# Graded Assignment 5.4

Name: Saad Sameer Khan
Employee#: 2303.KHI.DEG.034
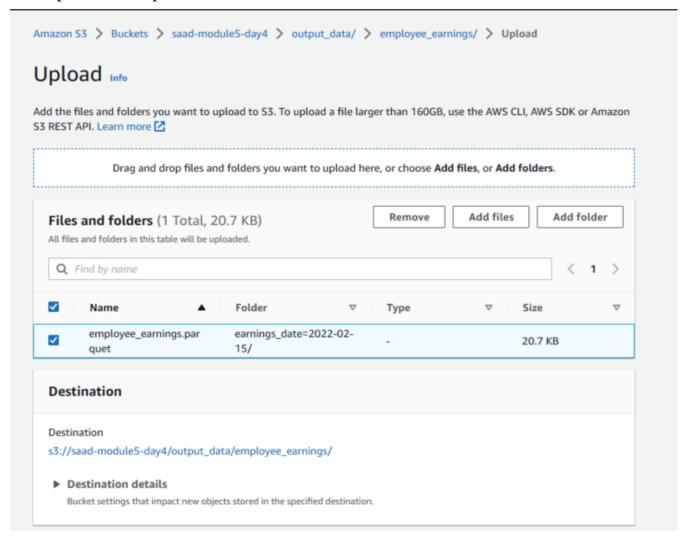Collaborated with: Mohammad Hamza Asim (2303.KHI.DEG.014)

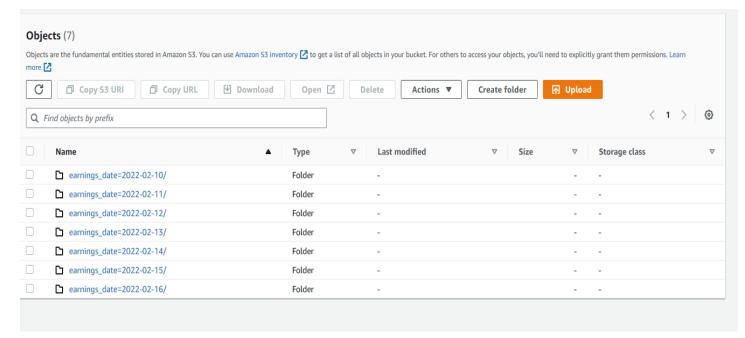## Creating data for 2 more days (day 6 & day 7)

We used numpy's random.randint method to generate random earnings data within the range of the maximum and minimum value of the 'earnings' column from one of the day's data. *ipynb file has been uploaded to git*

# Uploading new folders to S3 bucket

Next, we uploaded two new folders containing the data created in the previous step, in the output_data folder in the S3 bucket.

# Running crawler

After the new folders were uploaded, we ran the crawler so that it could fetch the new data from the bucket, and update its data catalog, so that we could use the new data in Athena.

**Crawler properties**

| | | | |
|---|---|---|---|
| Name | IAM role | Database | State |
| saad_combined_employee_earnings_ crawler | saad-glue-role ↗ | saad_glue_db | READY |
| Description | Security configuration | Lake Formation configuration | Table prefix |
| - | - | - | saad_ |
| Maximum table threshold | | | |
| - | | | |

▶ Advanced settings

**Crawler runs**    Schedule    Data sources    Classifiers    Tags

**Crawler runs** (2)
The list of crawler runs for this crawler.

Stop run    View CloudWatch logs ↗    View run details

| | Start time (UTC) ▲ | End time (UTC) ▽ | Current/last duration ▽ | Status ▽ | DPU hours ▽ | Table changes ▽ |
|---|---|---|---|---|---|---|
| ○ | May 18, 2023 at 12:02:29 | May 18, 2023 at 12:07:15 | 04 min 45 s | ⊘ Completed | 0.072 | 1 table change, 2 partition changes |
| ○ | May 18, 2023 at 08:06:43 | May 18, 2023 at 08:07:47 | 01 min 04 s | ⊘ Completed | 0.036 | 1 table change, 5 partition changes |

# Re-running queries on updated data

Now we re-ran the queries that were previously run, and observe what changes took place in the results.

## Query 1

*Before:*

| | Completed | | Time in queue: 172 ms | Run time: 958 ms | Data scanned: 19.04 KB |
|---|---|---|---|---|---|

**Results** (46)  — Copy — Download results

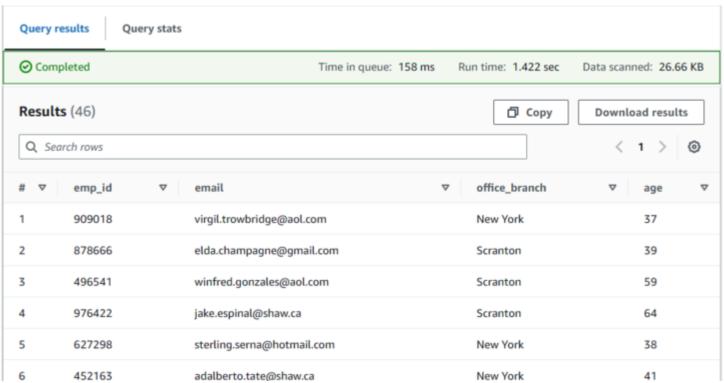| # | emp_id | email | office_branch | age |
|---|---|---|---|---|
| 1 | 900756 | benjamin.doss@gmail.com | Scranton | 38 |
| 2 | 215719 | brent.carrillo@aol.com | New York | 50 |
| 3 | 530134 | mathew.whitfield@gmail.com | New York | 36 |
| 4 | 597741 | tonya.wilson@aol.com | New York | 43 |
| 5 | 391837 | cory.hayden@gmail.com | New York | 56 |
| 6 | 622405 | harrison.hawk@hotmail.co.uk | Scranton | 60 |

*After*

| | Completed | | Time in queue: 158 ms | Run time: 1.422 sec | Data scanned: 26.66 KB |
|---|---|---|---|---|---|

**Results** (46)  — Copy — Download results

| # | emp_id | email | office_branch | age |
|---|---|---|---|---|
| 1 | 909018 | virgil.trowbridge@aol.com | New York | 37 |
| 2 | 878666 | elda.champagne@gmail.com | Scranton | 39 |
| 3 | 496541 | winfred.gonzales@aol.com | Scranton | 59 |
| 4 | 976422 | jake.espinal@shaw.ca | Scranton | 64 |
| 5 | 627298 | sterling.serna@hotmail.com | New York | 38 |
| 6 | 452163 | adalberto.tate@shaw.ca | New York | 41 |

## Query 2
*Before:*



*After:*

# Query 3

*Before:*



*After:*

# Calculating percentage change in earnings

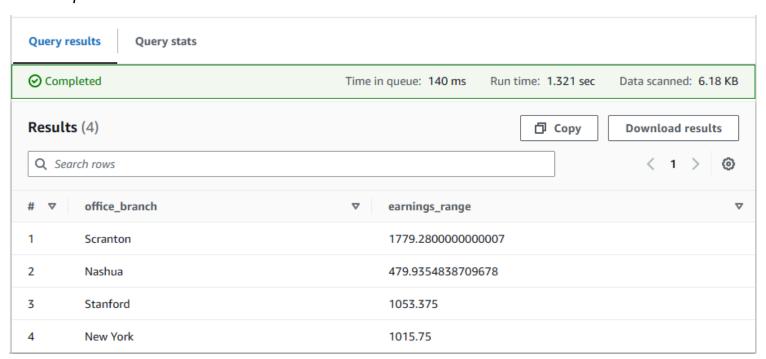Now running a new query in Athena that takes in an input day (day 16) and calculates the % change in earnings for every employee from compared to the previous day.

```sql
 1  WITH earnings_data AS (
 2    SELECT
 3      emp_id,
 4      earnings,
 5      earnings_date,
 6      LAG(earnings) OVER (ORDER BY earnings_date) AS previous_earnings
 7    FROM "saad_glue_db"."saad_employee_earnings"
 8    WHERE earnings_date IN ('2022-02-11', '2022-02-12', '2022-02-13', '2022-02-14', '2022-02-15', '2022
        -02-16')
 9  )
10  SELECT
11    emp_id,
12    earnings_date,
13    earnings AS current_earnings,
14    previous_earnings,
15    ((earnings - previous_earnings) / CAST(previous_earnings AS double)) * 100 AS percentage_change
16  FROM
17    earnings_data
18  WHERE
19    earnings_date = '2022-02-16';
```

Here are the results:

**Results (100)**

| # | emp_id | earnings_date | current_earnings | previous_earnings | percentage_change |
|---|--------|---------------|------------------|-------------------|-------------------|
| 1 | 289172 | 2022-02-16 | 9283 | 4707 | 97.21691098364138 |
| 2 | 915991 | 2022-02-16 | 9588 | 9283 | 3.2855757836906174 |
| 3 | 203380 | 2022-02-16 | 5985 | 9588 | -37.57822277847309 |
| 4 | 466832 | 2022-02-16 | 3031 | 5985 | -49.35672514619883 |
| 5 | 549389 | 2022-02-16 | 2425 | 3031 | -19.99340151765094 |
| 6 | 174955 | 2022-02-16 | 2969 | 2425 | 22.432989690721648 |
| 7 | 149972 | 2022-02-16 | 8636 | 2969 | 190.87234759178176 |
| 8 | 856379 | 2022-02-16 | 2304 | 8636 | -73.32098193608152 |
| 9 | 962291 | 2022-02-16 | 9248 | 2304 | 301.38888888888886 |
| 10 | 819367 | 2022-02-16 | 3492 | 9248 | -62.240484429065745 |
| 11 | 242388 | 2022-02-16 | 6296 | 3492 | 80.29782359679267 |
| 12 | 402180 | 2022-02-16 | 4198 | 6296 | -33.32274459974587 |