# Advanced Data Engineering in Cloud
## ASSIGNMENT-3

**Hem Prakash Dev**                                    **Roll-G23AI1054**

1. **Configure AWS Redshift or Amazon Athena to aggregate the processed data and enable fast querying and analysis.**

Configure AWS Redshift
**Create a Redshift Cluster:**

aws  Services  Q Search  [Alt+S]  Ohio ▾  Hem_Prakash_Dev ▾

Amazon Redshift > Clusters > my-redshift-cluster

# my-redshift-cluster

Actions ▾ | Edit | Add partner integration | Query data ▾

## General information  Info

| | | | |
|---|---|---|---|
| **Cluster identifier** | **Status** | **Node type** | **Endpoint** |
| my-redshift-cluster | ✓ Available | ra3.4xlarge | my-redshift-cluster.cr5987udnnd3.us-east-2.redshift.amazonaws.com:5439/dev |
| **Custom domain name** | **Date created** | **Number of nodes** | |
| - | July 17, 2024, 00:07 (UTC+05:30) | 2 | **JDBC URL** |
| | | | jdbc:redshift://my-redshift-cluster.cr5987udnnd3.us-east-2.redshift.amazonaws.com:5439/dev |
| **Cluster namespace ARN** | **Storage used** | **Patch version** | |
| arn:aws:redshift:us-east-2:339713036364:namespace:ba669921-6189-47f9-be16-48fcea11dff4 | - | Patch 182 ⬚ | **ODBC URL** |
| | **Multi-AZ** | | Driver={Amazon Redshift (x64)}; Server=my-redshift-cluster.cr5987udnnd3.us-east- |
| **Cluster configuration** | No | | |
| Production | | | |

CloudShell  Feedback  © 2024, Amazon Web Services, Inc. or its affiliates.  Privacy  Terms  Cookie preferences

---

aws  Services  Q Search  [Alt+S]  Ohio ▾  Hem_Prakash_Dev ▾

Database=dev

< Query monitoring | Databases | Datashares | Zero-ETL integrations | Resource Policy | Schedules | Maintenance | **Properties** >

## Database configurations  Info

Edit admin credentials | Rotate encryption keys | Edit ▾

| | | | |
|---|---|---|---|
| **Database name** | **Parameter group** | **Encryption** | **Audit logging** |
| dev | Defines database parameter and query queues for all the databases. | Disabled | Disabled |
| **Port** | default.redshift-1.0 | **AWS KMS key ID** | |
| 5439 | | - | |
| | **SSH ingestion setting (cluster public key)** | | |
| **Admin user name** | ssh-rsa | | |
| awsuser | AAAAB3NzaC1yc2EAAAADAQABAAA BAQCd/2THZ1Nc3BSf8WArL3FlDhlH GfOBsESBPwZzkGmkjYpnkSlE4y8uU iirt0cuuvK1g2x/jUSZDMJ9DNxOp42 YuIRO9TFbe4FQJMGw88g1BgeYxT1 HvLWasvsRticeot8OMLWoip69ZSWR lwViXcpGh86Dv3O3MSbCt8Ki3aYyzs | | |

CloudShell  Feedback  © 2024, Amazon Web Services, Inc. or its affiliates.  Privacy  Terms  Cookie preferences

---

## Successfully runed the query on Redshift query

aws  Services  Q Search  [Alt+S]  Ohio ▾  Hem_Prakash_Dev ▾

Status  ✓ Connected  database  dev  user  awsuser  | Change connection

Query 1 +

Resources Info

**Select database** Info
To view schemas, select a database.
dev ▾

**Select schema** Info
To view tables, select a schema.
public ▾
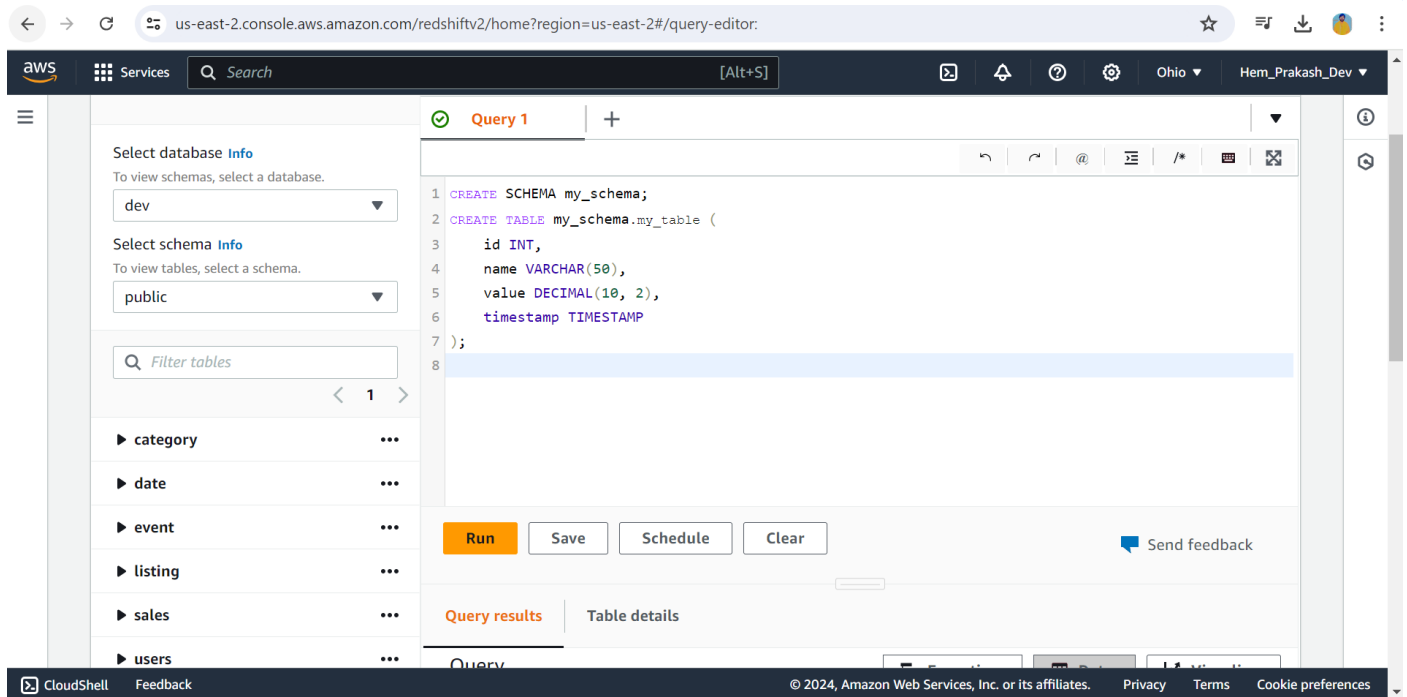
Q Filter tables
< 1 >

▶ category  •••
▶ date  •••
▶ event  •••
▶ listing  •••
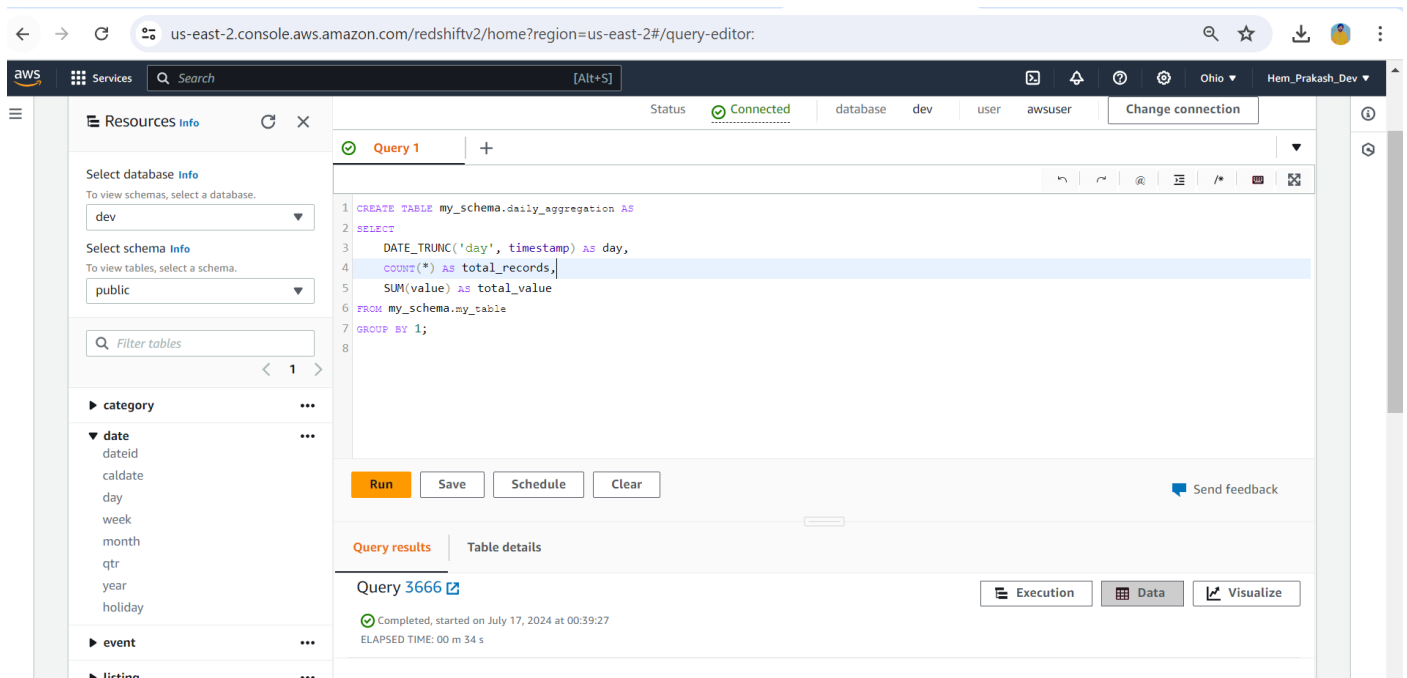
Run | Save | Schedule | Clear  | Send feedback

CloudShell  Feedback  © 2024, Amazon Web Services, Inc. or its affiliates.  Privacy  Terms  Cookie preferences

**2. Design and implement the data aggregation queries and optimize them for performance.**

**Example Aggregation Queries for Redshift**

Daily Aggregation:



Optimizing queries in a database involves creating appropriate indexes, using **SORTKEY** and **DISTKEY** in databases like Amazon Redshift, and performing regular maintenance like analyzing and vacuuming tables.

aws | Services | Search [Alt+S] | Ohio ▼ | Hem_Prakash_Dev ▼

Amazon Redshift > Query editor

Editor | Query history | Saved queries | Scheduled queries

Resources Info

Select database Info
To view schemas, select a database.
dev

Select schema Info
To view tables, select a schema.
public

Filter tables

◁ 1 ▷

▶ category ...
▼ date ...
  dateid
  caldate
  day
  week
  month
  qtr

Status ⊘ Connected | database dev | user awsuser | Change connection

⊗ Query 1 +

```
1  -- Drop the schema if it exists
2  DROP SCHEMA IF EXISTS my_schema CASCADE;
3
4  -- Create schema
5  CREATE SCHEMA my_schema;
6
7  -- Create table with DISTKEY and SORTKEY
8  CREATE TABLE my_schema.my_table (
9      id INT PRIMARY KEY,
10     name VARCHAR(50),
11     value DECIMAL(10, 2),
12     timestamp TIMESTAMP DEFAULT CURRENT_TIMESTAMP
13  )
```

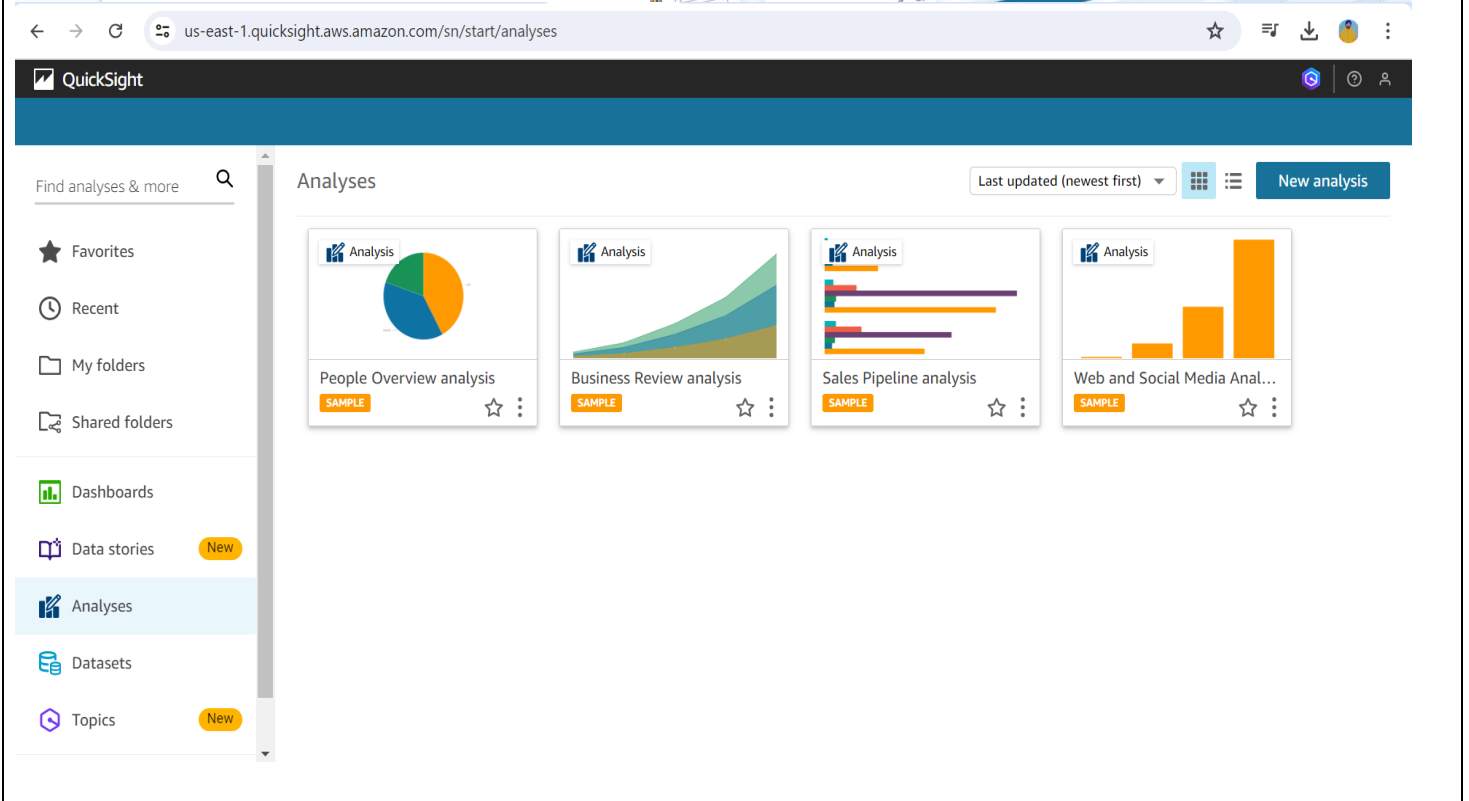Run | Save | Schedule | Clear | 💬 Send feedback

Query results | Table details

CloudShell | Feedback | © 2024, Amazon Web Services, Inc. or its affiliates. | Privacy | Terms | Cookie preferences

---

aws | Services | Search [Alt+S] | Ohio ▼ | Hem_Prakash_Dev ▼

Amazon Redshift > Query editor

Editor | Query history | Saved queries | Scheduled queries

Resources Info

Select database Info
To view schemas, select a database.
dev

Select schema Info
To view tables, select a schema.
public

Filter tables

◁ 1 ▷

▶ category ...
▼ date ...
  dateid
  caldate
  day
  week
  month
  qtr

Status ⊘ Connected | database dev | user awsuser | Change connection

⊗ Query 1 +

```
13  )
14  DISTKEY(id)
15  SORTKEY(timestamp);
16
17  -- Insert data into the table
18  INSERT INTO my_schema.my_table (id, name, value) VALUES
19  (1, 'Item A', 10.00),
20  (2, 'Item B', 20.00),
21  (3, 'Item C', 30.00);
22
23  -- Query the data
24  SELECT * FROM my_schema.my_table;
25
26  -- Update a record
```
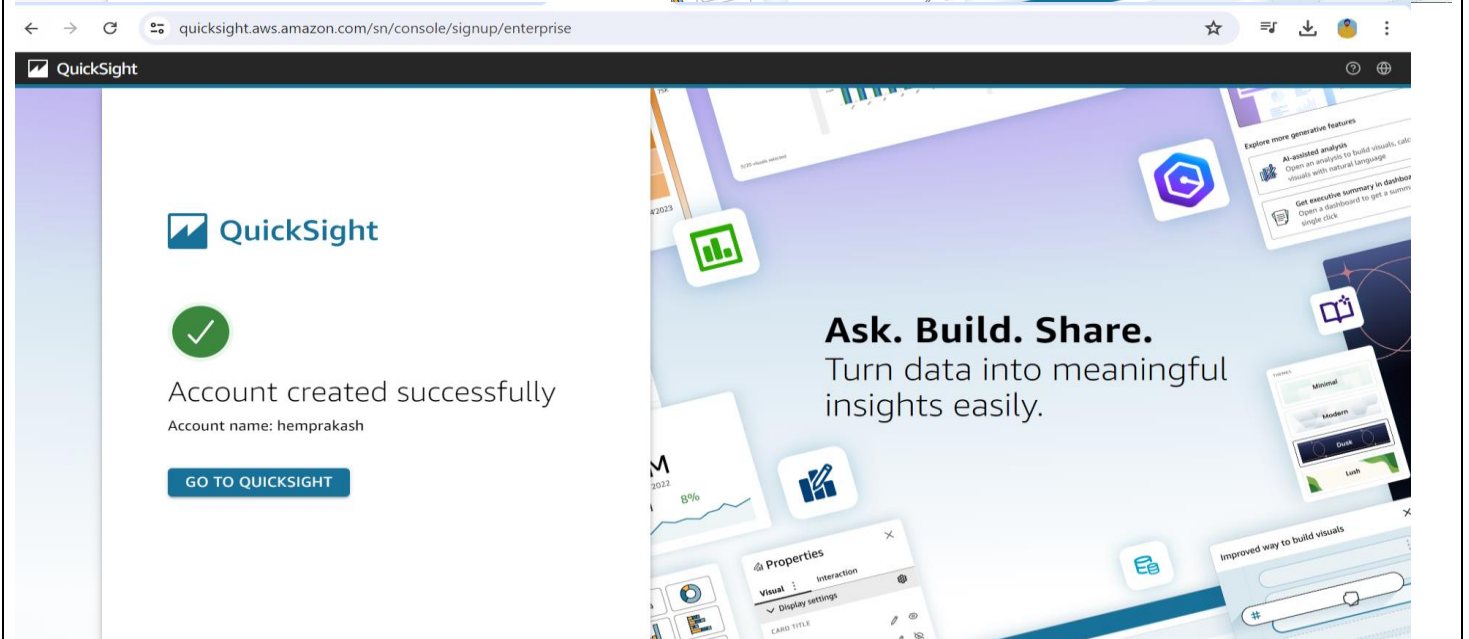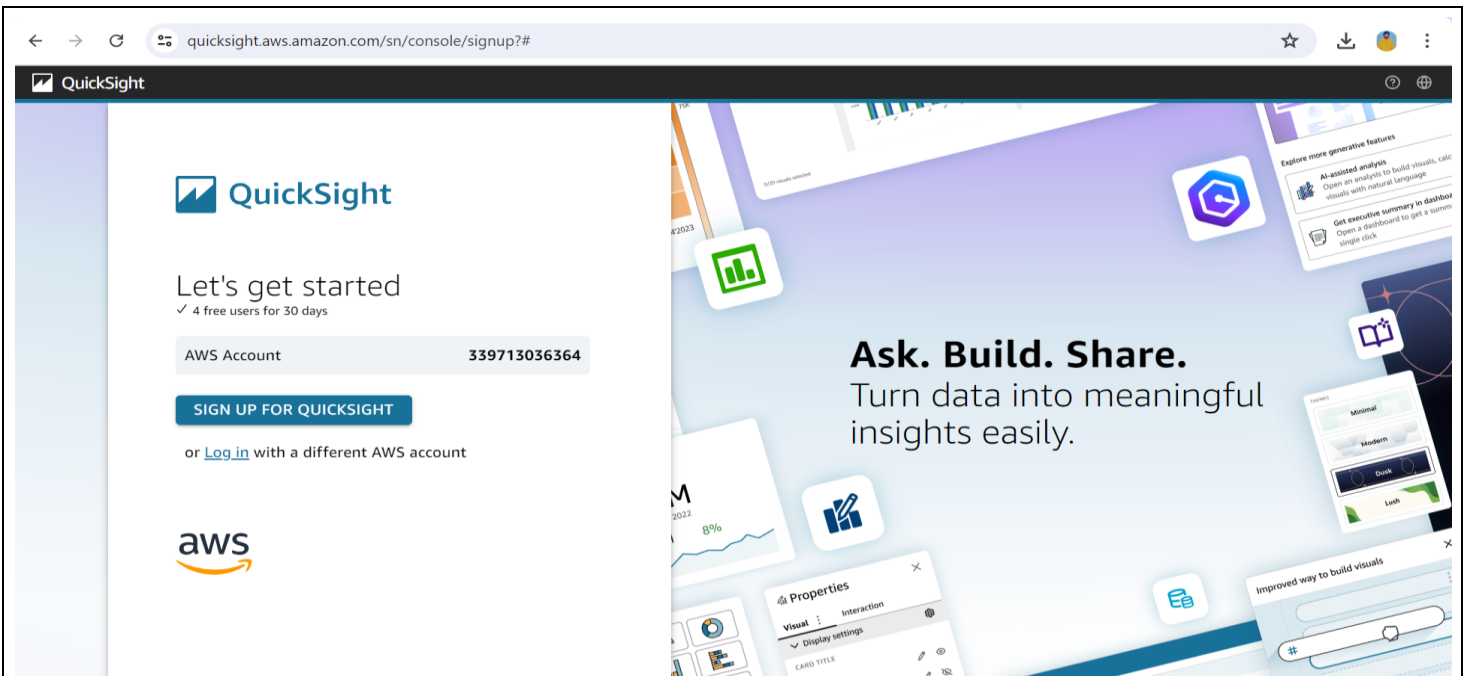
Run | Save | Schedule | Clear | 💬 Send feedback

Query results | Table details

CloudShell | Feedback | © 2024, Amazon Web Services, Inc. or its affiliates. | Privacy | Terms | Cookie preferences

aws Services | Q Search [Alt+S] | Ohio ▾ | Hem_Prakash_Dev ▾

Amazon Redshift > Query editor

**Editor** | Query history | Saved queries | Scheduled queries

Status ⊘ Connected | database dev | user awsuser | **Change connection**

⊗ Query 1 +

```
27  UPDATE my_schema.my_table
28  SET value = 25.00
29  WHERE id = 2;
30
31  -- Delete a record
32  DELETE FROM my_schema.my_table
33  WHERE id = 3;
34
35  -- Query the data again to see the changes
36  SELECT * FROM my_schema.my_table;
37
38  -- Generate and execute VACUUM and ANALYZE commands for each table in my_schema
39  DECLARE
```

**Run** | Save | Schedule | Clear | 💬 Send feedback

**Query results** | Table details

CloudShell | Feedback | © 2024, Amazon Web Services, Inc. or its affiliates. | Privacy | Terms | Cookie preferences

---

aws Services | Q Search [Alt+S] | Ohio ▾ | Hem_Prakash_Dev ▾

Amazon Redshift > Query editor

**Editor** | Query history | Saved queries | Scheduled queries

Status ⊘ Connected | database dev | user awsuser | **Change connection**

⊗ Query 1 +

```
40      sql_command TEXT;
41      table_name TEXT;
42
43  -- Cursor to fetch table names from information schema
44  FOR table_name IN
45      SELECT tablename
46      FROM pg_tables
47      WHERE schemaname = 'my_schema'
48  LOOP
49      -- Generate VACUUM ANALYZE command
50      sql_command := 'VACUUM ANALYZE my_schema.' || table_name;
51
52      -- Execute the command
```

**Run** | Save | Schedule | Clear | 💬 Send feedback

**Query results** | Table details

CloudShell | Feedback | © 2024, Amazon Web Services, Inc. or its affiliates. | Privacy | Terms | Cookie preferences

**3.  Integrate Amazon QuickSight or a third-party BI tool for data visualization.**

## QuickSight

### Let's get started
✓ 4 free users for 30 days

| AWS Account | 339713036364 |
|---|---|

**SIGN UP FOR QUICKSIGHT**

or Log in with a different AWS account

**aws**

# Ask. Build. Share.
Turn data into meaningful insights easily.

---

## QuickSight

✓

## Account created successfully
Account name: hemprakash

**GO TO QUICKSIGHT**

# Ask. Build. Share.
Turn data into meaningful insights easily.

---

## QuickSight

Find analyses & more

- ⭐ Favorites
- 🕐 Recent
- 📁 My folders
- 📁 Shared folders
- 📊 Dashboards
- 📖 Data stories  `New`
- 📊 Analyses
- 🗄 Datasets
- 🔵 Topics  `New`

### Analyses

Last updated (newest first) ▾   | **New analysis**

| Analysis | Analysis | Analysis | Analysis |
|---|---|---|---|
| People Overview analysis | Business Review analysis | Sales Pipeline analysis | Web and Social Media Anal... |
| SAMPLE | SAMPLE | SAMPLE | SAMPLE |

# Uploaded CSV file (amazon –bestseller-dataset) in amazon s3 bucket
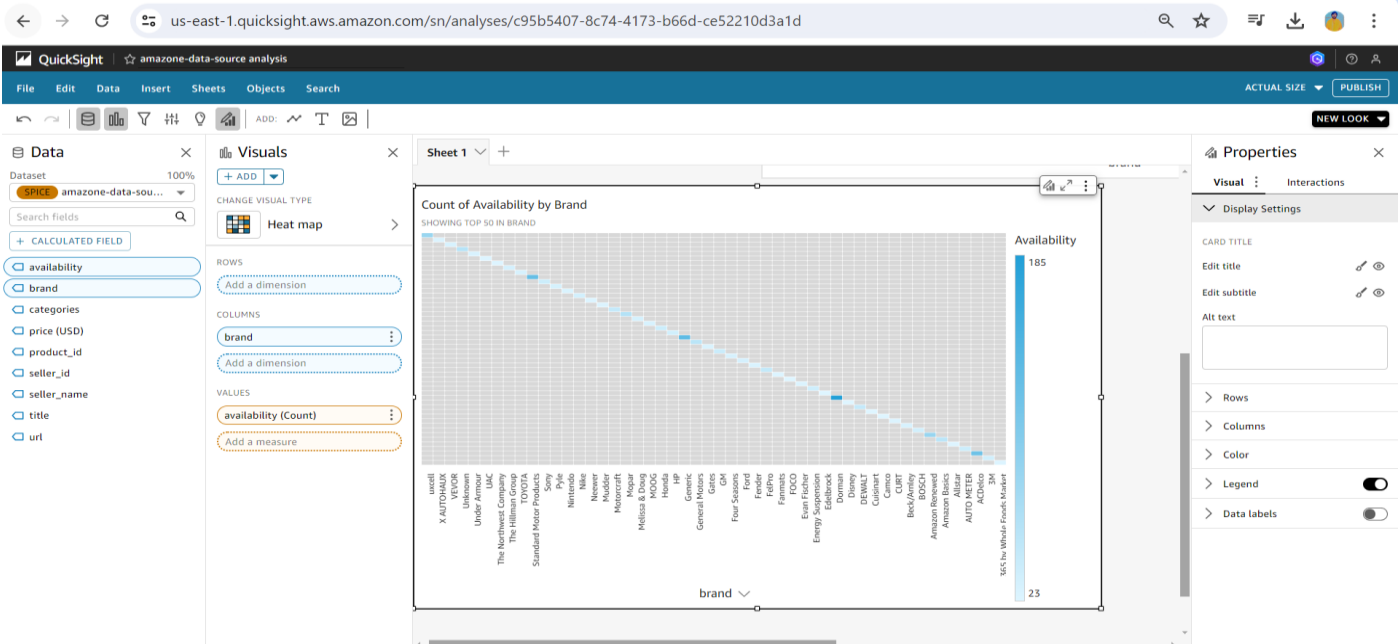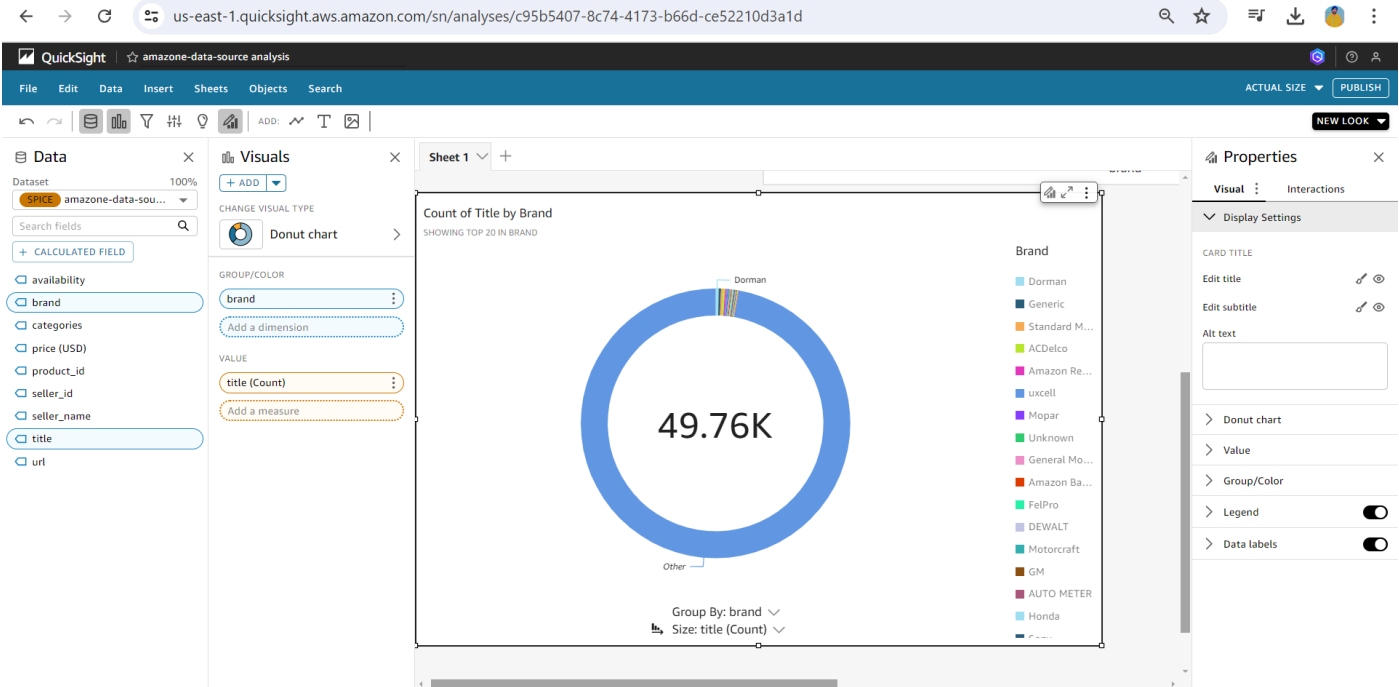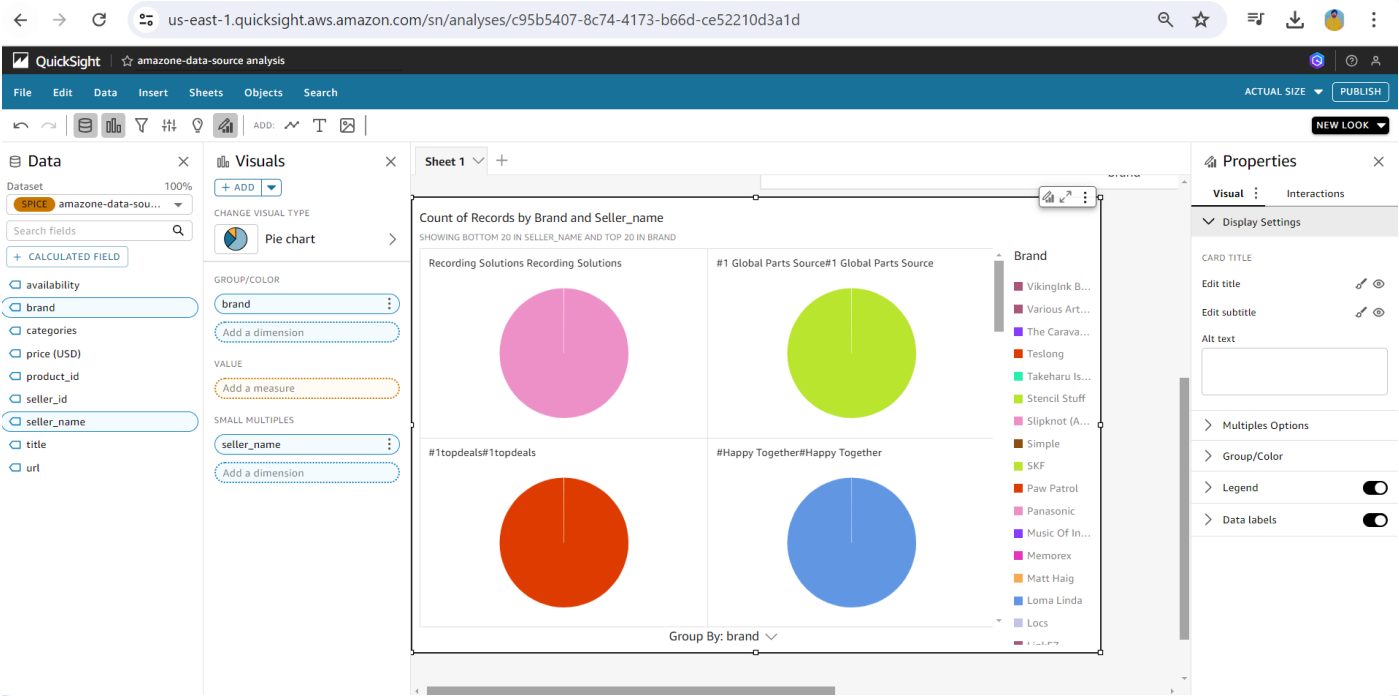
**4. Create sample dashboards and reports to showcase the insights derived from the aggregated data.**

5. **Update the GitHub repository with the code and configuration files for data aggregation and visualization**

  **git init**

  **git remote add origin https://github.com/Hem-Prakash-Dev-Bharadwaj/data-engineering**

  **git add .**

  **git commit -m "Initial commit with Redshift setup, aggregation queries, and QuickSight configuration"**

  **git push -u origin master**