

Name : Saurabh ShivPrakash Tripathi

Prn No. : 045

Q.1

Q.1) You have an S3 bucket student-data-lake containing multiple CSV files under the prefix logs/2025/11/. Each file has columns: user_id, login_time, region, device_type. Task: 1. Create an Athena external table to query CSV files directly from S3. 2. Write a query to return number of logins per region. 3. Save the query results to S3 folder athena-results

us-east-1.console.aws.amazon.com/s3/buckets/student-data-lake-045?region=us-east-1&prefix=logs/2025/11/&showversions=false

Amazon S3

General purpose buckets

Directory buckets

Table buckets

Vector buckets

Access Grants

Access Points (General Purpose Buckets, FSx file systems)

Access Points (Directory Buckets)

Object Lambda Access Points

Multi-Region Access Points

Batch Operations

IAM Access Analyzer for S3

Block Public Access settings for this account

Storage Lens

Dashboards

Storage Lens groups

AWS Organizations settings

11/

Objects

Properties

Objects (1)

Copy S3 URI

Copy URL

Download

Open

Delete

Actions

Create folder

Upload

Find objects by prefix

Name	Type	Last modified	Size	Storage class
s3_athena_logs.csv	csv	November 10, 2025, 08:53:47 (UTC+05:30)	2.3 KB	Standard

CloudShell

Feedback

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

28°C Sunny

us-east-1.console.aws.amazon.com/athena/home?region=us-east-1#/query-editor/history/7220ce8f-382a-496c-9e02-d99:0b5819c4

Amazon Athena

Query editor tabs

Data source

AwsDataCatalog

Catalogue

None

Database

athena_logs_database

Tables and views

Create

Filter tables and views

Tables (2)

athena_log

athena_logs

user_id string

login_time timestamp

region string

device_type string

Views (0)

```
2 CREATE EXTERNAL TABLE IF NOT EXISTS athena_logs_database.athena_log (  
3   user_id string,  
4   login_time timestamp,  
5   region string,  
6   device_type string  
7 )  
8 ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'  
9 WITH SERDEPROPERTIES (  
10   'separatorChar' = ',',  
11   'quoteChar' = '"',  
12   'escapeChar' = '\\',  
13   'skip.header.line.count' = '1'  
14 )  
15 STORED AS TEXTFILE  
16 LOCATION 's3://student-data-lake-045/logs/2025/11/'
```

SQL Ln 30, Col 22

Run again

Explain

Cancel

Clear

Create

Reuse query results up to 60 minutes ago

Query results

Query status

Completed

Time in queue: 123 ms

Run time: 539 ms

Data scanned: 2.25 KB

Results (4)

Copy

Download results CSV

Search rows

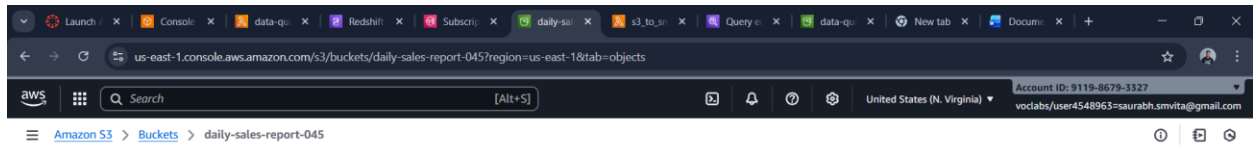
The image shows two screenshots from the AWS Management Console. The top screenshot is the Amazon Athena Query Editor. The left sidebar shows the 'Catalogue' with 'Database' set to 'athena_logs_database'. The 'Tables and views' section lists 'athena_log' and 'athena_logs'. The 'athena_logs' table has columns: 'user_id' (string), 'login_time' (timestamp), 'region' (string), and 'device_type' (string). The main editor shows a SQL query:

```
18 'classification'='csv',
19 'skip.header.line.count'='1',
20 'has_encrypted_data'='false'
21 );
22 SELECT * FROM athena_logs_database.athena_log LIMIT 10;
23
24 SELECT
25   region,
26   COUNT(*) AS logins
27 FROM athena_logs_database.athena_log
28 WHERE region IS NOT NULL
29 GROUP BY region
30 ORDER BY logins DESC;
```

 The query is executed, showing a 'Completed' status with 'Time in queue: 123 ms', 'Run time: 539 ms', and 'Data scanned: 2.25 KB'. The results are displayed in a table with columns 'region' and 'logins'. The bottom screenshot is the Amazon S3 console, showing the 'athena-results/' folder. The 'Objects' tab shows a single object 'athena-results/' of type 'Folder'. The 'Properties' tab is also visible. The console shows the path 'us-east-1.console.aws.amazon.com/s3/buckets/student-data-lake-045?region=us-east-1&prefix=athena-results/&showversions=false'.

Save the query results to S3 folder athena-results

Q.2



daily-sales-report-045 [Info](#)

[Objects](#) [Metadata](#) [Properties](#) [Permissions](#) [Metrics](#) [Management](#) [Access Points](#)

Objects (1) [Copy S3 URI](#) [Copy URL](#) [Download](#) [Open in new tab](#) [Delete](#) [Actions](#) [Create folder](#) [Upload](#)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	reports/	Folder	-	-	-

CloudShell Feedback

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

NIFTY +0.28%

us-east-1.console.aws.amazon.com/sns/v3/home?region=us-east-1#/topic/arn:aws:sns:us-east-1:911986793327:daily-report-alerts

Account ID: 9119-8679-3327 voclabs/user4548963=saurabh.smvita@gmail.com

[Amazon SNS](#) [Topics](#) [daily-report-alerts](#)

[Edit](#) [Delete](#) [Publish message](#)

daily-report-alerts

Details

Name daily-report-alerts	Display name -
ARN arn:aws:sns:us-east-1:911986793327:daily-report-alerts	Topic owner 911986793327
Type Standard	

[Subscriptions](#) [Access policy](#) [Data protection policy](#) [Delivery policy \(HTTP/S\)](#) [Delivery status logging](#) [Encryption](#) [Tags](#)

Subscriptions (1) [Edit](#) [Delete](#) [Request confirmation](#) [Confirm subscription](#) [Create subscription](#)

<input type="checkbox"/>	ID	Endpoint	Status	Protocol
<input type="radio"/>	a52c1596-2110-4d2b-806f-806e1...	saurabh6373@gmail.com	Confirmed	EMAIL

us-east-1.console.aws.amazon.com/lambda/home?region=us-east-1#/functions/s3_to_sns_metadata?subtab=triggers&tab=code

Account ID: 9119-8679-3327
voclabs/user4548963=saurabh.smvita@gmail.com

Explorer: s3_TO_SNS_METADATA, lambda_function.py

Deploy: Deploy (Ctrl+Shift+U), Test (Ctrl+Shift+I)

TEST EVENTS (NONE SELECTED)
+ Create new test event

ENVIRONMENT VARIABLES

```
1 import json
2 import boto3
3 import logging
4
5 logger = logging.getLogger()
6 logger.setLevel(logging.INFO)
7
8 sns = boto3.client('sns')
9 s3 = boto3.client('s3')
10
11 # Replace with your SNS topic ARN
12 SNS_TOPIC_ARN = 'arn:aws:sns:us-east-1:YOUR_ACCOUNT_ID:daily-report-alerts'
13
14 def lambda_handler(event, context):
15     try:
16         record = event['Records'][0]
17         bucket = record['s3']['bucket']['name']
18         key = record['s3']['object']['key']
19
20         response = s3.head_object(Bucket=bucket, Key=key)
21         size = response['contentLength']
22         upload_time = response['LastModified'].strftime('%Y-%m-%d %H:%M:%S')
23
24         message = {
25             'Filename': key,
26             'Bucket': bucket,
27             'Size': response['contentLength']
```

us-east-1.console.aws.amazon.com/cloudwatch/home?region=us-east-1#/logs/groups/log-group/\$252faws\$252flambda\$252f\$3_to_sns_metadata/log-events/2025\$252f11\$252f10\$252f\$2558\$...

Account ID: 9119-8679-3327
voclabs/user4548963=saurabh.smvita@gmail.com

CloudWatch > Log groups > /aws/lambda/s3_to_sns_metadata > 2025/11/10/[LATEST]af5faf39b20f4f9991defee0c42e7593

Log events

You can use the filter bar below to search for and match terms, phrases, or values in your log events. [Learn more about filter patterns](#)

Filter events - press enter to search

Clear 1m 30m 1h 12h Custom Local timezone

Display

Timestamp	Message
No older events at this moment. Retry	
2025-11-10T09:44:33.368+05:30	INIT_START Runtime Version: python:3.12.v93 Runtime Version ARN: arn:aws:lambda:us-east-1::runtime:c010741a32d362a2a9a45adf0a7418c...
2025-11-10T09:44:33.666+05:30	[INFO] 2025-11-10T04:14:33.666Z Found credentials in environment variables.
2025-11-10T09:44:33.873+05:30	START RequestId: f810ac77-6b3b-4528-a1e0-75e9300a3415 Version: \$LATEST
2025-11-10T09:44:34.383+05:30	[ERROR] 2025-11-10T04:14:34.383Z f810ac77-6b3b-4528-a1e0-75e9300a3415 Error processing file: An error occurred (InvalidParameter) w...
2025-11-10T09:44:34.423+05:30	[ERROR] InvalidParameterException: An error occurred (InvalidParameter) when calling the Publish operation: Invalid parameter: Topi...
2025-11-10T09:44:34.443+05:30	END RequestId: f810ac77-6b3b-4528-a1e0-75e9300a3415
2025-11-10T09:44:34.443+05:30	REPORT RequestId: f810ac77-6b3b-4528-a1e0-75e9300a3415 Duration: 569.41 ms Billed Duration: 1071 ms Memory Size: 128 MB Max Memory ...
2025-11-10T09:45:33.612+05:30	START RequestId: f810ac77-6b3b-4528-a1e0-75e9300a3415 Version: \$LATEST
2025-11-10T09:45:34.082+05:30	[ERROR] 2025-11-10T04:15:34.082Z f810ac77-6b3b-4528-a1e0-75e9300a3415 Error processing file: An error occurred (InvalidParameter) w...
2025-11-10T09:45:34.082+05:30	[ERROR] InvalidParameterException: An error occurred (InvalidParameter) when calling the Publish operation: Invalid parameter: Topi...

Q.3

The screenshot shows the AWS S3 console interface. The breadcrumb navigation is **Amazon S3 > Buckets > data-quality-alert-045 > incoming/**. The left sidebar shows the **Amazon S3** menu with options like **General purpose buckets**, **Directory buckets**, **Table buckets**, **Vector buckets**, **Access Grants**, **Access Points (General Purpose Buckets, FSx file systems)**, **Access Points (Directory Buckets)**, **Object Lambda Access Points**, **Multi-Region Access Points**, **Batch Operations**, and **IAM Access Analyzer for S3**. The main content area shows the **incoming/** bucket with a **Copy S3 URI** button. Below the bucket name, there are tabs for **Objects** and **Properties**. The **Objects** tab shows a list of objects with columns: **Name**, **Type**, **Last modified**, **Size**, and **Storage class**. The list contains two items: **pipeline_data_quality.csv** (Type: csv, Last modified: November 10, 2025, 10:14:39 (UTC+05:30), Size: 1.4 KB, Storage class: Standard) and **Unsaved/** (Type: Folder, Last modified: -, Size: -, Storage class: -). The **pipeline_data_quality.csv** file is selected.

The screenshot shows the AWS Athena Query Editor interface. The breadcrumb navigation is **Amazon Athena > Query editor tabs**. The left sidebar shows the **Data** section with **Data source** (AwsDataCatalog), **Catalogue** (None), **Database** (default), and **Tables and views** (Create button). The main content area shows a SQL query in the **Editor** tab. The query is:

```
1 CREATE EXTERNAL TABLE IF NOT EXISTS order_data (
2   order_id STRING,
3   product_id STRING,
4   quantity INT,
5   price DOUBLE,
6   order_date STRING
7 )
8 ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.lazy.SimpleSerDe'
9 WITH SERDEPROPERTIES (
10   'serialization.format' = ',',
11   'field.delim' = ','
12 )
13 LOCATION 's3://data-quality-check-bucket/incoming/'
14 TBLPROPERTIES ('skip.header.line.count'='1');
```

 The query is labeled **Query 10**. Below the query, there are buttons for **Run again**, **Explain**, **Cancel**, **Clear**, and **Create**. The **Run again** button is highlighted. The **Query results** tab is selected, showing a green bar indicating the query is **Completed**. The bottom status bar shows **Time in query: 01 ms**, **Run Memory: 420 ms**, and **Data scanned: 0 B**.

us-east-1.console.aws.amazon.com/sns/v3/home?region=us-east-1#/subscription/arn:aws:sns:us-east-1:911986793327:data-quality-alerts:0f91e419-23f9-4371-bbca-b94161ce8ff0

Amazon SNS > Topics > data-quality-alerts > Subscription: 0f91e419-23f9-4371-bbca-b94161ce8ff0

Subscription: 0f91e419-23f9-4371-bbca-b94161ce8ff0

[Edit](#) [Delete](#)

Details

ARN
arn:aws:sns:us-east-1:911986793327:data-quality-alerts:0f91e419-23f9-4371-bbca-b94161ce8ff0

Endpoint
saurabh6373@gmail.com

Topic
[data-quality-alerts](#)

Subscription Principal
arn:aws:iam::911986793327:role/voclabs

Status
Pending confirmation

Protocol
EMAIL

[Subscription filter policy](#) [Redrive policy \(dead-letter queue\)](#)

Subscription filter policy [Info](#)

This policy filters the messages that a subscriber receives.

No filter policy configured for this subscription.
To apply a filter policy, edit this subscription.

CloudShell Feedback

Air: Moderate Now

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

ENG IN 10:25 10-11-2025

us-east-1.console.aws.amazon.com/lambda/home?region=us-east-1#/functions/data-quality-check-lambda?newfunction=true&tab=code

AWS > Lambda > Functions > data-quality-check-lambda

data-quality-check-lambda

EXPLORER

- DATA-QUALITY-CHECK-LAMB...
- lambda_function.py

DEPLOY [UNDEPLOYED CHANGES]

Deploy (Ctrl+Shift+U)

Test (Ctrl+Shift+T)

TEST EVENTS [NONE SELECTED]

Create new test event

ENVIRONMENT VARIABLES

Amazon Q

```
def lambda_handler(event, context):  
    # Log metric  
    cloudwatch.put_metric_data(  
        Namespace='DataQuality',  
        MetricData=[  
            {'MetricName': 'BadRecordCount',  
             'Value': bad_count,  
             'Unit': 'Count'}  
        ])  
  
    # Send SNS alert  
    if bad_count > 0:  
        sns.publish(  
            TopicArn=SNES_TOPIC_ARN,  
            Subject='Data Quality Alert',  
            Message=f'{bad_count} bad records found in incoming data.'  
        )  
  
    return {'status': 'done', 'bad_count': bad_count}
```

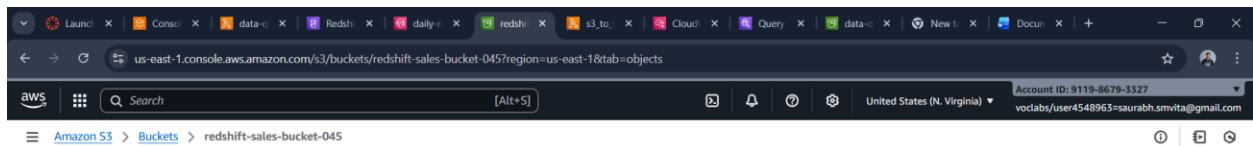
Ln 53, Col 1 Spaces: 4 UTF-8 LF Python Lambda Layout: US

CloudShell Feedback

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

ENG IN 10:27 10-11-2025

Q.4



redshift-sales-bucket-045 Info

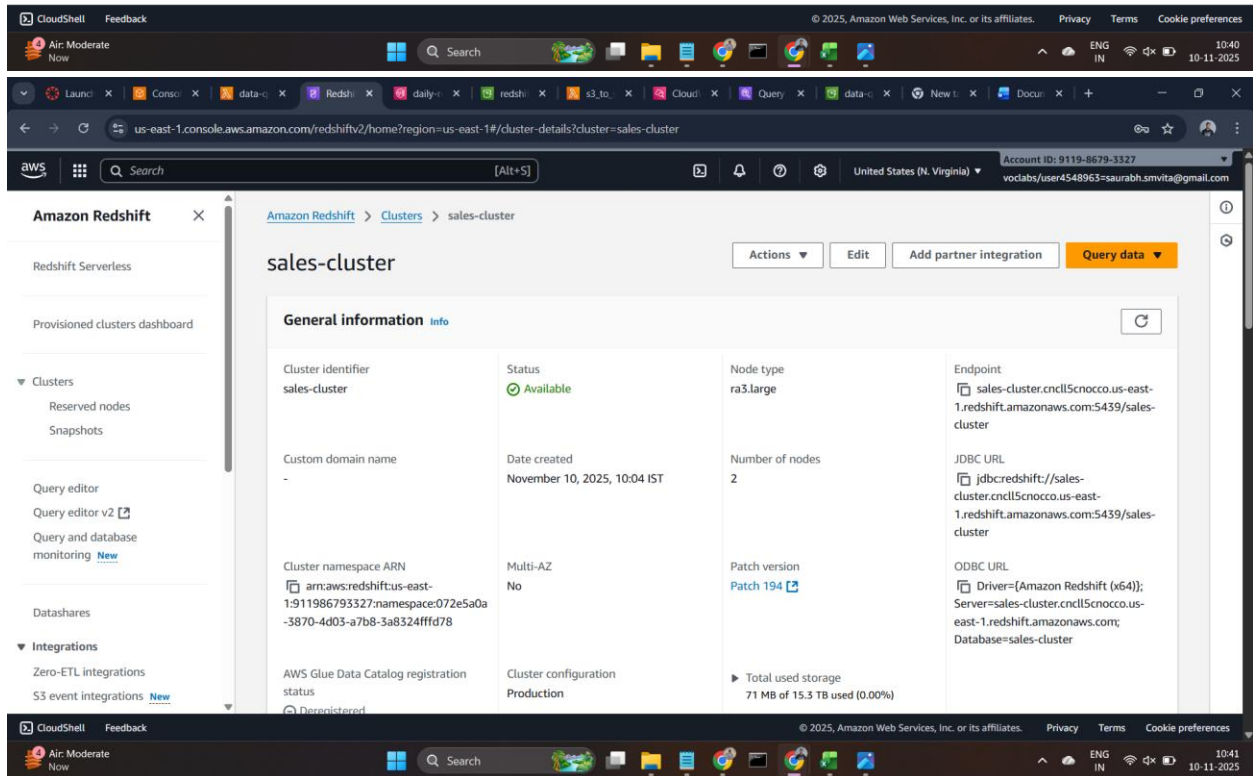
[Objects](#) [Metadata](#) [Properties](#) [Permissions](#) [Metrics](#) [Management](#) [Access Points](#)

Objects (1)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Find objects by prefix

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	redshift_sales_data.csv	csv	November 10, 2025, 10:11:43 (UTC+05:30)	2.6 KB	Standard



Launch

Conso

data

Redsh

daily

redsh

s3_to

Cloud

Query

data

New t

Docu

us-east-1.console.aws.amazon.com/redshiftv2/home?region=us-east-1#/query-editor

aws

Search

[Alt+S]

United States (N. Virginia)

Account ID: 9119-8679-3327
voclabr/user4548963=saurabh.smvita@gmail.com

Amazon Redshift

Redshift Serverless

Provisioned clusters dashboard

Clusters

Reserved nodes

Snapshots

Query editor

Query editor v2

Query and database monitoring

Datashares

Integrations

Zero-ETL integrations

S3 event integrations

Editor

Query history

Saved queries

Scheduled queries

Resources info

Select database info

sales-cluster

Select schema info

public

Filter tables

fact_sales

sale_id

store_id

product_id

quantity

unit_price

sale_timestamp

channel

Status Connected

database

sales-cluster

user

admin

Change connection

Query 1

CREATE TABLE public.fact_sales (
sale_id INT,
store_id INT,
product_id INT,
quantity INT,
unit_price DECIMAL(10,2),
sale_timestamp TIMESTAMP,
channel VARCHAR(10)
)
DISTKEY(store_id)
SORTKEY(sale_timestamp);

Run

Save

Schedule

Clear

Send feedback

Query results

Table details

CloudShell

Feedback

© 2025, Amazon Web Services, Inc. or its affiliates.

Privacy

Terms

Cookie preferences

Air: Moderate
Now

Search

10:41
10-11-2025