## Problem Statement:

Using the data manipulation tool of your choice (eg. Python) simulate the earnings predictions for 2 more days. Load it to the Data Lake that you've created today (Task 1-2).
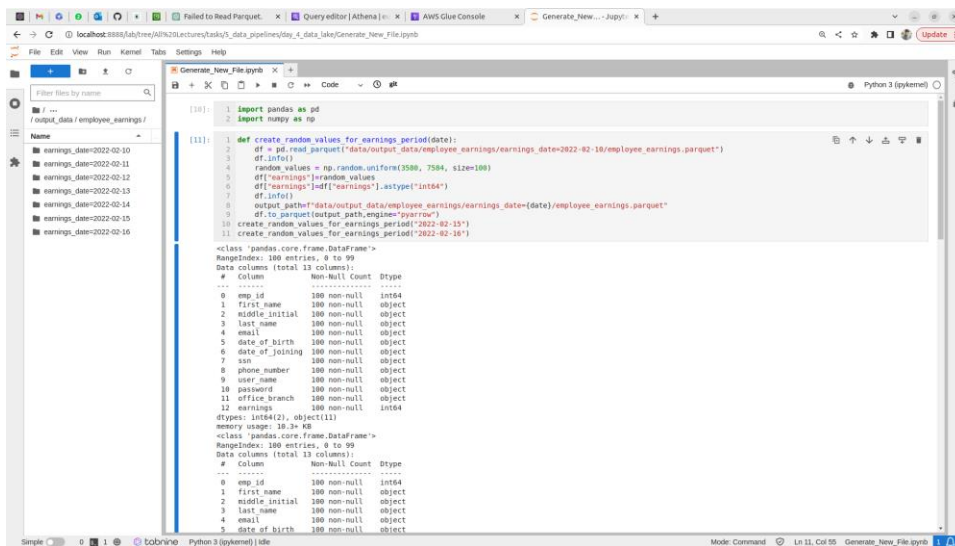 Rerun queries from Task 3 and Task 4 and see how the results change with this new data.
 Create a new query in Athena that calculates the % change in earnings for every employee from a given day compared to the previous day.

## Solution:

# Creating the new data using this script

We firstly analyzed the given data and based on that data we created the new data using a suitable range of random data.



# Task 1: Loading the Data into S3:



# Task2: Running the Glue Crawler on the new data:

# Task 3: Rerunning the previous query to see the change:

These are the previous query result

**Screenshot 1 — Results (46)**

| # | emp_id | email | office_branch | age |
|---|--------|-------|---------------|-----|
| 1 | 896517 | jenell.almanza@yahoo.com | New York | 64 |
| 2 | 633636 | bertram.carlisle@aol.com | Scranton | 40 |
| 3 | 495667 | cory.clarke@shell.com | New York | 42 |
| 4 | 500905 | harris.beavers@shell.com | Scranton | 52 |
| 5 | 492527 | hilton.mcgehee@microsoft.com | New York | 36 |
| 6 | 932773 | clair.harwell@ibp.com | Scranton | 54 |
| 7 | 405534 | clement.hidalgo@gmail.com | New York | 63 |
| 8 | 909018 | virgil.trowbridge@aol.com | New York | 57 |
| 9 | 878666 | elda.champagne@gmail.com | Scranton | 39 |
| 10 | 496541 | winfred.gonzales@aol.com | Scranton | 59 |
| 11 | 976422 | jake.espinal@shaw.ca | Scranton | 64 |
| 12 | 627298 | sterling.serna@hotmail.com | New York | 38 |
| 13 | 452163 | adalberto.tate@shaw.ca | New York | 41 |
| 14 | 147133 | tommie.weller@cox.net | Scranton | 63 |
| 15 | 403207 | michal.maurer@yahoo.com | Scranton | 46 |
| 16 | 397283 | rex.ng@yahoo.com | New York | 40 |
| 17 | 754455 | anastasia.childers@hotmail.com | New York | 34 |

**Screenshot 2 — Results (20)**

| # | office_branch | min_earnings | max_earnings | avg_earnings | total_earnings | earnings_date |
|---|---------------|--------------|--------------|--------------|----------------|---------------|
| 1 | Nashua | 2098 | 9728 | 6099.8387096774195 | 189095 | 2022-02-14 |
| 2 | Nashua | 2005 | 9786 | 6049.451612903225 | 187533 | 2022-02-13 |
| 3 | Nashua | 2006 | 9605 | 5997.967741935484 | 185937 | 2022-02-11 |
| 4 | New York | 2295 | 9889 | 6631.285714285715 | 185676 | 2022-02-12 |
| 5 | Nashua | 2124 | 9978 | 5764.5161290322585 | 178700 | 2022-02-12 |
| 6 | Nashua | 2066 | 9801 | 5619.903225806452 | 174217 | 2022-02-10 |
| 7 | New York | 2040 | 9954 | 6109.035714285715 | 171053 | 2022-02-14 |
| 8 | Scranton | 2788 | 9916 | 6850.6 | 170765 | 2022-02-13 |
| 9 | New York | 2141 | 9462 | 5998.170571428572 | 167949 | 2022-02-11 |
| 10 | New York | 2576 | 9972 | 5991.321428571428 | 167757 | 2022-02-10 |
| 11 | New York | 2195 | 9734 | 5615.535714285715 | 157235 | 2022-02-13 |
| 12 | Scranton | 2465 | 9027 | 6149.72 | 153743 | 2022-02-14 |
| 13 | Scranton | 2025 | 9846 | 6063.44 | 151586 | 2022-02-12 |

**Screenshot 3 — Query and Results (4)**

```sql
SELECT DISTINCT office_branch, (MAX(avg_earnings.value) - MIN(avg_earnings.value)) as earnings_range
FROM (
    SELECT office_branch as ob, AVG(earnings) AS value FROM "muzammilmehmood_glue_database"."muzammilmehmoodoutput_data" GROUP BY office_branch,
        earnings_date
) avg_earnings, "muzammilmehmood_glue_database"."muzammilmehmoodoutput_data"
WHERE office_branch = avg_earnings.ob
GROUP BY office_branch;
```

| # | office_branch | earnings_range |
|---|---------------|----------------|
| 1 | New York | 1015.75 |
| 2 | Stanford | 1053.375 |
| 3 | Nashua | 479.9554858709678 |
| 4 | Scranton | 1779.2800000000007 |

Now let me show you the result of the previous query on new data



SQL query editor showing:
```
SELECT * FROM "mtalha_glue_database"."mtalhaemployee_earnings" ;
```

**Results (700)**

| # | emp_id | first_name | middle_initial | last_name | email | date_of_birth | date_of_joining | ssn | phone_number |
|---|--------|-----------|----------------|-----------|-------|---------------|-----------------|-----|--------------|
| 1 | 526540 | Angelique | K | Goodwin | angelique.goodwin@gmail.com | 1964-05-15 | 2001-03-24 | 471-57-0559 | 212-884-7146 |
| 2 | 859527 | Jeni | S | Shaffer | jeni.shaffer@gmail.com | 1962-01-13 | 2015-12-10 | 624-05-4146 | 205-665-7020 |
| 3 | 887387 | Donald | T | Farris | donald.farris@bellsouth.net | 1958-04-11 | 1979-11-12 | 097-02-3315 | 205-959-7879 |
| 4 | 779497 | Steven | D | Rendon | steven.rendon@gmail.com | 1982-04-04 | 2008-09-18 | 134-98-6566 | 217-858-0054 |

**Results (46)**

| # | emp_id | email | office_branch | age |
|---|--------|-------|---------------|-----|
| 1 | 909018 | virgil.trowbridge@aol.com | New York | 37 |
| 2 | 878666 | elda.champagne@gmail.com | Scranton | 39 |
| 3 | 391857 | cory.hayden@gmail.com | New York | 56 |
| 4 | 496541 | winfred.gonzales@aol.com | Scranton | 59 |
| 5 | 976422 | jake.espinal@shaw.ca | Scranton | 64 |
| 6 | 627298 | sterling.serna@hotmail.com | New York | 38 |
| 7 | 622405 | harrison.hawk@hotmail.co.uk | Scranton | 60 |
| 8 | 452163 | adalberto.tate@shaw.ca | New York | 41 |
| 9 | 147133 | tommie.weller@cox.net | Scranton | 63 |
| 10 | 595558 | denisha.fitch@msn.com | Scranton | 32 |
| 11 | 314661 | charles.quintero@gmail.com | New York | 65 |
| 12 | 403207 | michal.maurer@yahoo.com | Scranton | 46 |
| 13 | 896517 | jenell.almanza@yahoo.com | New York | 64 |
| 14 | 633656 | bertram.carlisle@aol.com | Scranton | 40 |
| 15 | 823898 | carlton.leclair@cox.net | Scranton | 37 |
| 16 | 495667 | cory.clarke@shell.com | New York | 42 |
| 17 | 500905 | harris.beavers@shell.com | Scranton | 52 |
| 18 | 492527 | hilton.mcgehee@microsoft.com | New York | 36 |
| 19 | 952773 | clair.harwell@bp.com | Scranton | 54 |

**Results (28)**

| # | office_branch | min_earnings | max_earnings | avg_earnings | total_earnings | earnings_date |
|---|---------------|--------------|--------------|--------------|----------------|---------------|
| 1 | Nashua | 2098 | 9728 | 6099.8387096774195 | 189095 | 2022-02-14 |
| 2 | Nashua | 2005 | 9786 | 6049.451612903225 | 187533 | 2022-02-13 |
| 3 | Nashua | 2006 | 9605 | 5997.967741935484 | 185937 | 2022-02-11 |
| 4 | New York | 2295 | 9889 | 6631.285714285715 | 185676 | 2022-02-12 |
| 5 | Nashua | 2124 | 9978 | 5764.5161290322585 | 178700 | 2022-02-12 |
| 6 | Nashua | 2066 | 9801 | 5619.903225806452 | 174217 | 2022-02-10 |
| 7 | New York | 2040 | 9954 | 6109.035714285715 | 171053 | 2022-02-14 |
| 8 | Scranton | 2788 | 9916 | 6830.6 | 170765 | 2022-02-13 |
| 9 | New York | 2141 | 9462 | 5998.178571428572 | 167949 | 2022-02-11 |
| 10 | Nashua | 3729 | 7541 | 5413.387096774193 | 167815 | 2022-02-16 |
| 11 | New York | 2376 | 9972 | 5991.321428571428 | 167757 | 2022-02-10 |
| 12 | Nashua | 3740 | 7453 | 5373.387096774193 | 166575 | 2022-02-15 |
| 13 | New York | 2195 | 9734 | 5615.535714285715 | 157235 | 2022-02-13 |
| 14 | New York | 5688 | 7567 | 5611.571428571428 | 157124 | 2022-02-15 |
| 15 | Scranton | 2465 | 9827 | 6149.72 | 153743 | 2022-02-14 |
| 16 | Scranton | 2023 | 9846 | 6063.44 | 151586 | 2022-02-12 |
| 17 | Scranton | 2033 | 9888 | 6005.56 | 150139 | 2022-02-10 |
| 18 | New York | 5671 | 7530 | 5349.821428571428 | 149795 | 2022-02-16 |
| 19 | Scranton | 5618 | 7191 | 5557.84 | 138446 | 2022-02-16 |
| 20 | Scranton | 3749 | 7558 | 5267.56 | 131684 | 2022-02-15 |

## Task4: Querying the data using S3 Select

Due to the limitation of s3 select and the complexity of your queries, we were unable to run the queries directly on s3 select as it operates on the single object at a time, while Athena allows you to run queries across multiple objects and supports more complex queries and it would need multiple subsets of data, which cause errors in s3 select.

Create a new query in Athena that calculates the % change in earnings for every employee from a given day compared to the previous day.

Data

Data source
AwsDataCatalog ▼

Database
mtalha_glue_database ▼

Tables and views | Create ▼

Q Filter tables and views

▼ Tables (5)  ‹ 1 ›
⊞ employee_earnings        Partitioned ⋮
⊞ mtalha_earnings          Partitioned ⋮
⊞ mtalha_employees_earnings Partitioned ⋮
⊞ mtalha_location_location             ⋮
⊞ mtalhaemployee_earnings  Partitioned ⋮
▸ Views (0)   ‹ 1 ›

Query 7 ⋮ ✕ | Query 8 ⋮ ✕ | Query 9 ⋮ ✕ | ⊘ Query 10 ⋮ ✕ | ⊘ Query 11 ⋮ ✕

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

SQL   Ln 17, Col 19

Run again | Explain ⧉ | Cancel | Clear | Create ▼

Reuse query results
up to 60 minutes ago ✎

Query results | Query stats

⊘ Completed          Time in queue: 110 ms    Run time: 761 ms    Data scanned: 22.75 KB

Results (100)                                          Copy    Download results

Q Search rows                                          ‹ 1 … › ⚙

| # ▼ | emp_id ▼ | first_name ▼ | last_name ▼ | earnings_date ▼ | previous_earnings_date ▼ | earnings ▼ | previous_earnings ▼ | percentage_change ▼ |
|---|---|---|---|---|---|---|---|---|
| 1 | 220965 | Almeta | Brookins | 2022-02-15 | 2022-02-14 | 3859 | 9378 | -58.85050117295799 |
| 2 | 235295 | Yevette | Mullis | 2022-02-15 | 2022-02-14 | 7221 | 5760 | 25.364583333333336 |
| 3 | 312726 | Celine | Lumpkin | 2022-02-15 | 2022-02-14 | 7305 | 6055 | 20.644095788604456 |
| 4 | 314661 | Charles | Quintero | 2022-02-15 | 2022-02-14 | 3688 | 8483 | -56.52481433455145 |
| 5 | 316572 | Alexander | Goad | 2022-02-15 | 2022-02-14 | 3980 | 6686 | -40.47262937481304 |
| 6 | 366431 | Sade | Shay | 2022-02-15 | 2022-02-14 | 3882 | 9018 | -56.9527611445779 |
| 7 | 403534 | Clement | Hidalgo | 2022-02-15 | 2022-02-14 | 5206 | 5530 | -5.8589511754068715 |
| 8 | 549389 | Clemente | Gould | 2022-02-15 | 2022-02-14 | 7581 | 7944 | -4.569486404853037 |
| 9 | 597741 | Tonya | Wilson | 2022-02-15 | 2022-02-14 | 7010 | 9094 | -22.9162084891137 |
| 10 | 622405 | Harrison | Hawk | 2022-02-15 | 2022-02-14 | 4308 | 3974 | 8.40463009562154 |
| 11 | 728053 | Leigh | Fite | 2022-02-15 | 2022-02-14 | 6416 | 3605 | 78.07382736608582 |
| 12 | 819367 | Yolande | Piper | 2022-02-15 | 2022-02-14 | 4135 | 8535 | -51.55243116578794 |
| 13 | 878666 | Elda | Champagne | 2022-02-15 | 2022-02-14 | 5367 | 7755 | -30.79303675048356 |
| 14 | 885395 | Tayna | Poston | 2022-02-15 | 2022-02-14 | 7567 | 4478 | 68.98168825368468 |
| 15 | 896517 | Jenell | Almanza | 2022-02-15 | 2022-02-14 | 6283 | 2057 | 205.44482255712202 |
| 16 | 936158 | Sofia | Poole | 2022-02-15 | 2022-02-14 | 7428 | 9493 | -21.752870536184556 |
| 17 | 962291 | Whitney | Shipman | 2022-02-15 | 2022-02-14 | 6377 | 5286 | 94.06575341448569 |
| 18 | 976422 | Jake | Espinal | 2022-02-15 | 2022-02-14 | 3802 | 6532 | -41.79424372320815 |