to resource
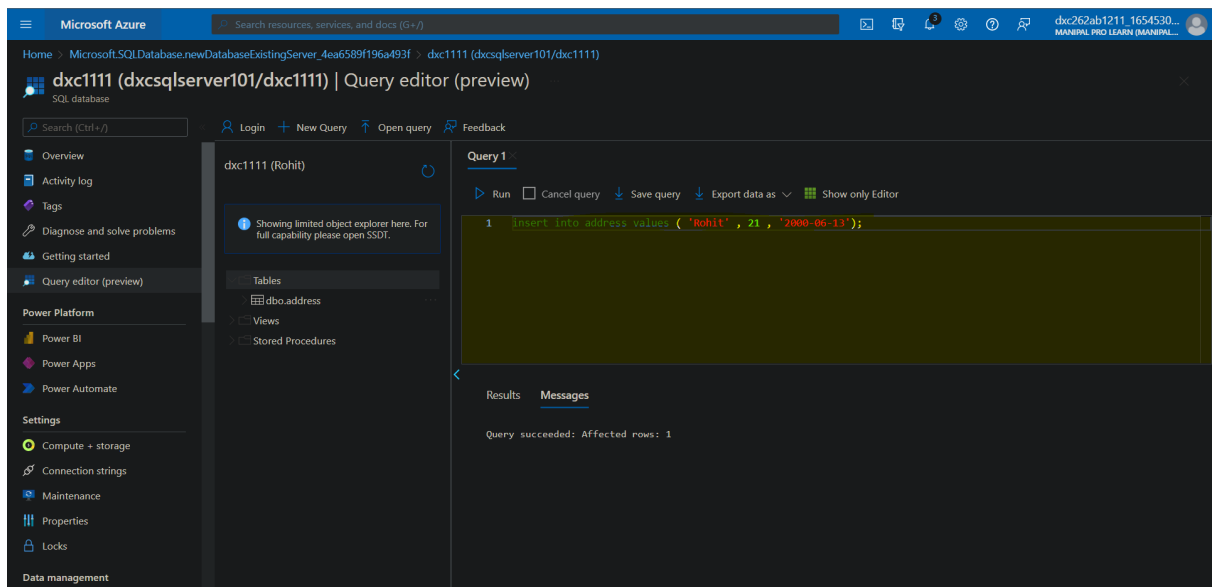
Proceed to Query Editor



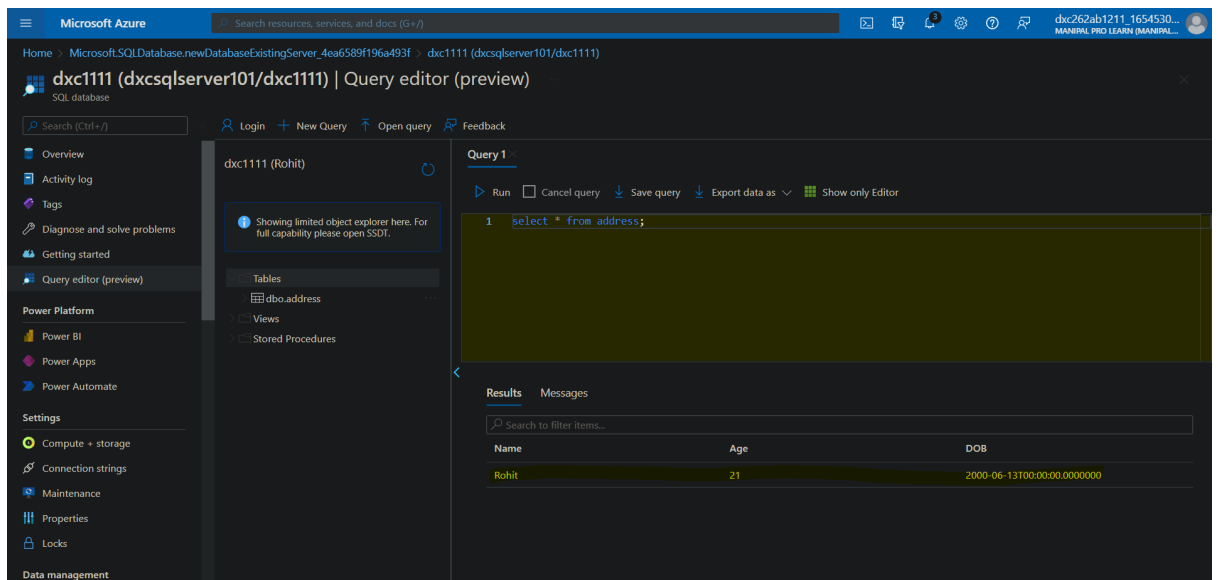Type in your server name and password

Allow the access for the IP and click on OK and you will now see a screen where you can now create table and insert data into it. Create a sample table



Click on run and insert a sample data in to the table address

So now the data is inserted we can view the inserted data using select command



Tada! We have successfully created and inserted data into the Azure Cloud platform SQL database