# Assignment 5.4 Data Lake

## Peer Members:

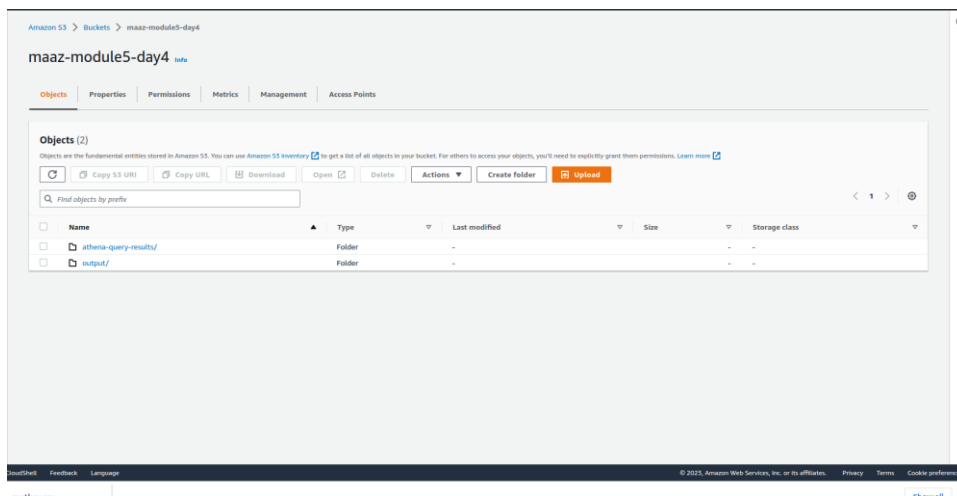- Syed Muhammad Raqim Ali Shah (2303.KHI.DEG.008)
- Maaz Javaid Siddique (2303.KHI.DEG.004)
- Qadeer Hussain (2303.KHI.DEG.006)

Answer:

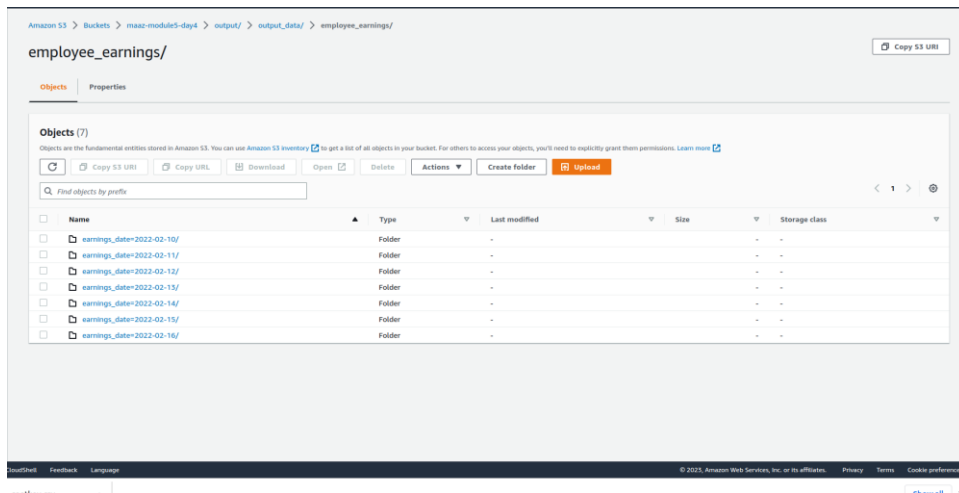S3 Bucket:



Folders create in a bucket:



Data we add in the output folder:

## employee_earnings/

Copy S3 URI

**Objects** | Properties

### Objects (7)

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

| | Name | Type | Last modified | Size | Storage class |
|---|---|---|---|---|---|
| ☐ | 🗀 earnings_date=2022-02-10/ | Folder | - | - | - |
| ☐ | 🗀 earnings_date=2022-02-11/ | Folder | - | - | - |
| ☐ | 🗀 earnings_date=2022-02-12/ | Folder | - | - | - |
| ☐ | 🗀 earnings_date=2022-02-13/ | Folder | - | - | - |
| ☐ | 🗀 earnings_date=2022-02-14/ | Folder | - | - | - |
| ☐ | 🗀 earnings_date=2022-02-15/ | Folder | - | - | - |
| ☐ | 🗀 earnings_date=2022-02-16/ | Folder | - | - | - |

Here i create a crawler:

### Crawlers

A crawler connects to a data store, progresses through a prioritized list of classifiers to determine the schema for your data, and then creates metadata tables in your data catalog.

**Crawlers (4)** Info
View and manage all available crawlers.

| | Name | State | Schedule | Last run | Last run timestamp | Log | Table changes from last... |
|---|---|---|---|---|---|---|---|
| ☐ | maaz_combined_employ... | ⊘ Ready | | ⊘ Succeeded | May 19, 2023 at 04:00:22 | View log | 1 updated |
| ☐ | maazjavaid_rds_employe... | ⊘ Ready | | ⊗ Failed | May 17, 2023 at 05:11:51 | View log | - |
| ☐ | maazjavaid_s3_earnings_... | ⊘ Ready | | ⊘ Succeeded | May 17, 2023 at 05:11:51 | View log | 1 updated |
| ☐ | maazjavaid_s3_office_loc... | ⊘ Ready | | ⊘ Succeeded | May 17, 2023 at 05:29:01 | View log | 1 created |

Athena Query editor:

```
1   WITH earnings_table AS (
2       SELECT
3           emp_id,
4           earnings_date,
5           earnings,
6           LAG(earnings) OVER (PARTITION BY emp_id ORDER BY earnings_date) AS previous_earnings,
7           LAG(earnings_date) OVER (PARTITION BY emp_id ORDER BY earnings_date) AS previous_earnings_date
8       FROM
9           "day4_glue_database"."employee_earnings"
10  )
11  SELECT
12      emp_id,
13      earnings_date,
14      previous_earnings_date,
15      earnings,
16      previous_earnings,
17      (earnings - previous_earnings) / CAST(previous_earnings AS double) * 100 AS percentage_change
18  FROM
19      earnings_table
20  WHERE
21      earnings_date = '2022-02-15';
```

Output Table:

**Results** (100)                                                      🗇 Copy          **Download results**

🔍 Search rows                                                                     ‹  1  …  ›  ⚙

| # ▽ | emp_id ▽ | earnings_date ▽ | previous_earnings_date ▽ | earnings ▽ | previous_earnings ▽ | percentage_change ▽ |
|---|---|---|---|---|---|---|
| 1 | 220965 | 2022-02-15 | 2022-02-14 | 3485 | 9378 | -62.8385583280017 |
| 2 | 235295 | 2022-02-15 | 2022-02-14 | 5082 | 5760 | -11.770833333333334 |
| 3 | 312726 | 2022-02-15 | 2022-02-14 | 3725 | 6055 | -38.48059454995871 |
| 4 | 314661 | 2022-02-15 | 2022-02-14 | 3289 | 8483 | -61.22833903100319 |
| 5 | 316372 | 2022-02-15 | 2022-02-14 | 5601 | 6686 | -16.227938976966797 |
| 6 | 366431 | 2022-02-15 | 2022-02-14 | 3883 | 9018 | -56.94167221113329 |
| 7 | 403534 | 2022-02-15 | 2022-02-14 | 5596 | 5530 | 1.193490054249548 |
| 8 | 549389 | 2022-02-15 | 2022-02-14 | 4317 | 7944 | -45.657099697885194 |
| 9 | 597741 | 2022-02-15 | 2022-02-14 | 4987 | 9094 | -45.16164504068617 |
| 10 | 622405 | 2022-02-15 | 2022-02-14 | 4142 | 3974 | 4.227478610971314 |
| 11 | 728053 | 2022-02-15 | 2022-02-14 | 3970 | 3603 | 10.185956147654732 |
| 12 | 819367 | 2022-02-15 | 2022-02-14 | 5080 | 8535 | -40.480374926772114 |
| 13 | 878666 | 2022-02-15 | 2022-02-14 | 5644 | 7755 | -27.22114764667956 |
| 14 | 885395 | 2022-02-15 | 2022-02-14 | 5089 | 4478 | 13.644484144707459 |
| 15 | 896517 | 2022-02-15 | 2022-02-14 | 4220 | 2057 | 105.15313563441906 |
| 16 | 936158 | 2022-02-15 | 2022-02-14 | 5123 | 9493 | -46.03391973032761 |
| 17 | 962291 | 2022-02-15 | 2022-02-14 | 5605 | 3286 | 70.57212416311624 |

Rerun queries from Task 3 and Task 4 and see how the results change with this new data.



Query is not run in the S3 bucket because s3 is a file storage service it's not a query editor. It accepts a very basic query. If I apply multi functions query it fails the query.