

WORK PROGRESS

By

NANDHINI R

ETL PROCESS

STEP 1: Extract the data from the CSV file

STEP 2: Transform the data [Ex: Convert the Column name into lower Case]

STEP 3: Load the data into Postgresql Database and Hadoop [HDFS]

ETL

etl.py - Sample - Visual Studio Code

File Edit Selection View Go Run Terminal Help

EXPLORER SAMPLE > hadoop > postgres checkeach.py docker-compose.yml etl.py

```
etl.py x checkeach.py
etl.py > ...
1 import psycopg2
2 from hdfs import InsecureClient
3 import pandas as pd
4 import io
5
6 # HDFS Configuration
7 hdfs_url = "http://localhost:9870"
8 client = InsecureClient(hdfs_url, user='nandhumidhun')
9
10 # PostgreSQL Configuration
11 pg_config = {
12     'dbname': 'project1',
13     'user': 'postgres',
14     'password': 'nandhu01',
15     'host': '127.0.0.1',
16     'port': '5432'
17 }
18
19 def extract_data(source_file):
20     # Example: Read CSV file
21     data = pd.read_csv(source_file)
22     return data
23
24 def transform_data(data):
25     # Example transformation: Convert column names to lowercase
```

PROBLEMS OUTPUT TERMINAL PORTS bash

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/ETLTTEST/Sample\$ python3 etl.py
Successful
Data successfully loaded to PostgreSQL
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/ETLTTEST/Sample\$

OUTLINE TIMELINE DEBUG CONSOLE

Filter (e.g. text, \exclude, \escape)

x 0 △ 0 ⌂ 0

Ln 80, Col 1 Spaces: 4 UTF-8 LF ↵ Python 3.10.12 64-bit

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $ hdfs dfs -ls /dockerhadoop  
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $ psql -h 127.0.0.1 -U postgres -d project1 -p 5432  
Password for user postgres:  
psql (16.3 (Ubuntu 16.3-1.pgdg22.04+1))  
SSL connection (protocol: TLSv1.3, cipher: TLS_AES_256_GCM_SHA384, compression: off)  
Type "help" for help.
```

```
project1=# SELECT * FROM sample;  
| id | name | email  
-----+-----+-----  
(0 rows)
```

```
project1=# SELECT * FROM sample;  
| id | name | email  
-----+-----+-----  
| 1 | Nandhu | nandhu@gmail.com  
| 2 | Midhun | midhun@gmail.com  
| 3 | xxxx | xx@gmail.com  
(3 rows)
```

```
project1=# exit;  
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $ hdfs dfs -ls /dockerhadoop
```

```
Found 1 items  
-rw-r--r-- 1 nandhumidhun supergroup 86 2024-06-20 15:05 /dockerhadoop/file.csv  
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $ hdfs dfs -cat /dockerhadoop/file.csv  
id,name,email  
1,Nandhu,nandhu@gmail.com  
2,Midhun,midhun@gmail.com  
3,xxxx,xx@gmail.com  
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $
```

DATA PIPELINE PROCESS

STEP 1: Insert the data from one database into another database

STEP 2: Fetch the data from API URL and insert into database

STEP 3: Read the CSV file and insert into the database

DATA PIPELINE

The screenshot shows a Visual Studio Code window with two terminal panes and a code editor.

Terminal 1: A PostgreSQL shell session connected to a local database on port 5432. It lists all users from the 'users' table.

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~$ psql -h 127.0.0.1 -U postgres -d project -p 5432
Password for user postgres:
psql (16.3 (Ubuntu 16.3-1.pgdg22.04+1))
SSL connection (protocol: TLSv1.3, cipher: TLS_AES_256_GCM_SHA384, compression: off)
Type "help" for help.

project=# SELECT * FROM users;
 id | name | email
----+-----+
(0 rows)

project=# SELECT * FROM users;
 id | name | email
----+-----+
 1 | Alice | alice@example.com
 2 | Nandhu | nandhu@example.com
 3 | Midhun | midhun@example.com
 4 | xxxxxx | xxxxx@example.com
 5 | yyyy | yyyy@gmail.com
(5 rows)

project=#[/]
```

Terminal 2: The same PostgreSQL session, showing the results of a SELECT query and the execution of an INSERT statement.

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~$ psql -h 127.0.0.1 -U postgres -d project -p 5432
SSL connection (protocol: TLSv1.3, cipher: TLS_AES_256_GCM_SHA384, compression: off)
Type "help" for help.

project1=# SELECT * FROM users;
 id | name | email
----+-----+
 1 | Alice | alice@example.com
 2 | Nandhu | nandhu@example.com
 3 | Midhun | midhun@example.com
 4 | xxxxxx | xxxxx@example.com
 5 | yyyy | yyyy@gmail.com
(5 rows)

project1=# INSERT INTO users(id,name,email) VALUES(5,'yyyy','yyyyy@gmail.com');
INSERT 0 1
project1=# SELECT * FROM users;
 id | name | email
----+-----+
 1 | Alice | alice@example.com
 2 | Nandhu | nandhu@example.com
 3 | Midhun | midhun@example.com
 4 | xxxxxx | xxxxx@example.com
 5 | yyyy | yyyy@gmail.com
(5 rows)[/]
```

Code Editor: The file `dbcopy.py` is open in the code editor. It contains Python code for a data pipeline, specifically for copying data from a source database to a destination database.

```
destination_cursor = destination_connection.cursor()
insert_query = "INSERT INTO users (id, name, email) VALUES (%s, %s, %s)"
for row in data:
    destination_cursor.execute(insert_query, row)

destination_connection.commit()
print("Success copy of data")
destination_cursor.close()
destination_connection.close()
except Exception as error:
    print("Error inserting data into destination:", error)

# Main function to execute the data transfer
if __name__ == "__main__":
    data = fetch_data_from_source()
    if data:
        insert_data_to_destination(data)
```

Bottom Status Bar: Shows the current file is `bash-python`, and the terminal output pane is active. Other tabs include PROBLEMS, OUTPUT, DEBUG CONSOLE, TERMINAL, and PORTS.

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/postgresql/python$ python3 dbcopy.py
Success copy of data
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/postgresql/python$ [/]
```

Bottom Right: Includes status icons for battery, signal, and network, along with the Python 3.10.12 64-bit interpreter information.

Ln 40, Col 40 Spaces:4 UTF-8 LF (Python 3.10.12 64-bit)

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~/nifi-2.0.0-M4

userid	id	title	completed
1	1	delectus aut autem	false
1	2	quis ut nam facilis et officia qui	false
1	3	fugiat veniam minus	false
1	4	et porro tempora	true
1	5	laboriosam mollitia et enim quasi adipisci quia provident illum	false
1	6	qui ullam ratione quibusdam voluptatem quia omnis	false
1	7	illo expedita consequatur quia in	false
1	8	quo adipisci enim quam ut ab	true
1	9	molestiae perspiciatis ipsa	false
1	10	illo est ratione doloremque quia maiores aut	true
1	11	vero rerum temporibus dolor	true
1	12	ipsa repellendus fugit nisi	true
1	13	et doloremque nulla	false
1	14	repellendus sunt dolores architecto voluptatum	true
1	15	ab voluptatum amet voluptas	true
1	16	accusamus eos facilis sint et aut voluptatem	true
1	17	quo laboriosam deleniti aut qui	true
1	18	dolorum est consequatur ea mollitia in culpa	false
1	19	molestiae ipsa aut voluptatibus pariatur dolor nihil	true
1	20	ullam nobis libero sapiente ad optio sint	true
2	21	suscipit repellat esse quibusdam voluptatem incident	false
2	22	distinctio vitae autem nihil ut molestias quo	true
2	23	et itaque necessitatibus maxime molestiae qui quas velit	false
2	24	adipisci non ad dicta qui amet quaerat doloribus ea	false
2	25	voluptas quo tenetur perspiciatis explicabo natus	true
2	26	aliquam aut quasi	true
2	27	veritatis pariatur delectus	true
2	28	nesciunt totam sit blanditiis sit	false
2	29	laborum aut in quam	false
2	30	nemo perspiciatis repellat ut dolor libero commodi blanditiis omnis	true
2	31	repudiandae totam in est sint facere fuga	false
2	32	earum doloribus ea doloremque quis	false
2	33	sint sit aut vero	false
2	34	porro aut necessitatibus eaque distinctio	false
2	35	repellendus veritatis molestias dicta incident	true
2	36	excepturi deleniti adipisci voluptatem et neque optio illum ad	true
2	37	sunt cum tempora	false
2	38	totam quia non	false
2	39	doloremque quibusdam asperiores libero corrupti illum qui omnis	false
2	40	totam atque quo nesciunt	true
3	41	aliquid amet impedit consequatur aspernatur placeat eaque fugiat suscipit	false
3	42	rerum perferendis error quia ut eveniet	false
3	43	tempore ut sint quis recusandae	true
3	44	cum debitis quis accusamus doloremque ipsa natus sapiente omnis	true
3	45	velit soluta adipisci molestias reiciendis harum	false
3	46	vel voluptatem repellat nihil placeat corporis	false
3	47	nam qui rerum fugiat accusamus	false
3	48	sit reprehenderit omnis quia	false
3	49	ut necessitatibus aut maiores debitis officia blanditiis velit et	false
3	50	cupiditate necessitatibus ullam aut quis dolor voluptate	true

File Edit Selection View Go Run Terminal Help

EXPLORER

POSTGRESQL

- python
 - _pycache_
 - api_data.csv
 - api.py
 - dbcopy.py
 - docker-compose.yml
 - dockerfile
 - file.csv
 - filesystem.py
 - main.py
 - sampleapi.py

```
python > file.csv
1 1,Alice,100
2 2,Bob,200
3 3,Charlie,300
4 4,Nandhu,400
```

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $ psql -h 127.0.0.1 -U postgres -d samples -p 5432
Password for user postgres:
psql (16.3 (Ubuntu 16.3-1.pgdg22.04+1))
SSL connection (protocol: TLSv1.3, cipher: TLS_AES_256_GCM_SHA384, compression: off)
Type "help" for help.

samples=# SELECT * FROM system_data;
 id | name   | value
----+-----+-----
(0 rows)

samples=# SELECT * FROM system_data;
 id | name   | value
----+-----+-----
 1 | Alice  | 100
 2 | Bob    | 200
 3 | Charlie| 300
 4 | Nandhu | 400
(4 rows)

samples=#

```

OUTLINE

No symbols found in document 'file.csv'

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

bash - python + □ ☰ ⌂ ⌂ ⌂ ⌂ ⌂

Success copy of data

- nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/postgresql/python\$ python3 filesystem.py
- nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/postgresql/python\$

Ln 4, Col 13 Spaces: 4 UTF-8 LF Plain Text

0 △ 0 ⌂ 0

SALES DATA ANALYSIS WITH HADOOP AND HIVE PROCESS

STEP 1: Load the sales data into HDFS

STEP 2: Create an external Hive table to map the CSV file

STEP 3: Load the data into the Hive table.

STEP 4: Save the query results to HDFS for further processing or reporting.

STEP 5: Optionally, clean up the Hive table and HDFS files if no longer needed.

SALES DATA ANALYSIS WITH HADOOP AND HIVE

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~          nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~          nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $ hdfs dfs -put '/home/nandhumidhun/Downloads/Sales_data.csv' /hadoop/sales_data/
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $ hdfs dfs -cat /hadoop/sales_analysis_results/top_products/*
Widget D600.0
Widget E550.0
Widget A500.0
Widget C450.0
Widget B350.0
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $ hdfs dfs -cat /hadoop/sales_analysis_results/total_sales/*
2450.0
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $ hdfs dfs -cat /hadoop/sales_analysis_results/average_sales/*
49.0
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $
```

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

sales_data					
sale_id	product_id	product_name	sale_date	sale_amount	customer_id
1	101	Widget A	2024-01-01	50.0	1001
2	102	Widget B	2024-01-02	35.0	1002
3	103	Widget C	2024-01-03	45.0	1003
4	104	Widget D	2024-01-04	60.0	1004
5	105	Widget E	2024-01-05	55.0	1005
6	101	Widget A	2024-01-06	50.0	1006
7	102	Widget B	2024-01-07	35.0	1007
8	103	Widget C	2024-01-08	45.0	1008
9	104	Widget D	2024-01-09	60.0	1009
10	105	Widget E	2024-01-10	55.0	1010
11	101	Widget A	2024-01-11	50.0	1011
12	102	Widget B	2024-01-12	35.0	1012
13	103	Widget C	2024-01-13	45.0	1013
14	104	Widget D	2024-01-14	60.0	1014
15	105	Widget E	2024-01-15	55.0	1015
16	101	Widget A	2024-01-16	50.0	1016
17	102	Widget B	2024-01-17	35.0	1017
18	103	Widget C	2024-01-18	45.0	1018
19	104	Widget D	2024-01-19	60.0	1019
20	105	Widget E	2024-01-20	55.0	1020
21	101	Widget A	2024-01-21	50.0	1021
22	102	Widget B	2024-01-22	35.0	1022
23	103	Widget C	2024-01-23	45.0	1023
24	104	Widget D	2024-01-24	60.0	1024
25	105	Widget E	2024-01-25	55.0	1025
26	101	Widget A	2024-01-26	50.0	1026
27	102	Widget B	2024-01-27	35.0	1027
28	103	Widget C	2024-01-28	45.0	1028
29	104	Widget D	2024-01-29	60.0	1029
30	105	Widget E	2024-01-30	55.0	1030
31	101	Widget A	2024-01-31	50.0	1031
32	102	Widget B	2024-02-01	35.0	1032
33	103	Widget C	2024-02-02	45.0	1033
34	104	Widget D	2024-02-03	60.0	1034
35	105	Widget E	2024-02-04	55.0	1035
36	101	Widget A	2024-02-05	50.0	1036
37	102	Widget B	2024-02-06	35.0	1037
38	103	Widget C	2024-02-07	45.0	1038
39	104	Widget D	2024-02-08	60.0	1039

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

```

INFO : Job running in-process (local Hadoop)
INFO : 2024-06-28 16:57:26,211 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 0.21 sec
INFO : MapReduce Total cumulative CPU time: 210 msec
INFO : Ended Job = job_local822845722_0010
INFO : MapReduce Jobs Launched:
INFO : Stage-Stage-1: Cumulative CPU: 0.21 sec HDFS Read: 46724 HDFS Write: 162 HDFS EC Read: 0 SUCCESS
INFO : Total MapReduce CPU Time Spent: 210 msec
INFO : Completed executing command(queryId=nandhumidhun_20240628165724_62740a7a-072e-45d9-91c5-932865291d2e); Time taken: 1.403 seconds
+-----+
| total_sales |
+-----+
| 2450.0 |
+-----+
1 row selected (1.601 seconds)

0: jdbc:hive2://localhost:10000> SELECT AVG(sale_amount) AS average_sales FROM sales_data;
INFO : Compiling command(queryId=nandhumidhun_20240628165745_67b656ee-f617-4099-b59a-edeb0997b4ce): SELECT AVG(sale_amount) AS average_sales FROM sales_data
INFO : No Stats for hadoop@sales_data, Columns: sale_amount
INFO : Semantic Analysis Completed (retrying = false)
INFO : Created Hive schema: Schema(fieldSchemas:[FieldSchema(name:average_sales, type:double, comment:null)], properties:null)
INFO : Completed compiling command(queryId=nandhumidhun_20240628165745_67b656ee-f617-4099-b59a-edeb0997b4ce); Time taken: 0.194 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=nandhumidhun_20240628165745_67b656ee-f617-4099-b59a-edeb0997b4ce): SELECT AVG(sale_amount) AS average_sales FROM sales_data
WARN : Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. tez) or using Hive 1.X releases.
INFO : Query ID = nandhumidhun_20240628165745_67b656ee-f617-4099-b59a-edeb0997b4ce
INFO : Total jobs = 1
INFO : Launching Job 1 out of 1
INFO : Starting task [Stage-1:MAPRED] in serial mode
INFO : Number of reduce tasks determined at compile time: 1
INFO : In order to change the average load for a reducer (in bytes):
INFO :   set hive.exec.reducers.bytes.per.reducer=<number>
INFO : In order to limit the maximum number of reducers:
INFO :   set hive.exec.reducers.max=<number>
INFO : In order to set a constant number of reducers:
INFO :   set mapreduce.job.reduces=<number>
INFO : number of splits:1
INFO : Submitting tokens for job: job_local1785308139_0011
INFO : Executing with tokens: []
INFO : The url to track the job: http://localhost:8080/
INFO : Job running in-process (local Hadoop)
INFO : 2024-06-28 16:57:47,173 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 0.19 sec
INFO : MapReduce Total cumulative CPU time: 190 msec
INFO : Ended Job = job_local1785308139_0011
INFO : MapReduce Jobs Launched:
INFO : Stage-Stage-1: Cumulative CPU: 0.19 sec HDFS Read: 50506 HDFS Write: 162 HDFS EC Read: 0 SUCCESS
INFO : Total MapReduce CPU Time Spent: 190 msec
INFO : Completed executing command(queryId=nandhumidhun_20240628165745_67b656ee-f617-4099-b59a-edeb0997b4ce); Time taken: 1.421 seconds
+-----+
| average_sales |
+-----+
| 49.0 |
+-----+
1 row selected (1.653 seconds)

0: jdbc:hive2://localhost:10000>
```

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

WARN : Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. tez) or using Hive 1.X releases.

INFO : Query ID = nandhumidhun_20240628165858_d9e8389c-ae92-4da7-a911-4104db755d21
 INFO : Total jobs = 2
 INFO : Launching Job 1 out of 2
 INFO : Starting task [Stage-1:MAPRED] in serial mode
 INFO : Number of reduce tasks not specified. Estimated from input data size: 1
 INFO : In order to change the average load for a reducer (in bytes):
 INFO : set hive.exec.reducers.bytes.per.reducer=<number>
 INFO : In order to limit the maximum number of reducers:
 INFO : set hive.exec.reducers.max=<number>
 INFO : In order to set a constant number of reducers:
 INFO : set mapreduce.job.reduces=<number>
 INFO : number of splits:1
 INFO : Submitting tokens for job: job_local1275533349_0012
 INFO : Executing with tokens: []
 INFO : The url to track the job: http://localhost:8080/
 INFO : Job running in-process (local Hadoop)
 INFO : 2024-06-28 16:58:59,747 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 0.21 sec
 INFO : MapReduce Total cumulative CPU time: 210 msec
 INFO : Ended Job = job_local1275533349_0012
 INFO : Launching Job 2 out of 2
 INFO : Starting task [Stage-2:MAPRED] in serial mode
 INFO : Number of reduce tasks determined at compile time: 1
 INFO : In order to change the average load for a reducer (in bytes):
 INFO : set hive.exec.reducers.bytes.per.reducer=<number>
 INFO : In order to limit the maximum number of reducers:
 INFO : set hive.exec.reducers.max=<number>
 INFO : In order to set a constant number of reducers:
 INFO : set mapreduce.job.reduces=<number>
 INFO : number of splits:1
 INFO : Submitting tokens for job: job_local1032948224_0013
 INFO : Executing with tokens: []
 INFO : The url to track the job: http://localhost:8080/
 INFO : Job running in-process (local Hadoop)
 INFO : 2024-06-28 16:59:01,193 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 0.12 sec
 INFO : MapReduce Total cumulative CPU time: 120 msec
 INFO : Ended Job = job_local1032948224_0013
 INFO : MapReduce Jobs Launched:
 INFO : Stage-Stage-1: Cumulative CPU: 0.21 sec HDFS Read: 54288 HDFS Write: 162 HDFS EC Read: 0 SUCCESS
 INFO : Stage-Stage-2: Cumulative CPU: 0.12 sec HDFS Read: 54288 HDFS Write: 162 HDFS EC Read: 0 SUCCESS
 INFO : Total MapReduce CPU Time Spent: 330 msec
 INFO : Completed executing command(queryId=nandhumidhun_20240628165858_d9e8389c-ae92-4da7-a911-4104db755d21); Time taken: 2.839 seconds

product_name	total_sales
Widget D	600.0
Widget E	550.0
Widget A	500.0
Widget C	450.0
Widget B	350.0

5 rows selected (3.052 seconds)

0: jdbc:hive2://localhost:10000>

DATA ANALYTICS PLATFORM PROCESS

STEP 1: Upload Data to HDFS

STEP 2: Create Hive Tables

STEP 3: Load Data into Hive Tables

STEP 4: Join Tables to Create a Unified View

STEP 5: Sample Queries for Analytics

STEP 6: Data Visualization

DATA ANALYTICS PLATFORM

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~$ hdfs dfs -cat /output/*
```

Category,avg(Discount)

Furniture,0.1

Office Supplies,0.05

Technology,0.0999999999999999

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~$ hdfs dfs -cat /output1/*
```

ProductID,Category,Discount,Amount,Profit,total_amount

P001,Office Supplies,0.05,100.5,20.25,2035.125

P002,Technology,0.1,1500.0,300.0,450000.0

P003,Furniture,0.08,250.75,50.15,12575.1125

P004,Office Supplies,0.02,20.25,4.5,91.125

P005,Technology,0.15,1200.0,180.0,216000.0

P006,Office Supplies,0.05,50.8,10.2,518.16

P007,Office Supplies,0.03,15.5,2.1,32.550000000000004

P008,Furniture,0.12,900.0,120.0,108000.0

P009,Technology,0.05,35.0,5.25,183.75

P010,Office Supplies,0.1,200.0,30.0,6000.0

Activities Terminal Jun 29 15:00

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~
```

INFO : Compiling command(queryId=nandhumidhun_20240629145743_5835f7f2-ee31-4720-95d2-ee2771c4ad3b); SELECT * FROM users
INFO : No Stats for hadoopusers, Columns: user_id, name, location, age
INFO : Semantic Analysis Completed (retrial = false)
INFO : Created Hive schema: Schema(fieldschemas:[FieldSchema(name:users.user_id, type:int, comment:null), FieldSchema(name:users.name, type:string, comment:null), FieldSchema(name:users.age, type:int, comment:null), FieldSchema(name:users.location, type:string, comment:null)], properties:null)
INFO : Completed compiling command(queryId=nandhumidhun_20240629145743_5835f7f2-ee31-4720-95d2-ee2771c4ad3b); Time taken: 0.112 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=nandhumidhun_20240629145743_5835f7f2-ee31-4720-95d2-ee2771c4ad3b); SELECT * FROM users
INFO : Completed executing command(queryId=nandhumidhun_20240629145743_5835f7f2-ee31-4720-95d2-ee2771c4ad3b); Time taken: 0.0 seconds

	users.user_id	users.name	users.age	users.location
1	User A	40	California	
2	User B	55	Washington	
3	User C	29	Washington	
4	User D	33	California	
5	User E	38	New Jersey	
6	User F	22	Texas	
7	User G	65	California	
8	User H	43	California	
9	User I	44	Florida	
10	User J	36	Florida	
11	User K	19	Washington	
12	User L	19	New Jersey	
13	User M	40	Washington	
14	User N	52	California	
15	User O	24	California	
16	User P	63	New York	
17	User Q	46	California	
18	User R	31	California	
19	User S	59	New Jersey	
20	User T	52	California	
21	User U	33	Florida	
22	User V	44	Florida	
23	User W	35	Washington	
24	User X	42	Texas	
25	User Y	33	Texas	
26	User Z	64	New Jersey	
27	User [42	New York	
28	User \	55	Florida	
29	User]	66	Texas	
30	User ^	40	Washington	
31	User _	52	New Jersey	
32	User -	60	Florida	
33	User a	29	California	
34	User b	52	Florida	
35	User c	24	New York	
36	User d	40	Texas	
37	User e	23	Washington	
38	User f	50	California	
39	User g	62	New Jersey	
40	User h	58	New Jersey	
41	User i	68	Florida	
42	User j	64	Florida	

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~



nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~ x

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~ x

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~ x

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~ x

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $ psql -h 127.0.0.1 -p 5432 -U postgres -d project
Password for user postgres:
psql (16.3 (Ubuntu 16.3.1.pgdg22.04+1))
SSL connection (protocol: TLSv1.3, cipher: TLS_AES_256_GCM_SHA384, compression: off)
Type "help" for help.
```

```
project=# SELECT * FROM transactional_data;
transaction_id | product_id | user_id | quantity | price |      timestamp
-----+-----+-----+-----+-----+-----+
 1 | 144 | 48 | 2 | 47.03 | 2023-06-01 11:00:00
 2 | 115 | 43 | 5 | 40.47 | 2023-06-01 12:00:00
 3 | 140 | 5 | 5 | 41.04 | 2023-06-01 13:00:00
 4 | 101 | 28 | 3 | 46.53 | 2023-06-01 14:00:00
 5 | 130 | 3 | 1 | 28.55 | 2023-06-01 15:00:00
 6 | 130 | 41 | 2 | 24.74 | 2023-06-01 16:00:00
 7 | 145 | 40 | 5 | 38.94 | 2023-06-01 17:00:00
 8 | 150 | 50 | 5 | 25.32 | 2023-06-01 18:00:00
 9 | 102 | 48 | 5 | 25.34 | 2023-06-01 19:00:00
10 | 106 | 43 | 4 | 29.71 | 2023-06-01 20:00:00
11 | 150 | 14 | 4 | 48.42 | 2023-06-01 21:00:00
12 | 144 | 37 | 2 | 40.56 | 2023-06-01 22:00:00
13 | 115 | 25 | 4 | 42.06 | 2023-06-01 23:00:00
14 | 138 | 47 | 1 | 34.13 | 2023-06-02 00:00:00
15 | 135 | 34 | 4 | 26.19 | 2023-06-02 01:00:00
16 | 149 | 29 | 5 | 14.69 | 2023-06-02 02:00:00
17 | 125 | 50 | 1 | 11.06 | 2023-06-02 03:00:00
18 | 103 | 1 | 2 | 30.46 | 2023-06-02 04:00:00
19 | 137 | 6 | 1 | 18.44 | 2023-06-02 05:00:00
20 | 115 | 46 | 5 | 14.46 | 2023-06-02 06:00:00
21 | 125 | 24 | 3 | 14.53 | 2023-06-02 07:00:00
22 | 123 | 2 | 3 | 27.28 | 2023-06-02 08:00:00
23 | 124 | 20 | 3 | 35.73 | 2023-06-02 09:00:00
24 | 132 | 34 | 5 | 32.12 | 2023-06-02 10:00:00
25 | 135 | 31 | 2 | 11.98 | 2023-06-02 11:00:00
26 | 124 | 29 | 3 | 35.93 | 2023-06-02 12:00:00
27 | 141 | 19 | 5 | 45.49 | 2023-06-02 13:00:00
28 | 111 | 5 | 2 | 31.96 | 2023-06-02 14:00:00
29 | 102 | 17 | 5 | 11.88 | 2023-06-02 15:00:00
30 | 119 | 47 | 3 | 30.39 | 2023-06-02 16:00:00
31 | 120 | 32 | 3 | 14.6 | 2023-06-02 17:00:00
32 | 106 | 29 | 5 | 19.5 | 2023-06-02 18:00:00
33 | 101 | 20 | 1 | 21.37 | 2023-06-02 19:00:00
34 | 108 | 20 | 2 | 27.57 | 2023-06-02 20:00:00
35 | 106 | 8 | 4 | 31.43 | 2023-06-02 21:00:00
36 | 134 | 39 | 2 | 28.49 | 2023-06-02 22:00:00
37 | 123 | 37 | 2 | 44.61 | 2023-06-02 23:00:00
38 | 120 | 35 | 5 | 39.35 | 2023-06-03 00:00:00
39 | 142 | 50 | 1 | 14.03 | 2023-06-03 01:00:00
40 | 110 | 42 | 5 | 18.9 | 2023-06-03 02:00:00
41 | 121 | 50 | 5 | 19.85 | 2023-06-03 03:00:00
42 | 105 | 33 | 4 | 19.71 | 2023-06-03 04:00:00
43 | 130 | 21 | 1 | 44.68 | 2023-06-03 05:00:00
```

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~ x nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~ x nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~ x nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~ x
```

ORDER BY total_sales DESC
WARN : Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. tez) or using Hive 1.X releases.
INFO : Query ID = nandhumidhun_20240701161819_967810f3-103b-4d21-bf39-355be7f42a89
INFO : Total jobs = 2

INFO : Launching Job 1 out of 2
INFO : Starting task [Stage-1:MAPRED] in serial mode
INFO : Number of reduce tasks not specified. Estimated from input data size: 1
INFO : In order to change the average load for a reducer (in bytes):
INFO : set hive.exec.reducers.bytes.per.reducer=<number>
INFO : In order to limit the maximum number of reducers:
INFO : set hive.exec.reducers.max=<number>
INFO : In order to set a constant number of reducers:
INFO : set mapreduce.job.reduces=<number>
INFO : number of splits:1
INFO : Submitting tokens for job: job_local1577149079_0002
INFO : Executing with tokens: []
INFO : The url to track the job: http://localhost:8080/
INFO : Job running in-process (local Hadoop)
INFO : 2024-07-01 16:18:20,908 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 0.49 sec
INFO : MapReduce Total cumulative CPU time: 490 msec
INFO : Ended Job = job_local1577149079_0002

INFO : Launching Job 2 out of 2
INFO : Starting task [Stage-2:MAPRED] in serial mode
INFO : Number of reduce tasks determined at compile time: 1
INFO : In order to change the average load for a reducer (in bytes):
INFO : set hive.exec.reducers.bytes.per.reducer=<number>
INFO : In order to limit the maximum number of reducers:
INFO : set hive.exec.reducers.max=<number>
INFO : In order to set a constant number of reducers:
INFO : set mapreduce.job.reduces=<number>
INFO : number of splits:1
INFO : Submitting tokens for job: job_local1793540612_0003
INFO : Executing with tokens: []
INFO : The url to track the job: http://localhost:8080/
INFO : Job running in-process (local Hadoop)
INFO : 2024-07-01 16:18:22,336 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 0.31 sec
INFO : MapReduce Total cumulative CPU time: 310 msec
INFO : Ended Job = job_local1793540612_0003

INFO : MapReduce Jobs Launched:
INFO : Stage-Stage-1: Cumulative CPU: 0.49 sec HDFS Read: 15284 HDFS Write: 0 HDFS EC Read: 0 SUCCESS
INFO : Stage-Stage-2: Cumulative CPU: 0.31 sec HDFS Read: 15284 HDFS Write: 0 HDFS EC Read: 0 SUCCESS
INFO : Total MapReduce CPU Time Spent: 800 msec
INFO : Completed executing command(queryId=nandhumidhun_20240701161819_967810f3-103b-4d21-bf39-355be7f42a89); Time taken: 2.961 seconds

category	total_sales
Hardware	1534.0
Apparel	1435.0
Electronics	958.0
Home	856.0

4 rows selected (3.33 seconds)

0: jdbc:hive2://localhost:10000> ■

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~      nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~      nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~      nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~
INFO : Number of reduce tasks not specified. Estimated from input data size: 1
INFO : In order to change the average load for a reducer (in bytes):
INFO :   set hive.exec.reducers.bytes.per.reducer=<number>
INFO : In order to limit the maximum number of reducers:
INFO :   set hive.exec.reducers.max=<number>
INFO : In order to set a constant number of reducers:
INFO :   set mapreduce.job.reduces=<number>
INFO : number of splits:1
INFO : Submitting tokens for job: job_local988817042_0004
INFO : Executing with tokens: []
INFO : The url to track the job: http://localhost:8080/
INFO : Job running in-process (local Hadoop)
INFO : 2024-07-01 16:18:44,321 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 0.46 sec
INFO : MapReduce Total cumulative CPU time: 460 msec
INFO : Ended Job = job_local988817042_0004
INFO : Launching Job 2 out of 2
```

```
INFO : Starting task [Stage-2:MAPRED] in serial mode
INFO : Number of reduce tasks determined at compile time: 1
INFO : In order to change the average load for a reducer (in bytes):
INFO :   set hive.exec.reducers.bytes.per.reducer=<number>
INFO : In order to limit the maximum number of reducers:
INFO :   set hive.exec.reducers.max=<number>
INFO : In order to set a constant number of reducers:
INFO :   set mapreduce.job.reduces=<number>
INFO : number of splits:1
INFO : Submitting tokens for job: job_local931053625_0005
INFO : Executing with tokens: []
INFO : The url to track the job: http://localhost:8080/
INFO : Job running in-process (local Hadoop)
INFO : 2024-07-01 16:18:45,760 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 0.21 sec
INFO : MapReduce Total cumulative CPU time: 210 msec
INFO : Ended Job = job_local931053625_0005
INFO : MapReduce Jobs Launched:
INFO : Stage-Stage-1: Cumulative CPU: 0.46 sec  HDFS Read: 22926 HDFS Write: 0 HDFS EC Read: 0 SUCCESS
INFO : Stage-Stage-2: Cumulative CPU: 0.21 sec  HDFS Read: 22926 HDFS Write: 0 HDFS EC Read: 0 SUCCESS
INFO : Total MapReduce CPU Time Spent: 670 msec
INFO : Completed executing command(queryId=nandhumidhun_20240701161842_f39fb62a-f0d5-4174-a904-bd54ddda5e5b); Time taken: 2.913 seconds
```

product_name	total_sales
"Product i"	569.0
"Product O"	438.0
"Product F"	344.0
"Product r"	317.0
"Product B"	273.0
"Product T"	240.0
"Product X"	216.0
"Product h"	210.0
"Product m"	195.0
"Product l"	176.0

10 rows selected (3.25 seconds)
0: jdbc:hive2://localhost:10000>

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

```

INFO : Ended Job = job_local334021857_0006
INFO : Launching Job 2 out of 2
INFO : Starting task [Stage-2:MAPRED] in serial mode
INFO : Number of reduce tasks determined at compile time: 1
INFO : In order to change the average load for a reducer (in bytes):
INFO :   set hive.exec.reducers.bytes.per.reducer=<number>
INFO : In order to limit the maximum number of reducers:
INFO :   set hive.exec.reducers.max=<number>
INFO : In order to set a constant number of reducers:
INFO :   set mapreduce.job.reduces=<number>
INFO : number of splits:1
INFO : Submitting tokens for job: job_local1849434121_0007
INFO : Executing with tokens: []
INFO : The url to track the job: http://localhost:8080/
INFO : Job running in-process (local Hadoop)
INFO : 2024-07-01 16:21:00,142 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 0.32 sec
INFO : MapReduce Total cumulative CPU time: 320 msec
INFO : Ended Job = job_local1849434121_0007
INFO : MapReduce Jobs Launched:
INFO : Stage-Stage-1: Cumulative CPU: 0.33 sec   HDFS Read: 30568 HDFS Write: 0 HDFS EC Read: 0 SUCCESS
INFO : Stage-Stage-2: Cumulative CPU: 0.32 sec   HDFS Read: 30568 HDFS Write: 0 HDFS EC Read: 0 SUCCESS
INFO : Total MapReduce CPU Time Spent: 650 msec
INFO : Completed executing command(queryId=nandhumidhun_20240701162056_15ee8806-5d17-496b-acf2-50eab4cb97a4); Time taken: 2.904 seconds
+-----+
| age | user_count |
+-----+
| 22  | 1
| 23  | 1
| 24  | 1
| 29  | 2
| 30  | 1
| 33  | 4
| 38  | 1
| 40  | 1
| 42  | 1
| 43  | 1
| 44  | 1
| 46  | 1
| 50  | 1
| 52  | 4
| 54  | 1
| 55  | 2
| 58  | 2
| 59  | 1
| 60  | 1
| 62  | 1
| 64  | 1
| 65  | 1
| 66  | 1
| 68  | 1
+-----+
24 rows selected (3.219 seconds)
0: jdbc:hive2://localhost:10000>
```

Home Data_Visualizati... data_json userId Analysis Dash...

Premium Trial - 14 days left
UPGRADE

userId Analysis Dashboard

A complete analysis on userId

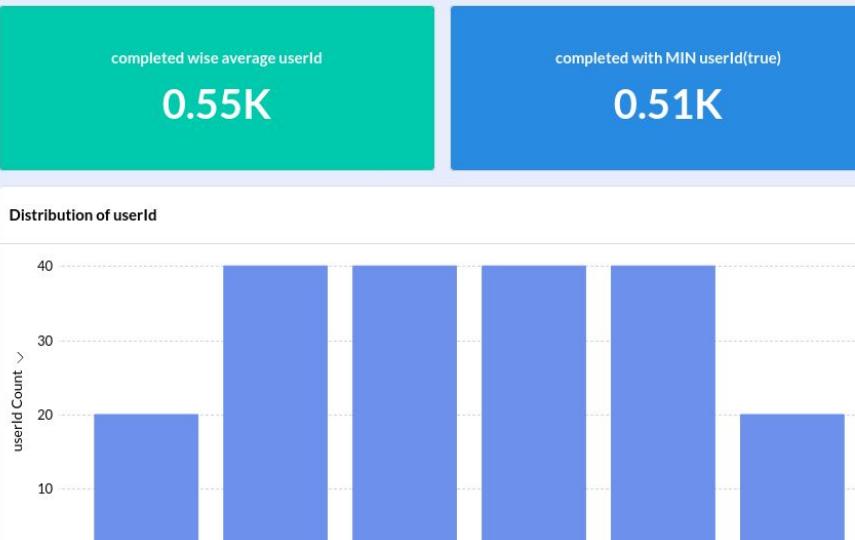
Completed wise average userId: 0.55K

Completed with MIN userId(true): 0.51K

Completed with MAX userId(false): 0.59K

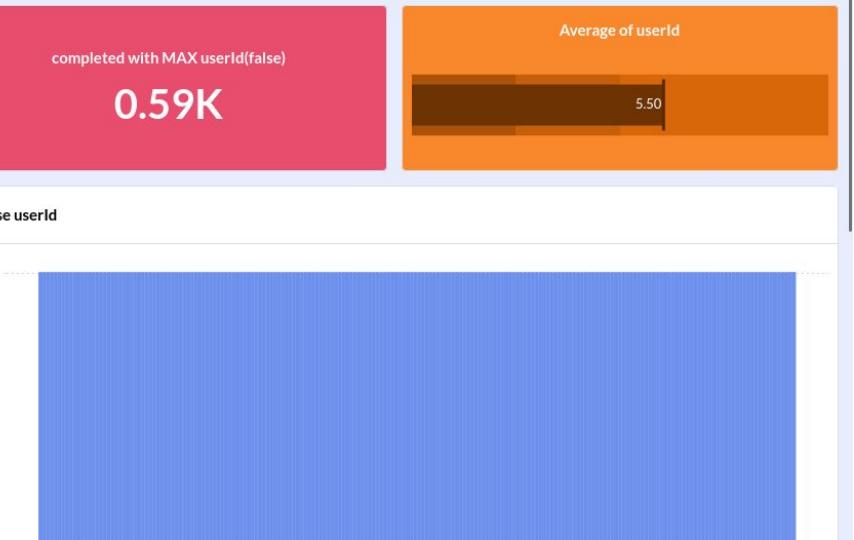
Average of userId: 5.50

Distribution of userId



userId Range	Userid Count
0 to 2	~20
2 to 4	~40
4 to 6	~40
6 to 8	~40
8 to 10	~40
10 to 12	~20

id-wise userId



id	Userid Count
1	1
2	1
3	1
4	1
5	1
6	1
7	1
8	1
9	1
10	1
11	1
12	1
13	1
14	1
15	1
16	1
17	1
18	1
19	1
20	1
21	1
22	1
23	1
24	1
25	1
26	1
27	1
28	1
29	1
30	1
31	1
32	1
33	1
34	1
35	1
36	1
37	1
38	1
39	1
40	1
41	1
42	1
43	1
44	1
45	1
46	1
47	1
48	1
49	1
50	1
51	1
52	1
53	1
54	1
55	1
56	1
57	1
58	1
59	1
60	1
61	1
62	1
63	1
64	1
65	1
66	1
67	1
68	1
69	1
70	1
71	1
72	1
73	1
74	1
75	1
76	1
77	1
78	1
79	1
80	1
81	1
82	1
83	1
84	1
85	1
86	1
87	1
88	1
89	1
90	1
91	1
92	1
93	1
94	1
95	1
96	1
97	1
98	1
99	1
100	1
101	1
102	1
103	1
104	1
105	1
106	1
107	1
108	1
109	1
110	1
111	1
112	1
113	1
114	1
115	1
116	1
117	1
118	1
119	1
120	1
121	1
122	1
123	1
124	1
125	1
126	1
127	1
128	1
129	1
130	1
131	1
132	1
133	1
134	1
135	1
136	1
137	1
138	1
139	1
140	1
141	1
142	1
143	1
144	1
145	1
146	1
147	1
148	1
149	1
150	1
151	1
152	1
153	1
154	1
155	1
156	1
157	1
158	1
159	1
160	1
161	1
162	1
163	1
164	1
165	1
166	1
167	1
168	1
169	1
170	1
171	1

completed-wise userId

Top 10 title by userId

Here is your Smart Chat (Ctrl+Space)

SEMI-STRUCTURED DATA PIPELINE PROCESS

STEP 1: Create a pipeline that ingests semi-structured data

STEP 2: Transform and normalize the data into a unified format.

STEP 3: Load the processed data into a HDFS

SEMI-STRUCTURED DATA PIPELINE

Activities Visual Studio Code

File Edit Selection View Go Run Terminal Help

data.json - Semi_Data_Pipeline - Visual Studio Code

Jul 2 16:15

EXPLORER

SEMI_DATA_PIPELINE

- automate_data.csv
- csvfile.py
- Data_process.py
- data.json
- Ingestion.py
- Quation.txt

datafile.py data.json automate_data.csv Data_process.py

```
[{"id": 1, "userId": 1, "title": "delectus aut autem", "completed": false}, {"id": 2, "userId": 1, "title": "quis ut nam facilis et officia qui", "completed": false}, {"id": 3, "userId": 1, "title": "fugiat veniam minus", "completed": false}, {"id": 4, "userId": 1, "title": "et porro tempora", "completed": true}, {"id": 5, "userId": 1, "title": "laboriosam mollitia et enim quasi adipisci quia provident illum", "completed": false}, {"id": 6, "userId": 1, "title": "qui ullam ratione quibusdam voluptatem quia omnis", "completed": false}, {"id": 7, "userId": 1, "title": "illo expedita consequatur quia in", "completed": false}, {"id": 8, "userId": 1, "title": "quo adipisci enim quam ut ab", "completed": true}]
```

OUTLINE

TIMELINE

DEBUG CONSOLE

Filter (e.g. text, exclude, \escape)

Ln 20, Col 6 Spaces: 4 UTF-8 JSON

Browse Directory

/user/hive/warehouse

Show 25 entries Search:

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	<input type="checkbox"/>
<input type="checkbox"/>	-rw-r--r--	nandhumidhun	supergroup	27.65 KB	Jul 02 12:00	1	128 MB	data.json	<input type="checkbox"/>
<input type="checkbox"/>	drwxrwxr-x	anonymous	supergroup	0 B	Jul 01 11:36	0	0 B	hadoop.db	<input type="checkbox"/>
<input type="checkbox"/>	drwxrwxr-x	anonymous	supergroup	0 B	Jun 27 13:46	0	0 B	hive.db	<input type="checkbox"/>
<input type="checkbox"/>	drwxrwxr-x	anonymous	supergroup	0 B	Jun 30 13:21	0	0 B	temp_products	<input type="checkbox"/>
<input type="checkbox"/>	drwxrwxr-x	anonymous	supergroup	0 B	Jun 30 13:21	0	0 B	temp_transactions	<input type="checkbox"/>
<input type="checkbox"/>	drwxrwxr-x	anonymous	supergroup	0 B	Jun 30 19:31	0	0 B	transaction_details_output	<input type="checkbox"/>

Showing 1 to 6 of 6 entries

AUTOMATED DOCUMENT PROCESSING PROCESS

STEP 1: Build an AI-driven solution that extracts relevant information from unstructured documents

STEP 2: Use natural language processing (NLP) techniques to identify key entities

STEP 3: Showcase how this automation reduces manual effort, improves accuracy, and accelerates decision-making

AUTOMATED DOCUMENT PROCESSING

The screenshot shows a code editor interface with the following details:

- EXPLORER** sidebar:
 - SEMI_DATA_PIPELINE** folder:
 - __pycache__
 - kafka
 - load_kafka_hadoop.py
 - seemessage.py
 - Unstructured_data
 - Post_Semi_data
 - Unstructured_data
 - automate_data.csv
 - automate_data.json
 - automate.csv
 - convertjson.py
 - Data_process.py
 - data.json
 - demo1.pdf
 - extracted_entities.csv
 - final_pdf.py
 - pdf.py
 - unstructured_pdf.py
- py** tab selected in the top bar.
- Code Editor Content:**

```
19 def load_to_hdfs(data, hdfs_path):
20     output.seek(0)
21     writer.write(output.getvalue())
22
23     def etl_process(source_file, hdfs_path):
24         # Extract
25         data = extract_data(source_file)
26
27         # Transform
28         data = transform_data(data)
29
30         # Load to HDFS
31         load_to_hdfs(data, hdfs_path)
32
33         # Load to PostgreSQL
34
35     print("Successfull")
```
- PROBLEMS**, **OUTPUT**, **TERMINAL**, **PORTS** tabs at the bottom.
- TERMINAL** tab content:

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/Semi_Data_Pipeline/Unstructured_data$ python3 pdf.py
[('January 1, 2024', 'DATE'), ('100.00', 'MONEY'), ('December 31, 2023', 'DATE'), ('50.50', 'MONEY'), ('December 30, 2023', 'DATE'), ('150.50', 'MONEY'),
 ('December 39, 2023', 'DATE'), ('52.50', 'MONEY')]
o nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/Semi_Data_Pipeline/Unstructured_data$
```
- PDF Preview**: A PDF viewer window titled "sample.pdf" showing the extracted content from the document.
- Page 1**:

/home/nandhumidhun/Semi_Data_Pipeline
Sample Document with Dates and Amounts:

Date: January 1, 2024
Amount: \$100.00

Another Date: December 31, 2023
Another Amount: \$50.50

Another Date: December 30, 2023
Another Amount: \$150.50

Another Date: December 39, 2023
Another Amount: \$52.50
- Bottom Status Bar**:

Ln 41, Col 22 Spaces: 4 UTF-8 LF Python 3.10.12 64-bit

The screenshot shows a Python development environment in Visual Studio Code (VS Code) with the following details:

- EXPLORER**: Shows the project structure under `SEMI_DATA_PIPELINE`. The `Unstructured_data` folder contains files: `final_pdf.py`, `pdf.py`, `seemessage.py`, `load_kafka_hadoop.py`, `gengpt.py`, `loadhadoop.py`, `unstructured_pdf.py`, and `Final_pdf.py`.
- CODE**: The `final_pdf.py` file is open, containing code to handle exceptions and write to HDFS:

```
25
26     except Exception as e:
27         print(f"Error writing to HDFS: {e}")
28
29
30
31     # Example: Write back to HDFS
32     df_transformed.write.csv("hdfs://localhost:9000/pdf_Output/transformed_data.csv", mode="overwrite")
```

- TERMINAL**: Shows command-line logs from running the `python3 pdf.py` and `python3 final_pdf.py` scripts. It also shows a native code loader warning and a successful `hdfs dfs -cat` command outputting transformed data.
- STATUS BAR**: Shows the current line (Ln 32), column (Col 63), and other system information like spaces, tabs, and file type (Python).

The screenshot shows a dark-themed Python development environment, likely PyCharm, with the following interface elements:

- Left Sidebar (EXPLORER):** Shows the project structure under "SEMI_DATA_PIPELINE".
- Top Bar:** Contains tabs for "csvfile.py", "automate_data.json", "Ingestion.py", "automate_data.csv", "kafka_replication.py", "kafka_postgre.py", "kafka.py", "seemessage.py", and "load_kafka_hadoop".
- Code Editor:** Displays the "csvfile.py" script. The code performs data cleaning, aggregation, and saves the results to HDFS.
- Terminal:** Shows the command "python3 csvfile.py" being run and its output. It includes several WARN messages from Spark and HDFS.
- Bottom Status Bar:** Shows the current line (Ln 43), column (Col 65), and other system information like "Spaces: 4", "UTF-8", "LF", "Python 3.10.12 64-bit".

```
16 # Clean and transform data
17 df_cleaned = df.dropna().select("ProductID", "Category", "Discount", "Amount", "Profit", "Quantity")
18
19 # Example: Aggregation and computation
20 df_transformed1 = df_cleaned.withColumn("total_amount", df_cleaned["Amount"] * df_cleaned["Profit"])
21 df_transformed = df_cleaned.groupBy("Category").agg({"Discount": 'avg'}).orderBy("Category")
22 # Write processed data back to HDFS
23 try:
24     df_transformed1.write.format("csv").mode("overwrite").option("header", "true").save("hdfs://localhost:9000/output1")
25     df_transformed.write.format("csv").mode("overwrite").option("header", "true").save("hdfs://localhost:9000/output")
26 except Exception as e:
27     print(f"Error writing to HDFS: {e}")
28
29
```

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/Semi_Data_Pipeline/Post_Semi_data$ python3 csvfile.py
24/07/05 10:50:42 WARN Utils: Your hostname, nandhumidhun-HP-Laptop-15-bs0xx resolves to a loopback address: 127.0.1.1; using 192.168.1.36 instead (on interface wlo1)
24/07/05 10:50:42 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
24/07/05 10:50:43 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
24/07/05 10:51:38 WARN Instrumentation: [8ede5b19] regParam is zero, which might cause numerical instability and overfitting.
24/07/05 10:51:40 WARN InstanceBuilder: Failed to load implementation from:dev.ludovic.netlib.blas.JNIBLAS
24/07/05 10:51:40 WARN InstanceBuilder: Failed to load implementation from:dev.ludovic.netlib.lapack.JNILAPACK
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/Semi_Data_Pipeline/Post_Semi_data$
```

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $ hdfs dfs -cat /output/transformed_data.csv/*
Furniture,0.1
Office Supplies,0.05
Technology,0.0999999999999999
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $ hdfs dfs -cat /output1/transformed_data.csv/*
P001,Office Supplies,0.05,100.5,20.25,3,2035.125
P002,Technology,0.1,1500.0,300.0,1,450000.0
P003,Furniture,0.08,250.75,50.15,2,12575.1125
P004,Office Supplies,0.02,20.25,4.5,5,91.125
P005,Technology,0.15,1200.0,180.0,1,216000.0
P006,Office Supplies,0.05,50.8,10.2,4,518.16
P007,Office Supplies,0.03,15.5,2.1,2,32.550000000000004
P008,Furniture,0.12,900.0,120.0,1,108000.0
P009,Technology,0.05,35.0,5.25,3,183.75
P010,Office Supplies,0.1,200.0,30.0,1,6000.0
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $
```

Ln 43, Col 65 Spaces: 4 UTF-8 LF Python 3.10.12 64-bit

RELATIONAL DATA PIPELINE PROCESS

STEP 1: Set up a pipeline to replicate data between relational databases

STEP 2: Highlight the benefits of real-time data replication for reporting, analytics, and business continuity

RELATIONAL DATA PIPELINE

The screenshot shows a terminal window with two tabs. The left tab is titled "nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~/spark" and the right tab is titled "nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~". The left tab contains a PostgreSQL command-line interface (psql) session connected to a database named "project". The user has entered the password for the "postgres" user and is now executing a query to select all columns from the "users" table. The output shows 8 rows of data with columns "id", "name", and "email". The right tab is currently empty.

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/spark
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~$ psql -h 127.0.0.1 -p 5432 -U postgres -d project
Password for user postgres:
psql (16.3 (Ubuntu 16.3-1.pgdg22.04+1))
SSL connection (protocol: TLSv1.3, cipher: TLS_AES_256_GCM_SHA384, compression: off)
Type "help" for help.

project=# SELECT * FROM users;
 id | name  |      email
----+-----+-----
 1 | Alice  | alice@example.com
 2 | Nandhu | nandhu@example.com
 3 | Midhun | midhun@example.com
 4 | xxxxxx | xxxx@example.com
 5 | yyyy   | yyyy@gmail.com
 6 | zzzz   | zzzz@gamil.com
 7 | zzzz   | zzzz@gamil.com
 8 | wwwww | ww@gamil.com
(8 rows)

project=#
```

EXPLORER

- SEMI_DATA_PIPELINE
 - __pycache__
 - kafka
 - __pycache__
 - kafka_postgre.py
 - kafka_replication.py
 - kafka.py
 - load_kafka_hadoop.py
 - seemessage.py**
- Post_Semi_data
 - convertcsvtojson.py
 - csvfile.py
 - loadhadoop.py
- Unstructured_data
 - final_pdf.py
 - pdf.py
 - unstructured_pdf.py
- automate_data.csv
- {} automate_data.json
- automate.csv
- convertjson.py
- Data_process.py
- {} data.json
- demo1.pdf
- extracted_entities.csv
- gengpt.py
- Ingestion.py
- Quetion.txt

OUTLINE

TIMELINE

DEBUG CONSOLE

Filter (e.g. text, \exclude, \escape)

kafka > seemessage.py > ...

```

18 consumer.subscribe([kafka_topic])
19
20 # Consume messages
21 try:
22     i=0
23     while i<16:
24         msg = consumer.poll(timeout=1.0)
25         if msg is None:
26             continue
27         if msg.error():
28             if msg.error().code() == KafkaError._PARTITION_EOF:
29                 continue
30             else:
31                 print(msg.error())
32         print(f"Received message: {msg.value().decode('utf-8')}")

```

PROBLEMS OUTPUT TERMINAL PORTS

● nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/Semi_Data_Pipeline/kafka\$ python3 kafka_postgre.py
 ○ nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/Semi_Data_Pipeline/kafka\$ python3 seemessage.py

```

Received message: [2, "Nandhu", "nandhu@example.com"]
Received message: [3, "Midhun", "midhun@example.com"]
Received message: [4, "xxxxxx", "xxxxx@example.com"]
Received message: [5, "yyyy", "yyyy@gmail.com"]
Received message: [6, "zzzz", "zzzz@gmail.com"]
Received message: [7, "zzzz", "zzzz@gmail.com"]
Received message: [8, "www", "ww@gmail.com"]
Received message: [1, "Alice", "alice@example.com"]
Received message: [3, "Midhun", "midhun@example.com"]
Received message: [4, "xxxxxx", "xxxxx@example.com"]
Received message: [5, "yyyy", "yyyy@gmail.com"]
Received message: [2, "Nandhu", "nandhu@example.com"]
Received message: [6, "zzzz", "zzzz@gmail.com"]
Received message: [7, "zzzz", "zzzz@gmail.com"]
Received message: [8, "www", "ww@gmail.com"]

```

REAL-TIME EVENT MONITORING AND ANALYTICS SYSTEM PROCESS

STEP 1: Create Data Producers

STEP 2: Data Processing with Kafka Streams

STEP 3: Store Processed Data

STEP 4: Visualization

REAL-TIME EVENT MONITORING AND ANALYTICS SYSTEM

Production_kafka.py - Kafka_POC - Visual Studio Code

File Edit Selection View Go Run Terminal Help

EXPLORER > OUTPUT KAFKA_POC > __pycache__ > consume_data.py & kafka_json.py {} kafka_messages.json Production_kafka.py

OUTLINE > TIMELINE DEBUG CONSOLE Filter (e.g. text, | exclude, \escape)

PROBLEMS 14 TERMINAL PORTS

```
1 import json
2 import time
3 import random
4
5 producer = KafkaProducer(bootstrap_servers='localhost:9092', value_serializer=lambda v: json.dumps(v).encode('utf-8'))
6
7 try:
8     while True:
9         data = {
10             'sensor_id': random.randint(1, 10),
11             'temperature': random.uniform(20.0, 30.0),
12             'timestamp': time.time()
13         }
14         producer.send('processed_sensor_data', data)
15         print(f"Sent: {data}") # Optional: for monitoring what data is being sent
16         time.sleep(1)
17     except KeyboardInterrupt:
18         print("Data production stopped.")
19     finally:
20         producer.close()
21
22
```

Sent: {'sensor_id': 7, 'temperature': 24.7987910888684, 'timestamp': 1720174282.97844}
Sent: {'sensor_id': 1, 'temperature': 22.986402925501388, 'timestamp': 1720174283.9805424}
Sent: {'sensor_id': 3, 'temperature': 27.683645853699956, 'timestamp': 1720174284.981764}
Sent: {'sensor_id': 5, 'temperature': 28.06945446796101, 'timestamp': 1720174285.9832044}
Sent: {'sensor_id': 6, 'temperature': 23.52127952073205, 'timestamp': 1720174286.9867878}
Sent: {'sensor_id': 7, 'temperature': 21.630011250513213, 'timestamp': 1720174287.9888284}
Sent: {'sensor_id': 8, 'temperature': 24.430602432300322, 'timestamp': 1720174288.9925873}
Sent: {'sensor_id': 9, 'temperature': 21.03931136268142, 'timestamp': 1720174289.9938836}
Sent: {'sensor_id': 10, 'temperature': 28.070882146994517, 'timestamp': 1720174290.9959443}
Sent: {'sensor_id': 6, 'temperature': 26.17139905758973, 'timestamp': 1720174291.9985301}
Sent: {'sensor_id': 8, 'temperature': 29.071654506008663, 'timestamp': 1720174293.0005314}
Sent: {'sensor_id': 2, 'temperature': 22.993665489635195, 'timestamp': 1720174294.002528}
Sent: {'sensor_id': 1, 'temperature': 24.75543865368347, 'timestamp': 1720174295.0070255}
Sent: {'sensor_id': 1, 'temperature': 22.963364083823873, 'timestamp': 1720174296.0085807}
Sent: {'sensor_id': 6, 'temperature': 26.745373661948616, 'timestamp': 1720174297.012767}
Sent: {'sensor_id': 7, 'temperature': 28.67574981534133, 'timestamp': 1720174298.0169497}
Sent: {'sensor_id': 5, 'temperature': 28.553610575614883, 'timestamp': 1720174299.0193534}
Sent: {'sensor_id': 5, 'temperature': 29.040654512707093, 'timestamp': 1720174300.0204828}
Sent: {'sensor_id': 7, 'temperature': 25.379216492968148, 'timestamp': 1720174301.0233214}
Sent: {'sensor_id': 4, 'temperature': 20.766815659148858, 'timestamp': 1720174302.0286348}
Sent: {'sensor_id': 1, 'temperature': 29.58385503108417, 'timestamp': 1720174303.03297}
Sent: {'sensor_id': 8, 'temperature': 20.526557353738824, 'timestamp': 1720174304.0366783}
Sent: {'sensor_id': 4, 'temperature': 28.525512702369253, 'timestamp': 1720174305.0376632}
Sent: {'sensor_id': 3, 'temperature': 27.478198234573128, 'timestamp': 1720174306.0397494}
Sent: {'sensor_id': 3, 'temperature': 28.05446611671858, 'timestamp': 1720174307.044295}
Sent: {'sensor_id': 4, 'temperature': 25.59728567883666, 'timestamp': 1720174308.0488935}

L 15, Col 45 Spaces: 4 UTF-8 LF Python 3.10.12 64-bit

EXPLORER

- > OUTPUT
- KAFKA_POC
 - > __pycache__
 - consume_data.py
 - kafka_json.py
 - { kafka_messages.json
 - Production_kafka.py
- > OUTLINE
- > TIMELINE
- DEBUG CONSOLE

Filter (e.g. text, | exclude, \escape)

Production_kafka.py StreamsPOC.java 9+ Settings kafka_json.py kafka_messages.json StreamsPOC.class consume_data.py x seemessage.py

```

30         timestamp_iso = datetime.now().isoformat()
31         cur.execute("INSERT INTO sensor_data (sensor_id, temperature, timestamp) VALUES (%s, %s, %s)", (data['sensor_id'], data['temperature'], timestamp_iso))
32         conn.commit()
33         logger.info("Inserted data into PostgreSQL")
34     except Exception as e:
35         logger.error(f"Error inserting data: {e}")
36     finally:
37         # Closing connections
38         if cur:
39             cur.close()
40         if conn:
41             conn.close()
42         logger.info("PostgreSQL connection closed")
43
44
45
46
47
48
49

```

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~ nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

sensor_id	temperature	timestamp
1	22.40933920985446	2024-07-05 10:16:05.802685
5	27.42124842409339	2024-07-05 10:16:06.804544
7	21.7975174413984	2024-07-05 10:16:07.806637
3	20.77522279137072	2024-07-05 10:16:08.811042
2	25.92513589507562	2024-07-05 10:16:09.814896
8	24.334581623227784	2024-07-05 10:16:10.816644
5	27.62020879378973	2024-07-05 10:16:11.818464
1	25.32967979818165	2024-07-05 10:16:12.819935
8	28.897675280818063	2024-07-05 10:16:13.82325
9	20.970045815421422	2024-07-05 10:16:14.824566
5	27.239797155282197	2024-07-05 10:16:15.828567
1	27.683286133676397	2024-07-05 10:16:16.830034
1	24.123884129670635	2024-07-05 10:16:17.832505
2	22.273384673899457	2024-07-05 10:16:18.838801

(14 rows)

PROBLEMS 14 TERMINAL PORTS

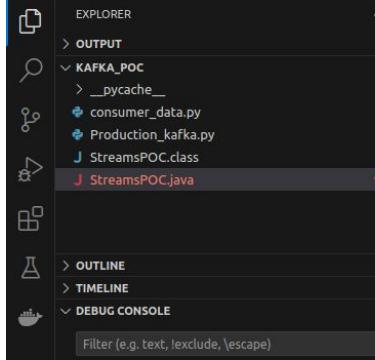
```

INFO: main :Received message: {'sensor_id': 1, 'temperature': 22.40933920985446, 'timestamp': 1720174575.828567}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 5, 'temperature': 27.42124842409339, 'timestamp': 1720174576.804544}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 7, 'temperature': 21.7975174413984, 'timestamp': 1720174577.806637}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 3, 'temperature': 20.77522279137072, 'timestamp': 1720174578.811042}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 2, 'temperature': 25.92513589507562, 'timestamp': 1720174579.814896}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 8, 'temperature': 24.334581623227784, 'timestamp': 1720174580.816644}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 5, 'temperature': 27.62020879378973, 'timestamp': 1720174581.818464}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 1, 'temperature': 25.32967979818165, 'timestamp': 1720174582.819935}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 8, 'temperature': 28.897675280818063, 'timestamp': 1720174583.82325}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 9, 'temperature': 20.970045815421422, 'timestamp': 1720174584.824566}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 5, 'temperature': 27.239797155282197, 'timestamp': 1720174585.828567}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 1, 'temperature': 27.683286133676397, 'timestamp': 1720174586.830034}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 1, 'temperature': 24.123884129670635, 'timestamp': 1720174587.832505}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 2, 'temperature': 22.273384673899457, 'timestamp': 1720174588.838801}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 4, 'temperature': 21.43756297991297, 'timestamp': 1720174589.8427944}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 1, 'temperature': 27.540560739417643, 'timestamp': 1720174590.8446128}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 9, 'temperature': 24.204914864014974, 'timestamp': 1720174591.8466108}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 2, 'temperature': 28.006405744759235, 'timestamp': 1720174592.8485348}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 8, 'temperature': 23.393432394055168, 'timestamp': 1720174593.8505433}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 3, 'temperature': 20.35990163472058, 'timestamp': 1720174594.8539484}
INFO: main :Inserted data into PostgreSQL
INFO: main :Received message: {'sensor_id': 10, 'temperature': 28.09579049645501, 'timestamp': 1720174595.8558493}
INFO: main :Inserted data into PostgreSQL

```

+ ... ^ x

- python3
- bash
- bash
- java
- java
- bash
- bash
- python3



Production_kafka.py StreamsPOC.java 9+ Settings StreamsPOC.class consumer_data.py

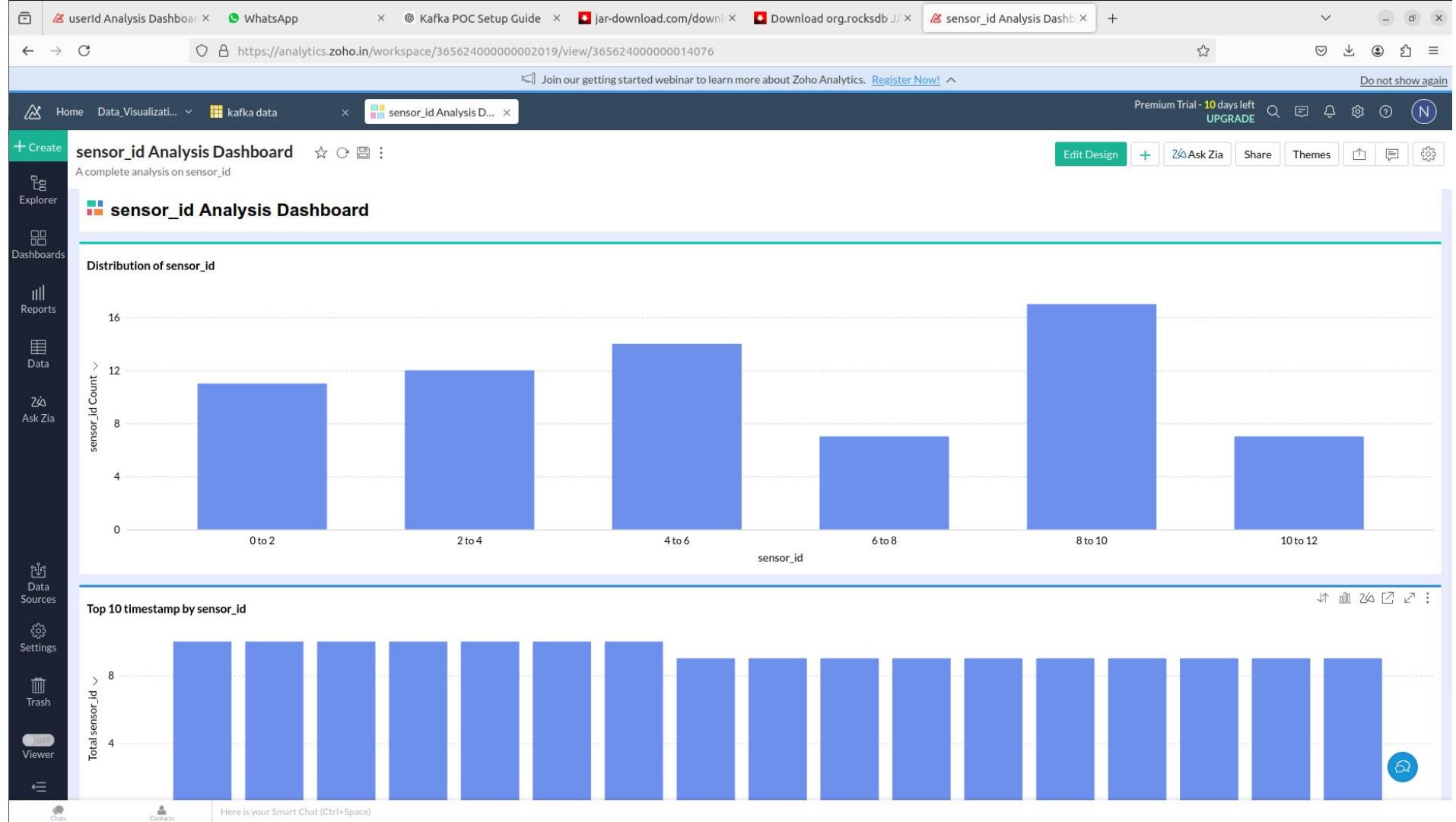
```

J StreamsPOC.java > StreamsPOC > main(String[])
7  public class StreamsPOC {
9    public static void main(String[] args) {
14      props.put(StreamsConfig.DEFAULT_VALUE_SERDE_CLASS_CONFIG, Serdes.String().getClass().getName());
15
16      StreamsBuilder builder = new StreamsBuilder();
17      KStream<String, String> sourceStream = builder.stream("processed_sensor_data");
18
19      // Process your stream or define transformations here
20      sourceStream.foreach((key, value) -> {
21        System.out.println("Key: " + key + ", Value: " + value);
22      });
23
24      KafkaStreams streams = new KafkaStreams(builder.build(), props);
25      streams.start();
26

```

[streams-poc-2-105d75c2-6050-4bc8-afb8-bdbc61ae8e55-StreamThread-1] INFO org.apache.kafka.clients.consumer.internals.ConsumerCoordinator - [Consumer clientId=streams-poc-2-105d75c2-6050-4bc8-afb8-bdbc61ae8e55-StreamThread-1-consumer, groupId=streams-poc-2] Found no committed offset for partition processed_sensor_data-0
[streams-poc-2-105d75c2-6050-4bc8-afb8-bdbc61ae8e55-StreamThread-1] INFO org.apache.kafka.clients.consumer.internals.SubscriptionState - [Consumer clientId=streams-poc-2-105d75c2-6050-4bc8-afb8-bdbc61ae8e55-StreamThread-1-consumer, groupId=streams-poc-2] Resetting offset for partition processed_sensor_data-0 to position FetchPosition{offset=0, offsetEpoch=Optional.empty, currentLeader=LeaderAndEpoch{leader=Optional[null], epoch=0}}.
Key: null, Value: {"sensor_id": 4, "temperature": 23.19852952294821, "timestamp": 1720241569.3474495}
Key: null, Value: {"sensor_id": 7, "temperature": 22.60240145854672, "timestamp": 1720241570.3576665}
Key: null, Value: {"sensor_id": 6, "temperature": 26.285091007789205, "timestamp": 1720241571.359596}
Key: null, Value: {"sensor_id": 7, "temperature": 29.06402170076111, "timestamp": 1720241572.3628755}
Key: null, Value: {"sensor_id": 8, "temperature": 27.698797717029294, "timestamp": 1720241573.366271}
Key: null, Value: {"sensor_id": 4, "temperature": 21.27913630499162, "timestamp": 1720241574.3673408}
Key: null, Value: {"sensor_id": 9, "temperature": 28.144526804550857, "timestamp": 1720241575.3713806}
Key: null, Value: {"sensor_id": 8, "temperature": 26.943628802194596, "timestamp": 1720241576.3753252}
Key: null, Value: {"sensor_id": 2, "temperature": 28.553616731767313, "timestamp": 1720241577.3774042}
Key: null, Value: {"sensor_id": 4, "temperature": 29.16208078654908, "timestamp": 1720241578.3793068}
Key: null, Value: {"sensor_id": 9, "temperature": 21.631640244999563, "timestamp": 1720241579.3808105}
Key: null, Value: {"sensor_id": 5, "temperature": 21.896282465271987, "timestamp": 1720241580.38255}
Key: null, Value: {"sensor_id": 4, "temperature": 25.55688583862497, "timestamp": 1720241581.3847117}
Key: null, Value: {"sensor_id": 4, "temperature": 29.37446271879176, "timestamp": 1720241582.386458}
Key: null, Value: {"sensor_id": 8, "temperature": 25.857415568051238, "timestamp": 1720241583.3874972}
Key: null, Value: {"sensor_id": 6, "temperature": 29.263712100359893, "timestamp": 1720241584.3893502}
Key: null, Value: {"sensor_id": 6, "temperature": 25.859430505798844, "timestamp": 1720241585.3911266}
Key: null, Value: {"sensor_id": 4, "temperature": 24.716697737499338, "timestamp": 1720241586.3935878}
Key: null, Value: {"sensor_id": 2, "temperature": 23.2595969873049, "timestamp": 1720241587.394147}
Key: null, Value: {"sensor_id": 8, "temperature": 25.137590893318645, "timestamp": 1720241588.3954144}
Key: null, Value: {"sensor_id": 1, "temperature": 20.881868362280766, "timestamp": 1720241589.3970525}
Key: null, Value: {"sensor_id": 6, "temperature": 24.844951993004745, "timestamp": 1720241590.398852}
Key: null, Value: {"sensor_id": 8, "temperature": 21.95386135713793, "timestamp": 1720241591.4009237}
Key: null, Value: {"sensor_id": 4, "temperature": 29.81146163668466, "timestamp": 1720241592.403235}
Key: null, Value: {"sensor_id": 10, "temperature": 25.21964599641193, "timestamp": 1720241593.4085279}
Key: null, Value: {"sensor_id": 8, "temperature": 27.2423285302754, "timestamp": 1720241594.4150567}
Key: null, Value: {"sensor_id": 7, "temperature": 27.983232449215933, "timestamp": 1720241595.4168944}





REAL-TIME STOCK MARKET ANALYTICS USING APACHE KAFKA AND HADOOP PROCESS

STEP 1: Fetch Data from NSE API

[<https://www.nseindia.com/api/equity-stockIndices?index=NIFTY%2050>]

STEP 2: Produce Data to Kafka

STEP 3: Consume Data from Kafka and Store in HDFS

STEP 4: Analyze Data using Spark

STEP 5: Visualize Results

REAL-TIME STOCK MARKET ANALYTICS USING APACHE KAFKA AND HADOOP

fetch_data.py - Real-time Stock Market - Visual Studio Code

File Edit Selection View Go Run Terminal Help

EXPLORER

OUTPUT

REAL-TIME STOCK MARKET

- > __pycache__
- > .venv
- copy_into_local.py
- fetch_data.py
- kafka_producer.py
- kafka_to_hdfs.py
- stock_prices.csv

PROBLEMS TERMINAL PORTS

nandhumidhun@nandhumidhun-MP-Laptop-15-bs0xx:~/Real-time Stock Market\$ python3 fetch_data.py

```
{'name': 'NIFTY 50', 'advance': {'declines': '31', 'advances': '19', 'unchanged': '0'}, 'timestamp': '10-Jul-2024 12:59:26', 'data': [{"priority": 1, "symbol": "NIFTY 50", "identifier": "NIFTY 50", "open": 24459.85, "dayHigh": 24461.05, "dayLow": 24141.8, "lastPrice": 24283.55, "previousClose": 24433.2, "change": -149.65000000000146, "pChange": -0.61, "ffmc": 1073554809.33, "yearHigh": 24443.6, "yearLow": 18837.85, "totalTradedVolume": 172904663, "totalTradedValue": 209233037589.77, "lastUpdateTime": '10-Jul-2024 12:59:26', "nearWKL": 0.6547726194177588, "nearWKL": -28.90828836624138, "perChange365d": 26.23, "date365dAgo": '10-Jul-2023', "chart365dPath": 'https://nsearchives.nseindia.com/365d/NIFTY-50.svg', "chartTodayPath": 'https://nsearchives.nseindia.com/today/NIFTY-50.svg'}, {"priority": 0, "symbol": "HDFCLIFE", "identifier": "HDFCLIFEON", "series": "EQ", "open": 624, "dayHigh": 639, "dayLow": 615.6, "lastPrice": 638.35, "previousClose": 623.65, "change": 14.7, "pChange": 2.36, "totalTradedValue": 3623630, "lastUpdateTime": '10-Jul-2024 13:00:25', "yearHigh": 710.6, "ffmc": 672333765471.14, "yearLow": 511.4, "nearWKH": 10.16746414832535, "nearWKL": -24.824012514665633, "perChange365d": -6.67, "date365dAgo": '10-Jul-2023', "chart365dPath": 'https://nsearchives.nseindia.com/365d/HDFCLIFEEQ-SVG', "date30dAgo": '07-Jun-2024', "perChange30d": 10.78, "chart30dPath": 'https://nsearchives.nseindia.com/30d/HDFCLIFEEQ-SVG', "chartTodayPath": 'https://nsearchives.nseindia.com/today/HDFCLIFEEQ-SVG', "meta": {"symbol": "HDFCLIFE", "companyName": "HDFC Life Insurance Company Limited", "activeSeries": ["EQ"], "debtSeries": [], "isFNOsec": True, "isCAsec": False, "isSLBsec": False, "isSuspended": False, "tempSuspendedSeries": [], "isETFSec": False, "isDelisted": False, "isin": 'INE795G01014', "isMunicipalBond": False, "quotePreopenstatus": {"equityTime": '10-Jul-2024 12:59:26', "preOpenTime": '10-Jul-2024 09:08:00', "QuotePreOpenFlag": False}}, {"priority": 0, "symbol": "SBILIFEON", "identifier": "SBILIFEON", "series": "EQ", "open": 1524.75, "dayHigh": 1557, "dayLow": 1518.6, "lastPrice": 1556.95, "previousClose": 1524.75, "change": 32.2, "pChange": 2.11, "totalTradedVolume": 710493, "totalTradedValue": 1090187564.13, "lastUpdateTime": '10-Jul-2024 13:00:25', "yearHigh": 1569.4, "ffmc": 700943546203.2, "yearLow": 1251.65, "nearWKH": 0.7932968013253501, "nearWKL": -24.391802820277228, "perChange365d": 17.9, "date365dAgo": '10-Jul-2023', "chart365dPath": 'https://nsearchives.nseindia.com/365d/SBILIFE-EQ-SVG', "date30dAgo": '07-Jun-2024', "perChange30d": 6.94, "chart30dPath": 'https://nsearchives.nseindia.com/30d/SBILIFE-EQ-SVG', "chartTodayPath": 'https://nsearchives.nseindia.com/today/SBILIFE-EQ-SVG', "meta": {"symbol": "SBILIFE", "companyName": "SBI Life Insurance Company Limited", "activeSeries": ["EQ"], "debtSeries": [], "isFNOsec": True, "isCAsec": False, "isSLBsec": True, "isDebtSec": False, "isSuspended": False, "tempSuspendedSeries": [], "isETFSec": False, "isDelisted": False, "isin": 'INE123W01016', "isMunicipalBond": False, "quotePreopenstatus": {"equityTime": '10-Jul-2024 12:59:26', "preOpenTime": '10-Jul-2024 09:08:00', "QuotePreOpenFlag": False}}, {"priority": 0, "symbol": "BRITANNIAEQON", "identifier": "BRITANNIAEQON", "series": "EQ", "open": 5670.05, "dayHigh": 5800, "dayLow": 5670.05, "lastPrice": 5768.4, "previousClose": 5668.85, "change": 99.55, "pChange": 1.76, "totalTradedVolume": 341442, "totalTradedValue": 1961495515.08, "lastUpdateTime": '10-Jul-2024 13:00:26', "yearHigh": 5800, "ffmc": 680800388716.98, "yearLow": 4347.7, "nearWKH": 0.544827582069029, "nearWKL": -32.6770476343819, "perChange365d": 12.66, "date365dAgo": '10-Jul-2023', "chart365dPath": 'https://nsearchives.nseindia.com/365d/BRITANNIA-EQ-SVG', "date30dAgo": '07-Jun-2024', "perChange30d": 3.76, "chart30dPath": 'https://nsearchives.nseindia.com/30d/BRITANNIA-EQ-SVG', "chartTodayPath": 'https://nsearchives.nseindia.com/today/BRITANNIAEQON-SVG', "meta": {"symbol": "BRITANNIA", "companyName": "Britannia Industries Limited", "activeSeries": ["EQ"], "debtSeries": [], "isFNOsec": True, "isCAsec": False, "isSLBsec": True, "isDebtSec": False, "isSuspended": False, "tempSuspendedSeries": [], "isETFSec": False, "isDelisted": False, "isin": 'INE216A07052', "isMunicipalBond": False, "quotePreopenstatus": {"equityTime": '10-Jul-2024 12:59:26', "preOpenTime": '10-Jul-2024 09:08:00', "QuotePreOpenFlag": False}}, {"priority": 0, "symbol": "DIVISLAB", "identifier": "DIVISLABON", "series": "EQ", "open": 4572, "dayHigh": 4635.55, "dayLow": 4536.8, "lastPrice": 4613.85, "previousClose": 4551.95, "change": 61.9, "pChange": 1.36, "totalTradedVolume": 325838, "totalTradedValue": 1496381689.58, "lastUpdateTime": '10-Jul-2024 13:00:25', "yearHigh": 4670, "ffmc": 587811148579.2, "yearLow": 3293.5, "nearWKH": 0.20352554603854312, "nearWKL": -40.01304887809916, "perChange365d": 25.06, "date365dAgo": '10-Jul-2023', "chart365dPath": 'https://nsearchives.nseindia.com/365d/DIVISLAB-EQ-SVG', "date30dAgo": '07-Jun-2024', "perChange30d": 0.62, "chart30dPath": 'https://nsearchives.nseindia.com/30d/DIVISLAB-EQ-SVG', "chartTodayPath": 'https://nsearchives.nseindia.com/today/DIVISLABEQ-SVG', "meta": {"symbol": "DIVISLAB", "companyName": "Pharmaceuticals", "activeSeries": ["EQ"], "debtSeries": [], "isFNOsec": True, "isCAsec": False, "isSLBsec": True, "isDebtSec": False, "isSuspended": False, "tempSuspendedSeries": [], "isETFSec": False, "isDelisted": False, "isin": 'INE361B01024', "isMunicipalBond": False, "quotePreopenstatus": {"equityTime": '10-Jul-2024 09:08:00', "QuotePreOpenFlag": False}}, {"priority": 0, "symbol": "FEO", "identifier": "FEO", "series": "T0", "open": 343, "dayHigh": 347.8, "dayLow": 335.5, "lastPrice": 345.55, "previousClose": 341.15, "change": 4.4, "pChange": 1.2, "totalTradedValue": 2422314625.3700004, "lastUpdateTime": '10-Jul-2024 13:00:23', "yearHigh": 348.7, "ffmc": 15743197590988, "nearWKL": -94.64315890272067, "perChange365d": 36.46, "date365dAgo": '10-Jul-2023', "chart365dPath": 'https://nsearchives.nseindia.com/365d/FEO-SVG', "date30dAgo": '07-Jun-2024', "perChange30d": 0.62, "chart30dPath": 'https://nsearchives.nseindia.com/30d/FEO-SVG', "chartTodayPath": 'https://nsearchives.nseindia.com/today/FEO-SVG', "meta": {"symbol": "FEO", "companyName": "Fertilizers", "activeSeries": ["T0"], "debtSeries": [], "isFNOsec": True, "isCAsec": False, "isSLBsec": True, "isDebtSec": False, "isSuspended": False, "tempSuspendedSeries": [], "isETFSec": False, "isDelisted": False, "isin": 'INE216A07052', "isMunicipalBond": False, "quotePreopenstatus": {"equityTime": '10-Jul-2024 09:08:00', "QuotePreOpenFlag": False}}}, {"priority": 0, "symbol": "PHARMACEUTICALS", "identifier": "PHARMACEUTICALS", "series": "EQ", "open": 343, "dayHigh": 347.8, "dayLow": 335.5, "lastPrice": 345.55, "previousClose": 341.15, "change": 4.4, "pChange": 1.2, "totalTradedValue": 2422314625.3700004, "lastUpdateTime": '10-Jul-2024 13:00:23', "yearHigh": 348.7, "ffmc": 15743197590988, "nearWKL": -94.64315890272067, "perChange365d": 36.46, "date365dAgo": '10-Jul-2023', "chart365dPath": 'https://nsearchives.nseindia.com/365d/PHARMACEUTICALS-SVG', "date30dAgo": '07-Jun-2024', "perChange30d": 0.62, "chart30dPath": 'https://nsearchives.nseindia.com/30d/PHARMACEUTICALS-SVG', "chartTodayPath": 'https://nsearchives.nseindia.com/today/PHARMACEUTICALS-SVG', "meta": {"symbol": "PHARMACEUTICALS", "companyName": "Pharmaceuticals", "activeSeries": ["EQ"], "debtSeries": [], "isFNOsec": True, "isCAsec": False, "isSLBsec": True, "isDebtSec": False, "isSuspended": False, "tempSuspendedSeries": [], "isETFSec": False, "isDelisted": False, "isin": 'INE361B01024', "isMunicipalBond": False, "quotePreopenstatus": {"equityTime": '10-Jul-2024 09:08:00', "QuotePreOpenFlag": False}}}, {"priority": 0, "symbol": "POWERGRID", "identifier": "POWERGRID", "series": "EQ", "open": 1000, "dayHigh": 1010, "dayLow": 990, "lastPrice": 1005, "previousClose": 995, "change": 10, "pChange": 1.0, "totalTradedValue": 1000000000.0, "lastUpdateTime": '10-Jul-2024 13:00:25', "yearHigh": 1010, "ffmc": 1000000000.0, "nearWKL": -100.0, "perChange365d": 10.0, "date365dAgo": '10-Jul-2023', "chart365dPath": 'https://nsearchives.nseindia.com/365d/POWERGRID-SVG', "date30dAgo": '07-Jun-2024', "perChange30d": 0.0, "chart30dPath": 'https://nsearchives.nseindia.com/30d/POWERGRID-SVG', "chartTodayPath": 'https://nsearchives.nseindia.com/today/POWERGRID-SVG', "meta": {"symbol": "POWERGRID", "companyName": "Powergrid", "activeSeries": ["EQ"], "debtSeries": [], "isFNOsec": True, "isCAsec": False, "isSLBsec": True, "isDebtSec": False, "isSuspended": False, "tempSuspendedSeries": [], "isETFSec": False, "isDelisted": False, "isin": 'INE361B01024', "isMunicipalBond": False, "quotePreopenstatus": {"equityTime": '10-Jul-2024 09:08:00', "QuotePreOpenFlag": False}}}
```

Do you want to install the recommended 'Rainbow CSV' extension from mechatroner for stock_prices.csv? Install Show Recommendations

Ln9, Col1 Spaces:4 UTF-8 LF Python 3.10.12 (.venv:venv)

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/kafka \$ kafka-console-consumer.sh --bootstrap-server localhost:9092 --topic nse_stock_data --from-beginning

SLF4J: Class path contains multiple SLF4J bindings.

SLF4J: Found binding in [jar:file:/home/nandhumidhun/kafka/libs/slf4j-reload4j-1.7.36.jar!org/slf4j/impl/StaticLoggerBinder.class]

SLF4J: Found binding in [jar:file:/home/nandhumidhun/kafka/libs/slf4j-simple-1.7.36.jar!org/slf4j/impl/StaticLoggerBinder.class]

SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.

SLF4J: Actual binding is of type [org.slf4j.impl.Reload4jLoggerFactory]

```
{"name": "NIFTY 50", "advance": {"declines": "30", "advances": "20", "unchanged": "0"}, "timestamp": "10-Jul-2024 13:04:26", "data": [{"priority": 1, "symbol": "NIFTY 50", "identifier": "NIFTY 50", "open": 24459.85, "dayHigh": 24461.05, "dayLow": 24141.8, "lastPrice": 24283.65, "previousClose": 24433.2, "change": -149.54999999999927, "pChange": -0.61, "ffmc": 1073629150.29, "yearHigh": 24443.6, "yearLow": 18837.85, "totalTradedVolume": 174845350, "totalTradedValue": 211622411574.23, "lastUpdateTime": "10-Jul-2024 13:04:26", "nearWKh": 0.6543635143759393, "nearWKL": -28.908819212383595, "perChange365d": 26.23, "date365dAgo": "10-Jul-2023", "chart365dPath": "https://nsearchives.nseindia.com/365d/NIFTY-50.svg", "date30dAgo": "07-Jun-2024", "perChange30d": 4.91, "chart30dPath": "https://nsearchives.nseindia.com/30d/NIFTY-50.svg"}, {"priority": 0, "symbol": "HDFCLIFE", "identifier": "HDFCLIFEQN", "series": "EQ", "open": 624, "dayHigh": 640.75, "dayLow": 615.6, "lastPrice": 639.2, "previousClose": 623.65, "change": 15.55, "pChange": 2.49, "totalTradedVolume": 4025188, "totalTradedValue": 2533453327.2, "lastUpdateTime": "10-Jul-2024 13:04:26", "nearWKh": 710.6, "nearWKL": 10.04784688952151, "perChange365d": -6.67, "date365dAgo": "10-Jul-2023", "chart365dPath": "https://nsearchives.nseindia.com/365d/HDFCLIFE-EQ.svg", "date30dAgo": "07-Jun-2024", "perChange30d": 10.78, "chart30dPath": "https://nsearchives.nseindia.com/30d/HDFCLIFE-EQ.svg", "chartTodayPath": "https://nsearchives.nseindia.com/today/HDFCLIFEQN.svg", "meta": {"symbol": "HDFCLIFE", "companyName": "HDFC Life Insurance Company Limited", "activeSeries": ["EQ"], "debtSeries": [], "isFNOsec": true, "isCAsec": false, "isSLBsec": true, "isDebtSec": false, "isSuspended": false, "tempSuspendedSeries": [], "isETFSec": false, "isDelisted": false, "isin": "INE795G01014", "isMunicipalBond": false, "quotepreopenstatus": {"equityTime": "10-Jul-2024 13:04:26", "preOpenTime": "10-Jul-2024 09:08:00", "QuotePreOpenFlag": false}}, {"priority": 0, "symbol": "SBILIFE", "identifier": "SBILIFEQN", "series": "EQ", "open": 1524.75, "dayHigh": 1559.9, "dayLow": 1518.6, "lastPrice": 1557.5, "previousClose": 1524.75, "change": 32.75, "pChange": 2.15, "totalTradedVolume": 763503, "totalTradedValue": 1172786418.18, "lastUpdateTime": "10-Jul-2024 13:04:56", "nearWKh": 1569.4, "nearWKL": 0.758251561106161, "perChange365d": 17.9, "date365dAgo": "10-Jul-2023", "chart365dPath": "https://nsearchives.nseindia.com/365d/SBILIFE-EQ.svg", "date30dAgo": "07-Jun-2024", "perChange30d": 6.94, "chart30dPath": "https://nsearchives.nseindia.com/30d/SBILIFE-EQ.svg", "chartTodayPath": "https://nsearchives.nseindia.com/today/SBILIFEQN.svg", "meta": {"symbol": "SBILIFE", "companyName": "SBI Life Insurance Company Limited", "activeSeries": ["EQ"], "debtSeries": [], "isFNOsec": true, "isCAsec": false, "isSLBsec": true, "isDebtSec": false, "isSuspended": false, "tempSuspendedSeries": [], "isETFSec": false, "isDelisted": false, "isin": "INE123W01010", "isMunicipalBond": false, "quotepreopenstatus": {"equityTime": "10-Jul-2024 13:04:26", "preOpenTime": "10-Jul-2024 09:08:00", "QuotePreOpenFlag": false}}, {"priority": 0, "symbol": "BRITANNIA", "identifier": "BRITANNIAEQN", "series": "EQ", "open": 5670.05, "dayHigh": 5800, "dayLow": 5670.05, "lastPrice": 5760, "previousClose": 5668.85, "change": 91.15, "pChange": 1.61, "totalTradedVolume": 348441, "totalTradedValue": 2001179607.36, "lastUpdateTime": "10-Jul-2024 13:04:56", "nearWKh": 5800, "nearWKL": 0.6896551724137931, "perChange365d": 12.66, "date365dAgo": "10-Jul-2023", "chart365dPath": "https://nsearchives.nseindia.com/365d/BRITANNIA-EQ.svg", "date30dAgo": "07-Jun-2024", "perChange30d": 3.76, "chart30dPath": "https://nsearchives.nseindia.com/30d/BRITANNIA-EQ.svg", "chartTodayPath": "https://nsearchives.nseindia.com/today/BRITANNIAEQN.svg", "meta": {"symbol": "BRITANNIA", "companyName": "Britannia Industries Limited", "industry": "FOOD AND FOOD PROCESSING", "activeSeries": ["EQ"], "debtSeries": [], "isCAsec": false, "isDebtSec": false, "isSuspended": false, "tempSuspendedSeries": ["N1", "N2"], "isETFSec": false, "isDelisted": false, "isin": "INE216A07052", "isMunicipalBond": false, "quotepreopenstatus": {"equityTime": "10-Jul-2024 13:04:26", "preOpenTime": "10-Jul-2024 09:08:00", "QuotePreOpenFlag": false}}, {"priority": 0, "symbol": "POWERGRID", "identifier": "POWERGRIDEQN", "series": "EQ", "open": 343, "dayHigh": 347.8, "dayLow": 335.5, "lastPrice": 345.75, "previousClose": 341.15, "change": 4.6, "pChange": 1.35, "totalTradedVolume": 7133957, "totalTradedValue": 2443380272.5, "lastUpdateTime": "10-Jul-2024 13:04:55", "nearWKh": 348.7, "nearWKL": 0.8459994264461036, "perChange365d": -94.55818943633, "date365dAgo": 36.46, "chart365dPath": "https://nsearchives.nseindia.com/365d/POWERGRID-EQ.svg", "date30dAgo": "07-Jun-2024", "perChange30d": 10.28, "chart30dPath": "https://nsearchives.nseindia.com/30d/POWERGRID-EQ.svg", "chartTodayPath": "https://nsearchives.nseindia.com/today/POWERGRID-EQ.svg", "meta": {"symbol": "POWERGRID", "companyName": "Power Grid Corporation of India Limited", "industry": "POWER", "activeSeries": ["EQ"], "debtSeries": [], "isCAsec": false, "isDebtSec": true, "isSuspended": false, "tempSuspendedSeries": ["IL"], "isETFSec": false, "isDelisted": false, "isin": "INE752E01010", "isMunicipalBond": false, "quotepreopenstatus": {"equityTime": "10-Jul-2024 13:04:26", "preOpenTime": "10-Jul-2024 09:08:00", "QuotePreOpenFlag": false}}, {"priority": 0, "symbol": "GRASIM", "identifier": "GRASIMEQN", "series": "EQ", "open": 2781, "dayHigh": 2812.25, "dayLow": 2741, "lastPrice": 2794.7, "previousClose": 2761.8, "change": 32.9, "pChange": 1.19, "totalTradedValue": 774465, "totalTradedValue": 215602134.9, "lastUpdateTime": "10-Jul-2024 13:04:56", "nearWKh": 2812.25, "nearWKL": 0.6240444252127496, "perChange365d": 58.45, "date365dAgo": "10-Jul-2023", "chart365dPath": "https://nsearchives.nseindia.com/365d/GRASIM-EQ.svg", "date30dAgo": "07-Jun-2024", "perChange30d": 16.05, "chart30dPath": "https://nsearchives.nseindia.com/30d/GRASIM-EQ.svg", "chartTodayPath": "https://nsearchives.nseindia.com/today/GRASIMEQN.svg", "meta": {"symbol": "GRASIM", "companyName": "Grasim Industries Limited", "industry": "CEMENT AND CEMENT PRODUCTS", "activeSeries": ["EQ"], "debtSeries": [], "isCAsec": false, "isDebtSec": true, "isSuspended": false, "tempSuspendedSeries": ["IL"], "isETFSec": false, "isDelisted": false, "isin": "INE0472E010101", "isMunicipalBond": false, "quotepreopenstatus": {"equityTime": "10-Jul-2024 13:04:26", "preOpenTime": "10-Jul-2024 09:08:00", "QuotePreOpenFlag": false}}, {"priority": 0, "symbol": "DIVISLAB", "identifier": "DIVISLAQEQN", "series": "EQ", "open": 4572, "dayHigh": 4635.55, "dayLow": 4536.8, "lastPrice": 4605.3, "previousClose": 4551.95, "change": 53.35, "pChange": 1.17, "totalTradedValue": 328400, "totalTradedValue": 1508193420, "lastUpdateTime": "10-Jul-2024 13:04:55", "nearWKh": 5876837323660.8, "nearWKL": 3295.3, "perChange365d": 25.06, "date365dAgo": "10-Jul-2023", "chart365dPath": "https://nsearchives.nseindia.com/365d/DIVISLAB-EQ.svg", "date30dAgo": "07-Jun-2024", "perChange30d": 0.62, "chart30dPath": "https://nsearchives.nseindia.com/30d/DIVISLAB-EQ.svg", "chartTodayPath": "https://nsearchives.nseindia.com/today/DIVISLAQEQN.svg", "meta": {"symbol": "DIVISLAB", "companyName": "Divi's Laboratories Limited", "industry": "PHARMACEUTICALS", "activeSeries": ["EQ"], "debtSeries": [], "isCAsec": false, "isDebtSec": true, "isSuspended": false, "tempSuspendedSeries": [], "isETFSec": false, "isDelisted": false, "isin": "INE361B01024", "isMunicipalBond": false, "quotepreopenstatus": {"equityTime": "10-Jul-2024 13:04:26", "preOpenTime": "10-Jul-2024 09:08:00", "QuotePreOpenFlag": false}}, {"priority": 0, "symbol": "MARUTI", "identifier": "MARUTIEQN", "series": "EQ", "open": 12950.05, "dayHigh": 13300, "dayLow": 12828.05, "lastPrice": 12962.5, "previousClose": 12827.7, "change": 134.8, "pChange": 1.05, "totalTradedValue": 1604580, "totalTradedValue": 20986911560.39999, "lastUpdateTime": "10-Jul-2024 13:04:55", "nearWKh": 12300, "nearWKL": 2.5375398844147724, "perChange365d": 31.75, "date365dAgo": "10-Jul-2023", "chart365dPath": "https://nsearchives.nseindia.com/365d/MARUTI-EQ.svg", "date30dAgo": "07-Jun-2024", "perChange30d": 0.13, "chart30dPath": "https://nsearchives.nseindia.com/30d/MARUTI-EQ.svg", "chartTodayPath": "https://nsearchives.nseindia.com/today/MARUTIEQN.svg", "meta": {"symbol": "MARUTI", "companyName": "Maruti Suzuki India Limited", "industry": "AUTOMOBILES - 4 WHEELERS", "activeSeries": ["EQ"], "debtSeries": [], "isCAsec": true, "isDebtSec": false, "isSuspended": false, "tempSuspendedSeries": ["IL"], "isETFSec": false, "isDelisted": false, "isin": "INE585B01010", "isMunicipalBond": false, "quotepreopenstatus": {"equityTime": "10-Jul-2024 13:04:26", "preOpenTime": "10-Jul-2024 09:08:00", "QuotePreOpenFlag": false}}, {"priority": 0, "symbol": "APOLLOHOSP", "identifier": "APOLLOHOSEPN", "series": "EQ", "open": 6320.3, "dayHigh": 6417, "dayLow": 6301.15, "lastPrice": 6379.85, "previousClose": 6320.3, "change": 59.55, "pChange": 0.94, "totalTradedValue": 210819, "totalTradedValue": 1343570568.9, "lastUpdateTime": "10-Jul-2024 13:04:53", "nearWKh": 6874.45, "nearWKL": 7.1947573988773712, "perChange365d": -34.23, "date365dAgo": "10-Jul-2023", "chart365dPath": "https://nsearchives.nseindia.com/365d/APOLLOHOSP-EQ.svg", "date30dAgo": "07-Jun-2024", "perChange30d": 0.13, "chart30dPath": "https://nsearchives.nseindia.com/30d/APOLLOHOSP-EQ.svg", "chartTodayPath": "https://nsearchives.nseindia.com/today/APOLLOHOSPEPN.svg", "meta": {"symbol": "APOLLOHOSP", "companyName": "Apollo Hospitals Enterprise Limited", "industry": "HEALTHCARE", "activeSeries": ["EQ"], "debtSeries": [], "isCAsec": false, "isDebtSec": true, "isSuspended": false, "tempSuspendedSeries": ["IL"], "isETFSec": false, "isDelisted": false, "isin": "INE585B01010", "isMunicipalBond": false, "quotepreopenstatus": {"equityTime": "10-Jul-2024 13:04:26", "preOpenTime": "10-Jul-2024 09:08:00", "QuotePreOpenFlag": false}}}
```

File Edit Selection View Go Run Terminal Help



EXPLORER

...

fetch_data.py kafka_to_hdfs.py > ...

copy_into_local.py

stock_prices.csv

kafka.py

kafka_producer.py

D v I ...



OUTPUT

```

39     exploded_df.createOrReplaceTempView("stocks")
40
41     # Run the SQL query to calculate average of 'lastPrice'
42     result = spark.sql("SELECT data.symbol, AVG(CAST(data.lastPrice AS FLOAT)) AS avg_price FROM stocks GROUP BY data.symbol")
43     result.show()
44     consumer.close()
45     spark.stop()

46 except KeyboardInterrupt:
47     print("Interrupted by user")
48 except Exception as e:
49     print(f"Error: {e}")
50 finally:
51     consumer.close()
52     spark.stop() # Stop Spark session gracefully
53
54

```



REAL-TIME STOCK MARKET



> __pycache__

> .venv

copy_into_local.py

fetch_data.py

kafka_producer.py

kafka_to_hdfs.py

stock_prices.csv



PROBLEMS TERMINAL PORTS

24/07/10 12:47:36 ERROR Executor: Exception in task 2.0 in stage 73.0 (TID 168): Connection reset
 (.venv) nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/Real-time Stock Market\$ python3 kafka_to_hdfs.py
 24/07/10 12:49:33 WARN Utils: Your hostname, nandhumidhun-HP-Laptop-15-bs0xx resolves to a loopback address: 127.0.1.1; using 192.168.1.38 instead (on interface wl0)
 24/07/10 12:49:33 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
 Setting default log level to "WARN".
 To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
 24/07/10 12:49:35 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
 Starting to consume messages...
 24/07/10 12:49:46 WARN SparkStringUtils: Truncated the string representation of a plan since it was too large. This behavior can be adjusted by setting 'spark.sql.debug.maxToStringFields'.
 +-----+-----+
 | symbol| avg_price|
 +-----+-----+
LT	3638.550048828125
ICICIBANK	1238.75
TITAN	3234.300048828125
SHRI RAM FIN	2757.89990234375
NIFTY 50	24275.900390625
DIVISLAB	4622.89990234375
HDFCBANK	1626.8499755859375
CIPLA	1504.8499755859375
SBI LIFE	1543.5
HINDALCO	690.5
BAJAJ-AUTO	9515.7998046875
INFY	1646.75
MGM	2728.449951171875
POWERGRID	345.95001220703125
ONGC	294.0

> OUTLINE

> TIMELINE

> DEBUG CONSOLE

Filter (e.g. text, lexclude, \escape)



Do you want to install the recommended 'Rainbow CSV' extension from mechatroner for stock_prices.csv?

Install

Show Recommendations

Ln 45, Col 9 Spaces: 4 UTF-8 LF ↵ Python 3.10.12 ('.venv':venv) ↴

Home Data_Visualizati... Stock-price Stock-price-design area-chat New Dashboard Premium Trial - 5 days left UPGRADE

+ Create Drag and Drop the Views

Search :

New Dashboard

+ Add Tab

Widget Text Image Embed

View Mode Ask Zia Themes

Tables & Reports

- Age Total
- area-chat
- Bar-chat
- data_json
- kafka data
- Kafka data1
- Product
- Stock-price
- Stock-price-design
- total
- Total1
- transaction_details

product_id Analysis

user_id Analysis

age Analysis

product_id VS user_id

total Analysis

Distribution of category

age Analysis(1)

age-wise age

Distribution of age(1)

user_count Analysis

geographic user count

Chats

Contacts

Here is your Smart Chat (Ctrl+Space)

area-chat

Total avg_price > 20000

symbol: ADANIPOSTS, APOLLOHOSP, BAJA-AUTO, BAJFINANCE, CIRLA, DIVISLAB, HDFCBANK, HINDALCO, ICICIBANK, INFY, LITIM, MSM, NIFTY 50, ONCC, POWERGRID, SBLIFLIFE, SHIRAMAEIN, TATACONSUM, TITAN

Bar-chat

Total avg_price > 80000

symbol: NIFTY 50, BAJA-AUTO, BAJFINANCE, APOLLOHOSP, LTIM, DIVISLAB, LT

symbol	Total avg_price
NIFTY 50	24275.45
BAJA-AUTO	9545.00
BAJFINANCE	7056.20
APOLLOHOSP	6394.10
LTIM	5388.85
DIVISLAB	4602.55
LT	3636.50

Stock-price-design

symbol: NIFTY 50, BAJA-AUTO, BAJFINANCE, APOLLOHOSP, LTIM, DIVISLAB, LT, + 13 more...

symbol	Value
NIFTY 50	24275.45
BAJA-AUTO	9545.00
BAJFINANCE	7056.20
APOLLOHOSP	6394.10
LTIM	5388.85
DIVISLAB	4602.55
LT	3636.50

1239.95
2720.80
3247.40
3636.50
4602.55
6394.10
29.9%
4%
8.7%
11.8%
6.6%
7.9%
5.7%

Here is your Smart Chat (Ctrl+Space)

PREDICTIVE MAINTENANCE FOR MANUFACTURING EQUIPMENT PROCESS

STEP 1: Data Ingestion and ETL with Apache NiFi

STEP 2: Data Integration with Apache Kafka

STEP 3: Batch Processing with Apache Spark

STEP 4: Stream Processing with Apache Kafka Streams

PREDICTIVE MAINTENANCE FOR MANUFACTURING EQUIPMENT

The screenshot shows a Java IDE interface with the following details:

- File Bar:** File, Edit, Selection, View, Go, Run, Terminal, Help.
- EXPLORER:** Shows a project structure under CUBEISOLUTIONS:
 - HadoopDocker
 - Kafka_POC
 - postgresql
 - Predictive Maintenance
 - DataAnalysis.class
 - J DataAnalysis.java
 - J helloworld.class
 - J helloworld.java
 - kafka_consumer.py
 - kafka_producer.py (selected)
 - J KafkaStreamsConsumer.class
 - J KafkaStreamsConsumer.java
 - J KafkaStreamsRealTimeProcessor.class
 - J KafkaStreamsRealTimeProcessor.java
 - load_data.py
 - machine_data.json
 - machine.csv
 - ModelTraining.py
 - PredictiveMaintenanceStreamProcessor.class
 - PredictiveMaintenanceStreamProcessor.java
 - Question.txt
 - J SparkStreamingKafkaConsumer.class
 - J SparkStreamingKafkaConsumer.java
 - Real-time Stock Market
 - Semi_Data_Pipeline
- PROBLEMS:** 106 errors.
- TERMINAL:** Shows log output from a Python script sending data to Kafka:

```
INFO: main :Sent data to Kafka: {'timestamp': '2024-07-11T08:05:00Z', 'machine_id': 'Machine003', 'sensor_type': 'temperature', 'value': 35.8, 'location': 'Assembly Line C'}
```

and many more entries for machines 001-003 across floors A, B, and C.
- OUTPUT:** Shows log output for Java processes.
- OUTLINE:** Shows the project outline.
- TIMELINE:** Shows the project timeline.
- DEBUG CONSOLE:** Shows the debug console.
- Filter (e.g. text, lexclude, \escape):** A search/filter bar.
- Bottom Icons:** Help, Java Projects, and other standard IDE icons.

File Edit Selection View Go Run Terminal Help

EXPLORER

- CUBEASOLUTIONS
 - HadoopPocKer
 - Kafka_POC
 - _pycache_
 - spark_stream
 - consumer_data.py
 - KafkaStreamsRealTimeProcessor.class
 - Production_kafka.py
 - Question.txt
 - StreamsPOC.class
 - StreamsPOC.java
 - postgresql
 - python
 - docker-compose.yml
- Predictive Maintenance
 - hellworld.java
 - hellworld.class
 - kafka_consumer.py
 - kafka_producer.py
 - KafkaStreamsConsumer.class
 - KafkaStreamsConsumer.java
 - KafkaStreamsRealTimeProcessor.class
 - KafkaStreamsRealTimeProcessor.java
 - machine_data.json
 - PredictiveMaintenanceStreamProcessor.class
 - PredictiveMaintenanceStreamProcessor.java
 - Question.txt
- OUTPUT
- OUTLINE
- TIMELINE
- DEBUG CONSOLE

Filter (e.g. text, !exclude, \escape)

J KafkaStreamsConsumer.java 9+ J PredictiveMaintenanceStreamProcessor.java 9+ J KafkaStreamsRealTimeProcessor.java 9+ J StreamsPOC.java 9+ J hellworld.java 4 J KafkaStreamsRe...

7 public class KafkaStreamsConsumer {
8 public static void main(String[] args) {
9 Properties props = new Properties();
10 props.put(StreamsConfig.APPLICATION_ID_CONFIG, "kafka-streams-consumer");
11 props.put(StreamsConfig.BOOTSTRAP_SERVERS_CONFIG, "localhost:9092");
12 props.put(StreamsConfig.DEFAULT_KEY_SERDE_CLASS_CONFIG, Serdes.String().getClass());
13 props.put(StreamsConfig.DEFAULT_VALUE_SERDE_CLASS_CONFIG, Serdes.String().getClass());
14
15 StreamsBuilder builder = new StreamsBuilder();
16 Kstream<String, String> stream = builder.stream("machine_data");
17
18 stream.foreach((key, value) -> {
19 System.out.println("Key: " + key + ", Value: " + value);
20 // Add your processing logic here
21 });
22
23 KafkaStreams streams = new KafkaStreams(builder.build(), props);
24 streams.start();
25 }
26}

PROBLEMS 88 TERMINAL PORTS

[kafka-streams-consumer-67aa37ed-857e-4d0f-9be4-cb079da788b6-StreamThread-1] INFO org.apache.kafka.streams.processor.internals.StreamTask - stream-thread [kafka-streams-consumer-67aa37ed-857e-4d0f-9be4-cb079da788b6-StreamThread-1] task [0_0] Restored and ready to run
[kafka-streams-consumer-67aa37ed-857e-4d0f-9be4-cb079da788b6-StreamThread-1] INFO org.apache.kafka.streams.processor.internals.StreamThread - stream-thread [kafka-streams-consumer-67aa37ed-857e-4d0f-9be4-cb079da788b6-StreamThread-1] Restoration took 66 ms for all active tasks [0_0]
[kafka-streams-consumer-67aa37ed-857e-4d0f-9be4-cb079da788b6-StreamThread-1] INFO org.apache.kafka.streams.processor.internals.StreamThread - stream-thread [kafka-streams-consumer-67aa37ed-857e-4d0f-9be4-cb079da788b6-StreamThread-1] State transition from PARTITIONS_ASSIGNED to RUNNING
[kafka-streams-consumer-67aa37ed-857e-4d0f-9be4-cb079da788b6-StreamThread-1] INFO org.apache.kafka.streams.KafkaStreams - stream-client [kafka-streams-consumer-67aa37ed-857e-4d0f-9be4-cb079da788b6-StreamThread-1] State transition from REBALANCING to RUNNING
[kafka-streams-consumer-67aa37ed-857e-4d0f-9be4-cb079da788b6-StreamThread-1] INFO org.apache.kafka.clients.consumer.internals.LegacyKafkaConsumer - [Consumer clientId=kafka-streams-consumer-67aa37ed-857e-4d0f-9be4-cb079da788b6-StreamThread-1-consumer, groupId=kafka-streams-consumer] Requesting the log end offset for machine-data-0 in order to compute lag
Key: null, Value: {"timestamp": "2024-07-11T08:00:00Z", "machine_id": "Machine001", "sensor_type": "temperature", "value": 35.6, "location": "Factory Floor A"}
Key: null, Value: {"timestamp": "2024-07-11T08:01:00Z", "machine_id": "Machine002", "sensor_type": "temperature", "value": 36.2, "location": "Factory Floor B"}
Key: null, Value: {"timestamp": "2024-07-11T08:02:00Z", "machine_id": "Machine003", "sensor_type": "vibration", "value": 0.015, "location": "Assembly Line C"}
Key: null, Value: {"timestamp": "2024-07-11T08:03:00Z", "machine_id": "Machine001", "sensor_type": "pressure", "value": 1020.5, "location": "Factory Floor A"}
Key: null, Value: {"timestamp": "2024-07-11T08:04:00Z", "machine_id": "Machine002", "sensor_type": "vibration", "value": 0.018, "location": "Factory Floor B"}
Key: null, Value: {"timestamp": "2024-07-11T08:05:00Z", "machine_id": "Machine003", "sensor_type": "temperature", "value": 35.8, "location": "Assembly Line C"}
Key: null, Value: {"timestamp": "2024-07-11T08:06:00Z", "machine_id": "Machine001", "sensor_type": "temperature", "value": 35.9, "location": "Factory Floor A"}
Key: null, Value: {"timestamp": "2024-07-11T08:07:00Z", "machine_id": "Machine002", "sensor_type": "pressure", "value": 1021.2, "location": "Factory Floor B"}
Key: null, Value: {"timestamp": "2024-07-11T08:08:00Z", "machine_id": "Machine003", "sensor_type": "vibration", "value": 0.016, "location": "Assembly Line C"}
Key: null, Value: {"timestamp": "2024-07-11T08:09:00Z", "machine_id": "Machine001", "sensor_type": "vibration", "value": 0.014, "location": "Factory Floor A"}
Key: null, Value: {"timestamp": "2024-07-11T08:10:00Z", "machine_id": "Machine002", "sensor_type": "temperature", "value": 36.0, "location": "Factory Floor B"}
86 △ 2 ⚡ 0 Java: Ready Ln 37, Col 358 Spaces: 4 UTF-8 LF () Java

File Edit Selection View Go Run Terminal Help

EXPLORER

- CUBEASOLUTIONS
 - HadoopPOC
 - Kafka_POC
 - _pycache_
 - spark_stream
 - consumer_data.py
 - J KafkaStreamsRealTimeProcessor.class
 - J KafkaStreamsRealTimeProcessor.java 9+
 - Production_kafka.py
 - Question.txt
 - J StreamsPOC.class
 - J StreamsPOC.java
 - J StreamsPOC.java 9+
 - postgresql
 - python
 - docker-compose.yml
 - Predictive Maintenance
 - J helloworld.class
 - J helloworld.java 4
 - kafka_consumer.py
 - kafka_producer.py
 - J KafkaStreamsConsumer.class
 - J KafkaStreamsConsumer.java 9+
 - { machine_data.json
 - J PredictiveMaintenanceStreamProcessor.class
 - J PredictiveMaintenanceStreamProcessor.java 9+
 - Question.txt
 - J SparkStreamingKafkaConsumer.class
 - OUTPUT
- OUTLINE
- TIMELINE
- DEBUG CONSOLE

Filter (e.g. text, !exclude, \escape)

J Consumer.java 9+ J KafkaStreamsConsumer.java 9+ J PredictiveMaintenanceStreamProcessor.java 9+ J KafkaStreamsRealTimeProcessor.java 9+ X J StreamsPOC.java 9+ J helloworld.java 4

Kafka_POC > J KafkaStreamsRealTimeProcessor.java 9+

```
1 import org.apache.kafka.streams.*;
2 import org.apache.kafka.streams.kstream.*;
3 import org.apache.kafka.common.serialization.Serdes;
4
5
6 import java.util.Properties;
7
8 public class KafkaStreamsRealTimeProcessor {
9     Run|Debug
10    public static void main(String[] args) {
11        // Configure Kafka properties
12
13        ALERT: Potential equipment failure detected!
14        Key: null, Value: {"timestamp": "2024-07-11T08:20:00Z", "machine_id": "Machine003", "sensor_type": "temperature", "value": 35.8, "location": "Assembly Line C"}
15        Anomaly Score: 0.9071891688509176
16        Processing: Key = null, Value = {"timestamp": "2024-07-11T08:21:00Z", "machine_id": "Machine001", "sensor_type": "temperature", "value": 36.0, "location": "Factory Floor A"}
17        ALERT: Potential equipment failure detected!
18        Key: null, Value: {"timestamp": "2024-07-11T08:21:00Z", "machine_id": "Machine001", "sensor_type": "temperature", "value": 36.0, "location": "Factory Floor A"}
19        Anomaly Score: 0.9159941720257557
20        Processing: Key = null, Value = {"timestamp": "2024-07-11T08:22:00Z", "machine_id": "Machine002", "sensor_type": "pressure", "value": 1021.0, "location": "Factory Floor B"}
21        Processing: Key = null, Value = {"timestamp": "2024-07-11T08:23:00Z", "machine_id": "Machine003", "sensor_type": "vibration", "value": 0.017, "location": "Assembly Line C"}
22        ALERT: Potential equipment failure detected!
23        Key: null, Value: {"timestamp": "2024-07-11T08:23:00Z", "machine_id": "Machine003", "sensor_type": "vibration", "value": 0.017, "location": "Assembly Line C"}
24        Anomaly Score: 0.8657048602127781
25        Processing: Key = null, Value = {"timestamp": "2024-07-11T08:24:00Z", "machine_id": "Machine001", "sensor_type": "vibration", "value": 0.013, "location": "Factory Floor A"}
26        Processing: Key = null, Value = {"timestamp": "2024-07-11T08:25:00Z", "machine_id": "Machine002", "sensor_type": "temperature", "value": 35.7, "location": "Factory Floor B"}
27        Processing: Key = null, Value = {"timestamp": "2024-07-11T08:26:00Z", "machine_id": "Machine003", "sensor_type": "temperature", "value": 36.2, "location": "Assembly Line C"}
28        Processing: Key = null, Value = {"timestamp": "2024-07-11T08:27:00Z", "machine_id": "Machine001", "sensor_type": "pressure", "value": 1020.2, "location": "Factory Floor A"}
29        Processing: Key = null, Value = {"timestamp": "2024-07-11T08:28:00Z", "machine_id": "Machine002", "sensor_type": "vibration", "value": 0.018, "location": "Factory Floor B"}
30        Processing: Key = null, Value = {"timestamp": "2024-07-11T08:29:00Z", "machine_id": "Machine003", "sensor_type": "pressure", "value": 1020.8, "location": "Assembly Line C"}
31        Processing: Key = null, Value = {"timestamp": "2024-07-11T08:30:00Z", "machine_id": "Machine001", "sensor_type": "temperature", "value": 35.9, "location": "Factory Floor A"}
32        Processing: Key = null, Value = {"timestamp": "2024-07-11T08:31:00Z", "machine_id": "Machine002", "sensor_type": "temperature", "value": 35.8, "location": "Factory Floor B"}
33        Processing: Key = null, Value = {"timestamp": "2024-07-11T08:32:00Z", "machine_id": "Machine003", "sensor_type": "vibration", "value": 0.016, "location": "Assembly Line C"}
34        Processing: Key = null, Value = {"timestamp": "2024-07-11T08:33:00Z", "machine_id": "Machine001", "sensor_type": "vibration", "value": 0.014, "location": "Factory Floor A"}
```

PROBLEMS 88 TERMINAL PORTS

java Kafka_POC bash Predictive...

Ln 3, Col 3 Spaces: 4 UTF-8 LF () Java

File Edit Selection View Go Run Terminal Help

EXPLORER

- CUBEAI SOLUTIONS
 - Predictive Maintenance
 - J DataAnalysis.java 9+
 - J helloworld.class
 - J helloworld.java 4
 - ↳ kafka_producer.py
 - J KafkaStreamsConsumer.class
 - J KafkaStreamsConsumer.java 9+
 - J KafkaStreamsRealTimeProcessor.class
 - J KafkaStreamsRealTimeProcessor.java 9+
 - ↳ load_data.py
 - { machine_data.json
 - machine.csv
 - ↳ ModelTraining.py
 - J PredictiveMaintenanceStreamProcessor.class
 - J PredictiveMaintenanceStreamProcessor.java 9+
 - J PredictiveMaintenanceStreamProcessor\$Tempe...
 - J PredictiveMaintenanceStreamProcessor\$Tempe...
 - J PredictiveMaintenanceStreamProcessor\$Tempe...
 - J PredictiveMaintenanceStreamProcessor\$Tempe...
 - Question.txt
 - J SparkStreamingKafkaConsumer.class

OUTPUT

PROBLEMS 123 TERMINAL PORTS

```

java Predictive Main...
python3 Predictive ...
java Predictive M...
bash Predictive M...
java Predictive Main...
bash Predictive Mai...
  
```

Received record: key=null, value={"timestamp": "2024-07-11T08:27:00Z", "machine_id": "Machine001", "sensor_type": "vibration", "value": 0.018, "location": "Factory Floor A"}
Received record: key=null, value={"timestamp": "2024-07-11T08:28:00Z", "machine_id": "Machine002", "sensor_type": "vibration", "value": 0.018, "location": "Factory Floor B"}
Received record: key=null, value={"timestamp": "2024-07-11T08:29:00Z", "machine_id": "Machine003", "sensor_type": "pressure", "value": 1020.8, "location": "Assembly Line C"}
Received record: key=null, value={"timestamp": "2024-07-11T08:30:00Z", "machine_id": "Machine001", "sensor_type": "temperature", "value": 35.9, "location": "Factory Floor A"}
Received record: key=null, value={"timestamp": "2024-07-11T08:31:00Z", "machine_id": "Machine002", "sensor_type": "temperature", "value": 35.8, "location": "Factory Floor B"}
Received record: key=null, value={"timestamp": "2024-07-11T08:32:00Z", "machine_id": "Machine003", "sensor_type": "vibration", "value": 0.016, "location": "Assembly Line C"}
Received record: key=null, value={"timestamp": "2024-07-11T08:33:00Z", "machine_id": "Machine001", "sensor_type": "vibration", "value": 0.014, "location": "Factory Floor A"}
Received record: key=null, value={"timestamp": "2024-07-11T08:35:00Z", "machine_id": "Machine001", "sensor_type": "vibration", "value": 0.014, "location": "Factory Floor A"}
Received record: key=null, value={"timestamp": "2024-07-11T08:32:00Z", "machine_id": "Machine003", "sensor_type": "vibration", "value": 0.016, "location": "Assembly Line C"}
Received record: key=null, value={"timestamp": "2024-07-11T08:33:00Z", "machine_id": "Machine001", "sensor_type": "vibration", "value": 0.014, "location": "Factory Floor A"}
Received record: key=null, value={"timestamp": "2024-07-11T08:35:00Z", "machine_id": "Machine001", "sensor_type": "vibration", "value": 0.014, "location": "Factory Floor A"}

Ln 45, Col 9 Spaces: 4 UTF-8 LF () Java Q

File Edit Selection View Go Run Terminal Help

EXPLORER

CUBEASOLUTIONS

- E-commerce transaction
- ETL
- Hadoopdocker
- Kafka_POC
- postgresql
- Predictive Maintenance
- DataAnalysis.class
- J DataAnalysis.java 9+
- J helloworld.class
- J helloworld.java 4
- kafka_consumer.py
- kafka_producer.py
- J KafkaStreamsConsumer.class
- J KafkaStreamsConsumer.java 9+
- J KafkaStreamsRealTimeProcessor.class
- J KafkaStreamsRealTimeProcessor.java 9+
- load_data.py
- machine_data.json
- machine.csv
- J ModelTraining.java 9+
- J PredictiveMaintenanceStreamProcessor.class
- J PredictiveMaintenanceStreamProcessor.java 9+
- Question.txt
- J SparkStreamingKafkaConsumer.class
- J SparkStreamingKafkaConsumer.java 9+
- Real-time Stock Market

PROBLEMS 131

TUTORIAL PORTS

bash Kafka_POC

bash Predictiv...

bash Predictiv...

bash Predictiv...

bash Predictiv...

24/07/11 17:31:08 INFO TaskSchedulerImpl: Removed TaskSet 3.0, whose tasks have all completed, from pool

24/07/11 17:31:08 INFO DAGScheduler: ResultStage 3 (show at DataAnalysis.java:25) finished in 0.187 s

24/07/11 17:31:08 INFO DAGScheduler: Job 2 is finished. Cancelling potential speculative or zombie tasks for this job

24/07/11 17:31:08 INFO TaskSchedulerImpl: Killing all running tasks in stage 3: Stage finished

24/07/11 17:31:08 INFO DAGScheduler: Job 2 finished: show at DataAnalysis.java:25, took 0.222068 s

24/07/11 17:31:08 INFO CodeGenerator: Code generated in 18.278092 ms

machine_id	avg_temperature	avg_vibration
Machine003	198.57990909090913	0.0158
Machine001	327.71961538461534	0.01375
Machine002	201.98809090909097	0.01725

24/07/11 17:31:08 INFO SparkContext: SparkContext is stopping with exitCode 0.

24/07/11 17:31:08 INFO SparkUI: Stopped Spark web UI at http://192.168.1.38:4040

24/07/11 17:31:09 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!

24/07/11 17:31:09 INFO MemoryStore: MemoryStore cleared

Ln 35, Col 107 Spaces: 4 UTF-8 LF () Java

File Edit Selection View Go Run Terminal Help

EXPLORER

- CUBEASOLUTIONS
 - / HadoopPocKer
 - > Kafka_POC
 - > postgresql
 - > Predictive Maintenance
 - J DataAnalysis.class
 - J DataAnalysis.java 9+
 - J helloworld.class
 - J helloworld.java 4
 - kafka_consumer.py
 - kafka_producer.py
 - J KafkaStreamsConsumer.class
 - J KafkaStreamsConsumer.java 9+
 - J KafkaStreamsRealTimeProcessor.class
 - J KafkaStreamsRealTimeProcessor.java 9+
 - load_data.py
 - machine_data.json
 - machine.csv
 - ModelTraining.py
 - J PredictiveMaintenanceStreamProcessor.class
 - J PredictiveMaintenanceStreamProcessor.java 9+
 - Question.txt
 - J SparkStreamingKafkaConsumer.class
 - J SparkStreamingKafkaConsumer.java 9+
 - > Real-time Stock Market
 - > Semi_Data_Pipeline

PREDICTIVE MAINTENANCE > ModelTraining.py > ...

```
1 import pandas as pd
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.model_selection import train_test_split
4
5 # Load data exported from Spark (example: CSV file)
6 data = pd.read_csv('machine.csv')
7
8 # Convert timestamp to numeric (Unix timestamp)
9 data['timestamp'] = pd.to_datetime(data['timestamp']).astype(int) // 10**9
10
11 # One-hot encode categorical variables (assuming 'sensor_type' and 'location' are categorical)
12 data_encoded = pd.get_dummies(data, columns=['sensor_type', 'location'])
13
14 # Prepare data for training
15 X = data_encoded.drop(['machine_id'], axis=1) # Assuming 'machine_id' is not the target
16 y = data_encoded['machine_id'] # Assuming 'machine_id' is the target variable
17
18 # Split data into training and testing sets
19 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
20
21 # Initialize and train a machine learning model (example: Random Forest)
22 model = RandomForestClassifier()
23 model.fit(X_train, y_train)
24
25 # Evaluate model performance
26 accuracy = model.score(X_test, y_test)
27 print(f"Model accuracy: {accuracy}")
28
```

OUTPUT

PROBLEMS 101 TERMINAL PORTS

- nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/CUBEASOLUTIONS/Predictive Maintenance\$ python3 ModelTraining.py
- Model accuracy: 1.0
- nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/CUBEASOLUTIONS/Predictive Maintenance\$

OUTLINE

TIMELINE

DEBUG CONSOLE

Filter (e.g. text, !exclude, \escape)

Java: Ready

Ln 28, Col 1 Spaces: 8 UTF-8 LF Python 3.10.12 64-bit

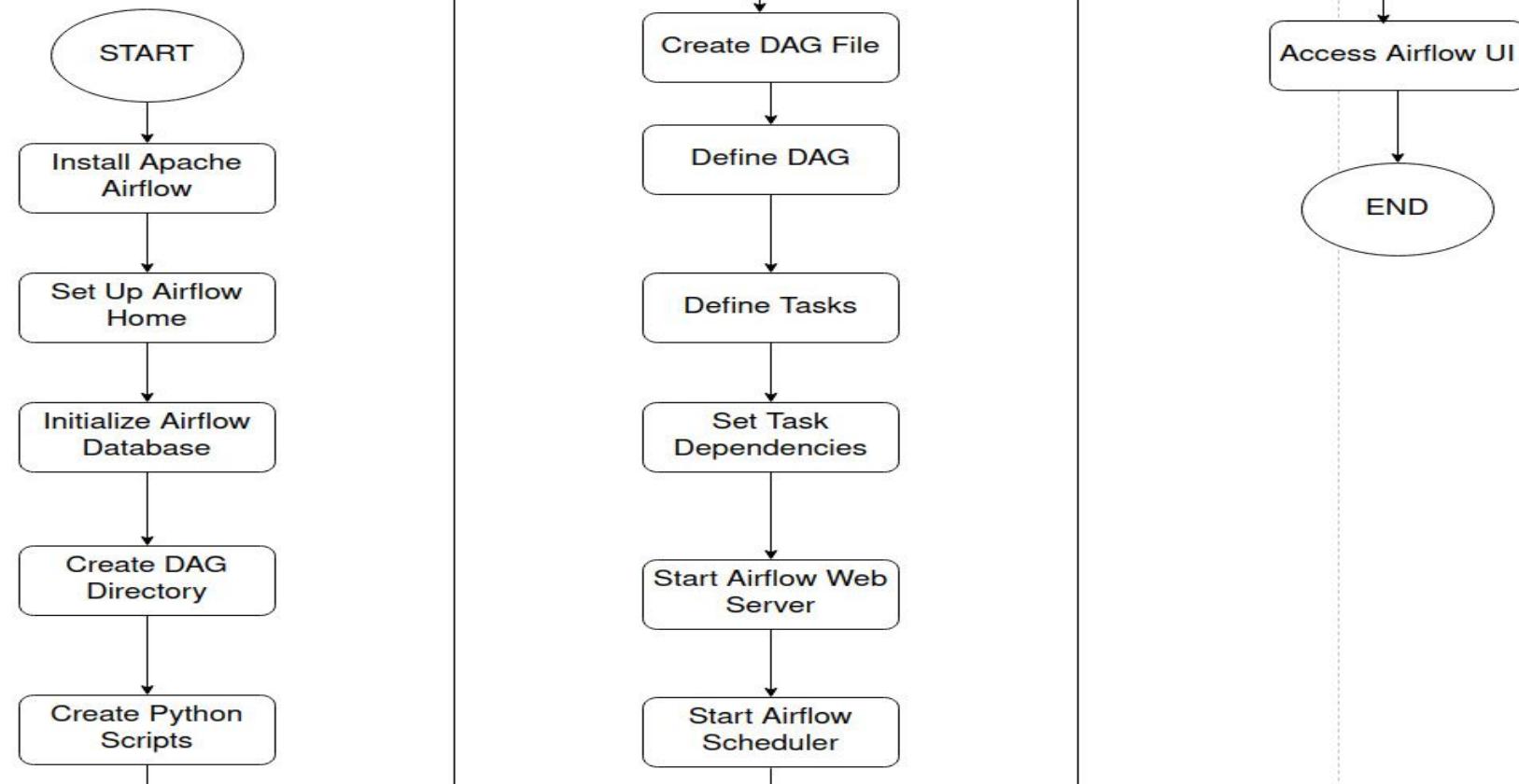
EXECUTE DAILY PYTHON SCRIPTS PROCESS

STEP 1: Create the Python Scripts

STEP 2: Create the DAG File

STEP 3: Start the Airflow Web Server

EXECUTE DAILY PYTHON SCRIPTS



localhost:9080/dags/daily_python_scripts/grid?tab=graph&dag_run_id=scheduled_2024-07-11T00%3A00%3A00%2B00%3A00&task_id=run_script2

Airflow DAGs Cluster Activity Datasets Security Browse Admin Docs 11:24 UTC AU

Triggered daily_python_scripts, it should start any moment now.

DAG: daily_python_scripts Run Python scripts daily Schedule: 1 day, 0:00:00 Next Run ID: 2024-07-12, 00:00:00 UTC Auto-refresh 25

12/07/2024 11:23:01 am All Run Types All Run States Clear Filters

Press shift + / for Shortcuts

Duration: 00:01:50

DAG: daily_python_scripts Run: 2024-07-11, 00:00:00 UTC / Task: run_script2

Details Graph Gantt Audit Log Logs XCom Task Duration

Clear task Mark state as... Filter DAG by task

Layout: Left > Right

run_script1 success PythonOperator

run_script2 success PythonOperator

React Flow

```
graph LR; run_script1[run_script1] --> run_script2[run_script2]
```

MEDALLION TASK

Bronze Layer

- **Purpose:** Raw data ingestion.
- **Characteristics:** This layer contains raw data from various sources, often stored in its original format without any transformations or cleaning. It acts as the landing zone for all incoming data.
- **Examples:** Logs, sensor data, transactional data from databases.

Silver Layer

- **Purpose:** Data cleansing and enrichment.
- **Characteristics:** This layer processes and cleans the data from the bronze layer. It involves removing duplicates, handling missing values, and transforming the data into a more usable format. This layer often includes normalized and enriched data.
- **Examples:** Processed sensor readings, cleaned and standardized transactional data.

Gold Layer

- **Purpose:** High-quality, aggregated data for analytics and reporting.
- **Characteristics:** This layer contains highly refined data ready for consumption by business intelligence tools, machine learning models, and other analytical applications. It often involves aggregations, calculations, and business logic applied to the silver

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/kafka
```

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/kafka$ cd
```

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/kafka$ hdfs dfs -cat /Kafka/kafka_etl/bronze/*
```

```
userId,id,title,completed
```

```
1,1,delectus aut autem,False
```

```
1,2,quis ut nam facilis et officia qui,False
```

```
1,3,fugiat veniam minus,False
```

```
1,4,et porro tempora,True
```

```
1,5,laboriosam mollitia et enim quasi adipisci quia provident illum,False
```

```
1,6,qui ullam ratione quibusdam voluptatem quia omnis,False
```

```
1,7,illo expedita consequatur quia in,False
```

```
1,8,quo adipisci enim quam ut ab,True
```

```
1,9,molestiae perspiciat is ipsa,False
```

```
1,10,illo est ratione doloremque quia maiores aut,True
```

```
1,11,vero rerum temporibus dolor,True
```

```
1,12,ipsa repellendus fugit nisi,True
```

```
1,13,et doloremque nulla,False
```

```
1,14,repellendus sunt dolores architecto voluptatum,True
```

```
1,15,ab voluptatum amet voluptas,True
```

```
1,16,accusamus eos facilis sint et aut voluptatem,True
```

```
1,17,quo laboriosam deleniti aut qui,True
```

```
1,18,dolorum est consequatur ea mollitia in culpa,False
```

```
1,19,molestiae ipsa aut voluptatibus pariatur dolor nihil,True
```

```
1,20,ullam nobis libero sapiente ad optio sint,True
```

```
2,21,suscipit repellat esse quibusdam voluptatem incident, False
```

```
2,22,distinctio vitae autem nihil ut molestias quo,True
```

```
2,23,et itaque necessitatibus maxime molestiae qui quas velit,False
```

```
2,24,adipisci non ad dicta qui amet quaerat doloribus ea,False
```

```
2,25,voluptas quo tenetur perspiciat is explicabo natus,True
```

```
2,26,aliquam aut quasi,True
```

```
2,27,veritatis pariatur delectus,True
```

```
2,28,nesciunt totam sit blanditiis sit,False
```

```
2,29,laborum aut in quan,False
```

```
2,30,nemo perspiciat is repellat ut dolor libero commodi blanditiis omnis,True
```

```
2,31,repudiandae totam in est sint facere fuga,False
```

```
2,32,earum doloribus ea doloremque quis,False
```

```
2,33,sint sit aut vero,False
```

```
2,34,porro aut necessitatibus eaque distinctio,False
```

```
2,35,repellendus veritatis molestias dicta incident, True
```

```
2,36,excepturi deleniti adipisci voluptatem et neque optio illum ad,True
```

```
2,37,sunt cum tempora,False
```

```
2,38,totam quia non,False
```

```
2,39,doloremque quibusdam asperiores libero corrupti illum qui omnis,False
```

```
2,40,totam atque quo nesciunt,True
```

```
3,41,aliquid amet impedit consequatur aspernatur placeat eaque fugiat suscipit,False
```

```
3,42,rerum perferendis error quia ut eveniet,False
```

```
3,43,tempore ut sint quis recusandae,True
```

```
3,44,cum debitis quis accusamus doloremque ipsa natus sapiente omnis,True
```

```
3,45,velit soluta adipisci molestias reiciendis harum,False
```

```
3,46,vel voluptatem repellat nihil placeat corporis,False
```

```
3,47,nam qui rerum fugiat accusamus,False
```

```
3,48,sit reprehenderit omnis quia,False
```

```
3,49,ut necessitatibus aut maiores debitis officia blanditiis velit et,False
```

```
3,50,cupiditate necessitatibus ullam aut quis dolor voluntate,True
```

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~
```

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/kafka

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~

9,161,ex hic consequuntur earum omnis alias ut occaecati culpa,True
9,162,omnis laboriosam molestias animi sunt dolore,True
9,163,natus corrupti maxime laudantium et voluptatem laboriosam odit,False
9,164,reprehenderit quos aut aut consequatur est sed,False
9,165,fugiat perferendis sed aut quidem,False
9,166,quos quo possimus suscipit minima ut,False
9,167,et quis minus quo a asperiores molestiae,False
9,168,recusandae quia qui sunt libero,False
9,169,ea odio perferendis officitis,True
9,170,quisquam aliquam quia doloribus aut,False
9,171,fugiat aut voluptatibus corrupti delecti velit iste odio,True
9,172,et provident amet rerum consectetur et voluptatum,False
9,173,harum ad aperiam quis,False
9,174,similique aut quo,False
9,175,laudantium eius officia perferendis provident perspiciatibus asperiores,True
9,176,magni soluta corrupti ut maiores rem quidem,False
9,177,et placeat temporibus voluptas est tempora quos quibusdam,False
9,178,nesciunt itaque commodi tempore,True
9,179,omnis consequuntur cupiditate impedit itaque ipsam quo,True
9,180,debitis nisi et dolorem repellat et,True
10,181,ut cupiditate sequi aliquam fuga maiores,False
10,182,inventore saepe cunque et aut illum enim,True
10,183,omnis nulla eum aliquam distinctio,True
10,184,molestias modi perferendis perspiciatibus,False
10,185,voluptates dignissimos sed doloribus animi quaerat aut,False
10,186,explicabo odio est et,False
10,187,consequuntur animi possimus,False
10,188,vel non beatae est,True
10,189,culpa eius et voluptatem et,True
10,190,accusamus sint iusto et voluptatem exercitationem,True
10,191,temporibus atque distinctio omnis eius impedit tempore molestias pariatur,True
10,192,ut quas possimus exercitationem sint voluptates,False
10,193,rerum debitum voluptatem qui eveniet tempora distinctio a,True
10,194,sed ut vero sit molestiae,False
10,195,rerum ex veniam mollitia voluptatibus pariatur,True
10,196,consequuntur aut ut fugit similique,True
10,197,dignissimos quo nobis earum saepe,True
10,198,quis eius est sint explicabo,True
10,199,numquam repellendus a magnam,True
10,200,ipsam aperiam voluptates qui,False
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~\$ hdfs dfs -cat /Kafka/kafka_etl/gold/*
userId,id
1,20
2,20
3,20
4,20
5,20
6,20
7,20
8,20
9,20
10,20
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~\$

Banking Applications

STEP 1: Data Ingestion

STEP 2: Data Storage

STEP 3: Data Processing

STEP 4: Data Transformation

STEP 5: Fraud Detection Analytics

STEP 6: Real-time Alerts

STEP 7: Data Visualization and Reporting



File Edit Selection View Go Run Terminal Help

EXPLORER



- TEST
- ATM
 - ATM_transaction.csv
 - ATM_transaction.json
 - convert_json_csv.py
 - kafka_consumer_atm.py
 - kafka_consumer_online.py**
 - kafka_hadoop_atm.py
 - kafka_hadoop_mobile.py
 - kafka_hadoop_online.py
 - kafka_producer_atm.py
 - kafka_producer_mobile.py
 - kafka_producer_online.py
 - mobile_banking.csv
 - mobile_banking.json
 - online_banking.csv
 - online_banking.json
- > checkpoint
- > kafka_topic
- > spark_log
- csv_stream.py

OUTPUT

PROBLEMS

TERMINAL

PORTS

bash

bash

bash ATM

bash ATM

bash ATM

- nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/test/ATM\$ python3 kafka_producer_atm.py
All messages have been sent to Kafka topic: atm-transactions
- nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/test/ATM\$ python3 kafka_producer_online.py
All messages have been sent to Kafka topic: online-banking-transactions
- nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/test/ATM\$ python3 kafka_producer_mobile.py
All messages have been sent to Kafka topic: mobile-banking-transactions
- nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/test/ATM\$

OUTLINE

TIMELINE

DEBUG CONSOLE

Filter (e.g. text, |exclude, \escape)

0 0 0 % 0

Ln 18, Col 1 Spaces: 4 UTF-8 LF Python 3.10.12 64-bit

```
{"transaction_id": "ATM17", "user_id": "U7", "amount": 65.0, "location": "ATM159", "timestamp": "2024-07-01T11:20:00Z"}, {"transaction_id": "ATM18", "user_id": "U8", "amount": 75.0, "location": "ATM140", "timestamp": "2024-07-01T11:25:00Z"}, {"transaction_id": "ATM19", "user_id": "U9", "amount": 65.0, "location": "ATM141", "timestamp": "2024-07-01T11:30:00Z"}, {"transaction_id": "ATM20", "user_id": "U10", "amount": 55.0, "location": "ATM142", "timestamp": "2024-07-01T11:35:00Z"}, {"transaction_id": "ATM21", "user_id": "U11", "amount": 45.0, "location": "ATM143", "timestamp": "2024-07-01T11:40:00Z"}, {"transaction_id": "ATM22", "user_id": "U2", "amount": 35.0, "location": "ATM144", "timestamp": "2024-07-01T11:45:00Z"}, {"transaction_id": "ATM23", "user_id": "U3", "amount": 25.0, "location": "ATM145", "timestamp": "2024-07-01T11:50:00Z"}, {"transaction_id": "ATM24", "user_id": "U4", "amount": 20.0, "location": "ATM146", "timestamp": "2024-07-01T11:55:00Z"}, {"transaction_id": "ATM25", "user_id": "U5", "amount": 15.0, "location": "ATM147", "timestamp": "2024-07-01T12:00:00Z"}, {"transaction_id": "ATM26", "user_id": "U6", "amount": 10.0, "location": "ATM148", "timestamp": "2024-07-01T12:05:00Z"}, {"transaction_id": "ATM27", "user_id": "U7", "amount": 5.0, "location": "ATM149", "timestamp": "2024-07-01T12:10:00Z"}, {"transaction_id": "ATM28", "user_id": "U8", "amount": 25.0, "location": "ATM150", "timestamp": "2024-07-01T12:15:00Z"}, {"transaction_id": "ATM29", "user_id": "U9", "amount": 35.0, "location": "ATM151", "timestamp": "2024-07-01T12:20:00Z"}, {"transaction_id": "ATM30", "user_id": "U10", "amount": 45.0, "location": "ATM152", "timestamp": "2024-07-01T12:25:00Z"}]
```

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $ hdfs dfs -cat /kafka/ATM/mobile_transaction
{"transaction_id": "U005", "user_id": "U005", "amount": 900.0, "account": "AC005", "timestamp": "2024-07-01T10:20:00Z"}, {"transaction_id": "U006", "user_id": "U036", "amount": 1000.0, "account": "AC006", "timestamp": "2024-07-01T10:25:00Z"}, {"transaction_id": "U010", "user_id": "U040", "amount": 1400.0, "account": "AC010", "timestamp": "2024-07-01T10:45:00Z"}, {"transaction_id": "U011", "user_id": "U041", "amount": 1500.0, "account": "AC011", "timestamp": "2024-07-01T10:50:00Z"}, {"transaction_id": "U015", "user_id": "U045", "amount": 1900.0, "account": "AC015", "timestamp": "2024-07-01T11:00:00Z"}, {"transaction_id": "U018", "user_id": "U048", "amount": 2200.0, "account": "AC018", "timestamp": "2024-07-01T11:25:00Z"}, {"transaction_id": "U019", "user_id": "U049", "amount": 2300.0, "account": "AC019", "timestamp": "2024-07-01T11:30:00Z"}, {"transaction_id": "U023", "user_id": "U053", "amount": 2700.0, "account": "AC023", "timestamp": "2024-07-01T11:50:00Z"}, {"transaction_id": "U026", "user_id": "U056", "amount": 3000.0, "account": "AC026", "timestamp": "2024-07-01T12:05:00Z"}, {"transaction_id": "U001", "user_id": "U031", "amount": 500.0, "account": "AC001", "timestamp": "2024-07-01T10:00:00Z"}, {"transaction_id": "U003", "user_id": "U033", "amount": 700.0, "account": "AC003", "timestamp": "2024-07-01T10:10:00Z"}, {"transaction_id": "U008", "user_id": "U038", "amount": 1200.0, "account": "AC008", "timestamp": "2024-07-01T10:35:00Z"}, {"transaction_id": "U012", "user_id": "U042", "amount": 1600.0, "account": "AC012", "timestamp": "2024-07-01T10:55:00Z"}, {"transaction_id": "U014", "user_id": "U044", "amount": 1800.0, "account": "AC014", "timestamp": "2024-07-01T11:05:00Z"}, {"transaction_id": "U016", "user_id": "U046", "amount": 2000.0, "account": "AC016", "timestamp": "2024-07-01T11:15:00Z"}, {"transaction_id": "U017", "user_id": "U047", "amount": 2100.0, "account": "AC017", "timestamp": "2024-07-01T11:20:00Z"}, {"transaction_id": "U0822", "user_id": "U052", "amount": 2600.0, "account": "AC022", "timestamp": "2024-07-01T11:45:00Z"}, {"transaction_id": "U0825", "user_id": "U055", "amount": 2900.0, "account": "AC025", "timestamp": "2024-07-01T12:00:00Z"}, {"transaction_id": "U0830", "user_id": "U060", "amount": 3400.0, "account": "AC030", "timestamp": "2024-07-01T12:25:00Z"}, {"transaction_id": "U0802", "user_id": "U032", "amount": 600.0, "account": "AC002", "timestamp": "2024-07-01T10:05:00Z"}, {"transaction_id": "U0804", "user_id": "U034", "amount": 800.0, "account": "AC004", "timestamp": "2024-07-01T10:15:00Z"}, {"transaction_id": "U0807", "user_id": "U037", "amount": 1100.0, "account": "AC007", "timestamp": "2024-07-01T10:30:00Z"}, {"transaction_id": "U0809", "user_id": "U039", "amount": 1300.0, "account": "AC009", "timestamp": "2024-07-01T10:40:00Z"}, {"transaction_id": "U0813", "user_id": "U043", "amount": 1700.0, "account": "AC013", "timestamp": "2024-07-01T11:00:00Z"}, {"transaction_id": "U0820", "user_id": "U050", "amount": 2400.0, "account": "AC020", "timestamp": "2024-07-01T11:35:00Z"}, {"transaction_id": "U0821", "user_id": "U051", "amount": 2500.0, "account": "AC021", "timestamp": "2024-07-01T11:40:00Z"}, {"transaction_id": "U0824", "user_id": "U054", "amount": 2800.0, "account": "AC024", "timestamp": "2024-07-01T11:55:00Z"}, {"transaction_id": "U0827", "user_id": "U057", "amount": 3100.0, "account": "AC027", "timestamp": "2024-07-01T12:10:00Z"}, {"transaction_id": "U0828", "user_id": "U058", "amount": 3200.0, "account": "AC028", "timestamp": "2024-07-01T12:15:00Z"}, {"transaction_id": "U0829", "user_id": "U059", "amount": 3300.0, "account": "AC029", "timestamp": "2024-07-01T12:20:00Z"}]
```

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: $ hdfs dfs -cat /kafka/ATM/mobile_transaction
{"transaction_id": "MB003", "user_id": "U063", "amount": 350.0, "device_id": "D003", "timestamp": "2024-07-01T10:10:00Z"}, {"transaction_id": "MB006", "user_id": "U066", "amount": 650.0, "device_id": "D006", "timestamp": "2024-07-01T10:25:00Z"}, {"transaction_id": "MB009", "user_id": "U069", "amount": 950.0, "device_id": "D009", "timestamp": "2024-07-01T10:40:00Z"}, {"transaction_id": "MB015", "user_id": "U075", "amount": 1550.0, "device_id": "D015", "timestamp": "2024-07-01T11:10:00Z"}, {"transaction_id": "MB020", "user_id": "U080", "amount": 2050.0, "device_id": "D020", "timestamp": "2024-07-01T11:35:00Z"}, {"transaction_id": "MB025", "user_id": "U085", "amount": 2550.0, "device_id": "D025", "timestamp": "2024-07-01T12:00:00Z"}, {"transaction_id": "MB026", "user_id": "U086", "amount": 2650.0, "device_id": "D026", "timestamp": "2024-07-01T12:05:00Z"}, {"transaction_id": "MB028", "user_id": "U088", "amount": 2850.0, "device_id": "D028", "timestamp": "2024-07-01T12:15:00Z"}]
```

```

INFO : Completed executing command(queryId=nandhumidhun_20240724150540_40471927-6865-4072-b844-013c6050b492); Time taken: 0.073 seconds
0: jdbc:hive2://localhost:10000> SELECT * FROM atm_transactions;
INFO : Compiling command(queryId=nandhumidhun_20240724150552_a3c4b815-3137-44c6-b7db-4eea6accfed0): SELECT * FROM atm_transactions
INFO : No Stats for hadoop@atm_transactions, Columns: transaction_id, amount, user_id, location, timestamp
INFO : Semantic Analysis Completed (retrial = false)
INFO : Created Hive schema: Schema(fieldschemas:[FieldSchema(name:atm_transactions.transaction_id, type:string, comment:null), FieldSchema(name:atm_transactions.user_id, type:string, comment:null), FieldSchema(name:atm_transactions.amount, type:double, comment:null), FieldSchema(name:atm_transactions.location, type:string, comment:null), FieldSchema(name:atm_transactions.timestamp, type:string, comment:null)], properties:null)
INFO : Completed compiling command(queryId=nandhumidhun_20240724150552_a3c4b815-3137-44c6-b7db-4eea6accfed0); Time taken: 0.174 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=nandhumidhun_20240724150552_a3c4b815-3137-44c6-b7db-4eea6accfed0): SELECT * FROM atm_transactions
INFO : Completed executing command(queryId=nandhumidhun_20240724150552_a3c4b815-3137-44c6-b7db-4eea6accfed0); Time taken: 0.0 seconds
+-----+-----+-----+-----+-----+
| atm_transactions.transaction_id | atm_transactions.user_id | atm_transactions.amount | atm_transactions.location | atm_transactions.timestamp |
+-----+-----+-----+-----+-----+
| ATM1           | U1            | 100.0          | ATM123          | 2024-07-01 10:00:00+00:00
| ATM2           | U2            | 50.0           | ATM124          | 2024-07-01 10:05:00+00:00
| ATM3           | U3            | 200.0          | ATM125          | 2024-07-01 10:10:00+00:00
| ATM4           | U4            | 80.0           | ATM126          | 2024-07-01 10:15:00+00:00
| ATM5           | U5            | 150.0          | ATM127          | 2024-07-01 10:20:00+00:00
| ATM6           | U6            | 30.0           | ATM128          | 2024-07-01 10:25:00+00:00
| ATM7           | U7            | 60.0           | ATM129          | 2024-07-01 10:30:00+00:00
| ATM8           | U8            | 40.0           | ATM130          | 2024-07-01 10:35:00+00:00
| ATM9           | U9            | 120.0          | ATM131          | 2024-07-01 10:40:00+00:00
| ATM10          | U10           | 70.0           | ATM132          | 2024-07-01 10:45:00+00:00
| ATM11          | U1            | 90.0           | ATM133          | 2024-07-01 10:50:00+00:00
| ATM12          | U2            | 110.0          | ATM134          | 2024-07-01 10:55:00+00:00
| ATM13          | U3            | 140.0          | ATM135          | 2024-07-01 11:00:00+00:00
| ATM14          | U4            | 160.0          | ATM136          | 2024-07-01 11:05:00+00:00
| ATM15          | U5            | 130.0          | ATM137          | 2024-07-01 11:10:00+00:00
| ATM16          | U6            | 95.0           | ATM138          | 2024-07-01 11:15:00+00:00
| ATM17          | U7            | 85.0           | ATM139          | 2024-07-01 11:20:00+00:00
| ATM18          | U8            | 75.0           | ATM140          | 2024-07-01 11:25:00+00:00
| ATM19          | U9            | 65.0           | ATM141          | 2024-07-01 11:30:00+00:00
| ATM20          | U10           | 55.0           | ATM142          | 2024-07-01 11:35:00+00:00
| ATM21          | U1            | 45.0           | ATM143          | 2024-07-01 11:40:00+00:00
| ATM22          | U2            | 35.0           | ATM144          | 2024-07-01 11:45:00+00:00
| ATM23          | U3            | 25.0           | ATM145          | 2024-07-01 11:50:00+00:00
| ATM24          | U4            | 20.0           | ATM146          | 2024-07-01 11:55:00+00:00
| ATM25          | U5            | 15.0           | ATM147          | 2024-07-01 12:00:00+00:00
| ATM26          | U6            | 10.0           | ATM148          | 2024-07-01 12:05:00+00:00
| ATM27          | U7            | 5.0            | ATM149          | 2024-07-01 12:10:00+00:00
| ATM28          | U8            | 25.0           | ATM150          | 2024-07-01 12:15:00+00:00
| ATM29          | U9            | 35.0           | ATM151          | 2024-07-01 12:20:00+00:00
| ATM30          | U10           | 45.0           | ATM152          | 2024-07-01 12:25:00+00:00
+-----+-----+-----+-----+-----+
30 rows selected (0.208 seconds)
0: jdbc:hive2://localhost:10000> SELECT * FROM online_transactions;
INFO : Compiling command(queryId=nandhumidhun_20240724150559_f3879f78-677f-47ce-b224-e12d5b8eeee61): SELECT * FROM online_transactions
INFO : No Stats for hadoop@online_transactions, Columns: transaction_id, amount, user_id, location, timestamp
INFO : Semantic Analysis Completed (rettrial = false)
INFO : Created Hive schema: Schema(fieldschemas:[FieldSchema(name:online_transactions.transaction_id, type:string, comment:null), FieldSchema(name:online_transactions.user_id, type:string, comment:null),

```

```
INFO : No Stats for hivequeryonline_transactions, Columns: transaction_id, amount, user_id, location, timestamp
INFO : Semantic Analysis Completed (retrial = false)
INFO : Created Hive schema: SchemaFieldSchemas:[FieldSchema(name:online_transactions.transaction_id, type:string, comment:null), FieldSchema(name:online_transactions.user_id, type:string, comment:null), FieldSchema(name:online_transactions.amount, type:double, comment:null), FieldSchema(name:online_transactions.location, type:string, comment:null), FieldSchema(name:online_transactions.timestamp, type:string, comment:null)], properties:null
INFO : Completed compiling command(queryId=nandhumidhun_20240724150559_f3879f78-677f-47ce-b224-e12d5b8eee61); Time taken: 0.163 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=nandhumidhun_20240724150559_f3879f78-677f-47ce-b224-e12d5b8eee61): SELECT * FROM online_transactions
INFO : Completed executing command(queryId=nandhumidhun_20240724150559_f3879f78-677f-47ce-b224-e12d5b8eee61); Time taken: 0.0 seconds
```

online_transactions.transaction_id	online_transactions.user_id	online_transactions.amount	online_transactions.location	online_transactions.timestamp
08001	U031	500.0	AC001	2024-07-01 10:00:00+00:00
08002	U032	600.0	AC002	2024-07-01 10:05:00+00:00
08003	U033	700.0	AC003	2024-07-01 10:10:00+00:00
08004	U034	800.0	AC004	2024-07-01 10:15:00+00:00
08005	U035	900.0	AC005	2024-07-01 10:20:00+00:00
08006	U036	1000.0	AC006	2024-07-01 10:25:00+00:00
08007	U037	1100.0	AC007	2024-07-01 10:30:00+00:00
08008	U038	1200.0	AC008	2024-07-01 10:35:00+00:00
08009	U039	1300.0	AC009	2024-07-01 10:40:00+00:00
08010	U040	1400.0	AC010	2024-07-01 10:45:00+00:00
08011	U041	1500.0	AC011	2024-07-01 10:50:00+00:00
08012	U042	1600.0	AC012	2024-07-01 10:55:00+00:00
08013	U043	1700.0	AC013	2024-07-01 11:00:00+00:00
08014	U044	1800.0	AC014	2024-07-01 11:05:00+00:00
08015	U045	1900.0	AC015	2024-07-01 11:10:00+00:00
08016	U046	2000.0	AC016	2024-07-01 11:15:00+00:00
08017	U047	2100.0	AC017	2024-07-01 11:20:00+00:00
08018	U048	2200.0	AC018	2024-07-01 11:25:00+00:00
08019	U049	2300.0	AC019	2024-07-01 11:30:00+00:00
08020	U050	2400.0	AC020	2024-07-01 11:35:00+00:00
08021	U051	2500.0	AC021	2024-07-01 11:40:00+00:00
08022	U052	2600.0	AC022	2024-07-01 11:45:00+00:00
08023	U053	2700.0	AC023	2024-07-01 11:50:00+00:00
08024	U054	2800.0	AC024	2024-07-01 11:55:00+00:00
08025	U055	2900.0	AC025	2024-07-01 12:00:00+00:00
08026	U056	3000.0	AC026	2024-07-01 12:05:00+00:00
08027	U057	3100.0	AC027	2024-07-01 12:10:00+00:00
08028	U058	3200.0	AC028	2024-07-01 12:15:00+00:00
08029	U059	3300.0	AC029	2024-07-01 12:20:00+00:00
08030	U060	3400.0	AC030	2024-07-01 12:25:00+00:00

30 rows selected (0.207 seconds)

```
0: jdbc:hive2://localhost:10000> SELECT * FROM mobile_transactions;
INFO : Compiling command(queryId=nandhumidhun_20240724150666_bbdfe6fa-8eea-48e7-8bf6-84b83fb327e8): SELECT * FROM mobile_transactions
INFO : No Stats for hadoop@mobile_transactions, Columns: transaction_id, amount, user_id, location, timestamp
INFO : Semantic Analysis Completed (retrial = false)
INFO : Created Hive schema: SchemaFieldSchemas:[FieldSchema(name:mobile_transactions.transaction_id, type:string, comment:null), FieldSchema(name:mobile_transactions.user_id, type:string, comment:null), FieldSchema(name:mobile_transactions.amount, type:double, comment:null), FieldSchema(name:mobile_transactions.location, type:string, comment:null), FieldSchema(name:mobile_transactions.timestamp, type:string, comment:null)], properties:null
INFO : Completed compiling command(queryId=nandhumidhun_20240724150666_bbdfe6fa-8eea-48e7-8bf6-84b83fb327e8); Time taken: 0.168 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=nandhumidhun_20240724150666_bbdfe6fa-8eea-48e7-8bf6-84b83fb327e8): SELECT * FROM mobile_transactions
```

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~

Q E - x

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~ ...

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~ ...

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~ ...

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx: ~ ...

08027	U057	3100.0	AC027	2024-07-01 12:10:00+00:00
08028	U058	3200.0	AC028	2024-07-01 12:15:00+00:00
08029	U059	3300.0	AC029	2024-07-01 12:20:00+00:00
08030	U060	3400.0	AC030	2024-07-01 12:25:00+00:00

+-----+-----+-----+-----+-----+

30 rows selected (0.207 seconds)

0: jdbc:hive2://localhost:10000> SELECT * FROM mobile_transactions;

INFO : Compiling command(queryId=nandhumidhun_20240724150606_bbdfe6fa-8eea-48e7-8bf6-84b83fb327e8): SELECT * FROM mobile_transactions

INFO : No Stats for hadoop@mobile_transactions, Columns: transaction_id, amount, user_id, location, timestamp

INFO : Semantic Analysis Completed (retire = false)

INFO : Created Hive schema: Schema(fieldSchemas:[FieldSchema(name:mobile_transactions.transaction_id, type:string, comment:null), FieldSchema(name:mobile_transactions.user_id, type:string, comment:null), FieldSchema(name:mobile_transactions.amount, type:double, comment:null), FieldSchema(name:mobile_transactions.location, type:string, comment:null), FieldSchema(name:mobile_transactions.timestamp, type:string, comment:null)], properties:null)

INFO : Completed compiling command(queryId=nandhumidhun_20240724150606_bbdfe6fa-8eea-48e7-8bf6-84b83fb327e8); Time taken: 0.168 seconds

INFO : Concurrency mode is disabled, not creating a lock manager

INFO : Executing command(queryId=nandhumidhun_20240724150606_bbdfe6fa-8eea-48e7-8bf6-84b83fb327e8): SELECT * FROM mobile_transactions

INFO : Completed executing command(queryId=nandhumidhun_20240724150606_bbdfe6fa-8eea-48e7-8bf6-84b83fb327e8); Time taken: 0.001 seconds

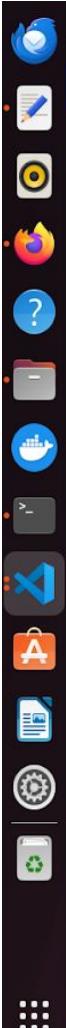
mobile_transactions.transaction_id	mobile_transactions.user_id	mobile_transactions.amount	mobile_transactions.location	mobile_transactions.timestamp
------------------------------------	-----------------------------	----------------------------	------------------------------	-------------------------------

MB001	U061	150.0	D001	2024-07-01 10:00:00+00:00
MB002	U062	250.0	D002	2024-07-01 10:05:00+00:00
MB003	U063	350.0	D003	2024-07-01 10:10:00+00:00
MB004	U064	450.0	D004	2024-07-01 10:15:00+00:00
MB005	U065	550.0	D005	2024-07-01 10:20:00+00:00
MB006	U066	650.0	D006	2024-07-01 10:25:00+00:00
MB007	U067	750.0	D007	2024-07-01 10:30:00+00:00
MB008	U068	850.0	D008	2024-07-01 10:35:00+00:00
MB009	U069	950.0	D009	2024-07-01 10:40:00+00:00
MB010	U070	1050.0	D010	2024-07-01 10:45:00+00:00
MB011	U071	1150.0	D011	2024-07-01 10:50:00+00:00
MB012	U072	1250.0	D012	2024-07-01 10:55:00+00:00
MB013	U073	1350.0	D013	2024-07-01 11:00:00+00:00
MB014	U074	1450.0	D014	2024-07-01 11:05:00+00:00
MB015	U075	1550.0	D015	2024-07-01 11:10:00+00:00
MB016	U076	1650.0	D016	2024-07-01 11:15:00+00:00
MB017	U077	1750.0	D017	2024-07-01 11:20:00+00:00
MB018	U078	1850.0	D018	2024-07-01 11:25:00+00:00
MB019	U079	1950.0	D019	2024-07-01 11:30:00+00:00
MB020	U080	2050.0	D020	2024-07-01 11:35:00+00:00
MB021	U081	2150.0	D021	2024-07-01 11:40:00+00:00
MB022	U082	2250.0	D022	2024-07-01 11:45:00+00:00
MB023	U083	2350.0	D023	2024-07-01 11:50:00+00:00
MB024	U084	2450.0	D024	2024-07-01 11:55:00+00:00
MB025	U085	2550.0	D025	2024-07-01 12:00:00+00:00
MB026	U086	2650.0	D026	2024-07-01 12:05:00+00:00
MB027	U087	2750.0	D027	2024-07-01 12:10:00+00:00
MB028	U088	2850.0	D028	2024-07-01 12:15:00+00:00
MB029	U089	2950.0	D029	2024-07-01 12:20:00+00:00
MB030	U090	3050.0	D030	2024-07-01 12:25:00+00:00

+-----+-----+-----+-----+-----+

30 rows selected (0.204 seconds)

0: jdbc:hive2://localhost:10000>



EXPLORER

- TEST
- ATM
 - kafka_consumer_atm.py
 - kafka_consumer_mobile.py
 - kafka_consumer_online.py
 - kafka_hadoop_atm.py
 - kafka_hadoop_mobile.py
 - kafka_hadoop_online.py
 - kafka_producer_atm.py
 - kafka_producer_mobile.py
 - kafka_producer_online.py
- mobile_banking.csv
- (!) mobile_banking.json
- online_banking.csv
- (!) online_banking.json
- stream.py
- > checkpoint
- > kafka_topic
- > spark_log
- csv_stream.py
- data_processing.py
- file.csv

OUTPUT

- > OUTLINE
- > TIMELINE
- DEBUG CONSOLE

Filter (e.g. text, t(exclude), \escape)

```
line.py    kafka_consumer_mobile.py    kafka_consumer_online.py    kafka_stream.py    ATM_transaction.json    Batch_Processing.py    stream.py    J BatchProcessing.java 9+    D V I ...  
ATM > stream.py > ...  
31 |)  
32 |  
33 |  
34 # Convert the value column from bytes to string and parse JSON  
35 df = df.selectExpr("CAST(value AS STRING)") \\\n
```

PROBLEMS 01 TERMINAL PORTS

24/07/24 17:29:48 WARN AdminClientConfig: The configuration 'auto.offset.reset' was supplied but isn't a known config.

Batch: 10

	transaction_id	user_id	amount	location	timestamp
1	ATM1	U1	100.0	ATM123	2024-07-01T10:00:00Z
2	ATM2	U2	50.0	ATM124	2024-07-01T10:05:00Z
3	ATM3	U3	200.0	ATM125	2024-07-01T10:10:00Z
4	ATM4	U4	80.0	ATM126	2024-07-01T10:15:00Z
5	ATM5	U5	150.0	ATM127	2024-07-01T10:20:00Z
6	ATM6	U6	30.0	ATM128	2024-07-01T10:25:00Z
7	ATM7	U7	60.0	ATM129	2024-07-01T10:30:00Z
8	ATM8	U8	40.0	ATM130	2024-07-01T10:35:00Z
9	ATM9	U9	120.0	ATM131	2024-07-01T10:40:00Z
10	ATM10	U10	70.0	ATM132	2024-07-01T10:45:00Z
11	ATM11	U1	90.0	ATM133	2024-07-01T10:50:00Z
12	ATM12	U2	110.0	ATM134	2024-07-01T10:55:00Z
13	ATM13	U3	140.0	ATM135	2024-07-01T11:00:00Z
14	ATM14	U4	160.0	ATM136	2024-07-01T11:05:00Z
15	ATM15	U5	130.0	ATM137	2024-07-01T11:10:00Z
16	ATM16	U6	95.0	ATM138	2024-07-01T11:15:00Z
17	ATM17	U7	85.0	ATM139	2024-07-01T11:20:00Z
18	ATM18	U8	75.0	ATM140	2024-07-01T11:25:00Z
19	ATM19	U9	65.0	ATM141	2024-07-01T11:30:00Z
20	ATM20	U10	55.0	ATM142	2024-07-01T11:35:00Z

only showing top 20 rows

Batch: 11

	transaction_id	user_id	amount	location	timestamp
1	ATM1	U1	100.0	ATM123	2024-07-01T10:00:00Z

Batch: 12

	transaction_id	user_id	amount	location	timestamp
--	----------------	---------	--------	----------	-----------

File Edit Selection View Go Run Terminal Help



EXPLORER

ATM > Batch_Processing.py

- batch_output
- ATM_transaction.csv
- ATM_transaction.json
- Batch_Processing.py
- convert_json_csv.py
- fra.py
- fraud_detection_model.joblib
- fraud_detection_model.pkl
- fraud.py
- history.csv
- job.py
- kafka_consumer_atm.py
- kafka_consumer_mobile.py
- kafka_consumer_online.py
- kafka_hadoop_atm.py
- kafka_hadoop_mobile.py

OUTPUT

```
ATM > Batch_Processing.py ...
13     StructField("transaction_id", StringType()),
14     StructField("user_id", StringType()),
15     StructField("amount", StringType()),
16     StructField("location", StringType()),
17     StructField("timestamp", StringType())
18 }

19
20 # Read the data from a batch source (e.g., JSON files in local file system)
21 input_path = "file:///home/nandhumidhun/CUBEAI SOLUTIONS/ATM/ATM_transaction.json" # Ensure this path is correct
22 df = spark.read \
23     .format("json") \
24     .schema(schema) \
25     .load(input_path)

26
27 # Optionally perform transformations or actions on the DataFrame
28 df.show(truncate=False)

29
30 # Save the DataFrame to a file or database (optional)
31 output_path = "/home/nandhumidhun/CUBEAI SOLUTIONS/ATM/batch_output"
32 df.write \
33     .mode("overwrite") \
34     .format("parquet") \
35     .save(output_path)
```

PROBLEMS TERMINAL PORTS

bash ATM

bash ATM

To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).

24/07/25 14:33:57 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable

transaction_id	user_id	amount	location	timestamp
NULL	NULL	NULL	NULL	NULL
ATM1	U1	100.0	ATM123	2024-07-01T10:00:00Z
ATM2	U2	50.0	ATM124	2024-07-01T10:05:00Z
ATM3	U3	200.0	ATM125	2024-07-01T10:10:00Z
ATM4	U4	80.0	ATM126	2024-07-01T10:15:00Z
ATM5	U5	150.0	ATM127	2024-07-01T10:20:00Z
ATM6	U6	30.0	ATM128	2024-07-01T10:25:00Z
ATM7	U7	60.0	ATM129	2024-07-01T10:30:00Z
ATM8	U8	40.0	ATM130	2024-07-01T10:35:00Z
ATM9	U9	120.0	ATM131	2024-07-01T10:40:00Z
ATM10	U10	70.0	ATM132	2024-07-01T10:45:00Z
ATM11	U1	90.0	ATM133	2024-07-01T10:50:00Z
ATM12	U2	110.0	ATM134	2024-07-01T10:55:00Z
ATM13	U3	140.0	ATM135	2024-07-01T11:00:00Z
ATM14	U4	160.0	ATM136	2024-07-01T11:05:00Z
ATM15	U5	130.0	ATM137	2024-07-01T11:10:00Z
ATM16	U6	95.0	ATM138	2024-07-01T11:15:00Z
ATM17	U7	85.0	ATM139	2024-07-01T11:20:00Z
ATM18	U8	75.0	ATM140	2024-07-01T11:25:00Z
ATM19	U9	65.0	ATM141	2024-07-01T11:30:00Z

Real-time Stock Market > Real-time Stock Market

OUTPUT

> OUTLINE

> TIMELINE

DEBUG CONSOLE

Filter (e.g. text, |exclude, \escape)

ChatGPT (4) WhatsApp Untitled-1 jinja2.exceptions.Template

https://analytics.zoho.in/workspace/36562400000002019/view/365624000000024002

Premium Trial - 8 days left UPGRADE

Edit Design + 260 Insights

Home Data_Visualizati... fraud_amount fraud Untitled-1

+ Create

Explorer

Dashboards

Reports

Data

Ask Zia

Data Sources

Settings

Trash

Viewer

Clouds

Here is your Smart Chat (Ctrl+Space)

NOT Fraudulent
● nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/CUBEASOLUTIONS/ATM\$ python3 job_amount.py
Fraudulent
Alert: Fraud detected for transaction: {'amount': 3420}
● nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/CUBEASOLUTIONS/ATM\$ python3 job_amount.py
Not Fraudulent
● nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/CUBEASOLUTIONS/ATM\$ python3 user_amount.py
Not Fraudulent: {'amount': 100, 'user_id': 'U1'}
Not Fraudulent: {'amount': 50, 'user_id': 'U2'}
Not Fraudulent: {'amount': 200, 'user_id': 'U3'}
Not Fraudulent: {'amount': 80, 'user_id': 'U4'}
Not Fraudulent: {'amount': 150, 'user_id': 'U5'}
Not Fraudulent: {'amount': 30, 'user_id': 'U6'}
Not Fraudulent: {'amount': 60, 'user_id': 'U7'}
Not Fraudulent: {'amount': 40, 'user_id': 'U8'}
Not Fraudulent: {'amount': 120, 'user_id': 'U9'}
Not Fraudulent: {'amount': 70, 'user_id': 'U10'}
Not Fraudulent: {'amount': 90, 'user_id': 'U11'}
Not Fraudulent: {'amount': 110, 'user_id': 'U2'}
Not Fraudulent: {'amount': 140, 'user_id': 'U3'}
Not Fraudulent: {'amount': 160, 'user_id': 'U4'}

Ln 19, Col 1 Spaces: 4 UTF-8 LF Python 3.10.12 64-bit

user_id	Total amount	Fraud
U0	~100	No
U1	~150	No
U2	~200	No
U3	~400	No
U4	~350	No
U5	~350	No
U6	~150	No
U7	~350	No
U8	~350	No
U9	~250	No
U10	~250	No
U11	~250	No
U12	~350	No
U13	~350	No
U14	~350	No
U15	~350	No
U16	~350	No
U17	~350	No
U18	~350	No
U19	~350	No
U20	~350	No
U21	~350	No
U22	~350	No
U23	~350	No
U24	~350	No
U25	~350	No
U26	~350	No
U27	~350	No
U28	~350	No
U29	~350	No
U30	~350	No
U31	~550	No
U32	~650	No
U33	~750	No
U34	~850	No
U35	~950	No
U36	~1050	No
U37	~1150	No
U38	~1250	No
U39	~1350	No
U40	~1450	No
U41	~1550	No
U42	~1650	No
U43	~1750	No
U44	~1850	No
U45	~1950	No
U46	~2050	No
U47	~2150	No
U48	~2250	No
U49	~2350	No
U50	~2450	No
U51	~2550	No
U52	~2650	No
U53	~2750	No
U54	~2850	No
U55	~2950	No
U56	~3050	No
U57	~3150	No
U58	~3250	No
U59	~3400	Yes

Credit card fraud Detection

STEP 1: Create kafka topic

STEP 2: Verify the data are correctly send to kafka topic

STEP 3: Real-time Streaming Process

STEP 4: Detect the fraud by using the email_id

STEP 5: load the data into kinesis data stream by using lambda



EXPLORER

...

jit-card.json
credit_card1 > kafka_producer.py > ...

TEST

Credit_Card

└ kafka_consumer.py

└ kafka_producer_csv.py

└ kafka_producer1.py

└ KafkaStreamProcessing.class

└ KafkaStreamProcessing.java

└ LeaderElection.class

└ LeaderElection.java

└ path_to_threads.py

└ reassignment.json

└ stream.py

└ Stream.py

└ Transaction.py

credit_card1

└ credit-card.csv

└ credit-card.json

└ credit-card.zip

└ csv_json.py

└ kafka_consumer.py

└ kafka_consumer1.py

└ kafka_producer.py

└ kafka_streaming.py

└ kafka_producer.py

└ KafkaTransaction.class

└ KafkaTransaction.java

└ test.py

> E-Commerce Analysis

> kafka_topic

> TIMELINE

> DEBUG CONSOLE

Filter (e.g. text, lexclude, \escape)



> OUTLINE



> OUTPUT



> JAVA PROJECTS

```

6  # Read the CSV file
7  csv_file_path = 'credit-card.csv'
8  df = pd.read_csv(csv_file_path)
9
10 # Create a Kafka producer
11 producer = KafkaProducer(
12     bootstrap_servers='localhost:9092',
13     value_serializer=lambda v: json.dumps(v).encode('utf-8')
14 )
15
16 # Send each row in the CSV file to the Kafka topic
17 topic_name = 'credit-card-transactions'
18 for _, row in df.iterrows():
19     message = row.to_dict()
20     producer.send(topic_name, message)
21
22 # Ensure all messages are sent before exiting
23 producer.flush()
24
25 print("All messages sent to Kafka topic successfully.")
26

```

PROBLEMS 58 TERMINAL PORTS

bash

bash credit...

- nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/test/credit_card1\$ python3 kafka_producer.py
 All messages sent to Kafka topic successfully.
- nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/test/credit_card1\$ python3 kafka_producer.py
 All messages sent to Kafka topic successfully.
- nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/test/credit_card1\$



EXPLORER



TEST



Credit_Card



kafka_consumer.py



kafka_producer_csv.py



kafka_producer1.py



KafkaStreamProcessing.class



KafkaStreamProcessing.java



LeaderElection.class



LeaderElection.java



path_to_threads.py



reassignment.json



stream.py



Stream.py



Transaction.py



credit_card1



credit-card.csv



credit-card.json



credit-card.zip



csv_json.py



kafka_consumer.py



kafka_consumer1.py



kafka_producer.py



kafka_streaming.py



kafkaproducer.py



J KafkaTransaction.class



J KafkaTransaction.java



test.py



> E-Commerce Analysis

> kafka_topic

> TIMELINE

> DEBUG CONSOLE

Filter (e.g. text, lexclude, \escape)

```

1  from kafka import KafkaConsumer
2
3  def consume_data(topic_name):
4      consumer = KafkaConsumer(
5          topic_name,
6          bootstrap_servers='localhost:9092',
7          auto_offset_reset='earliest', # Start reading at the earliest available message
8          enable_auto_commit=True,
9          group_id='nandhumidhun', # Consumer group ID
10         value_deserializer=lambda x: x.decode('utf-8'))
11
12
13     for message in consumer:
14         print(f"Received message: {message.value}")
15
16     consume_data('credit-card-transactions')
17

```

PROBLEMS 58



bash



python3 cre...



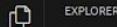
bash credit...



```

Received message: {"transaction_id": 1, "amount": 100.0, "timestamp": "2024-07-31T08:00:00Z", "email": "john.doe@example.com"}
Received message: {"transaction_id": 2, "amount": 5000.0, "timestamp": "2024-07-31T08:05:00Z", "email": "jane.smith@example.com"}
Received message: {"transaction_id": 3, "amount": 20.0, "timestamp": "2024-07-31T08:10:00Z", "email": "alice.l.jones@example.com"}
Received message: {"transaction_id": 4, "amount": 300.0, "timestamp": "2024-07-31T08:15:00Z", "email": "bob.brown@example.com"}
Received message: {"transaction_id": 5, "amount": 150.0, "timestamp": "2024-07-31T08:20:00Z", "email": "carol.white@example.com"}
Received message: {"transaction_id": 6, "amount": 5000.0, "timestamp": "2024-07-31T08:25:00Z", "email": "dan.wilson@example.com"}
Received message: {"transaction_id": 7, "amount": 200.0, "timestamp": "2024-07-31T08:30:00Z", "email": "emily.taylor@example.com"}
Received message: {"transaction_id": 8, "amount": 750.0, "timestamp": "2024-07-31T08:35:00Z", "email": "frank4.martin@example.com"}
Received message: {"transaction_id": 9, "amount": 1200.0, "timestamp": "2024-07-31T08:40:00Z", "email": "grace.moore@example.com"}
Received message: {"transaction_id": 10, "amount": 50.0, "timestamp": "2024-07-31T08:45:00Z", "email": "henry.jackson@example.com"}
Received message: {"transaction_id": 11, "amount": 100.0, "timestamp": "2024-07-31T08:50:00Z", "email": "iris.lee@example.com"}
Received message: {"transaction_id": 12, "amount": 5000.0, "timestamp": "2024-07-31T08:55:00Z", "email": "joan.lee@example.com"}
Received message: {"transaction_id": 13, "amount": 300.0, "timestamp": "2024-07-31T09:00:00Z", "email": "kelly.6.kim@example.com"}
Received message: {"transaction_id": 14, "amount": 600.0, "timestamp": "2024-07-31T09:05:00Z", "email": "lisa.white@example.com"}
Received message: {"transaction_id": 15, "amount": 50.0, "timestamp": "2024-07-31T09:10:00Z", "email": "michael.doe@example.com"}
Received message: {"transaction_id": 16, "amount": 5000.0, "timestamp": "2024-07-31T09:15:00Z", "email": "nina.johnson@example.com"}
Received message: {"transaction_id": 17, "amount": 80.0, "timestamp": "2024-07-31T09:20:00Z", "email": "oliver.thomas@example.com"}
Received message: {"transaction_id": 18, "amount": 400.0, "timestamp": "2024-07-31T09:25:00Z", "email": "peter.harris@example.com"}
Received message: {"transaction_id": 19, "amount": 90.0, "timestamp": "2024-07-31T09:30:00Z", "email": "quincy.roberts@example.com"}
Received message: {"transaction_id": 20, "amount": 5000.0, "timestamp": "2024-07-31T09:35:00Z", "email": "rachel.clark@example.com"}
Received message: {"transaction_id": 21, "amount": 60.0, "timestamp": "2024-07-31T09:40:00Z", "email": "samuel.lee@example.com"}
Received message: {"transaction_id": 22, "amount": 500.0, "timestamp": "2024-07-31T09:45:00Z", "email": "tina.morris@example.com"}
Received message: {"transaction_id": 23, "amount": 1000.0, "timestamp": "2024-07-31T09:50:00Z", "email": "uma.james@example.com"}
Received message: {"transaction_id": 24, "amount": 200.0, "timestamp": "2024-07-31T09:55:00Z", "email": "victor8.carter@example.com"}
Received message: {"transaction_id": 25, "amount": 50.0, "timestamp": "2024-07-31T10:00:00Z", "email": "wendy.rogers@example.com"}
Received message: {"transaction_id": 26, "amount": 5000.0, "timestamp": "2024-07-31T10:05:00Z", "email": "xander.baker@example.com"}
Received message: {"transaction_id": 27, "amount": 300.0, "timestamp": "2024-07-31T10:10:00Z", "email": "yvonne.perez@example.com"}
Received message: {"transaction_id": 28, "amount": 120.0, "timestamp": "2024-07-31T10:15:00Z", "email": "zachary.adams@example.com"}
Received message: {"transaction_id": 29, "amount": 700.0, "timestamp": "2024-07-31T10:20:00Z", "email": "anna.fisher@example.com"}

```



EXPLORER

TEST

Credit_Card

kafka_consumer.py

kafka_producer_csv.py

kafka_producer1.py

KafkaStreamProcessing.class

KafkaStreamProcessing.java

LeaderElection.class

LeaderElection.java

path_to_threads.py

reassignment.json

stream.py

Stream.py

Transaction.py

credit_card1

credit-card.csv

credit-card.json

credit-card.zip

csv_json.py

kafka_consumer.py

kafka_consumer1.py

kafka_producer.py

kafka_streaming.py

kafkaproducer.py

KafkaTransaction.class

KafkaTransaction.java

test.py

> E-Commerce Analysis

> kafka_topic

> TIMELINE

> DEBUG CONSOLE

Filter (e.g. text, \exclude, \escape)

mer1.py credit-card.csv kafka_producer.py kafka_consumer.py credit_card1 test.py J KafkaTransaction.java 9+ J KafkaTransaction.class kafka_streaming.py X ▶ ⓘ ...

```

30
31 # Fill null values
32 transaction_df = transaction_df.withColumn("amount", coalesce(col("amount"), lit(0.0)))
33
34 # Define fraud detection logic
35 # Example logic: Consider transactions with amount greater than a threshold as fraudulent
36 fraudulent_transactions = transaction_df.filter(col("amount") > 50)
37
38 # Define the query to print the fraudulent transactions to the console
39 query = fraudulent_transactions.writeStream \
40   .outputMode("append") \
41   .format("console") \
42   .option("failOnDataLoss", "false") \
43   .start()

```

PROBLEMS 58 TERMINAL PORTS

```

24/08/02 11:24:04 WARN AdminClientConfig: The configuration 'max.poll.records' was supplied but isn't a known config.
24/08/02 11:24:04 WARN AdminClientConfig: The configuration 'auto.offset.reset' was supplied but isn't a known config.

```

Batch: 0

transaction_id	amount	timestamp	email
1	100.0	2024-07-31 13:30:00	john.doe@example.com
2	5000.0	2024-07-31 13:35:00	jane.smith@example.com
4	300.0	2024-07-31 13:45:00	bob.brown@example.com
5	150.0	2024-07-31 13:50:00	carol.white@example.com
6	5000.0	2024-07-31 13:55:00	dan.wilson@example.com
7	200.0	2024-07-31 14:00:00	emily.taylor@example.com
8	750.0	2024-07-31 14:05:00	frank4.martin@example.com
9	1200.0	2024-07-31 14:10:00	grace.moore@example.com
11	100.0	2024-07-31 14:20:00	iris.lee@example.com
12	5000.0	2024-07-31 14:25:00	joan.lee@example.com
13	300.0	2024-07-31 14:30:00	kelly6.kim@example.com
14	600.0	2024-07-31 14:35:00	lisa.white@example.com
16	5000.0	2024-07-31 14:45:00	nina.johnson@example.com
17	80.0	2024-07-31 14:50:00	oliver.thomas@example.com
18	400.0	2024-07-31 14:55:00	peter.harris@example.com
19	90.0	2024-07-31 15:00:00	quincy.roberts@example.com
20	5000.0	2024-07-31 15:05:00	rachel.clark@example.com
21	60.0	2024-07-31 15:10:00	samuel.lee@example.com
22	500.0	2024-07-31 15:15:00	tina.morris@example.com
23	1000.0	2024-07-31 15:20:00	uma.james@example.com

Batch: 1

transaction_id	amount	timestamp	email
1	100.0	2024-07-31 13:30:00	john.doe@example.com
2	5000.0	2024-07-31 13:35:00	jane.smith@example.com
4	300.0	2024-07-31 13:45:00	bob.brown@example.com
5	150.0	2024-07-31 13:50:00	carol.white@example.com
6	5000.0	2024-07-31 13:55:00	dan.wilson@example.com
7	200.0	2024-07-31 14:00:00	emily.taylor@example.com
8	750.0	2024-07-31 14:05:00	frank4.martin@example.com
9	1200.0	2024-07-31 14:10:00	grace.moore@example.com
11	100.0	2024-07-31 14:20:00	iris.lee@example.com
12	5000.0	2024-07-31 14:25:00	joan.lee@example.com
13	300.0	2024-07-31 14:30:00	kelly6.kim@example.com
14	600.0	2024-07-31 14:35:00	lisa.white@example.com
16	5000.0	2024-07-31 14:45:00	nina.johnson@example.com
17	80.0	2024-07-31 14:50:00	oliver.thomas@example.com
18	400.0	2024-07-31 14:55:00	peter.harris@example.com
19	90.0	2024-07-31 15:00:00	quincy.roberts@example.com
20	5000.0	2024-07-31 15:05:00	rachel.clark@example.com
21	60.0	2024-07-31 15:10:00	samuel.lee@example.com
22	500.0	2024-07-31 15:15:00	tina.morris@example.com
23	1000.0	2024-07-31 15:20:00	uma.james@example.com

EXPLORER

TEST

- Credit_Card
 - kafka_consumer.py
 - kafka_producer_csv.py
 - kafka_producer1.py
 - KafkaStreamProcessing.class
 - KafkaStreamProcessing.java
 - LeaderElection.class
 - LeaderElection.java
 - path_to_threads.py
 - reassignment.json
 - stream.py
 - Stream.py
 - Transaction.py
- credit_card1
 - credit-card.csv
 - credit-card.json
 - credit-card.zip
 - csv_json.py
 - kafka_consumer.py
 - kafka_consumer1.py
 - kafka_producer.py
 - kafka_streaming.py
 - kafkaproducer.py
 - KafkaTransaction.class
 - KafkaTransaction.java
 - test.py
- E-Commerce Analysis
- kafka_topic

TIMELINE

DEBUG CONSOLE

Filter (e.g. text, !exclude, \escape)

mer1.py credit-card.csv kafka_producer.py kafka_consumer.py credit_card1 test.py KafkaTransaction.java 9+ KafkaTransaction.class kafka_streaming.py

```

11  public class KafkaTransaction {
12      public static void main(String[] args) {
13          // Consume and check messages
14          while (true) {
15              consumer.poll(Duration.ofMillis(100)).forEach(record -> {
16                  try {
17                      JsonNode jsonNode = mapper.readTree(record.value());
18                      String email = jsonNode.path("email").asText(null);
19                      int transactionId = jsonNode.path("transaction_id").asInt(-1);
20
21                      // Check if the email matches the pattern and if transaction_id meets criteria
22                      if (email != null && emailPattern.matcher(email).matches()) {
23                          System.out.println("Fraudulent Transaction: Email: " + email + ", Transaction ID: " + transactionId);
24                      } else {
25                          System.out.println("Normal Transaction: Email: " + email + ", Transaction ID: " + transactionId);
26                      }
27                  } catch (IOException e) {
28                      e.printStackTrace();
29                  }
30              }
31          }
32      }
33  }
34
35
36
37
38
39
40
41
42

```

PROBLEMS 58

TUTORIAL PORTS

bash

java credit_c...

bash credit_...

bash credit_...

```

[main] INFO org.apache.kafka.clients.consumer.internals.SubscriptionState - [Consumer clientId=consumer-transaction-consumer-group-1, groupId=transaction-consumer-group] Found no committed offset for partition credit-card-transactions-0
[main] INFO org.apache.kafka.clients.consumer.internals.SubscriptionState - [Consumer clientId=consumer-transaction-consumer-group-1, groupId=transaction-consumer-group] Resetting offset for partition credit-card-transactions-0 to position FetchPosition{offset=0, offsetEpoch=Optional.empty, currentLeader=LeaderAndEpoch[leader=Optional.empty, epoch=0]}.
Normal Transaction: Email: john.doe@example.com, Transaction ID: 1
Normal Transaction: Email: jane.smith@example.com, Transaction ID: 2
Fraudulent Transaction: Email: alice.l.jones@example.com, Transaction ID: 3
Normal Transaction: Email: bob.brown@example.com, Transaction ID: 4
Normal Transaction: Email: carol.white@example.com, Transaction ID: 5
Normal Transaction: Email: dan.wilson@example.com, Transaction ID: 6
Normal Transaction: Email: emily.taylor@example.com, Transaction ID: 7
Fraudulent Transaction: Email: frank4.martin@example.com, Transaction ID: 8
Normal Transaction: Email: grace.moore@example.com, Transaction ID: 9
Normal Transaction: Email: henry.jackson@example.com, Transaction ID: 10
Normal Transaction: Email: iris.lee@example.com, Transaction ID: 11
Normal Transaction: Email: joan.lee@example.com, Transaction ID: 12
Fraudulent Transaction: Email: kelly6.kim@example.com, Transaction ID: 13
Normal Transaction: Email: lisa.white@example.com, Transaction ID: 14
Normal Transaction: Email: michael.doe@example.com, Transaction ID: 15
Normal Transaction: Email: nina.johnson@example.com, Transaction ID: 16
Normal Transaction: Email: oliver.thomas@example.com, Transaction ID: 17
Normal Transaction: Email: peter.harris@example.com, Transaction ID: 18
Normal Transaction: Email: quincy.roberts@example.com, Transaction ID: 19
Normal Transaction: Email: rachel.clark@example.com, Transaction ID: 20
Normal Transaction: Email: samuel.lee@example.com, Transaction ID: 21
Normal Transaction: Email: tina.morris@example.com, Transaction ID: 22
Normal Transaction: Email: uma.james@example.com, Transaction ID: 23
Fraudulent Transaction: Email: victor8.carter@example.com, Transaction ID: 24
Normal Transaction: Email: wendy.rogers@example.com, Transaction ID: 25
Normal Transaction: Email: xander.baker@example.com, Transaction ID: 26
Normal Transaction: Email: yvonne.perez@example.com, Transaction ID: 27
Normal Transaction: Email: zachary.adams@example.com, Transaction ID: 28
Normal Transaction: Email: anna.fisher@example.com, Transaction ID: 29
Normal Transaction: Email: benjamin.green@example.com, Transaction ID: 30
Normal Transaction: Email: claire.nelson@example.com, Transaction ID: 31
Normal Transaction: Email: david.morris@example.com, Transaction ID: 32

```

LN 56, Col S7 Spaces: 4 UTF-8 LF () Java

Kinesis | us-west-2 DOCUMENTATION FOR X Untitled document - Go ChatGPT

https://us-west-2.console.aws.amazon.com/kinesis/home?region=us-west-2#/streams/create

Services Search [Alt+S]

Create data stream Info

Data stream configuration

Data stream name Acceptable characters are uppercase and lowercase letters, numbers, underscores, hyphens and periods.

Data stream capacity Info

Capacity mode

On-demand
Use this mode when your data stream's throughput requirements are unpredictable and variable. With on-demand mode, your data stream's capacity scales automatically.

Provisioned
Use provisioned mode when you can reliably estimate throughput requirements of your data stream. With provisioned mode, your data stream's capacity is fixed.

Total data stream capacity
By default, data streams with on-demand mode scale throughput automatically to accommodate traffic of up to 200 MiB per second and 200,000 records per second for the write capacity. If traffic exceeds capacity, your data stream will throttle. To request capacity increase up to 2GB per second write and 4GB per second read, submit a support ticket [↗](#)

Write capacity
Maximum
200 MiB/second and 200,000 records/second

Read capacity
Maximum (per consumer)
400 MiB/second
Up to 2 default consumers. Use Enhanced Fan-Out (EFO) for more consumers. EFO supports adding up to 20 consumers, each having a dedicated throughput.

ⓘ On-demand mode has a pay-per-throughput pricing model. See [Kinesis pricing for on-demand mode](#) ↗

Data stream settings
You can edit the settings after the data stream has been created and is in the active status.

CloudShell Feedback © 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

EXPLORER

✓ AWS

- > __pycache__
- aws_lambda.py
- check_dynamo.py
- check_lambda.py
- check.py
- dynamo.py
- dynamo1.py
- kafka_kinesis.py
- kafka_to_s3.py
- kinesis.py**
- knesis.py
- lambda.py
- load_file.py
- { response.json }

```

kinesis.py > ...
1 import boto3
2 import csv
3 import json
4 from time import sleep
5
6 # Define Kinesis stream details
7 my_stream_name = 'kinesisstream'
8
9 # Initialize Kinesis client without specifying credentials directly
10 kinesis_client = boto3.client('kinesis', region_name='us-west-2')
11
12 # Path to the CSV file
13 csv_file_path = '/home/nandhumidhun/test/credit_card1/credit-card.csv'
14
15 # Read data from CSV file and send to Kinesis
16 with open(csv_file_path, mode='r') as csv_file:
17     csv_reader = csv.DictReader(csv_file)
18     for row in csv_reader:
19         json_data = {
20             "transaction_id": row["transaction_id"],
21             "amount": row["amount"],
22             "timestamp": row["timestamp"],
23             "email": row["email"]
24         }
25         print(json_data) # Print the JSON data
26         kinesis_client.put_record(

```

> TIMELINE

PROBLEMS

TERMINAL PORTS

bash + ×

✓ DEBUG CONSOLE

Filter (e.g. text, lexclude, \escape)

```

● nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/AWS$ python3 kinesis.py
[{'transaction_id': '1', 'amount': '100.00', 'timestamp': '2024-07-31T08:00:00Z', 'email': 'john.doe@example.com'},
 {'transaction_id': '2', 'amount': '5000.00', 'timestamp': '2024-07-31T08:05:00Z', 'email': 'jane.smith@example.com'},
 {'transaction_id': '3', 'amount': '20.00', 'timestamp': '2024-07-31T08:10:00Z', 'email': 'alice1.jones@example.com'},
 {'transaction_id': '4', 'amount': '300.00', 'timestamp': '2024-07-31T08:15:00Z', 'email': 'bob.brown@example.com'},
 {'transaction_id': '5', 'amount': '150.00', 'timestamp': '2024-07-31T08:20:00Z', 'email': 'carol.white@example.com'},
 {'transaction_id': '6', 'amount': '5000.00', 'timestamp': '2024-07-31T08:25:00Z', 'email': 'dan.wilson@example.com'},
 {'transaction_id': '7', 'amount': '200.00', 'timestamp': '2024-07-31T08:30:00Z', 'email': 'emily.taylor@example.com'},
 {'transaction_id': '8', 'amount': '750.00', 'timestamp': '2024-07-31T08:35:00Z', 'email': 'frank4.martin@example.com'},
 {'transaction_id': '9', 'amount': '1200.00', 'timestamp': '2024-07-31T08:40:00Z', 'email': 'grace.moore@example.com'},
 {'transaction_id': '10', 'amount': '50.00', 'timestamp': '2024-07-31T08:45:00Z', 'email': 'henry.jackson@example.com'},
 {'transaction_id': '11', 'amount': '100.00', 'timestamp': '2024-07-31T08:50:00Z', 'email': 'iris.lee@example.com'},
 {'transaction_id': '12', 'amount': '5000.00', 'timestamp': '2024-07-31T08:55:00Z', 'email': 'joan.lee@example.com'},
 {'transaction_id': '13', 'amount': '300.00', 'timestamp': '2024-07-31T09:00:00Z', 'email': 'kelly6.kim@example.com'},
 {'transaction_id': '14', 'amount': '600.00', 'timestamp': '2024-07-31T09:05:00Z', 'email': 'lisa.white@example.com'},
 {'transaction_id': '15', 'amount': '50.00', 'timestamp': '2024-07-31T09:10:00Z', 'email': 'michael.doe@example.com'},
 {'transaction_id': '16', 'amount': '5000.00', 'timestamp': '2024-07-31T09:15:00Z', 'email': 'nina.johnson@example.com'},
 {'transaction_id': '17', 'amount': '80.00', 'timestamp': '2024-07-31T09:20:00Z', 'email': 'oliver.thomas@example.com'},
 {'transaction_id': '18', 'amount': '400.00', 'timestamp': '2024-07-31T09:25:00Z', 'email': 'peter.harris@example.com'},
 {'transaction_id': '19', 'amount': '90.00', 'timestamp': '2024-07-31T09:30:00Z', 'email': 'quincy.roberts@example.com'},
 {'transaction_id': '20', 'amount': '5000.00', 'timestamp': '2024-07-31T09:35:00Z', 'email': 'rachel.clark@example.com'},
 {'transaction_id': '21', 'amount': '60.00', 'timestamp': '2024-07-31T09:40:00Z', 'email': 'samuel.lee@example.com'}]
```



Kinesis | us-west-2 DOCUMENTATION FOR Untitled document - Google ChatGPT

https://us-west-2.console.aws.amazon.com/kinesis/home?region=us-west-2#/streams/details/kinesisstream/details

Amazon Kinesis Services Search [Alt+S]

Amazon Kinesis > Data streams > kinesisstream

kinesisstream [Info](#)

[Delete](#)

Data stream summary

Status Active	Capacity mode On-demand	ARN arn:aws:kinesis:us-west-2:637423373422:stream:kinesisstream	Creation time August 02, 2024 at 11:32 GMT+5:30
	Data retention period 1 day		

[Applications](#) [Monitoring](#) [Configuration](#) [Enhanced fan-out \(0\)](#) [Data viewer](#) [Data analytics - new](#) [Data stream sharing](#) [EventBridge Pipes](#)

Producers [Info](#)

Producers put records into Kinesis Data Streams.

Amazon Kinesis Agent Use a stand-alone Java software application to send data to the stream. Learn more	AWS SDK Use AWS SDK for Java to develop producers. Learn more	Amazon Kinesis Producer Library (KPL) Use KPL to develop producers. Learn more
--	--	---

[View in GitHub](#) [View in GitHub](#) [View in GitHub](#)

Consumers [Info](#)

Consumers get records from Kinesis Data Streams and process them.

Managed Apache Flink New Use an Amazon Managed Service for Apache Flink application to process and analyze using SQL or Java. Process data in real time	Amazon Data Firehose Use a Firehose stream to process and store records in a destination. Process with Firehose stream	Amazon Kinesis Client Library (KCL) Use Kinesis Client Library to develop consumers. Learn more View in GitHub
---	--	--

CloudShell Feedback © 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Screenshot of the AWS Kinesis Monitoring Dashboard.

The dashboard displays Stream metrics for a Kinesis stream. The time range is set to 1 day. The monitoring tab is selected.

Stream metrics:

- GetRecords - sum (MB/s):** A line chart showing GetRecords - sum (MB/s) over time. The Y-axis ranges from 0 to 400 MB/s. The X-axis shows 05:55, 06:00, and 06:05. A red line represents the data, and a blue line represents the Maximum GetRecords Limit (400 MB/s).
- GetRecords iterator age - maximum (Milliseconds):** A line chart showing GetRecords iterator age - maximum (Milliseconds) over time. The Y-axis ranges from 0 to 1 Millisecond. The X-axis shows 05:50, 05:55, 06:00, and 06:05. A blue line represents the data.
- GetRecords latency - average (Milliseconds):** A line chart showing GetRecords latency - average (Milliseconds) over time. The Y-axis ranges from 0 to 9 Milliseconds. The X-axis shows 05:50, 05:55, 06:00, and 06:05. A blue line represents the data.
- GetRecords - sum (Count):** A line chart showing GetRecords - sum (Count) over time. The Y-axis ranges from 0 to 50. The X-axis shows 05:50, 05:55, 06:00, and 06:05. A blue line represents the data.
- GetRecords success - average (Ratio):** A line chart showing GetRecords success - average (Ratio) over time. The Y-axis ranges from 0 to 2. The X-axis shows 05:50, 05:55, 06:00, and 06:05. A blue line represents the data.
- Incoming data - sum (MB/s):** A line chart showing Incoming data - sum (MB/s) over time. The Y-axis ranges from 0 to 200 MB/s. The X-axis shows 05:50, 05:55, 06:00, and 06:05. A single red dot is plotted at approximately 200 MB/s.

Navigation and Configuration:

- Time range: 1 day
- Metrics: Applications, Monitoring, Configuration, Enhanced fan-out (0), Data viewer, Data analytics - new, Data stream sharing, EventBridge Pipes
- Custom Metrics: Add to dashboard
- CloudWatch Metrics: 1h, 3h, 12h, 1d, 3d, 1w, Custom (15m), UTC timezone

Left Sidebar:

- Amazon Kinesis
- Dashboard
- Data streams
 - Amazon Data Firehose [New]
 - Managed Apache Flink [New]
- Resources
 - CloudFormation templates
 - AWS Glue Schema Registry [New]

Bottom Navigation:

- CloudShell
- Feedback
- © 2024, Amazon Web Services, Inc. or its affiliates.
- Privacy
- Terms
- Cookie preferences

The trigger kinesisstream was successfully added to function lambdakinesis. The trigger is in a disabled state.

Function overview

Diagram Template

lambdakinesis

Kinesis

+ Add destination

+ Add trigger

Description

Last modified 24 seconds ago

Function ARN arn:aws:lambda:us-west-2:637423373422:function:lambdakinesis

Function URL [Info](#)

Code Test Monitor Configuration Aliases Versions

Triggers (1) [Info](#)

Find triggers

Trigger

Kinesis: kinesisstream
arn:aws:kinesis:us-west-2:637423373422:stream/kinesisstream
state: Enabled

Details

Learn how to implement common use cases in AWS Lambda.

Create a simple web app

In this tutorial you will learn how to:

- Build a simple web app, consisting of a Lambda function with a function URL that outputs a webpage
- Invoke your function through its function URL

Learn more [\[\]](#)

Start tutorial

https://us-west-2.console.aws.amazon.com/lambda/home?region=us-west-2#/functions/lambdakinesis?tab=versions

© 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

ETL PROCESS IN AWS

STEP 1: Create an S3 Bucket

STEP 2: Set Up AWS Glue

STEP 3: Amazon Redshift

STEP 4: Data Loading

STEP 5: Schedule the ETL Job

STEP 6: Testing and Validation

Create S3 bucket | S3 | us X Untitled document - Google Sheets (1) WhatsApp ChatGPT

https://us-west-2.console.aws.amazon.com/s3/bucket/create?region=us-west-2&bucketType=general

Services Search [Alt+S]

Amazon S3 > Buckets > Create bucket

Create bucket Info

Buckets are containers for data stored in S3.

General configuration

AWS Region
US West (Oregon) us-west-2

Bucket type Info

General purpose
Recommended for most use cases and access patterns. General purpose buckets are the original S3 bucket type. They allow a mix of storage classes that redundantly store objects across multiple Availability Zones.

Directory - New
Recommended for low-latency use cases. These buckets use only the S3 Express One Zone storage class, which provides faster processing of data within a single Availability Zone.

Bucket name Info

Bucket name must be unique within the global namespace and follow the bucket naming rules. [See rules for bucket naming](#)

Copy settings from existing bucket - optional
Only the bucket settings in the following configuration are copied.

Format: s3://bucket/prefix

Object Ownership Info

Control ownership of objects written to this bucket from other AWS accounts and the use of access control lists (ACLs). Object ownership determines who can specify access to objects.

ACLs disabled (recommended)
All objects in this bucket are owned by this account. Access to this bucket and its objects is specified using only policies.

ACLs enabled
Objects in this bucket can be owned by other AWS accounts. Access to this bucket and its objects can be specified using ACLs.

Object Ownership

CloudShell Feedback © 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Visual - Editor - AWS Glue × Upload objects - S3 bucki × Untitled document - Googl × (1) WhatsApp × ChatGPT × +

https://us-west-2.console.aws.amazon.com/s3/upload/localdatainput?region=us-west-2&bucketType=general

AWS Services Search [Alt+S]

Amazon S3 > Buckets > localdatainput > Upload

Upload Info

Add the files and folders you want to upload to S3. To upload a file larger than 160GB, use the AWS CLI, AWS SDK or Amazon S3 REST API. [Learn more](#)

Drag and drop files and folders you want to upload here, or choose Add files or Add folder.

Files and folders (1 Total, 1.9 KB)

All files and folders in this table will be uploaded.

<input type="checkbox"/> Name	Folder
<input type="checkbox"/> Sales_data.csv	-

< 1 >

Destination Info

Destination
s3://localdatainput

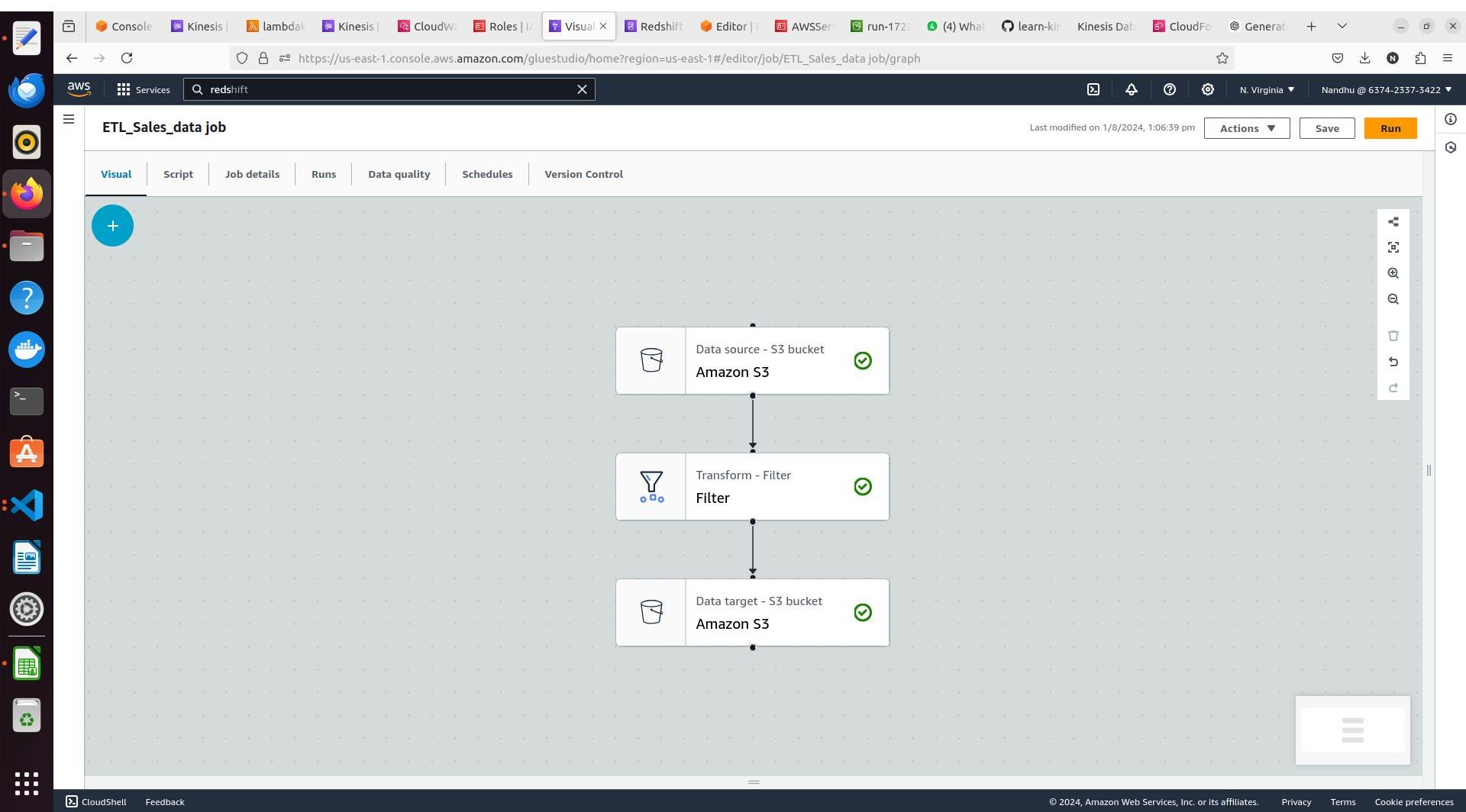
▶ Destination details
Bucket settings that impact new objects stored in the specified destination.

▶ Permissions
Grant public access and access to other AWS accounts.

▶ Properties
Specify storage class, encryption settings, tags, and more.

Cancel

CloudShell Feedback © 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences



Console Kinesis lambda Kinesis CloudWatch Roles Runs Redshift Editor AWS Services run-172 (4) What learn-kir Kinesis Data CloudFormation Generat... + N. Virginia Nandhu @ 6374-2337-3422

https://us-east-1.console.aws.amazon.com/sqlworkbench/home?region=us-east-1#/client

aws Services Search [Alt+S]

Redshift query editor v2

Create Load data Filter resources Serverless: default-workgroup... 1:1

Editor Queries Notebooks Charts History Scheduled queries

awsdatacatalog dev sample_data_dev

Run Limit 100 Explain Isolated session Serverless: de... dev

1 's:iam::637423373422:role/service-role/AmazonRedshift-CommandsAccessRole-20240729T163148' FORMAT AS CSV DELIMITER ',' QUOTE '\"' IGNOREHEADER 1 REGION AS 'us-east-1';
2

Row 1, Col 265, Chr 300

Result 1 Result 2 (20)

	product_id	date	product_name	customer_id	amount	id
□	104	2024-01-04	Widget D	1004	60	4
□	104	2024-01-09	Widget D	1009	60	9
□	104	2024-01-14	Widget D	1014	60	14
□	104	2024-01-19	Widget D	1019	60	19
□	104	2024-01-24	Widget D	1024	60	24
□	104	2024-01-29	Widget D	1029	60	29
□	104	2024-02-03	Widget D	1034	60	34
□	104	2024-02-08	Widget D	1039	60	39
□	104	2024-02-13	Widget D	1044	60	44
□	104	2024-02-18	Widget D	1049	60	49
□	104	2024-01-04	Widget D	1004	60	4
□	104	2024-01-09	Widget D	1009	60	9
□	104	2024-01-14	Widget D	1014	60	14

Query ID 3366 Elapsed time: 2847 ms Total rows: 20

CloudShell Feedback © 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Console Kinesis lambda Kinesis CloudWatch Roles | Runs - Redshift Editor AWS S3 bucket (5) What learn-kir Kinesis Data CloudFront General + N. Virginia Nandhu @ 6374-2337-3422

AWS Glue ETL_Sales_data job

Last modified on 1/8/2024, 1:06:39 pm Actions Save Run

Getting started ETL jobs Visual ETL Notebooks Job run monitoring Data Catalog tables Data connections Workflows (orchestration)

Data Catalog Databases Tables Stream schema registries Schemas Connections Crawlers Classifiers Catalog settings

Data Integration and ETL Legacy pages

What's New Documentation AWS Marketplace

Enable compact mode Enable new navigation

Job runs (1/5) Info Last updated (UTC) August 1, 2024 at 09:25:35

Filter job runs by property

Run status	Retries	Start time (Local)	End time (Local)	Duration	Capacity (DPU)	Worker type
Running	0	08/01/2024 14:51:53	-	1 m 34 s	10 DPU	G.1X
Succeeded	0	08/01/2024 14:31:21	08/01/2024 14:34:12	2 m 35 s	10 DPU	G.1X
Succeeded	0	08/01/2024 13:31:21	08/01/2024 13:34:55	3 m 15 s	10 DPU	G.1X

Run details Input arguments (10) Continuous logs Run insights Metrics Spark UI Stop job run

Continuous logs Info

Driver logs Driver and executor log streams

```
24/08/01 09:22:43 INFO SecurityManager: Changing view acls groups to:
24/08/01 09:22:43 INFO SecurityManager: Changing modify acls to: spark
24/08/01 09:22:43 INFO SecurityManager: Changing view acls to: spark
24/08/01 09:22:43 INFO ResourceProfileManager: Added ResourceProfile id: 0
24/08/01 09:22:43 INFO ResourceProfile: Limiting resource is cpus at 4 tasks per executor
24/08/01 09:22:43 INFO ResourceProfile: Default ResourceProfile created, executor resources: Map(cores -> name: cores, amount: 4, script: , venc
24/08/01 09:22:43 INFO SparkContext: Submitted application: nativespark-ETL_Sales_data_job-jr_d0b39efc0428a94787d5dad6cadde721767fb393af77fe2
24/08/01 09:22:43 INFO ResourceUtils: -----
24/08/01 09:22:43 INFO ResourceUtils: No custom resources configured for spark.driver.
24/08/01 09:22:43 INFO SparkContext: Running Spark version 3.3.0-amzn-1
24/08/01 09:22:35 INFO SafeLogging: Initializing logging subsystem
24/08/01 09:22:33 INFO PlatformInfo: Unable to read clusterId from /var/lib/info/job-flow.json, out of places to look
24/08/01 09:22:33 INFO PlatformInfo: Unable to read clusterId from /var/lib/instance-controller/extrainstanceData.json, trying EMR job-flow dat
24/08/01 09:22:33 INFO PlatformInfo: Unable to read clusterId from http://localhost:8321/configuration, trying extra instance data file: /var/li
24/08/01 09:22:32 INFO LogPusher: legacyLogging: true - logs will be written with spark application ID
24/08/01 09:22:32 INFO LogPusher: standardLogging: true - logs will be written with job run ID or session ID
24/08/01 09:22:32 INFO SparkUILogFileCleaner: SparkUILogFileCleanerThread started
```

This information isn't currently available.

REAL-TIME ANOMALY DETECTION IN NETWORK TRAFFIC

STEP 1: Data Ingestion Using kafka

STEP 2: Data Processing Spark Structured Streaming

STEP 3: Train an Anomaly Detection Model

STEP 4: Real-Time Anomaly Scoring

STEP 5: Detect the Anomaly

STEP 6: Visualize Data



EXPLORER

- NEWWTORK_TRAFFIC
 - __pycache__
 - anomalies_log.csv
 - anomalis_ext.py
 - anomalis.py
 - create_csv.py
 - isolation_forest_model.pkl
 - isolation_forest_model1.pkl
 - kafka_consumer.py
 - kafka_producer.py
 - network.csv
 - newtype.py
 - scaler.pkl
 - stream_kafka.py
 - trainset.py
 - visualization.csv
 - visualization.py

```

stream_kafka.py >...
54     query = df.writeStream \
55         .outputMode("update") \
56         .format("console") \
57         .option("checkpointLocation", "/home/nandhumidhun/test/spark_log") \
58         .start()
59
60     # Wait for the termination of the queries
61     try:
62         converted_data_query.awaitTermination()
63         query.awaitTermination()
64     except KeyboardInterrupt:
65         print("Streaming query terminated.")
66     finally:
67         # Stop the Spark session
68         spark.stop()
69

```

TIMELINE

DEBUG CONSOLE

Filter (e.g. text, !exclude, \escape)

OUTLINE

OUTPUT

PROBLEMS 50 TERMINAL PORTS

pyth... python3 bash

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/CUBEAI SOLUTIONS/Newtwork_traffic\$ python3 kafka_producer.py

All messages have been sent to Kafka topic: network-traffic

nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/CUBEAI SOLUTIONS/Newtwork_traffic\$ python3 kafka_consumer.py

Received message: {"source_ip": "192.168.0.1", "dest_ip": "10.0.0.1", "protocol": "TCP", "packet_size": 512, "timestamp": 1690797580}

Received message: {"source_ip": "192.168.0.2", "dest_ip": "10.0.0.2", "protocol": "UDP", "packet_size": 1024, "timestamp": 1690797581}

Received message: {"source_ip": "192.168.0.3", "dest_ip": "10.0.0.3", "protocol": "TCP", "packet_size": 256, "timestamp": 1690797582}

Received message: {"source_ip": "192.168.0.4", "dest_ip": "10.0.0.4", "protocol": "ICMP", "packet_size": 64, "timestamp": 1690797583}

Received message: {"source_ip": "192.168.0.5", "dest_ip": "10.0.0.5", "protocol": "TCP", "packet_size": 1280, "timestamp": 1690797584}

Received message: {"source_ip": "192.168.0.6", "dest_ip": "10.0.0.6", "protocol": "UDP", "packet_size": 512, "timestamp": 1690797585}

Received message: {"source_ip": "192.168.0.7", "dest_ip": "10.0.0.7", "protocol": "TCP", "packet_size": 256, "timestamp": 1690797586}

Received message: {"source_ip": "192.168.0.8", "dest_ip": "10.0.0.8", "protocol": "ICMP", "packet_size": 128, "timestamp": 1690797587}

Received message: {"source_ip": "192.168.0.9", "dest_ip": "10.0.0.9", "protocol": "TCP", "packet_size": 768, "timestamp": 1690797588}

Received message: {"source_ip": "192.168.0.10", "dest_ip": "10.0.0.10", "protocol": "UDP", "packet_size": 2048, "timestamp": 1690797589}

Received message: {"source_ip": "192.168.0.11", "dest_ip": "10.0.0.11", "protocol": "TCP", "packet_size": 1024, "timestamp": 1690797590}

Received message: {"source_ip": "192.168.0.12", "dest_ip": "10.0.0.12", "protocol": "ICMP", "packet_size": 256, "timestamp": 1690797591}

Received message: {"source_ip": "192.168.0.13", "dest_ip": "10.0.0.13", "protocol": "TCP", "packet_size": 512, "timestamp": 1690797592}

Received message: {"source_ip": "192.168.0.14", "dest_ip": "10.0.0.14", "protocol": "UDP", "packet_size": 1280, "timestamp": 1690797593}

Received message: {"source_ip": "192.168.0.15", "dest_ip": "10.0.0.15", "protocol": "TCP", "packet_size": 256, "timestamp": 1690797594}

Received message: {"source_ip": "192.168.0.16", "dest_ip": "10.0.0.16", "protocol": "ICMP", "packet_size": 64, "timestamp": 1690797595}

Received message: {"source_ip": "192.168.0.17", "dest_ip": "10.0.0.17", "protocol": "TCP", "packet_size": 768, "timestamp": 1690797596}

Received message: {"source_ip": "192.168.0.18", "dest_ip": "10.0.0.18", "protocol": "UDP", "packet_size": 1024, "timestamp": 1690797597}

Received message: {"source_ip": "192.168.0.19", "dest_ip": "10.0.0.19", "protocol": "TCP", "packet_size": 1280, "timestamp": 1690797598}

Received message: {"source_ip": "192.168.0.20", "dest_ip": "10.0.0.20", "protocol": "ICMP", "packet_size": 512, "timestamp": 1690797599}

Received message: {"source_ip": "192.168.0.21", "dest_ip": "10.0.0.21", "protocol": "TCP", "packet_size": 256, "timestamp": 1690797600}

Received message: {"source_ip": "192.168.0.22", "dest_ip": "10.0.0.22", "protocol": "UDP", "packet_size": 2048, "timestamp": 1690797601}

Received message: {"source_ip": "192.168.0.23", "dest_ip": "10.0.0.23", "protocol": "TCP", "packet_size": 512, "timestamp": 1690797602}

Received message: {"source_ip": "192.168.0.24", "dest_ip": "10.0.0.24", "protocol": "ICMP", "packet_size": 128, "timestamp": 1690797603}

File Edit Selection View Go Run Terminal Help

EXPLORER

NEWNTWORK_TRAFFIC

- __pycache__
- anomalies_log.csv
- anomalis_ext.py
- anomalis.py
- create_csv.py
- isolation_forest_model.pkl
- isolation_forest_model1.pkl
- kafka_consumer.py
- kafka_producer.py
- network.csv
- newtype.py
- scaler.pkl
- stream_kafka.py
- trainset.py
- visualization.csv
- visualization.py

stream_kafka.py >...

```
# Define the query to print the aggregated data to the console
54 query = df.writeStream \
55     .outputMode("update") \
56     .format("console") \
57     .option("checkpointLocation", "/home/nandhumidhun/test/spark_log") \
58     .start()
59
60 # Wait for the termination of the queries
61 try:
62     converted_data_query.awaitTermination()
63     query.awaitTermination()
64 except KeyboardInterrupt:
65     print("Streaming query terminated.")
66 finally:
67     # Stop the Spark session
68     spark.stop()
```

TIMELINE

DEBUG CONSOLE

Filter (e.g. text, !exclude, \escape)

PROBLEMS 50

TERMINAL

PORTS

bash	192.168.0.8	10.0.0.8	ICMP	128	2023-07-31 15:29:47		
python3	192.168.0.9	10.0.0.9	TCP	768	2023-07-31 15:29:48		
bash	192.168.0.10	10.0.0.10	UDP	2048	2023-07-31 15:29:49		
	192.168.0.11	10.0.0.11	TCP	1024	2023-07-31 15:29:50		
	192.168.0.12	10.0.0.12	ICMP	256	2023-07-31 15:29:51		
	192.168.0.13	10.0.0.13	TCP	512	2023-07-31 15:29:52		
	192.168.0.14	10.0.0.14	UDP	1280	2023-07-31 15:29:53		
	192.168.0.15	10.0.0.15	TCP	256	2023-07-31 15:29:54		
	192.168.0.16	10.0.0.16	ICMP	64	2023-07-31 15:29:55		
	192.168.0.17	10.0.0.17	TCP	768	2023-07-31 15:29:56		
	192.168.0.18	10.0.0.18	UDP	1024	2023-07-31 15:29:57		
	192.168.0.19	10.0.0.19	TCP	1280	2023-07-31 15:29:58		
	192.168.0.20	10.0.0.20	ICMP	512	2023-07-31 15:29:59		
	192.168.0.21	10.0.0.21	TCP	256	2023-07-31 15:30:00		
	192.168.0.22	10.0.0.22	UDP	2048	2023-07-31 15:30:01		
	192.168.0.23	10.0.0.23	TCP	512	2023-07-31 15:30:02		

only showing top 20 rows

Batch: 11

		window packet_size_avg	
		(2023-07-31 15:30:...	806.4
		(2023-07-31 15:25:...	691.2

Ln 50, Col 14 Spaces: 4 UTF-8 LF Python 3.10.12 64-bit

File Edit Selection View Go Run Terminal Help

EXPLORER

NEWTWORK_TRAFFIC

- > __pycache__
- anomalies_log.csv
- anomalis_ext.py
- anomalis.py
- create_csv.py
- isolation_forest_model.pkl
- isolation_forest_model1.pkl
- kafka_consumer.py
- kafka_producer.py
- network.csv
- newtype.py
- scaler.pkl
- stream_kafka.py
- trainset.py
- visualization.csv
- visualization.py

trainset.py > ...

```
6
7 # Features for anomaly detection
8 X = df[['packet_size']]
9
10 # Train Isolation Forest
11 model = IsolationForest(contamination=0.01)
12 model.fit(X)
13
14 # Save the model
15 import joblib
16 joblib.dump(model, "isolation_forest_model1.pkl")
17 joblib.dump(model, "isolation_forest_model.pkl")
18
```

TIMELINE

DEBUG CONSOLE

Filter (e.g. text, !exclude, \escape)

PROBLEMS 50

TERMINAL PORTS

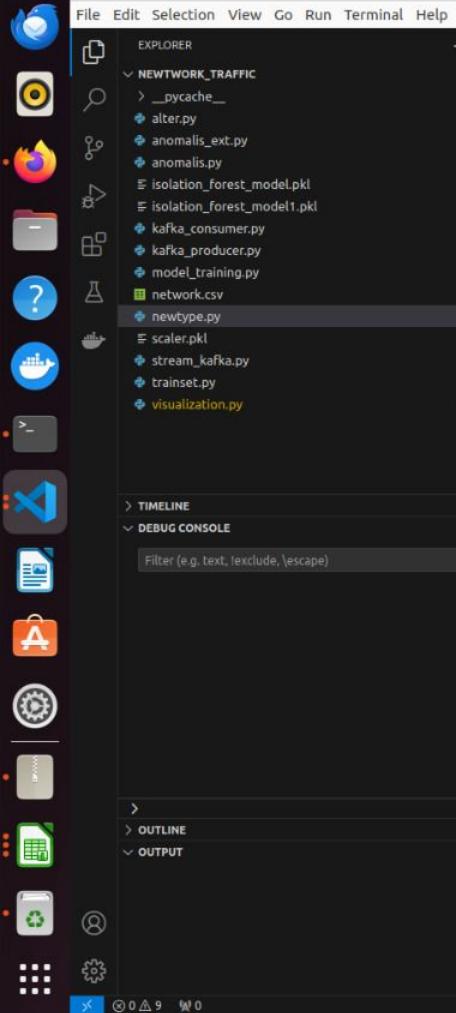
bash python3 bash

```
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/CUBEAI SOLUTIONS/Newtwork_traffic$ python3 trainset.py
nandhumidhun@nandhumidhun-HP-Laptop-15-bs0xx:~/CUBEAI SOLUTIONS/Newtwork_traffic$
```

OUTLINE

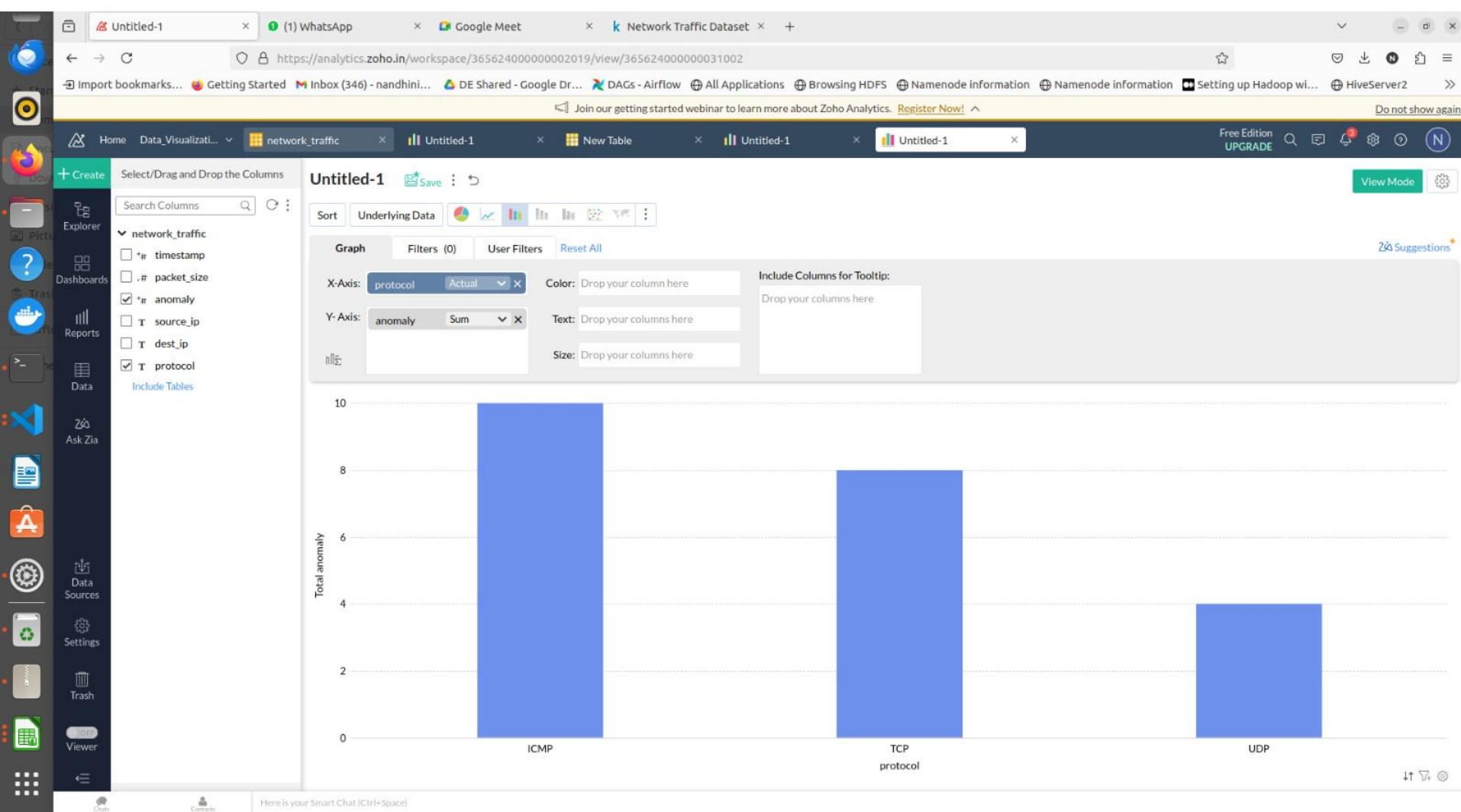
OUTPUT

Ln 16, Col 43 Spaces: 4 UTF-8 LF Python 3.10.12 64-bit



```
EXPLORER ... kafka_producer.py kafka_consumer.py stream_kafka.py trainset.py newtype.py x model_training.py anomalis.py anomalis_ext.py alter.py  
NEWNTWORK_TRAFFIC  
> __pycache__  
alter.py  
anomalis_ext.py  
anomalis.py  
isolation_forest_model.pkl  
isolation_forest_model1.pkl  
kafka_consumer.py  
kafka_producer.py  
model_training.py  
network.csv  
newtype.py  
scaler.pkl  
stream_kafka.py  
trainset.py  
visualization.py  
  
newtype.py > ...  
3 import pandas as pd  
4 import joblib  
5  
6 # Example training data  
7 data = pd.DataFrame({  
8     'packet_size': [1024, 2048, 64, 1280] # Ensure this is representative of normal and anomalous data  
9 })  
10  
11 # Scaling  
12 scaler = StandardScaler()  
13 data_scaled = scaler.fit_transform(data)  
14  
15 # Train the model  
16 model = IsolationForest()  
17 model.fit(data_scaled)  
18  
19 # Save the model and scaler  
20 joblib.dump(model, 'isolation_forest_model.pkl')  
21 joblib.dump(scaler, 'scaler.pkl')  
22
```

```
TIMELINE 9 PROBLEMS 9 TERMINAL PORTS  
python3 bash  
Received message: {'source_ip': '192.168.0.18', 'dest_ip': '10.0.0.18', 'protocol': 'UDP', 'packet_size': 1024, 'timestamp': 1690797597}  
Anomaly score: [1]  
Normal data: {'source_ip': '192.168.0.18', 'dest_ip': '10.0.0.18', 'protocol': 'UDP', 'packet_size': 1024, 'timestamp': 1690797597}  
Received message: {'source_ip': '192.168.0.19', 'dest_ip': '10.0.0.19', 'protocol': 'TCP', 'packet_size': 1280, 'timestamp': 1690797598}  
Anomaly score: [1]  
Normal data: {'source_ip': '192.168.0.19', 'dest_ip': '10.0.0.19', 'protocol': 'TCP', 'packet_size': 1280, 'timestamp': 1690797598}  
Received message: {'source_ip': '192.168.0.20', 'dest_ip': '10.0.0.20', 'protocol': 'ICMP', 'packet_size': 512, 'timestamp': 1690797599}  
Anomaly score: [1]  
Normal data: {'source_ip': '192.168.0.20', 'dest_ip': '10.0.0.20', 'protocol': 'ICMP', 'packet_size': 512, 'timestamp': 1690797599}  
Received message: {'source_ip': '192.168.0.21', 'dest_ip': '10.0.0.21', 'protocol': 'TCP', 'packet_size': 256, 'timestamp': 1690797600}  
Anomaly score: [-1]  
Anomaly detected: {'source_ip': '192.168.0.21', 'dest_ip': '10.0.0.21', 'protocol': 'TCP', 'packet_size': 256, 'timestamp': 1690797600}  
Received message: {'source_ip': '192.168.0.22', 'dest_ip': '10.0.0.22', 'protocol': 'UDP', 'packet_size': 2048, 'timestamp': 1690797601}  
Anomaly score: [-1]  
Anomaly detected: {'source_ip': '192.168.0.22', 'dest_ip': '10.0.0.22', 'protocol': 'UDP', 'packet_size': 2048, 'timestamp': 1690797601}  
Received message: {'source_ip': '192.168.0.23', 'dest_ip': '10.0.0.23', 'protocol': 'TCP', 'packet_size': 512, 'timestamp': 1690797602}  
Anomaly score: [1]  
Normal data: {'source_ip': '192.168.0.23', 'dest_ip': '10.0.0.23', 'protocol': 'TCP', 'packet_size': 512, 'timestamp': 1690797602}  
Received message: {'source_ip': '192.168.0.24', 'dest_ip': '10.0.0.24', 'protocol': 'ICMP', 'packet_size': 128, 'timestamp': 1690797603}  
Anomaly score: [-1]  
Anomaly detected: {'source_ip': '192.168.0.24', 'dest_ip': '10.0.0.24', 'protocol': 'ICMP', 'packet_size': 128, 'timestamp': 1690797603}  
Received message: {'source_ip': '192.168.0.25', 'dest_ip': '10.0.0.25', 'protocol': 'TCP', 'packet_size': 768, 'timestamp': 1690797604}  
Anomaly score: [1]  
Normal data: {'source_ip': '192.168.0.25', 'dest_ip': '10.0.0.25', 'protocol': 'TCP', 'packet_size': 768, 'timestamp': 1690797604}  
Received message: {'source_ip': '192.168.0.26', 'dest_ip': '10.0.0.26', 'protocol': 'UDP', 'packet_size': 1024, 'timestamp': 1690797605}  
Anomaly score: [1]  
Normal data: {'source_ip': '192.168.0.26', 'dest_ip': '10.0.0.26', 'protocol': 'UDP', 'packet_size': 1024, 'timestamp': 1690797605}  
Received message: {'source_ip': '192.168.0.27', 'dest_ip': '10.0.0.27', 'protocol': 'TCP', 'packet_size': 256, 'timestamp': 1690797606}
```



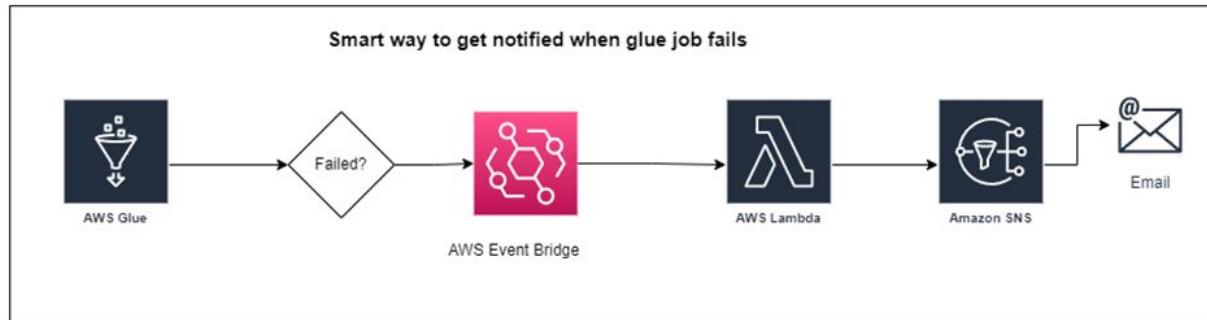
AWS GLUE JOB FAILURE NOTIFICATION WITH EVENTBRIDGE, LAMBDA, AND SNS

STEP 1: Open AWS Glue and Create the failed script

STEP 2: Create Event Bridge

STEP 3: Create lambda function and connect it to Event Bridge and SNS

STEP 4: Create SNS and create subscription



Activities Firefox Web Browser Aug 12 20:27

(2) WhatsApp x PLAYING Meet - apw-0 x An AWS Glue x Console Home x Rules | Amazon x gluealert | Fi x Subscription x Runs - Editor x ChatGPT x An AWS Glue x Subscription cor x + v https://us-east-1.console.aws.amazon.com/gluestudio/home?region=us-east-1#/editor/job/ETL_job/runs 110% ☆ Import bookmarks... Getting Started Inbox (346) - nandhini... DE Shared - Google Dr... DAGs - Airflow All Applications Browsing HDFS Namenode information Namenode information Setting up Hadoop wi... HiveServer2 N. Virginia NandhuCubeAI

aws Services Search [Alt+S] Actions Save Run

ETL_job

Last modified on 12/8/2024, 7:58:42 pm

Script Job details Runs Data quality Schedules Version Control

Job runs (1/6) Info Last updated (UTC) August 12, 2024 at 14:56:06 C View details Stop job run Table View Card View

Filter job runs by property

Run status	Retries	Start time (Local)	End time (Local)	Duration	Capacity (DPUs)	Worker type	Glue version
Failed	1	08/12/2024 20:25:24	08/12/2024 20:25:27	0 s	0.0625 DPUs	-	3.0
Failed	0	08/12/2024 20:24:50	08/12/2024 20:24:53	0 s	0.0625 DPUs	-	3.0
Failed	1	08/12/2024 20:23:51	08/12/2024 20:23:54	0 s	0.0625 DPUs	-	3.0
Failed	0	08/12/2024 20:23:17	08/12/2024 20:23:21	0 s	0.0625 DPUs	-	3.0
Failed	1	08/12/2024 19:59:24	08/12/2024 19:59:27	0 s	0.0625 DPUs	-	3.0
Failed	0	08/12/2024 19:59:50	08/12/2024 19:59:54	0 s	0.0625 DPUs	-	3.0

Run details Input arguments (6) Continuous logs Run insights Metrics Spark UI

JobName:ETL_job and JobRunId:jr_b20471131e344ba7857b6ca45252a2c642ecdff92d2cff266fbe49b4a69f4207 failed to execute with exception Role arn:aws:iam::637423373422:role/glueredshifts3 should be given assume role permissions for Glue Service. (Service: AWSGlueJobExecutor; Status Code: 400; Error Code: InvalidInputException; Request ID: 69a40afa-ed75-467a-9399-4b0d3f7c766b; Proxy: null)

Job name	Start time (Local)	Glue version	Last modified on (Local)
ETL_job	08/12/2024 19:58:50	3.0	08/12/2024 19:58:54
Id	End time (Local)	Worker type	Log group name
jr_b20471131e344ba7857b6ca45252a2c642ecdff92d2cff266fbe49b4a69f4207	08/12/2024 19:58:54	-	/aws-glue/python-jobs
Run status	Start-up time	Max capacity	Number of workers
Failed	0	0.0625 DPUs	-
Retry attempt number	Execution time	Execution class	Timeout
1	0 s	Standard	3000 ms

CloudShell Feedback © 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Activities Firefox Web Browser Aug 12 20:28

(2) WhatsApp X Google Meet X An AWS Glue X Console Home X Rules | Amazon X gluealert | Fx X Gluetopic | T X Jobs - AWS X ChatGPT X An AWS Glue X Subscription cor X +

Import bookmarks... Getting Started Inbox (346) - nandhini... DE Shared - Google Dr... DAGs - Airflow All Applications Browsing HDFS Namenode information Namenode information Setting up Hadoop wi... HiveServer2 N. Virginia NandhuCubeAI

aws Services Search [Alt+S] X ? X ⓘ X

Successfully updated the function gluealert.

Diagram Template

gluealert

Layers (0)

+ Add destination

EventBridge (CloudWatch Events)

+ Add trigger

Description

Last modified 37 minutes ago

Function ARN arn:aws:lambda:us-east-1:637423373422:function:gluealert

Function URL Info

Code Test Monitor Configuration Aliases Versions

Code source Info

Upload from

File Edit Find View Go Tools Window Test Deploy

Environment Var Execution results

Environment

lambda_function

lambda_function.py

```
1 # Import modules
2 import json
3 import logging
4 import boto3
5
6 # Set up logging
7 logger = logging.getLogger()
8 logger.setLevel(logging.INFO)
9
10 # Set up Boto3 client for SNS
11 client = boto3.client('sns')
12
13 # Variables for the SNS:
14 SNS_TOPIC_ARN = "arn:aws:sns:us-east-1:637423373422:Gluetopic"
```

Info Tutorials

Learn how to implement common use cases in AWS Lambda.

Create a simple web app

In this tutorial you will learn how to:

- Build a simple web app, consisting of a Lambda function with a function URL that outputs a webpage
- Invoke your function through its function URL

Learn more

Start tutorial

CloudShell Feedback © 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Activities Firefox Web Browser Aug 12 20:28 110%

CloudShell Feedback

https://us-east-1.console.aws.amazon.com/events/home?region=us-east-1#/rules

Import bookmarks... Getting Started Inbox (346) - nandhini... DE Shared - Google Dr... DAGs - Airflow All Applications Browsing HDFS Namenode information Namenode information Setting up Hadoop wi... HiveServer2

aws Services Search [Alt+S] Provide Feedback X

Amazon EventBridge X

Let us know what you think! We recently added a new dashboard in the console. We'd love to hear your feedback.

1 rule(s) deleted successfully

Amazon EventBridge > Rules

Rules

A rule watches for specific types of events. When a matching event occurs, the event is routed to the targets associated with the rule. A rule can be associated with one or more targets.

Select event bus

Event bus

Select or enter event bus name

default

Rules (1)

Find rules Any status

Create rule

Name	Status	Type	ARN	Description
glue_alert	Enabled	Standard	arn:aws:events:us-east-1:63742 3373422:rule/glue_alert	Alert for glue fail

CloudFormation Template

CloudShell Feedback

© 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Activities Firefox Web Browser Aug 12 20:28

(2) WhatsApp X Google Meet X An AWS Glue X Console Home X Rules | Amaz X glualert | Fi X Gluetopic | T C Jobs - AWS X ChatGPT X An AWS Glue X Subscription cor X +

Import bookmarks... Getting Started Inbox (346) - nandhini... DE Shared - Google Dr... DAGs - Airflow All Applications Browsing HDFS Namenode information Namenode information Setting up Hadoop wi... HiveServer2 N. Virginia NandhuCubeAI

aws Services Search [Alt+S] X

New Feature Amazon SNS now supports in-place message archiving and replay for FIFO topics. [Learn more](#)

Amazon SNS X

Dashboard Topics Subscriptions

Mobile Push notifications Text messaging (SMS) Origination numbers

Amazon SNS > Topics > Gluetopic

Gluetopic

Edit Delete Publish message

Details

Name	Display name
Gluetopic	glue alert
ARN	Topic owner
arn:aws:sns:us-east-1:637423373422:Gluetopic	637423373422
Type	
Standard	

Subscriptions Access policy Data protection policy Delivery policy (HTTP/S) Delivery status logging Encryption Tags Integrations

Subscriptions (2)

Search

ID	Endpoint	Status	Protocol
4c1d142d-b4a1-4e55-a489-4cb342d21...	nandhiriravi1402@gmail.com	Confirmed	EMAIL
c354bd8d-8974-473c-b405-c21a1cac3...	nandhiriravice01@gmail.com	Confirmed	EMAIL

Create subscription

CloudShell Feedback © 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Activities Firefox Web Browser Aug 12 20:30

https://mail.google.com/mail/u/1/?hl=en#inbox/ FMfcgzQVzFNjGngZNhdclRvrrLjtTss

Import bookmarks... Getting Started Inbox (346) - nandhini... DE Shared - Google Dr... DAGs - Airflow All Applications Browsing HDFS Namenode information Namenode information Setting up Hadoop wi... HiveServer2

Gmail Search mail

Compose

Inbox 665

Starred

Snoozed

Sent

Drafts 11

More

Labels +

An AWS Glue Job has failed Inbox

glue alert <no-reply@sns.amazonaws.com>
to me ▾
A Glue Job has failed after attempting to retry. JobName: ETL_job, JobRunID: jr_2bd2ea6e1feef8dd9b60da71b9331a6c6d1e2b4c1840e4ab3ea4dcde054b7509_attempt_1"

--
If you wish to stop receiving notifications from this topic, please click or visit the link below to unsubscribe:
<https://sns.us-east-1.amazonaws.com/unsubscribe.html?SubscriptionArn=arn:aws:sns:us-east-1:637423373422:Gluetopic:4c1d142d-b4a1-4e55-a489-4cb342d2171&Endpoint=nandhinirav1402@gmail.com>

Please do not reply directly to this email. If you have any questions or comments regarding this email, please contact us at <https://aws.amazon.com/support>

glue alert <no-reply@sns.amazonaws.com>
to me ▾
A Glue Job has failed after attempting to retry. JobName: ETL_job, JobRunID: jr_c171fc43b6a5179ff045536a01e1b507221973a7edc6694acfd7083e3bc11e1_attempt_1"

...
--
If you wish to stop receiving notifications from this topic, please click or visit the link below to unsubscribe:
<https://sns.us-east-1.amazonaws.com/unsubscribe.html?SubscriptionArn=arn:aws:sns:us-east-1:637423373422:Gluetopic:4c1d142d-b4a1-4e55-a489-4cb342d2171&Endpoint=nandhinirav1402@gmail.com>

Please do not reply directly to this email. If you have any questions or comments regarding this email, please contact us at <https://aws.amazon.com/support>

Reply Forward

THANK YOU

ANY QUERY