# NAME: Mohammad Hussam (2303.KHI.DEG.020)
## PARING WITH : MAVIA ALAM KHAN(2302.KHI.DEG.017)
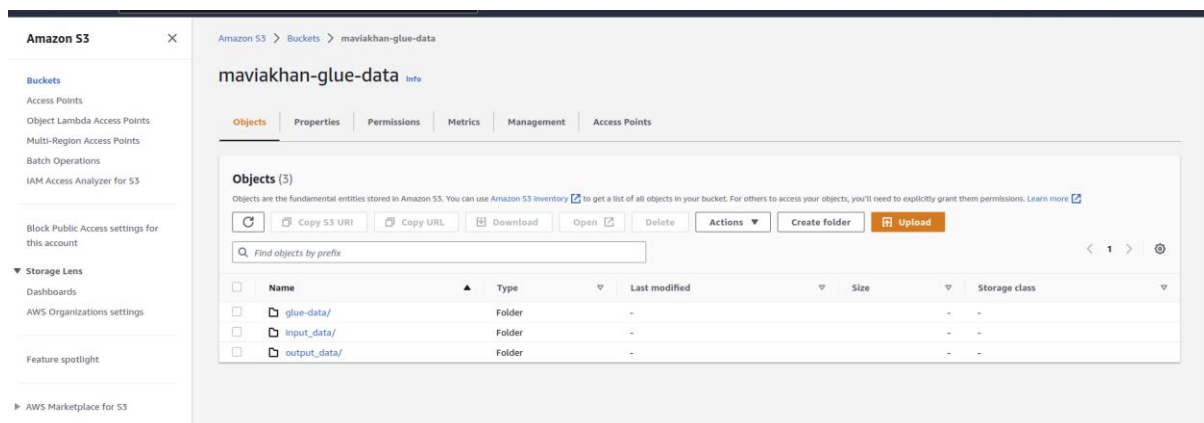## &
## AQSA TAUHEED(2303.KHI.DEG.011)

# ASSIGNMENT NO :5.2

Using the salary CSV as a base, prepare a new data file with employees' office locations. Make sure there are 5-6 distinct locations that are shared between employees.
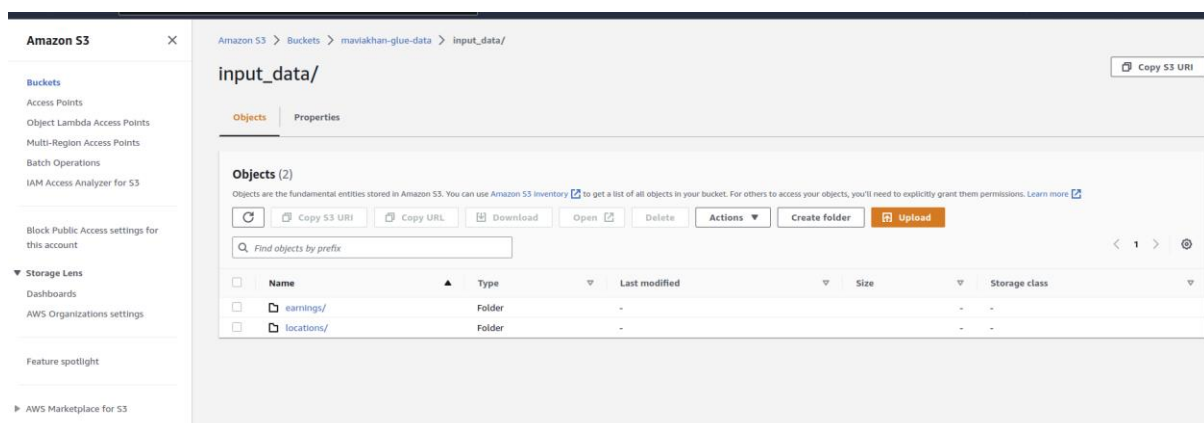Create a Glue job that aggregates the data based on the office location to calculate
average salaries and raise percentages for these locations.

## SOLUTION

First create a folder in s3 bucket name (input and output)



In input folder we store a two dataset earning.csv and location.csv

Now we create a crawler and extract the meta data



And after create the glue crawler we start creating the job

Now we create two s3 source one is employee earning data and location,we perform inner join on both data on emp_id and after that prepare for querying

Now here the queries based on the salaries and percentage of these locations.



| location | average_earnings | raise_percentage |
|---|---|---|
| B | 6286.75 | 155.14407467532467 |
| C | 5576.95 | 129.78780387309433 |
| A | 5926.05 | 191.49286768322676 |
| D | 5889.7 | 185.07744433688285 |
| E | 5599.2 | 158.74306839186693 |

Now finally we load the data and show in the output folder

# Amazon S3

- **Buckets**
- Access Points
- 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
- AWS Organizations settings

Feature spotlight

▶ AWS Marketplace for S3

## earningswithLocationTarget/

Copy S3 URI

**Objects** | Properties

### Objects (5)

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 ↗
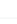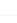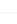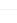
Copy S3 URI | Copy URL | Download | Open | Delete | Actions ▼ | Create folder | Upload

🔍 Find objects by prefix

< 1 >

| | Name | Type ▽ | Last modified ▽ | Size ▽ | Storage class ▽ |
|---|---|---|---|---|---|
| ☐ | 📄 run-1684517258119-part-block-0-r-00002-snappy.parquet | parquet | May 19, 2023, 22:27:46 (UTC+05:00) | 599.0 B | Standard |
| ☐ | 📄 run-1684517258119-part-block-0-r-00014-snappy.parquet | parquet | May 19, 2023, 22:27:45 (UTC+05:00) | 599.0 B | Standard |
| ☐ | 📄 run-1684517258119-part-block-0-r-00021-snappy.parquet | parquet | May 19, 2023, 22:27:45 (UTC+05:00) | 599.0 B | Standard |
| ☐ | 📄 run-1684517258119-part-block-0-r-00025-snappy.parquet | parquet | May 19, 2023, 22:27:45 (UTC+05:00) | 599.0 B | Standard |
| ☐ | 📄 run-1684517258119-part-block-0-r-00031-snappy.parquet | parquet | May 19, 2023, 22:27:44 (UTC+05:00) | 599.0 B | Standard |