**Azure Databricks and AWS S3 Storage**

Databricks is an integrated analytics environment powered by Apache Spark which let you connect and read from many data sources such as AWS S3, HDFS, MySQL, SQL Server, Cassandra etc. We are focusing on S3 , Since it is very easy to work with.
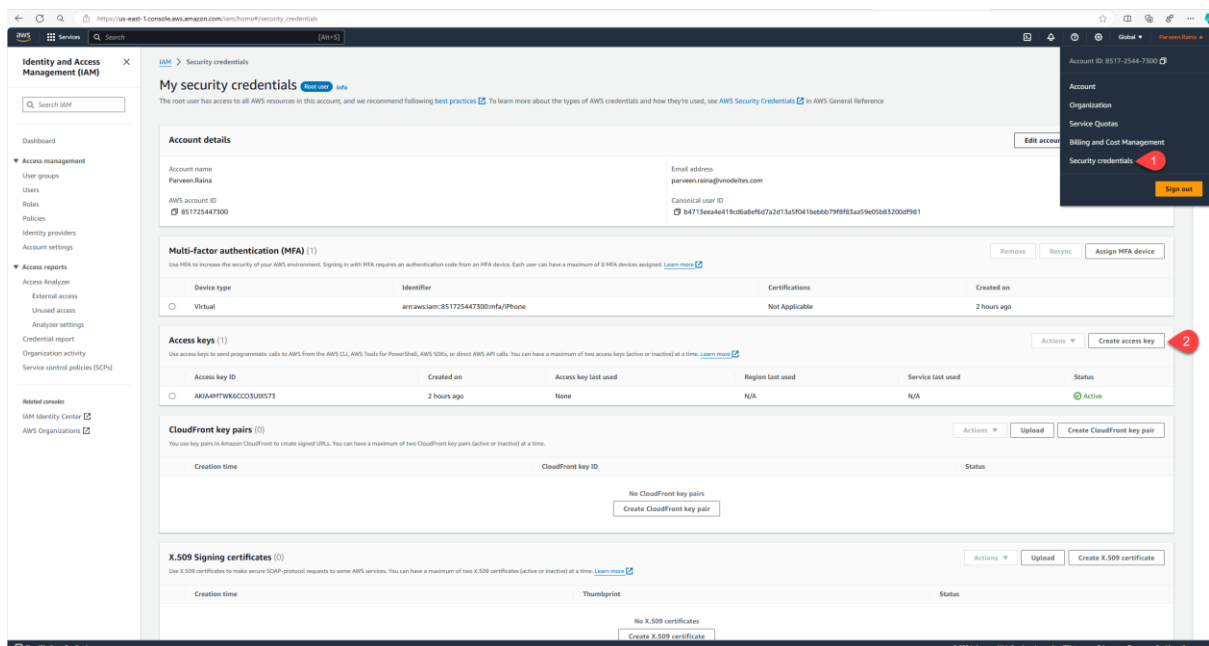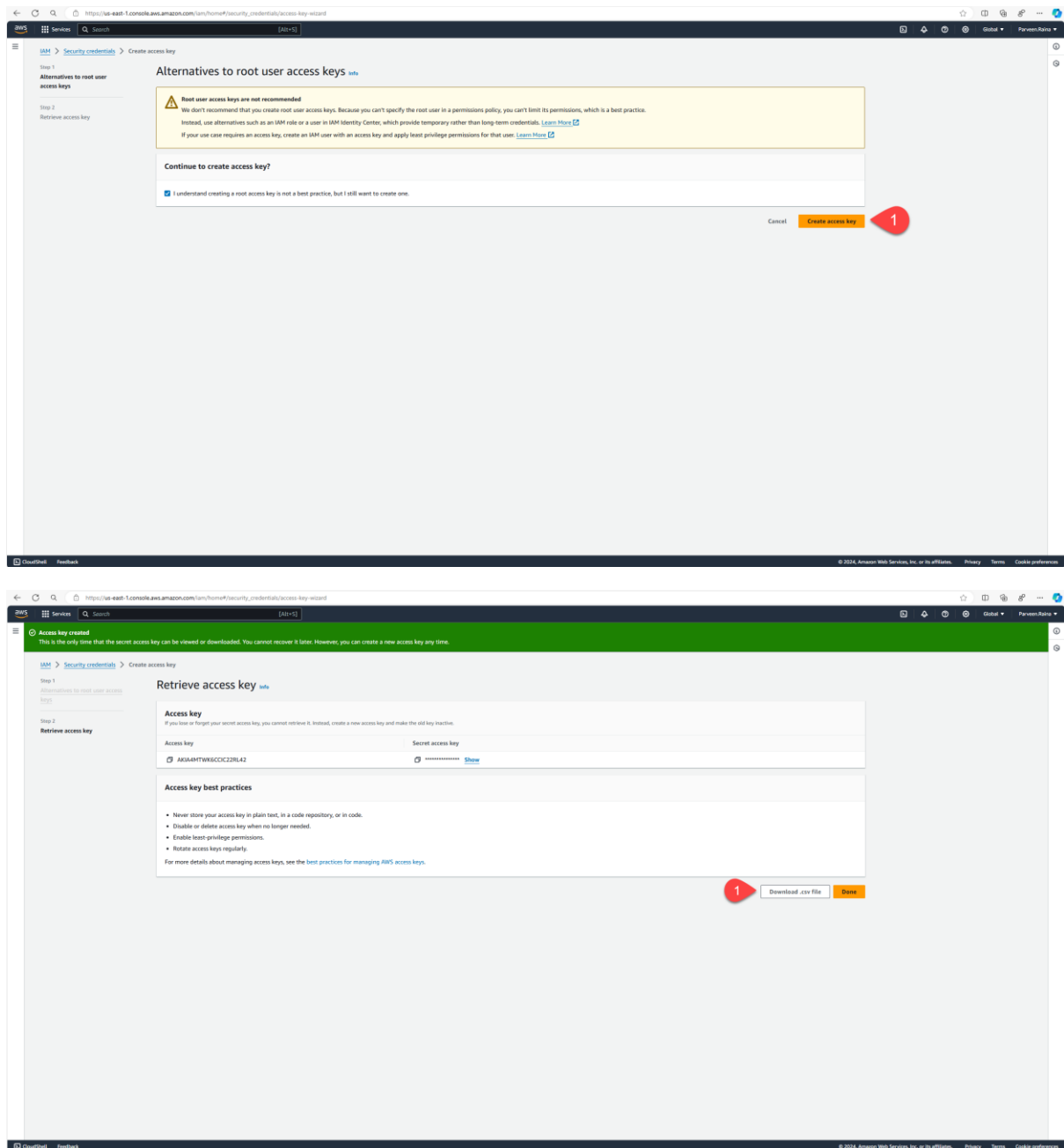
**Prerequisites:**

- *you must have AWS & Azure account with subscription*

- *S3 bucket created and contains some objects to work with*

Creating S3 storage bucket is beyond this section.
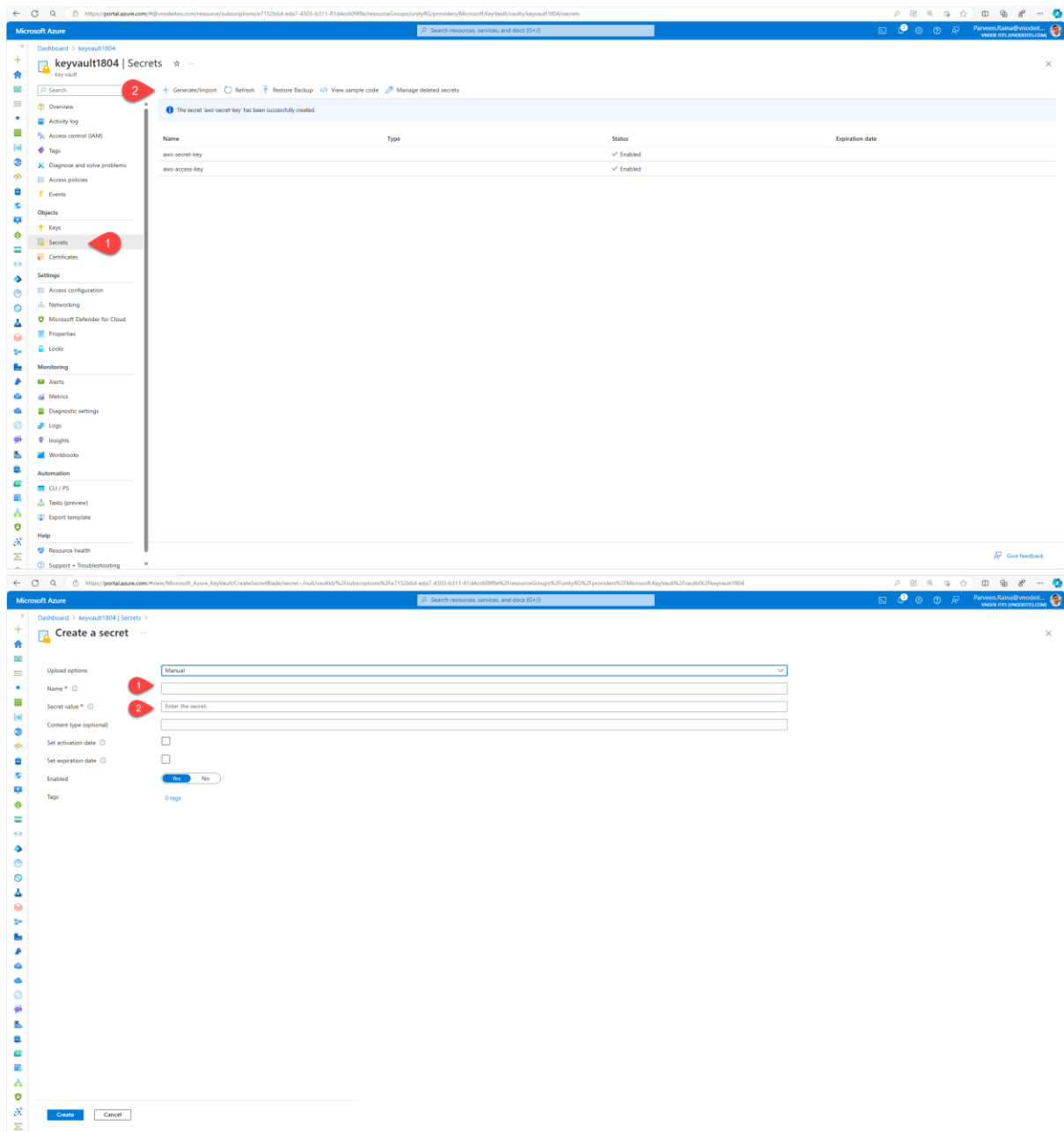
**Step 1.**

In AWS Management Console, go to 'Security Credentials' under your user account name. Then , you click 'Access keys (access key ID and secret access key)' section under My Security Credentials which will be provided with option called 'Create New Access Key'. once you click that option, you will be given **access key ID and secret access key.** please save it in notepad by pressing show access key or download key file for later use.
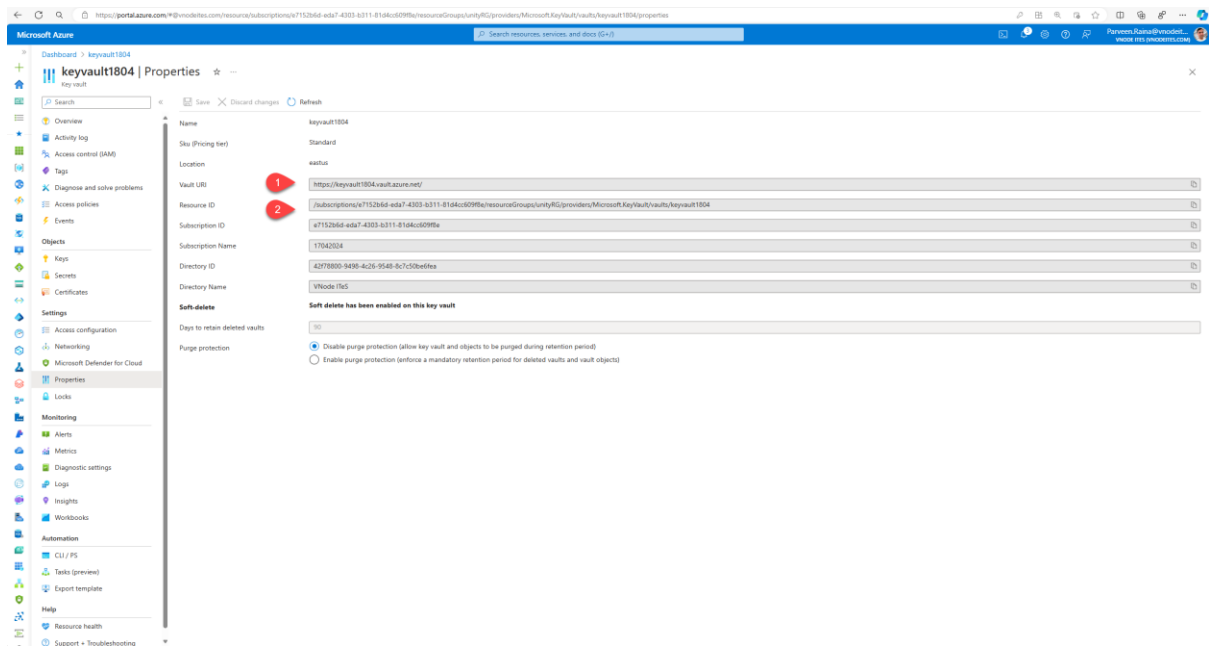
**Step 2.**

Since we access S3 bucket using databricks-backed scope, Secrets should be created by putting access key & secret key values in Azure key vault. Go to Azure Key Vault, in the resource menu, click secrets under Settings category. Then, click + sign (Generate/Import) in the command bar. you will be prompted with below window

Give *aws-access-key* and and *aws-secret-key* for name and paste copied *access key & secret access key* values in step 1 in the place of Value. By doing so two times, two secrets will be created. leave others default and click create. you are done with creating secret for accessing storage account.

## Step 3.

Navigate to properties under resource menu of Key Vault & copy DNS name and Resource ID and save it in notepad. This will be used whilst creating secret scope.
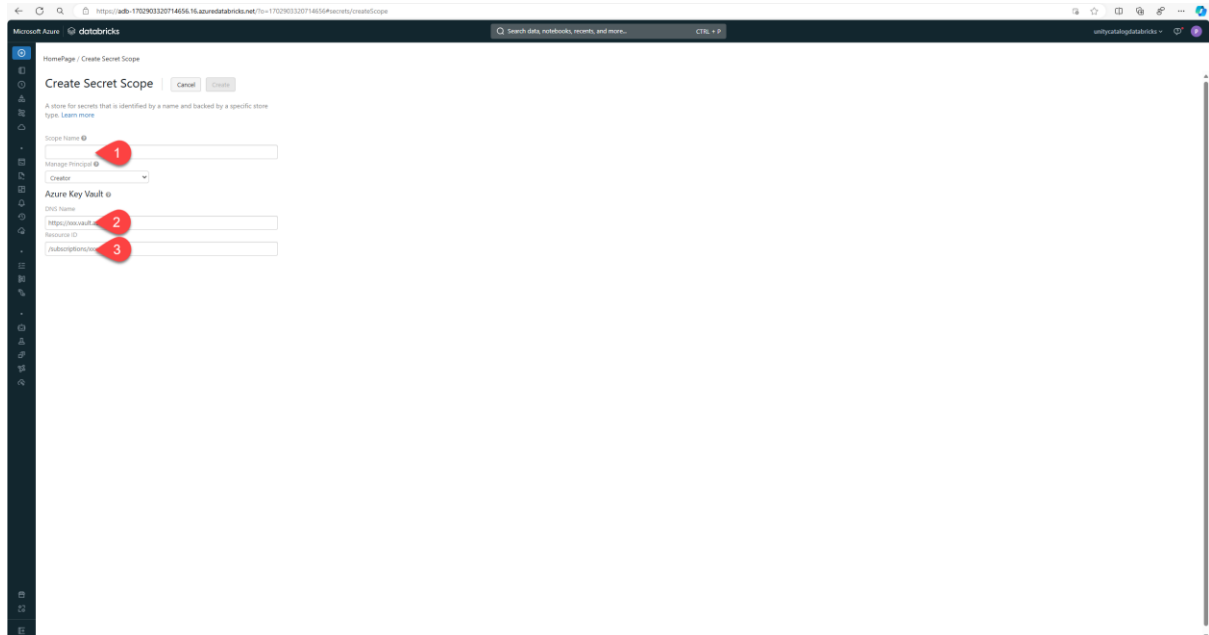
**Step 4.**

Go to https://<your_azure_databricks_url>#secrets/createScope

Ex- https://adb-1702903320714656.16.azuredatabricks.net/#secrets/createScope

you will be directed to below window



Give scope name *'aws'* which is uniquely identified in the database maintained by Databricks, leave Manage Principal as Creator and paste copied values of DNS name & resource ID in place of DNS Name and Resource ID fields. Click Create & you will be given with success message. Remember scope name or save it in file.

**Step 5.**

Accessing S3 bucketing directly by setting AWS keys in the spark context. Secrets utilities of Databricks utilities (DBUtils) will help us fetching sensitive credentials information without making them visible.

*ACCESS_KEY = dbutils.secrets.get(scope = "aws", key = "aws-access-key")*

*SECRET_KEY = dbutils.secrets.get(scope = "aws", key = "aws-secret-key")*

*sc._jsc.hadoopConfiguration().set("fs.s3n.awsAccessKeyId", ACCESS_KEY)*

*sc._jsc.hadoopConfiguration().set("fs.s3n.awsSecretAccessKey", SECRET_KEY)*



Spark context is now configured to connect to AWS S3 bucket by above diagrammatically.

Lets check whether we can be able to read below .csv file stored in my_bucket

**df = spark.read.csv("s3://unitycatalogbucket01/data/products.csv", header=True, inferSchema=True)**