

Name: - Muhammad Huzaifa Waseem (2303-KHI-DEG-021)

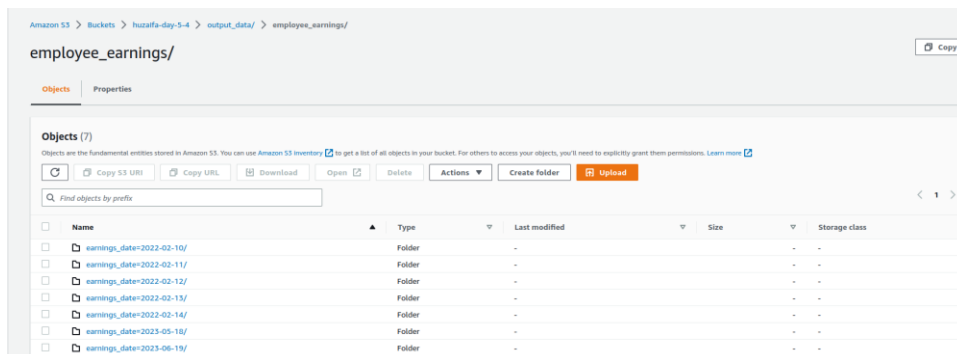
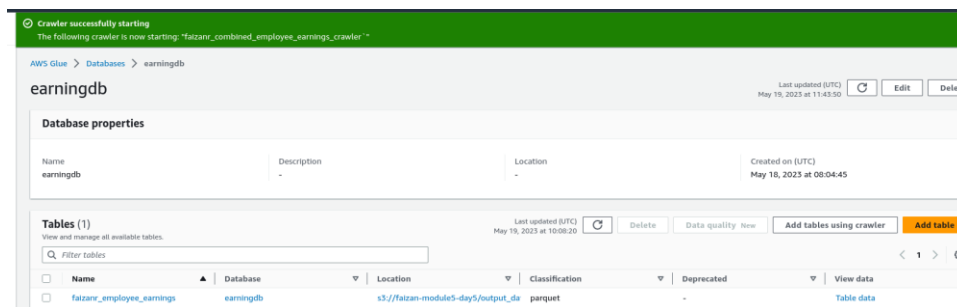
Pair Partner 1: - Muhammad Faizan Rafique (2303.005.KHI.DEG)

Pair Partner 2: - Syed Muhammad Hammad Irshad(2303.KHI.DEG.032)

UNIT 5.4:

Assignment

Use data from today's Daily Activities tasks/5_data_pipelines/day_4_data_lake/data/output_data/employee_earnings



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).

```

Name | sytchmuhammadhmdhmdh | Download | 1 | datafile syb | 1 | report pandas | 1 |
+ Code | 1 | Markdown | 1 | Run All | 1 | Clear All Outputs | 1 | 1 | Outline | 1 |

import pandas as pd
import numpy as np

MUHAMMAD FAIZAN RAFIQ 2303.005.DEG.KHI MUHAMMAD HUZAF A WASSEEM 2303.KHI.DEG-021 HAMMAD IRSHAD 2303.KHI.DEG-032

df1 = pd.read_parquet('/home/muhammadfaizanrafique/data_engineering_bootcamp_2303/tasks/5_data_pipeline/day_4_data lake/data/output_data/employee_earnings/earnings_date=2022-02-18/employee_earnings.parquet')

df1.head()

emp_id first_name middle_initial last_name email date_of_birth date_of_joining sex phone_number user_name password office_branch earnings
0 525456 Anygeline K Goodwin anygeline.goodwin@gmail.com 1964-05-15 2007-03-24 471-57-0231 212-884-7146 anygelinec d5m7nLzP Nashville 6221
1 809327 Jeno S Shaffer jeno.shaffer@gmail.com 1962-01-13 2015-12-10 428-85-4146 205-665-7020 joshaffer 705m7D Stanford 8145
2 887387 Donald T Farris donald.farris@earthlink.net 1958-06-11 1979-11-12 897-02-3315 205-959-7879 dfarris nLFG3m8AAx Stanford 5278
3 779497 Steven D Rowen steven.rowen@gmail.com 1962-04-04 2008-09-10 328-89-4345 212-858-0054 srowen02 Nashville 9647
4 696517 Jerald L Alvarado jerald.alvarado@yahoo.com 1958-07-01 1993-07-14 399-92-7345 314-893-2390 jalvarado OutRkgT New York 4138

95 149389 Clemente M Gould clemente.gould@hotmail.com 1981-12-31 1992-10-02 271-17-5467 228-485-0919 cmgould m7%uqg7VhJ Stanford 3402
96 466822 Chang K Roden chang.roden@yahoo.com 1988-09-07 2010-08-06 874-02-0202 316-256-7851 choden SBHG4W8Ry+5 Nashville 5373
97 213369 Marcin B Reade marcin.reak@gmail.com 1986-11-25 2012-10-06 552-99-3445 276-750-7760 mreade0 P7Qg_R Scotland 7284
98 913991 Ediva Y Triskie edivay.triskie@earthlink.net 1995-01-28 2016-10-17 763-12-2082 236-584-1914 eytriskie jzm0t584 Nashville 2810
99 289172 Azale L Layman azale.layman@hotmail.co.uk 1981-09-06 2004-03-26 637-28-1007 503-656-5889 alayman b7%u7m6_12V New York 2295

100 rows x 12 columns

random_numbers = np.random.randint(10000, 100000, size=len(df1))

df1.iloc[:,1:] = random_numbers

new_earning_col = df1.iloc[:, 'earnings']
new_earning_col

0 18902
1 33209
2 32309
3 50427
4 54208
...
95 45413
96 78823
97 62264

Name: earnings, Length: 100, dtype: int64

df1.head()

emp_id first_name middle_initial last_name email date_of_birth date_of_joining sex phone_number user_name password office_branch earnings
0 525456 Anygeline K Goodwin anygeline.goodwin@gmail.com 1964-05-15 2007-03-24 471-57-0231 212-884-7146 anygelinec d5m7nLzP Nashville 10171
1 809327 Jeno S Shaffer jeno.shaffer@gmail.com 1962-01-13 2015-12-10 428-85-4146 205-665-7020 joshaffer 705m7D Stanford 21660
2 887387 Donald T Farris donald.farris@earthlink.net 1958-06-11 1979-11-12 897-02-3315 205-959-7879 dfarris nLFG3m8AAx Stanford 60996
3 779497 Steven D Rowen steven.rowen@gmail.com 1962-04-04 2008-09-10 328-89-4345 212-858-0054 srowen02 Nashville 20000
4 696517 Jerald L Alvarado jerald.alvarado@yahoo.com 1958-07-01 1993-07-14 399-92-7345 314-893-2390 jalvarado OutRkgT New York 17811

path = '/home/muhammadfaizanrafique/data_engineering_bootcamp_2303/tasks/5_data_pipeline/day_4_data lake/data/output_data/employee_earnings/earnings_date=2022-05-18/employee_earnings.parquet'
df1.to_parquet(path)

df2 = pd.read_parquet('/home/muhammadfaizanrafique/data_engineering_bootcamp_2303/tasks/5_data_pipeline/day_4_data lake/data/output_data/employee_earnings/earnings_date=2022-02-12/employee_earnings.parquet')

df2.head()

emp_id first_name middle_initial last_name email date_of_birth date_of_joining sex phone_number user_name password office_branch earnings
0 525456 Anygeline K Goodwin anygeline.goodwin@gmail.com 1964-05-15 2007-03-24 471-57-0231 212-884-7146 anygelinec d5m7nLzP Nashville 6176
1 809327 Jeno S Shaffer jeno.shaffer@gmail.com 1962-01-13 2015-12-10 428-85-4146 205-665-7020 joshaffer 705m7D Stanford 8145
2 887387 Donald T Farris donald.farris@earthlink.net 1958-06-11 1979-11-12 897-02-3315 205-959-7879 dfarris nLFG3m8AAx Stanford 5278
3 779497 Steven D Rowen steven.rowen@gmail.com 1962-04-04 2008-09-10 328-89-4345 212-858-0054 srowen02 Nashville 9647
4 696517 Jerald L Alvarado jerald.alvarado@yahoo.com 1958-07-01 1993-07-14 399-92-7345 314-893-2390 jalvarado OutRkgT New York 4138

95 149389 Clemente M Gould clemente.gould@hotmail.com 1981-12-31 1992-10-02 271-17-5467 228-485-0919 cmgould m7%uqg7VhJ Stanford 3402
96 466822 Chang K Roden chang.roden@yahoo.com 1988-09-07 2010-08-06 874-02-0202 316-256-7851 choden SBHG4W8Ry+5 Nashville 5373
97 213369 Marcin B Reade marcin.reak@gmail.com 1986-11-25 2012-10-06 552-99-3445 276-750-7760 mreade0 P7Qg_R Scotland 7284
98 913991 Ediva Y Triskie edivay.triskie@earthlink.net 1995-01-28 2016-10-17 763-12-2082 236-584-1914 eytriskie jzm0t584 Nashville 2810
99 289172 Azale L Layman azale.layman@hotmail.co.uk 1981-09-06 2004-03-26 637-28-1007 503-656-5889 alayman b7%u7m6_12V New York 2295

100 rows x 12 columns

random_numbers = np.random.randint(10000, 100000, size=len(df2))
df2.iloc[:,1:] = random_numbers
new_earning_col = df2.iloc[:, 'earnings']
new_earning_col

0 88888
1 14329
2 76238
3 96335
4 54739
...
95 52549
96 77849
97 60889
98 58084
99 88028

Name: earnings, Length: 100, dtype: int64

df2.head()

emp_id first_name middle_initial last_name email date_of_birth date_of_joining sex phone_number user_name password office_branch earnings
0 525456 Anygeline K Goodwin anygeline.goodwin@gmail.com 1964-05-15 2007-03-24 471-57-0231 212-884-7146 anygelinec d5m7nLzP Nashville 8178
1 809327 Jeno S Shaffer jeno.shaffer@gmail.com 1962-01-13 2015-12-10 428-85-4146 205-665-7020 joshaffer 705m7D Stanford 8145
2 887387 Donald T Farris donald.farris@earthlink.net 1958-06-11 1979-11-12 897-02-3315 205-959-7879 dfarris nLFG3m8AAx Stanford 5278
3 779497 Steven D Rowen steven.rowen@gmail.com 1962-04-04 2008-09-10 328-89-4345 212-858-0054 srowen02 Nashville 9647
4 696517 Jerald L Alvarado jerald.alvarado@yahoo.com 1958-07-01 1993-07-14 399-92-7345 314-893-2390 jalvarado OutRkgT New York 4138

95 149389 Clemente M Gould clemente.gould@hotmail.com 1981-12-31 1992-10-02 271-17-5467 228-485-0919 cmgould m7%uqg7VhJ Stanford 3402
96 466822 Chang K Roden chang.roden@yahoo.com 1988-09-07 2010-08-06 874-02-0202 316-256-7851 choden SBHG4W8Ry+5 Nashville 5373
97 213369 Marcin B Reade marcin.reak@gmail.com 1986-11-25 2012-10-06 552-99-3445 276-750-7760 mreade0 P7Qg_R Scotland 7284
98 913991 Ediva Y Triskie edivay.triskie@earthlink.net 1995-01-28 2016-10-17 763-12-2082 236-584-1914 eytriskie jzm0t584 Nashville 2810
99 289172 Azale L Layman azale.layman@hotmail.co.uk 1981-09-06 2004-03-26 637-28-1007 503-656-5889 alayman b7%u7m6_12V New York 2295

100 rows x 12 columns

new_path = '/home/muhammadfaizanrafique/data_engineering_bootcamp_2303/tasks/5_data_pipeline/day_4_data lake/data/output_data/employee_earnings/earnings_date=2023-06-10/employee_earnings.parquet'
df2.to_parquet(new_path)


```

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

PREVIOUS DATA:

Data source

AwsDataCatalog

Database

huzafa-database

Tables and views

Create

Filter tables and views

Tables (1)

huzafaaw_employee_earnings

--SELECT * FROM "huzafa-database"."huzafaaw_employee_earnings" limit 10;

SELECT DISTINCT emp_id, email, office_branch, (date_diff('year', DATE(date_of_birth), current_date)) AS age

FROM "huzafa-database"."huzafaaw_employee_earnings"

WHERE office_branch IN ('New York', 'Scranton')

AND

(date_diff('year', DATE(date_of_birth), current_date)) > 30;

SQLLn 7, Col 64

Results (46)

CopyDownload results

Search rows

#	emp_id	email	office_branch	age
1	397283	rex.ng@yahoo.com	New York	40
2	734455	anastasia.children@hotmail.com	New York	34
3	961442	stephan.titus@hotmail.com	Scranton	55
4	155097	dorian.sage@shaw.ca	Scranton	52
5	160938	marcel.trotter@aol.com	Scranton	38
6	885395	layna.poston@yahoo.ca	New York	65
7	289172	azzie.layman@hotmail.co.uk	New York	61
8	900756	benjamin.doss@gmail.com	Scranton	38
9	215719	brent.carrillo@aol.com	New York	50
10	530134	mathew.whitfield@gmail.com	New York	36
11	597741	torya.wilson@aol.com	New York	43
12	220965	almeta.brookins@gmail.com	Scranton	38
13	537991	samuel.wendt@bellsouth.net	New York	46
14	767674	irena.dang@gmail.com	New York	51
15	505927	oswaldo.winchester@gmail.com	New York	64
16	432820	myron.marble@gmail.com	New York	42
17	317987	chastity.pineda@shaw.ca	New York	48
18	203580	marvin.nickel@ibm.com	Scranton	36
19	896517	jennell.almanza@yahoo.com	New York	64

Query 1Query 2

--SELECT * FROM "huzafa-database"."huzafaaw_employee_earnings" limit 10;

SELECT office_branch, MIN(earnings) as min_earnings, MAX(earnings) as max_earnings, AVG(earnings) as avg_earnings, SUM(earnings) as total_earnings,

earnings_date

FROM "huzafa-database"."huzafaaw_employee_earnings"

GROUP BY office_branch, earnings_date

ORDER BY SUM(earnings) desc;

Results (20)

CopyDownload results

Search rows

#	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	9603	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	New York	2376	9972	5991.321428571428	167757	2022-02-10
11	New York	2195	9734	5615.535714285715	157235	2022-02-13
12	Scranton	2465	9827	6149.72	153743	2022-02-14
13	Scranton	2023	9846	6063.44	151586	2022-02-12
14	Scranton	2033	9888	6005.56	150139	2022-02-10
15	Scranton	2199	9179	5051.32	126283	2022-02-11
16	Stanford	2921	9493	6617.75	105884	2022-02-14
17	Stanford	2210	9404	6237.6875	99803	2022-02-11
18	Stanford	2808	9850	6142.1875	98275	2022-02-12
19	Stanford	3642	9467	5841.25	93460	2022-02-10

Query 2 : X Query 3 : X Query 4 : X Query 5 : X

```

1 SELECT DISTINCT office_branch, (MAX(avg_earnings.value) - MIN(avg_earnings.value)) as earnings_range
2 , FROM (
3 SELECT office_branch as ob, AVG(earnings) AS value FROM "earningdb"."faizanr_employee_earnings" GROUP BY office_branch, earnings_date
4 ) avg_earnings, "earningdb"."faizanr_employee_earnings"
5 WHERE office_branch = avg_earnings.ob
6 GROUP BY office_branch]

```

SQL Ln 6, Col 24

Run again Explain Cancel Clear Create

Reuse query results
*Athena engine version 3 only

Query results Query stats

Completed Time in queue: 137 ms Run time: 1.678 sec Data scanned: 6.24 KB

Results (4)

Search rows

#	office_branch	earnings_range
1	Stanford	45506.6875
2	New York	46990.57142857143
3	Nashua	50212.83870967742
4	Scranton	51783.52

Query 2 : X Query 3 : X Query 4 : X Query 5 : X

```

1 SELECT office_branch, MIN(earnings) as min_earnings, MAX(earnings) as max_earnings, AVG(earnings) as avg_earnings, SUM(earnings) as total_earnings,
2 earnings_date
3 FROM "earningdb"."faizanr_employee_earnings"
4 GROUP BY office_branch, earnings_date
5 ORDER BY SUM(earnings) desc;

```

SQL Ln 2, Col 45

Run again Explain Cancel Clear Create

Reuse query results
*Athena engine version 3 only

Query results Query stats

Completed Time in queue: 106 ms Run time: 862 ms Data scanned: 5.30 KB

Results (28)

Search rows

#	office_branch	min_earnings	max_earnings	avg_earnings	total_earnings	earnings_date
1	Nashua	12752	98828	55832.74193548387	1730815	2023-05-18
2	New York	12098	97946	52606.107142857145	1472971	2023-05-18
3	Scranton	16667	97522	56834.84	1420871	2023-05-18
4	Stanford	12254	98733	51071.0625	817137	2023-05-18
5	Nashua	2098	9728	6099.8387096774195	189095	2022-02-14

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

Query 2

Query 3

Query 4

Query 5

Query 6

1 WITH previous_earnings 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 FROM
8 "earningsdb"."faizanr_employee_earnings"
9)
10 SELECT
11 emp_id,
12 earnings_date,
13 previous_earnings,
14 earnings,
15 (earnings - previous_earnings) / CAST(previous_earnings AS DOUBLE) * 100 AS percentage_change
16 FROM
17 previous_earnings
18 WHERE
19 earnings_date = '2023-05-18';

Completed

Time in queue: 157 msRun time: 763 msData scanned: 1 MB

Results (100)

CopyDownload

Search rows

#	emp_id	earnings_date	previous_earnings	earnings	percentage_change
1	138911	2023-05-18	6709	63631	848.4423908183038
2	391837	2023-05-18	4653	20161	333.2903505116269
3	413865	2023-05-18	7272	97946	1246.8921892189219
4	495667	2023-05-18	9462	27041	185.7852462481505
5	526540	2023-05-18	2716	15219	460.3460972017673
6	530134	2023-05-18	8483	19295	127.45490981963927
7	709884	2023-05-18	4495	50376	1020.7119021134594
8	713294	2023-05-18	8528	88995	943.5623827392121
9	779497	2023-05-18	8297	29060	250.24707725683984
10	932773	2023-05-18	7563	16667	120.37551236281898
11	984409	2023-05-18	2921	13727	369.9418007531667
12	160938	2023-05-18	3469	70028	1918.6797347938889
13	163409	2023-05-18	5323	38340	620.2705241405223
14	170637	2023-05-18	8950	82775	824.8603351955306
15	233136	2023-05-18	6499	51864	698.0304662255732
16	572204	2023-05-18	9168	33613	266.63394415357766
17	721091	2023-05-18	3557	32860	823.8122012932246
18	748190	2023-05-18	6157	58970	857.7716420334579