# ANALYTICS SYSTEMS ENGINEERING (MSDS 436)
## EXERCISE 2
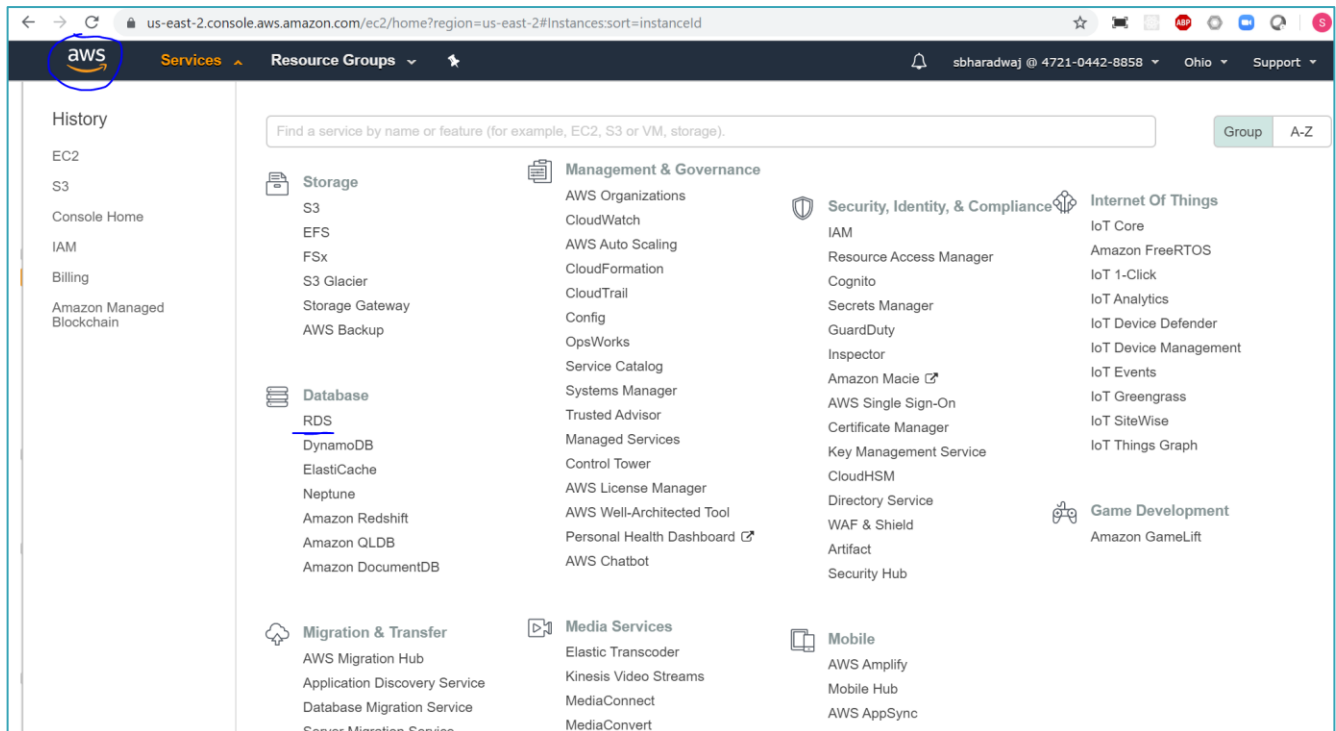
shreenidhi.bharadwaj@northwestern.edu | ChristopherFiore2015@u.northwestern.edu

## CONTENTS

CREATE POSTGRESQL DATABASE IN AWS

1. Go to AWS services and click on RDS (Relational Databases)

2. Click on Create Database and choose Standard Create > Amazon Aurora > Postgres Edition.



3. Scroll down, to select **One writer and multiple readers & Dev/Test**

**Database features**

○ **One writer and multiple readers**
Supports multiple reader instances connected to the same storage volume as a single writer instance. This is a good general-purpose option for most workloads.

○ **Serverless**
You specify the minimum and maximum amount of resources needed, and Aurora scales the capacity based on database load. This is a good option for intermittent or unpredictable workloads.

**Templates**
Choose a sample template to meet your use case.

○ **Production**
Use defaults for high availability and fast, consistent performance.

● **Dev/Test**
This instance is intended for development use outside of a production environment.

4. Scroll down to provide **database name (northwind), username (postgres) & password (rootroot)**

**Settings**

DB cluster identifier  Info
Type a name for your DB cluster. The name must be unique cross all DB clusters owned by your AWS account in the current AWS Region.

northwind

The DB cluster identifier is case-insensitive, but is stored as all lowercase (as in "mydbcluster"). Constraints: 1 to 60 alphanumeric characters or hyphens. First character must be a letter. Can't contain two consecutive hyphens. Can't end with a hyphen.

▼ **Credentials Settings**

Master username  Info
Type a login ID for the master user of your DB instance.

postgres

1 to 16 alphanumeric characters. First character must be a letter.

☐ Auto generate a password
Amazon RDS can generate a password for you, or you can specify your own password

Master password  Info

••••••••

Constraints: At least 8 printable ASCII characters. Can't contain any of the following: / (slash), "(double quote) and @ (at sign).

Confirm password  Info

••••••••

5. Scroll down to Select **'db.r4. large'** as instance size & **don't create an Aurora Replica**



**DB instance size**

DB instance class  **Info**
Choose a DB instance class that meets your processing power and memory requirements. The DB instance class options below are limited to those supported by the engine you selected above.

⦿ Memory Optimized classes (includes r and x classes)

◯ Burstable classes (includes t classes)

db.r4.large
2 vCPUs   15.25 GiB RAM   EBS: 400 Mbps

⬤ Include previous generation classes

**Availability & durability**

Multi-AZ deployment  **Info**

◯ Create an Aurora Replica/Reader node in a different AZ (recommended for scaled availability)
Creates an Aurora replica for fast failover and high availability.

⦿ Don't create an Aurora Replica

6. Scroll down and Select **Default VPC**, Your Existing **VPC Security group, and** Click on **Create database**

## Connectivity

**Virtual Private Cloud (VPC)** Info
VPC that defines the virtual networking environment for this DB cluster.

Default VPC (vpc-484b8d20) ▼

Only VPCs with a corresponding DB subnet group are listed.

ⓘ After a database is created, you can't change the VPC selection.

▼ **Additional connectivity configuration**

**Subnet group** Info
DB subnet group that defines which subnets and IP ranges the DB instance can use in the VPC you selected.

default ▼

**Publicly accessible** Info

🔘 Yes
Amazon EC2 instances and devices outside the VPC can connect to your database. Choose one or more VPC security groups that specify which EC2 instances and devices inside the VPC can connect to the database.

⚪ No
RDS will not assign a public IP address to the database. Only Amazon EC2 instances and devices inside the VPC can connect to your database.

**VPC security group**
Choose one or more RDS security groups to allow access to your database. Ensure that the security group rules allow incoming traffic from EC2 instances and devices outside your VPC. (Security groups are required for publicly accessible databases.)

🔘 Choose existing
Choose existing VPC security groups

⚪ Create new
Create new VPC security group

Existing VPC security groups

Choose VPC security groups ▼

launch-wizard-4 ✕

**Availability zone** Info

No preference ▼

**Database port** Info
TCP/IP port the database will use for application connections.

5432

▶ **Additional configuration**
Database options, encryption enabled, failover, backup enabled, backtrack disabled, Performance Insights enabled, Enhanced Monitoring enabled, maintenance, CloudWatch Logs, delete protection disabled

Cancel     **Create database**

7. It takes few moments to create the database, Once created, you get a success/Available message



8. Once your database status is **available** click on writer instance. Make a note of Endpoint/port (ex: northwind-instance-1.cv0pnsuqer81.us-east-2.rds.amazonaws.com/ 5432) & Click on VPC security Group (e.g. Launch-wizard-4)



9. Select **Inbound** tab. Click **Edit**.

10. Add new rule as shown below by clicking **Add Rule**. Once added, Click **Save** & Confirm the addition



11. Open DBeaver (client) installed on your laptop & Search for PostgreSQL, Select & click **Next**.

12. Provide hostname, port, database, username & password and click on Test Connection.



13. Connection confirmation popup will appear. Click on Next > Next > Finish

Note PostgreSQL installation is now complete. To generate ER diagram for a schema



## CREATE REDSHIFT CLUSTER IN AWS

For creating Redshift cluster and querying the underlying dataset from S3, we will need to follow the below mentioned steps

A. Add access roles to access Redshift
B. Add IAM policy for Redshift query editor

C. Create Redshift cluster
D. Execute sample query using query editor for verifying the install

A. Add access roles to access Redshift

1. Go to AWS services and click on IAM Console



2. Once in IAM dashboard, In the navigation pane, choose **Roles**



3. Choose **Create Role**

4. Choose **AWS service**, and then choose **Redshift.** Under **Select your use case**, choose **Redshift - Customizable** and then choose **Next: Permissions**.



5. The **Attach permissions policy** page appears. Search for S3 policies & select AmazonDMSRedshiftS3Role & AmazonS3FullAccess & then Click **Next: Tags**

6. The **Add tags** page appears. You can optionally add tags. Choose **Next: Review.**

7. For **Role name**, type a name for your role, for example **RedshiftDWS3FullAccess.** Choose **Create role**.



8. Confirm, New Role is successfully added.



B. Add IAM policy for Redshift Query Editor

1. In the navigation menu, click on policies and Search for **"**AmazonRedshiftQueryEditor**"** policy

2. Select "AmazonRedshiftQueryEditor" policy & Click on "**Policy actions**" & select **Attach**



3. In the Attach policy window, select the user ids & click **Attach policy**

4. Confirm the policy attachment to user



## C. Create Redshift Cluster

1. Go to AWS services and click on Amazon Redshift (Datawarehouse)

2. Click on **Quick launch cluster**



3. Provide the details as shown below. Note that this will only add 1 node/machine in the cluster. If you need more than one update the number of nodes accordingly. Select the role that was added in previous steps & click **Launch cluster.** Choose password as Root1234 or something simple.
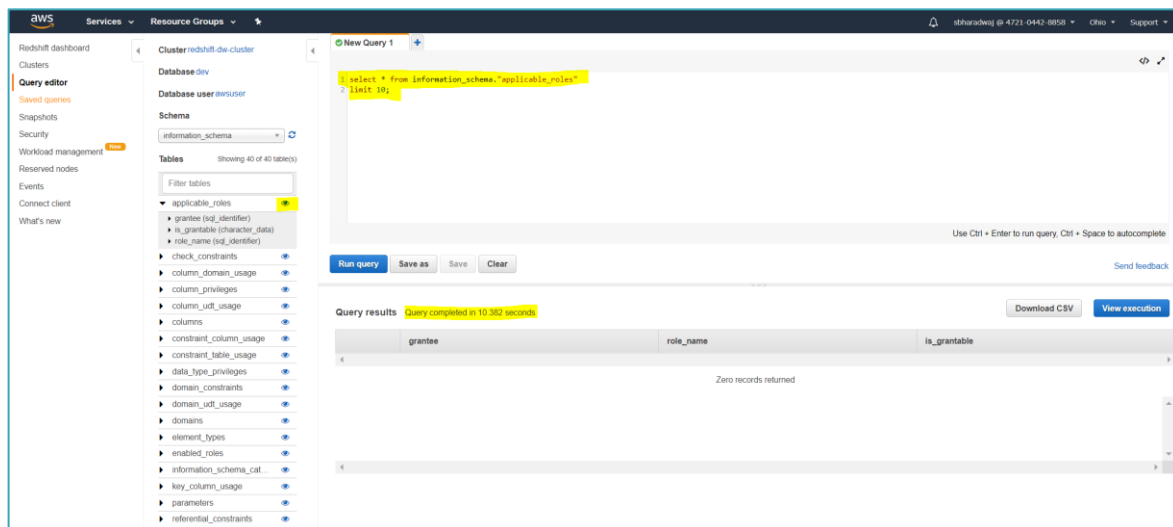
4. Confirm the cluster has been launched successfully.

## D. Query using Query Editor

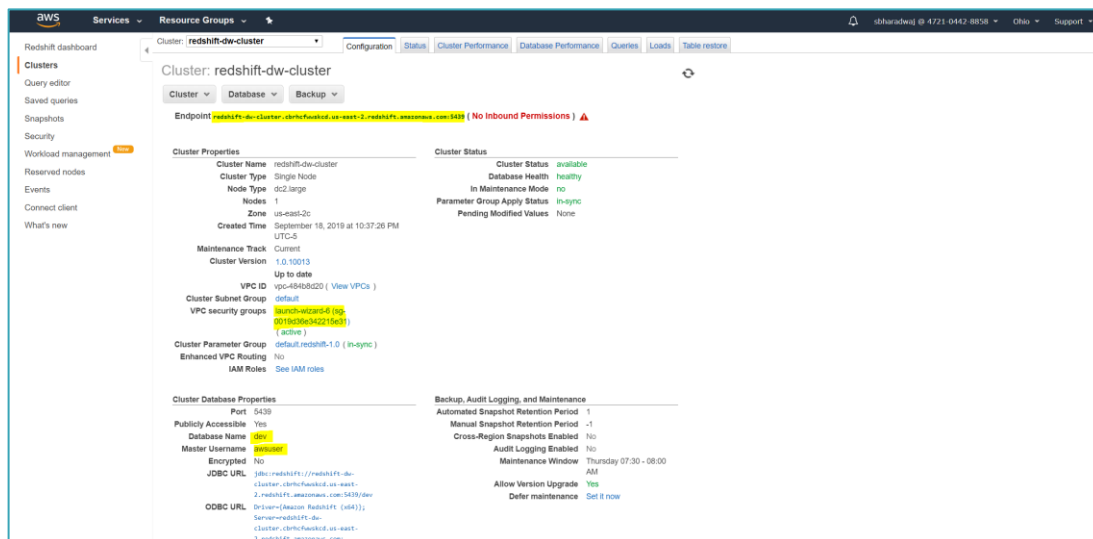1. Once the cluster gets created, click on **Open Query Editor** from the navigation menu
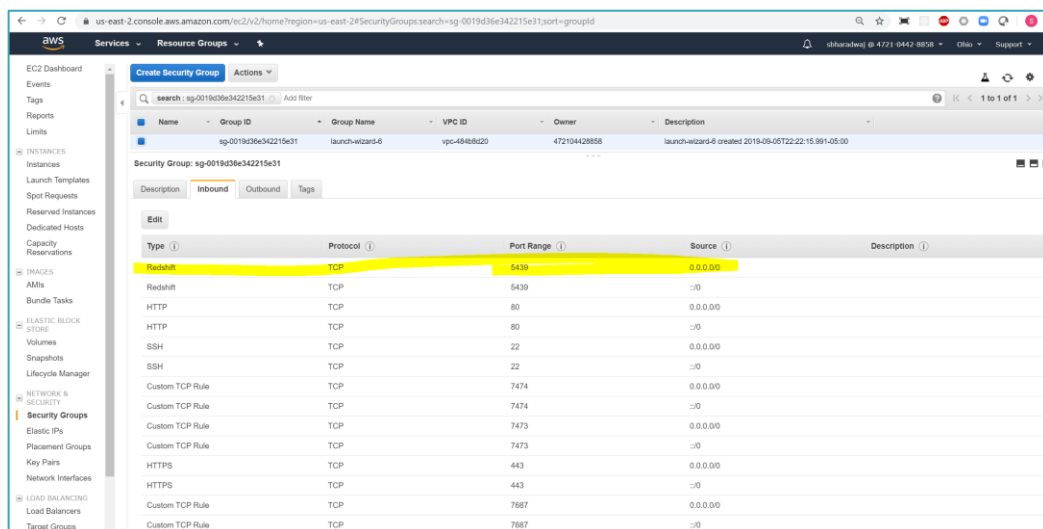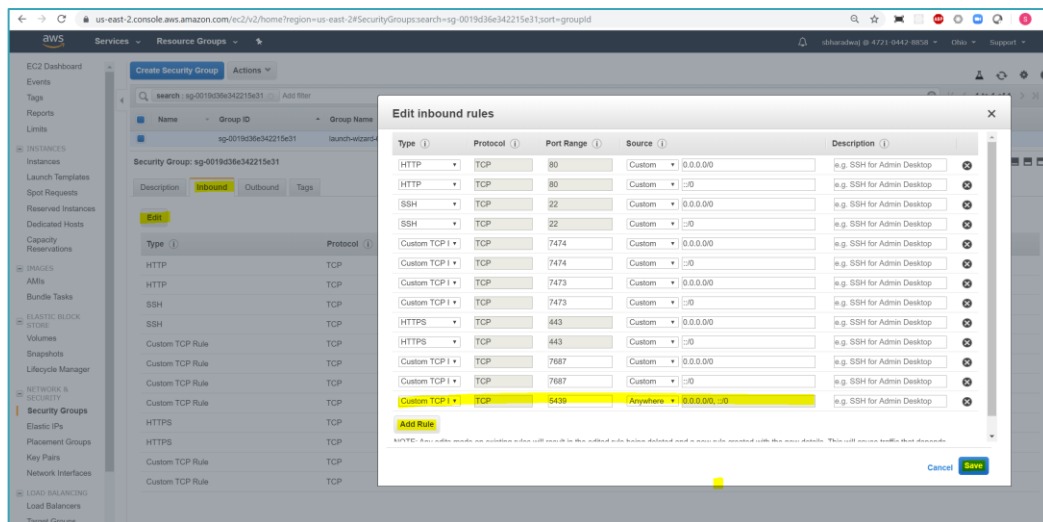
2. Run sample query

Note: Amazon Redshift installation is now complete.

CONNECT TABLEAU TO REDSHIFT

1. Login to AWS console and have 5439 port added to your VPC security group. Copy the Endpoint {redshift-dw-cluster.cbrhcfwwskcd.us-east-2.redshift.amazonaws.com} which will be added to Tableau as server. Make a note of the Database Name and Master Username

2. Click on the VPC Security Group and navigate to the Inbound tab and click on Edit and Add Rule and click Save

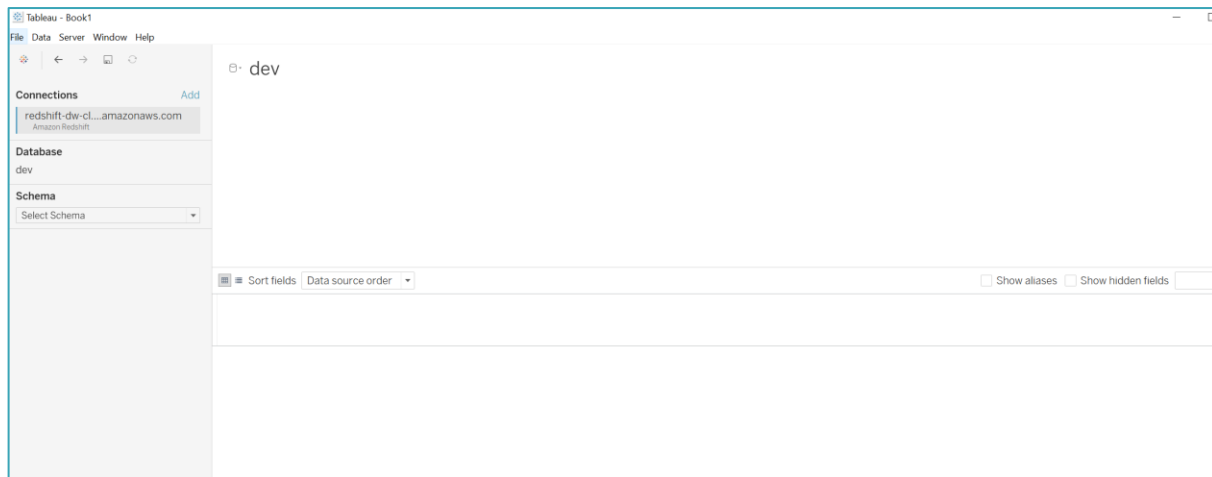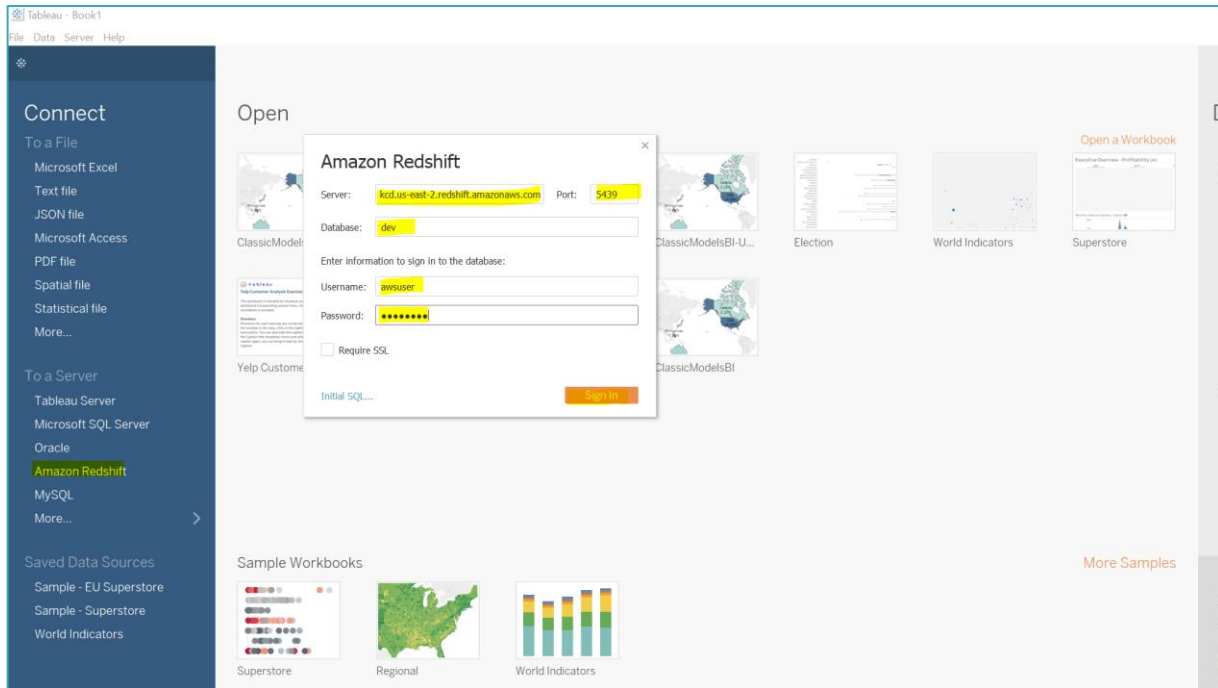3. Open Tableau click on 'Amazon Redshift' under server options. Add Server details to connect to AWS Redshift instance.





Tableau is now connected to Amazon Redshift.

End of Exercise 2.