

TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN
NGÀNH KHOA HỌC DỮ LIỆU



PYTHON CHO KHOA HỌC DỮ LIỆU

ĐỒ ÁN CUỐI KÌ
TAXI TRIP RECORD ETL-PIPELINE

Giảng viên phụ trách

Thầy Hà Văn Thảo

Thành phố Hồ Chí Minh, tháng 1, năm 2024

MỤC LỤC

1. GIỚI THIỆU	1
1.1 Thông tin nhóm	1
1.2 Nguyên tắc hoạt động nhóm	1
1.3 Mục tiêu project	1
1.4 Data sources	1
2. ARCHITECTURE	3
2.1 Dictionary tree	3
2.2 Database	4
2.3 ETL PIPELINE	5
2.3.1 resources	5
2.3.2 assets	9
3. CÀI ĐẶT	17
3.1 Ngoài thư mục lớn NYC-TRIPRECORD	17
3.2 Trong thư mục etl_pipeline	19
4. DEMO KẾT QUẢ VÀ VISUALIZE	20
4.1 DEMO	20
4.2 VISUALIZE	20
TÀI LIỆU THAM KHẢO	22

NỘI DUNG

1. GIỚI THIỆU

1.1 Thông tin nhóm

- Tên nhóm: Nhóm 7
- Thành viên:
 - 21280115 - Trần Đức Trung
 - 21280113 - Lê Quang Trung
 - 21280012 - Nguyễn Đông Hải
 - 21280018 - Trần Phi Hùng

1.2 Nguyên tắc hoạt động nhóm

- Ý tưởng được đóng góp từ tất cả thành viên của nhóm.
- Mọi người tôn trọng ý kiến của nhau.
- Giao và làm việc đúng thời hạn đưa ra. Nếu trễ thì sẽ chịu trách nhiệm về phần của mình.

1.3 Mục tiêu project

Với mục tiêu tìm hiểu và học tập về data engineering và big data, nhóm đã thực hiện project NYC Taxi Trip Record. Project tập trung vào quá trình ETL từ database (MySQL) đến datawarehouse (Psql), đồng thời tận dụng sức mạnh của Minio, Pyspark và Polars để thực hiện clean và transform. Cuối cùng sử dụng streamlit để visualize và analyze data Trong dự án này, nhóm mình sẽ minh họa rõ ràng và chi tiết các quy trình thực hiện ETL trên tập dữ liệu TLC Trip Record Data - một tập dữ liệu mở, phục vụ cho việc học tập và nghiên cứu.

1.4 Data sources

NYC Taxi Trip Record được chọn lọc từ tập dataset chính thức tại TLC Trip Record Data. Dữ liệu bao gồm:

- Yellow Taxi Trip Records: Đây là điều mọi người nghĩ đến khi nhắc về taxi ở New York. Chiếc taxi mang tính biểu tượng là tiêu chuẩn trong vận chuyển ô tô ở New York. Xe taxi màu vàng là phương tiện duy nhất được phép đón khách ở mọi nơi trong thành phố.

- Green Taxi Trip Records: Xe xanh có thể thả bạn đi bất cứ đâu nhưng chỉ được phép đón khách ở những khu vực nhất định.
- For Hire Vehicle Trip Records (FHV Trip Records): Là những taxi chủ yếu được sử dụng cho chuyển đi được sắp xếp trước.

Dữ liệu được sử dụng trong project:

- 1 tháng (January 2023) của TLC Trip Record Data

Link website: <https://www.nyc.gov/site/tlc/about/tlc-trip-record-data.page>

Với data dictionary như sau:

https://www.nyc.gov/assets/tlc/downloads/pdf/data_dictionary_trip_records_yellow.pdf

https://www.nyc.gov/assets/tlc/downloads/pdf/data_dictionary_trip_records_green.pdf

https://www.nyc.gov/assets/tlc/downloads/pdf/data_dictionary_trip_records_fhv.pdf

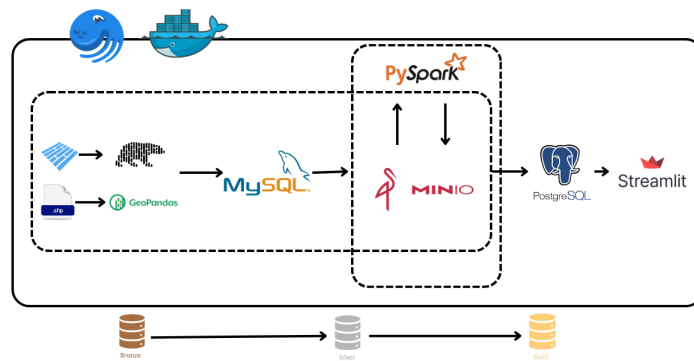
- Taxi Zone Shapefile: https://d37ci6vzurychx.cloudfront.net/misc/taxi_zones.zip

2. ARCHITECTURE

2.1 Dictionary tree

Chi tiết:

- **app**: The UI's application written with **streamlit**
- **dagster_home**: Dagit and dagster daemon's configurations
- **dockerimages**: self-built docker images, such as dagster (for dagit + daemon), spark master,...
- **etl_pipeline**: pipeline
- **load_dataset**: include files have .sql to create schema and load data to MySQL, Postgres
- **minio**: Docker container for MinIO
- **mysql**: Docker container for MySQL
- **postgresql**: Docker container for Psq
- **Test**: folder to test and EDA
- **.gitignore + .gitattributes**: Code versioning
- **docker-compose**: to compose docker containers
- **env**: Env variables.
- **makefile**: simplify terminal's commands
- **requirements.txt**: required library.



2.2 Database

<div>long_lat</div> <div> 123 LocationID 123 longitude 123 latitude </div>	<div>fhv_record</div> <div> abc dispatching_base_num 🕒 pickup_datetime 🕒 dropOff_datetime 123 PUlocationID 123 DOlocationID abc SR_Flag abc Affiliated_base_number 123 Column1 </div>	<div>yellow_record</div> <div> 123 VendorID 🕒 tpep_pickup_datetime 🕒 tpep_dropoff_datetime 123 passenger_count 123 trip_distance 123 RatecodeID abc store_and_fwd_flag 123 PUlocationID 123 DOLocationID 123 payment_type 123 fare_amount 123 extra 123 mta_tax 123 tip_amount 123 tolls_amount 123 improvement_surcharge 123 total_amount 123 congestion_surcharge 123 airport_fee 123 Column1 </div>	<div>green_record</div> <div> 123 VendorID 🕒 lpep_pickup_datetime 🕒 lpep_dropoff_datetime abc store_and_fwd_flag 123 RatecodeID 123 PUlocationID 123 DOLocationID 123 passenger_count 123 trip_distance 123 fare_amount 123 extra 123 mta_tax 123 tip_amount 123 tolls_amount abc ehail_fee 123 improvement_surcharge 123 total_amount 123 payment_type 123 trip_type 123 congestion_surcharge 123 Column1 </div>
--	--	---	---

Sau khi thực hiện load data vào database ở đây là MySQL ta được các table có như trên với:

- yellow_record, green_record, fhv_record: có data dictionary như đã mô tả
- long_lat: là 1 table chứa kinh độ, vĩ độ, tương ứng với các LocationID

Một vài xử lý nhỏ khi đọc file **taxizones.shp** để load các data về kinh độ, vĩ độ đưa vào trong database (MySQL). Sử dụng thư viện **GeoPandas** and **PyProj** và chuyển đổi từ **shapefile format** thành một **DataFrame**

```

Python
shapefile_path = f"{adr_data}/{zone_data}/{zone_data}.shp"

gdf = gpd.read_file(shapefile_path)
# Define
source_crs = gdf.crs # CRS of the shapefile
target_crs = 'EPSG:4326' # WGS84 - lat/lon CRS

# Create a PyProj transformer
transformer = pyproj.Transformer.from_crs(source_crs, target_crs, always_xy=True)

gdf['longitude'] = gdf.geometry.centroid.x
gdf['latitude'] = gdf.geometry.centroid.y
gdf['longitude'], gdf['latitude'] = transformer.transform(gdf['longitude'], gdf['latitude'])

df = pl.DataFrame(gdf[['LocationID', 'longitude', 'latitude']])

```

2.3 ETL PIPELINE

2.3.1 resources

Là folder chứa các input output manager

- **mysql_io_manager.py**

Là 1 file python chứa **class MySQLIOManager**, với hàm **extract_data** làm nhiệm vụ kết nối vào MySQL đọc và trả ra 1 polars dataframe

```

mysql_io_manager.py X
etl_pipeline > etl_pipeline > resources > mysql_io_manager.py > ...
1 from dagster import IOManager, InputContext, OutputContext
2 import polars as pl
3
4
5 def connect_mysql(config) -> str:
6     conn = (
7         f"mysql://{config['user']}:{config['password']}@{config['host']}:{config['port']}/{config['database']}"
8     )
9     return conn
10
11
12
13
14 class MySQLIOManager(IOManager):
15     def __init__(self, config):
16         self.config = config
17
18     def handle_output(self, context: "OutputContext", obj: pl.DataFrame):
19         pass
20
21     def load_input(self, context: "InputContext"):
22         pass
23
24     def extract_data(self, sql: str) -> pl.DataFrame:
25         conn = connect_mysql(self.config)
26         df_data = pl.read_database(query=sql, connection_uri=conn)
27         return df_data

```

- **minio_io_manager.py**

Là 1 file python chứa **class MinIOIOManager**, với các hàm chính như:

- **load_input** thực hiện lấy data từ MinIO đưa vào asset để xử lý
- **handle_output** thực hiện chuyển data sau khi đã được xử lý trong asset và đưa lên MinIO

```

minio_io_manager.py X
etl_pipeline > etl_pipeline > resources > minio_io_manager.py > ...
1  from typing import Union
2  import polars as pl
3  from dagster import IOManager, OutputContext, InputContext
4  from minio import Minio
5  import os
6  import pyarrow.parquet as pq
7
8
9  def make_bucket(client: Minio, bucket_name):
10     found = client.bucket_exists(bucket_name)
11     if not found:
12         client.make_bucket(bucket_name)
13     else:
14         print(f"Bucket {bucket_name} already exists.")
15
16
17  class MinIOIOManager(IOManager):
18     def __init__(self, config):
19         self.config = config
20         self.minio_client = Minio(
21             self.config["endpoint_url"],
22             access_key=self.config["aws_access_key_id"],
23             secret_key=self.config["aws_secret_access_key"],
24             secure=False
25         )
26
27     def _get_path(self, context: Union[OutputContext, InputContext]):
28         layer, schema, table = context.asset_key.path
29         key = "/".join([layer, schema, table.replace("{layer}_", "")])
30         tmp_file_path = "/tmp/file-{}.parquet".format(
31             "_".join(context.asset_key.path),
32         )
33         if context.has_partition_key:
34             partition_str = str(table) + "_" + context.asset_partition_key
35             return os.path.join(key, f"{partition_str}.parquet"), tmp_file_path
36         else:
37             return f"{key}.parquet", tmp_file_path
38
39
40
41     def handle_output(self, context: OutputContext, obj: pl.DataFrame):
42         key_name, tmp_file_path = self._get_path(context)
43         table = obj.to_arrow()
44         pq.write_table(table, tmp_file_path)
45
46         try:
47             bucket_name = self.config["bucket"]
48             make_bucket(self.minio_client, bucket_name)
49             self.minio_client.fput_object(
50                 bucket_name, key_name, tmp_file_path,
51             )
52             context.log.info(
53                 f"(MinIO handle_output) Number of rows and columns: {obj.shape}"
54             )
55             # Clean up tmp file
56             os.remove(tmp_file_path)
57         except Exception as e:
58             raise e
59
60     def load_input(self, context: InputContext) -> pl.DataFrame:
61         bucket_name = self.config["bucket"]
62         key_name, tmp_file_path = self._get_path(context)
63         try:
64             context.log.info(f"(MinIO load_input) from key_name is {key_name}")
65             self.minio_client.fget_object(
66                 bucket_name, key_name, tmp_file_path,
67             )
68             df = pl.read_parquet(tmp_file_path)
69             os.remove(tmp_file_path)
70             return df
71         except Exception as e:
72             raise e

```

- spark_io_manager.py

Là 1 file python chứa class **SparkIOManager** như sau:


```

spark_io_manager.py M X
etl_pipeline > etl_pipeline > resources > spark_io_manager.py > get_spark_session
1 from dagster import IOManager, OutputContext, InputContext
2 from pyspark.sql import SparkSession, DataFrame
3 from contextlib import contextmanager
4
5
6
7 @contextmanager
8 def get_spark_session(config, run_id="Spark IO Manager"):
9     executor_memory = "1g" if run_id != "Spark IO Manager" else "1500m"
10    try:
11        spark = (
12            SparkSession.builder.master("spark://spark-master:7077")
13            # SparkSession.builder.master("local[*]")
14            .appName(run_id)
15            .config("spark.driver.memory", "1g")
16            .config("spark.executor.memory", executor_memory)
17            .config("spark.cores.max", "4")
18            .config("spark.executor.cores", "2")
19            .config(
20                "spark.jars",
21                "/usr/local/spark/jars/delta-core_2.12-2.2.0.jar,/usr/local/spark/jars/hadoop-aws-3.3.2.jar,/usr/local/spark/jars/delta-storage-2.2.0.jar,/usr/l
22            )
23            .config("spark.driver.extraClassPath", "/usr/local/spark/jars/postgresql-42.7.1.jar")
24            .config("spark.executor.extraClassPath", "/usr/local/spark/jars/postgresql-42.7.1.jar")
25            .config(
26                "spark.sql.catalog.spark_catalog",
27                "org.apache.spark.sql.delta.catalog.DeltaCatalog",
28            )
29            .config("spark.sql.extensions", "io.delta.sql.DeltaSparkSessionExtension")
30            .config("spark.hadoop.fs.s3a.endpoint", f"http://{config['endpoint_url']}")
31            .config("spark.hadoop.fs.s3a.access.key", str(config["minio_access_key"]))
32            .config("spark.hadoop.fs.s3a.secret.key", str(config["minio_secret_key"]))
33            .config("spark.hadoop.fs.s3a.path.style.access", "true")
34            .config("spark.hadoop.fs.connection.ssl.enabled", "false")
35            .config(
36                "spark.hadoop.fs.s3a.impl", "org.apache.hadoop.fs.s3a.S3AFileSystem"
37            )
38            .config("spark.sql.execution.arrow.pyspark.enabled", "true")
39            .config("spark.sql.execution.arrow.pyspark.fallback.enabled", "true")
40            .getOrCreate()
41        )
42    except Exception as e:

```

- Hình trên là về **hàm get_spark_session**, hàm sẽ thực hiện tạo 1 spark session và config các thông tin về core, memory, cũng như các extension của Spark mà sẽ được sử dụng (ví dụ như ext pgJDBC không có sẵn trong spark, để kết nối được với psql)

```

spark_io_manager.py M X
etl_pipeline > etl_pipeline > resources > spark_io_manager.py > SparkIOManager > load_input

45
46 class SparkIOManager(IOManager):
47     def __init__(self, config):
48         self._config = config
49
50     def handle_output(self, context: OutputContext, obj: DataFrame):
51         # Write output to s3a (MinIO)
52         context.log.debug("(Spark handle_output) Writing output to MinIO ...")
53
54         file_path = "s3a://lakehouse/" + "/".join(context.asset_key.path)
55         if context.has_partition_key:
56             file_path += f"/{context.partition_key}"
57         context.log.debug(f"(Spark handle_output) File path: {file_path}")
58         file_name = str(context.asset_key.path[-1])
59         context.log.debug(f"(Spark handle_output) File name: {file_name}")
60
61         try:
62             obj.write.mode("overwrite").parquet(file_path)
63             context.log.debug(f"Saved {file_name} to {file_path}")
64         except Exception as e:
65             raise Exception(f"(Spark handle_output) Error while writing output: {e}")
66
67     def load_input(self, context: InputContext) -> DataFrame:
68         context.log.debug(f"Loading input from {context.asset_key.path} ...")
69         file_path = "s3a://lakehouse/" + "/".join(context.asset_key.path)
70         check_partition = (context.metadata or {}).get("partition", True)
71         if check_partition == True:
72             if context.has_partition_key:
73                 file_path += f"/{context.partition_key}"
74         full_load = (context.metadata or {}).get("full_load", False)
75
76         try:
77             with get_spark_session(self._config) as spark:
78                 df = None
79                 if full_load:
80                     df = (
81                         spark.read.format("parquet")
82                         .options(header=True, inferSchema=False)
83                         .load(file_path + "/*")
84                     )
85                 else:
86                     df = spark.read.parquet(file_path)
87
88             return df
89         except Exception as e:
90             raise Exception(f"Error while loading input: {e}")

```

- Kể tới là các hàm **load_input** và **handle_output** để đưa data lên s3a (đã được setting là MinIO) hoặc lấy data về asset
- **psql_io_manager.py** - Cuối cùng là file python chứa **class PSQLIOManager**, chứa **hàm handle_output** sẽ sử dụng spark và ext pgJDBC của nó để kết nối spark session với psql và load data vào psql

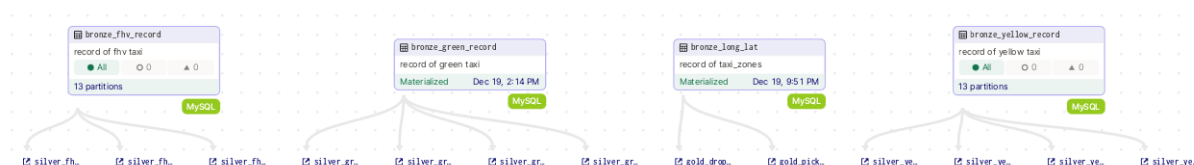
```

psql_io_manager.py X
etl_pipeline > etl_pipeline > resources > psql_io_manager.py > ...
1 from dagster import IOManager, InputContext, OutputContext
2 from pyspark.sql import DataFrame
3 import polars as pl
4
5
6 def connect_psql(config, table):
7     conn = {
8         "url": f"jdbc:postgresql://{config['host']}:{config['port']}/{config['database']}",
9         "dbtable": table,
10        "user": config["user"],
11        "password": config["password"],
12    }
13    return conn
14
15
16 class PostgreSQLIOManager(IOManager):
17     def __init__(self, config):
18         self._config = config
19
20     def handle_output(self, context: OutputContext, obj: DataFrame):
21         table = context.asset_key.path[-1]
22         schema_ = context.asset_key.path[-2]
23         full_table = f"{schema_}.{table}"
24         conn = connect_psql(self._config, full_table)
25         context.log.info(f"Now use pgjdbc load to postgres")
26         obj.write \
27             .mode("overwrite") \
28             .format("jdbc") \
29             .option("url", conn['url']) \
30             .option("dbtable", conn['dbtable']) \
31             .option("user", conn['user']) \
32             .option("password", conn['password']) \
33             .save()
34
35     def load_input(self, context: InputContext) -> pl.DataFrame:
36         pass

```

2.3.2 assets

- bronze_layer



Tại các asset trên layer này, ta sẽ sử dụng `mysql_io_manager` làm **required resource keys** và `minio_io_manager` là **io manager key** để thực hiện đọc data trên database (MySQL) và đưa vào folder `bronze_layer` trên lakehouse (MinIO)

Bao gồm các assets:

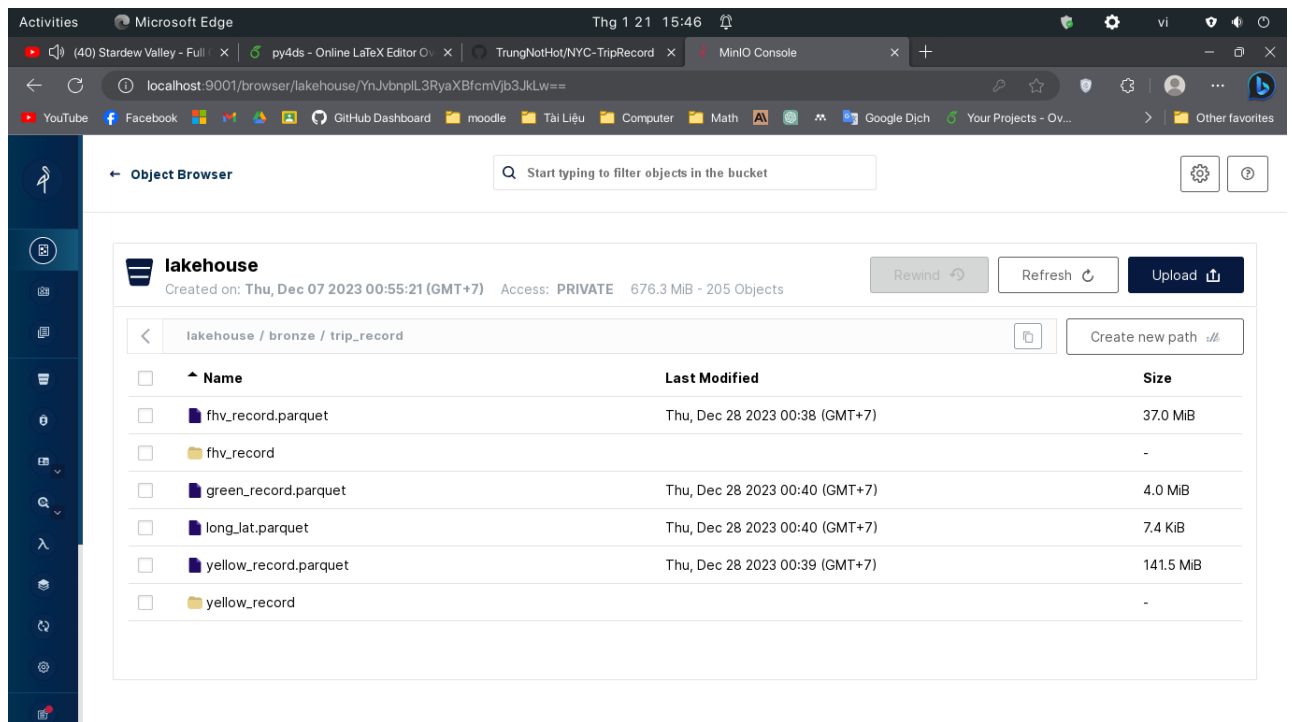
- `bronze_yellow_record`: là table `yellow_record` từ MySQL, bởi vì data quá lớn với số lượng dòng khoảng 3.5 triệu dòng, nên asset này sẽ được thực hiện partition ra theo tuần
- `bronze_green_record`: là table `green_record` từ MySQL, chỉ khoảng 65000 dòng
- `bronze_fhv_record`: là table `fhv_record` từ MySQL, khoảng 1000000 dòng, nên cũng sẽ partition tương tự `bronze_yellow_record`
- `bronze_long_lat`: chứa data về các kinh độ vĩ độ ở khu vực các taxi có thể đến

```

19
20 @asset(
21     name="bronze_yellow_record",
22     description="record of yellow taxi",
23     io_manager_key="minio_io_manager",
24     required_resource_keys={"mysql_io_manager"},
25     key_prefix=["bronze", "trip_record"],
26     compute_kind="MySQL",
27     group_name="bronze",
28     partitions_def=WEEKLY,
29 )
30 def bronze_yellow_record(context) → Output[pl.DataFrame]:
31     query = "SELECT * FROM yellow_record"
32     try:
33         partition = context.asset_partition_key_for_output()
34         partition_by = "tpep_pickup_datetime"
35         query += f" WHERE DATE({partition_by}) ≥ '{partition}' AND DATE({partition_by}) < DATE_ADD('{partition}', INTERVAL 1 WEEK);"
36         context.log.info(f"Partition by {partition_by}: {partition} to 1 week later")
37     except Exception:
38         context.log.info("No partition key found")
39     df_data = context.resources.mysql_io_manager.extract_data(query)
40     context.log.info(f"Table extracted with shape: {df_data.shape}")
41
42     return Output(
43         df_data,
44         metadata={
45             "table": "yellow_record",
46             "row_count": df_data.shape[0],
47             "column_count": df_data.shape[1],
48             "columns": str(df_data.columns),
49         },
50     )

```

Picture 2.1 1 asset ví dụ trong bronze

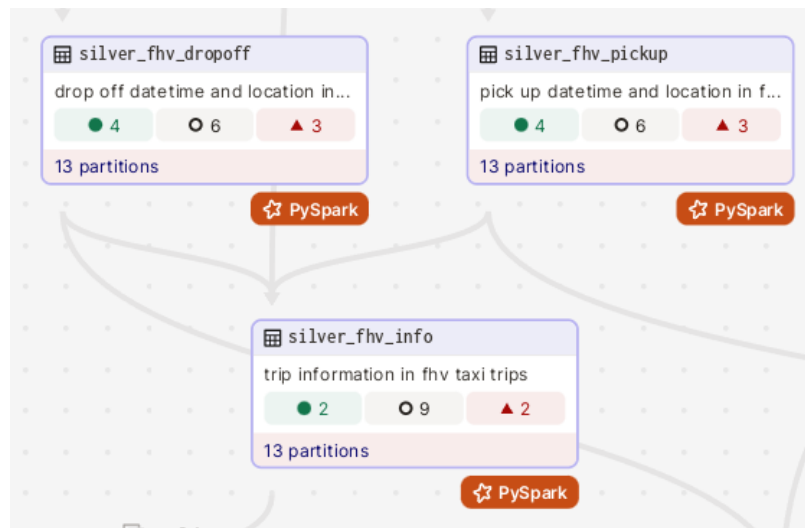


Picture 2.2 folder bronze trong minio sau khi chạy xong

• silver_layer

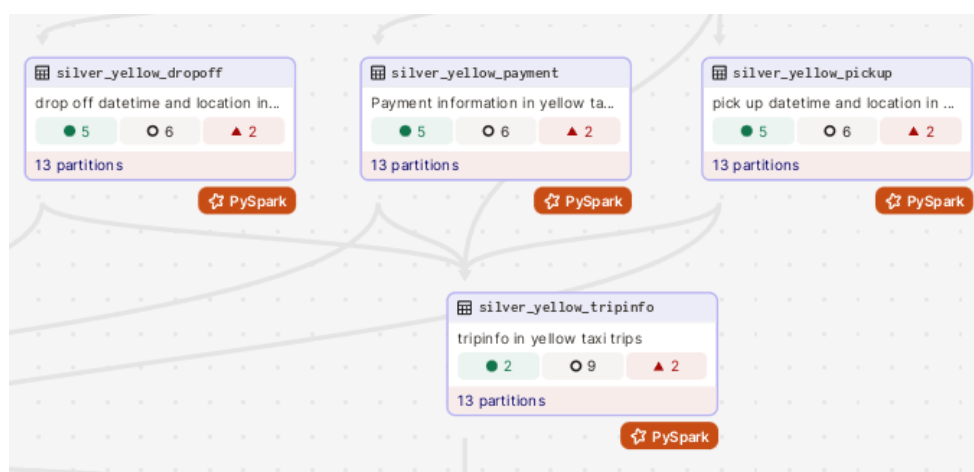
Silver FHV: dùng **spark_io_manager** làm **io manager key**, tiếp đó thực hiện data cleaning và transform từ upstream asset là `bronze_fhv_record`, tại các assets này thực hiện tách `bronze_fhv_record` ra thành các asset:

- `silver_fhv_dropoff`: chứa các thông tin về thời gian và địa điểm khách xuống xe của các fhv taxi
- `silver_fhv_pickup`: tương tự là các thông tin về thời gian địa điểm đón khách
- `silver_fhv_info`: là thông tin về chuyến đi như quãng đường, số hành khác, ...



Silver yellow: dùng **spark_io_manager** làm **io manager key**, data cleaning và transform từ upstream asset là **bronze_yellow_record**, tại các assets này thực hiện tách **bronze_yellow_record** ra thành các asset:

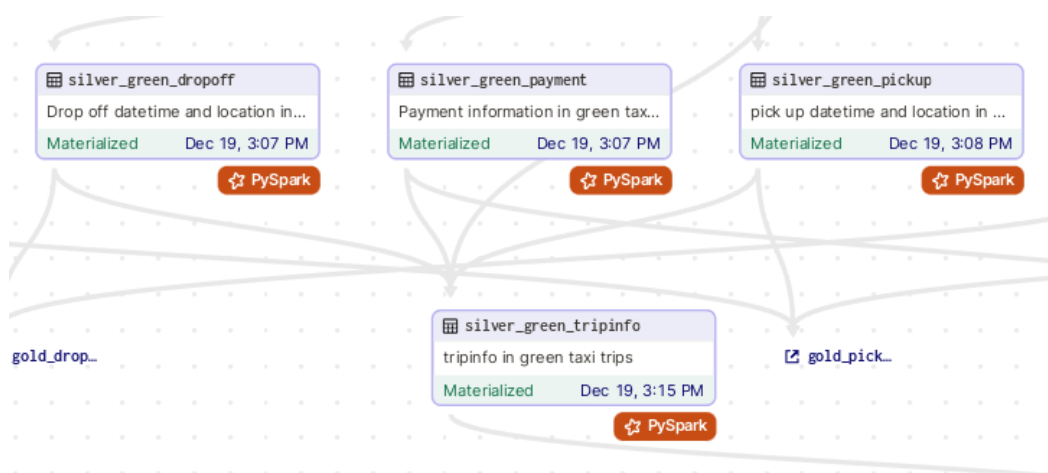
- **silver_yellow_dropoff:** chứa các thông tin về thời gian và địa điểm khách xuống xe của các taxi
- **silver_yellow_pickup:** tương tự là các thông tin về thời gian địa điểm đón khách
- **silver_yellow_payment:** Là thông tin về các chi phí như là phí xe, phụ phí, tips, ...
- **silver_yellow_info:** là thông tin về chuyến đi như quãng đường, số hành khách, ...



Silver green: dùng **spark_io_manager** làm **io manager key**, data cleaning và transform từ upstream asset là **bronze_green_record**, tại các assets này thực hiện tách **bronze_green_record** ra thành các asset:

- silver_green_dropoff: chứa các thông tin về thời gian và địa điểm khách xuống xe của các taxi
- silver_green_pickup: tương tự là các thông tin về thời gian địa điểm đón khách
- silver_green_payment: Là thông tin về các chi phí như là phí xe, phụ phí, tips, ...
- silver_green_info: là thông tin về chuyến đi như quãng đường, số hành khách, ...

Bên trong mỗi asset thực hiện cài đặt ID thích hợp, khử trùng lặp, xóa bớt null và các giá trị chưa thích hợp.



Và sau khi xử lý ta có được

Yellow

- silver_yellow_pickup: PickupID, Pickup_datetime, PULocationID
- silver_yellow_dropoff: DropOffID, Dropoff_datetime, DOLocationID
- silver_yellow_payment: PaymentID, Fare_amount, MTA_tax, Improvement_surcharge, Payment_type, RateCodeID, Extra, Tip_amount, Tolls_amount, Total_amount, Congestion_Surcharge, Airport_fee
- silver_yellow_tripinfo: VendorID, PickupID, DropOffID, PaymentID, Passenger_count, Trip_distance, Store_and_fwd_flag

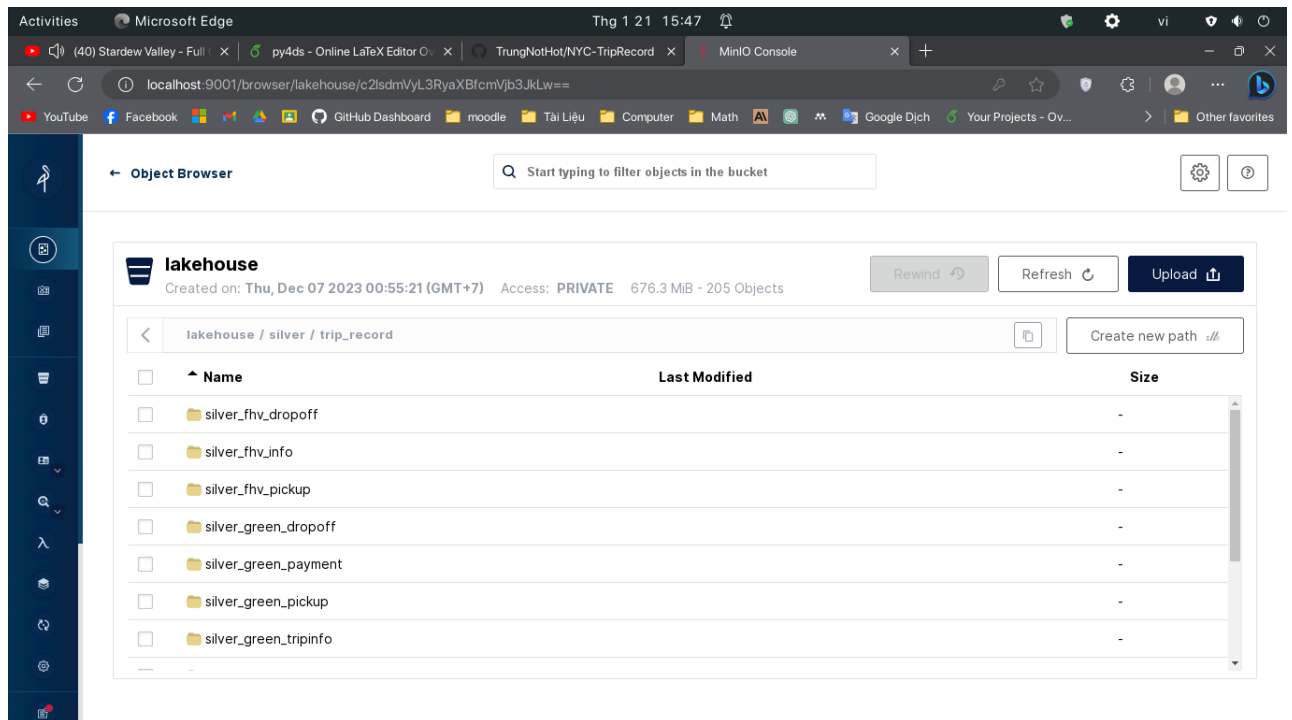
Green

- silver_green_pickup: PickupID, Pickup_datetime, PULocation
- silver_green_dropoff: DropOffID, Dropoff_datetime, DOLocation
- silver_green_payment: PaymentID, Fare_amount, MTA_tax, Improvement_surcharge, Payment_type, RateCodeID, Extra, Tip_amount, Tolls_amount, Total_amount

- silver_green_tripinfo: VendorID, PickUpID, DropOffID, PaymentID, Passenger_count, Trip_distance, Store_and_fwd_flag, Trip_type

Fhv

- silver_fhv_pickup: PickUpID, Pickup_datetime, PULocationID
 - silver_fhv_dropoff: DropOffID, Dropoff_datetime, DOLocationID
 - silver_fhvinfo: PickUpID, DropOffID, Dispatch_base_num, SR_Flag, Affiliated_base_number



Picture 2.3 folder silver trong minio sau khi chạy xong

• gold_layer

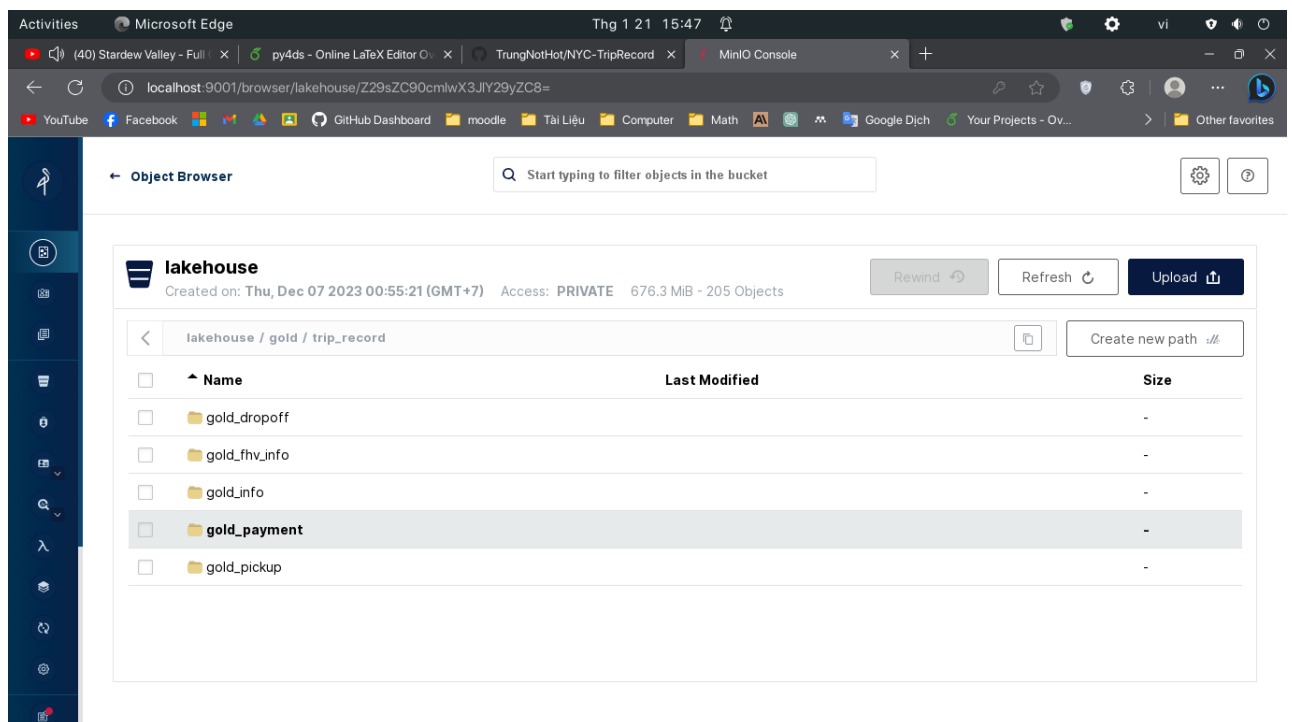


Gold layer dùng **spark_io_manager** làm **io manager key**, tiếp đó thực hiện transform từ upstream asset

Bao gồm các assets được gộp và transform từ các upstream như sau:

- gold_pickup = silver_yellow_pickup + silver_green_pickup + silver_fhv_pickup + bronze_long_lat

- $\text{gold_dropoff} = \text{silver_yellow_dropoff} + \text{silver_green_dropoff} + \text{silver_fhv_dropoff} + \text{bronze_long_lat}$
- $\text{gold_payment} = \text{silver_yellow_payment} + \text{silver_green_payment}$
- $\text{gold_tripinfo} = \text{silver_yellow_tripinfo} + \text{silver_green_tripinfo}$: vì cấu trúc table của `silver_yellow_tripinfo`, `silver_green_tripinfo` khác với `silver_fhv_tripinfo` nên nhóm chỉ thực hiện gom của yellow và green và giữ nguyên `fhv_info`
- $\text{gold_info} = \text{silver_fhv_tripinfo}$



Picture 2.4 folder gold trong minio sau khi chạy xong

- **warehouse_layer**

Tại layer này sử dụng **psql_io_manager** làm **io manager key**, và thực hiện lấy data của gold layer đưa vào psql

```

4
5 @asset(
6     name = 'warehouse_pickup',
7     description="Load gold_pickup data from spark to postgres",
8     ins={
9         "gold_pickup": AssetIn(
10             key_prefix=["gold", "trip_record"],
11         ),
12     },
13     metadata={
14         "primary_keys": ["PickUpID"],
15         "columns": ["PickUpID", "PULocationID", "pickup_datetime", "longitude", "latitude"],
16     },
17     io_manager_key="psql_io_manager",
18     key_prefix=["warehouse"],
19     compute_kind="Postgres",
20     group_name="warehouse",
21 )
22 def warehouse_pickup(context, gold_pickup: DataFrame):
23
24     context.log.info("Got spark DataFrame, loading to postgres")
25
26     return Output(
27         gold_pickup,
28         metadata={
29             "database": "trip_record",
30             "schema": "warehouse",
31             "table": "warehouse_pickup",
32             "primary_keys": ["PickUpID"],
33             "row_count": gold_pickup.count(),
34         },
35     )
36

```

Picture 2.5 1 asset ví dụ trong warehouse layer

Dữ liệu từ gold_layer sẽ được tải lên data warehouse (Postgresql) để lưu trữ và sử dụng sau khi đã được xử lý trong pipeline

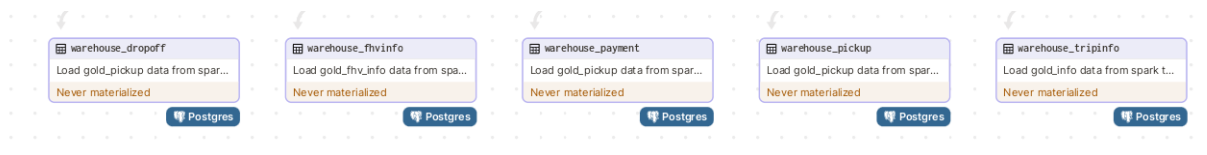
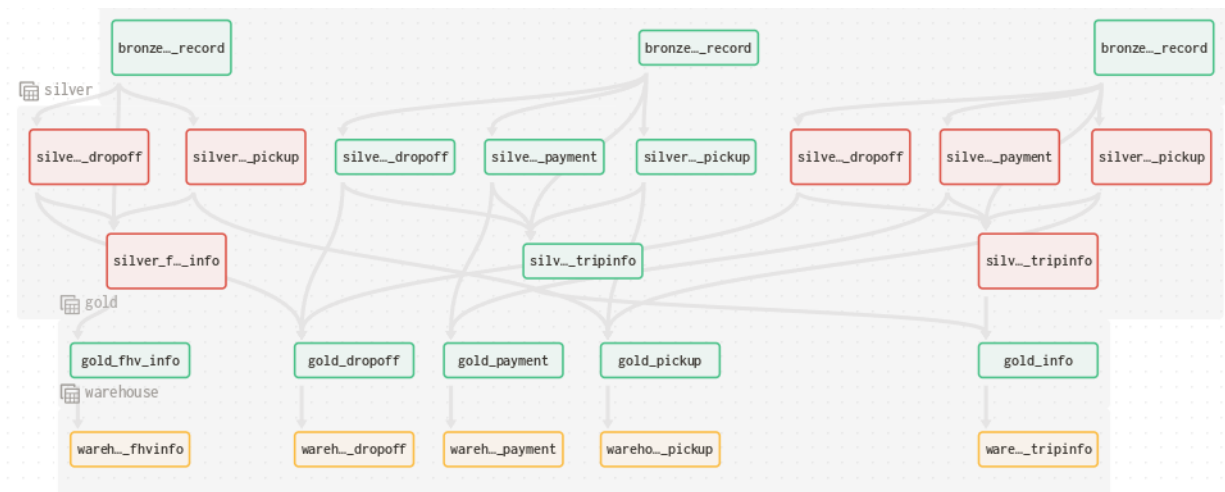


Table có được trong Postgresql sau khi đi qua pipeline

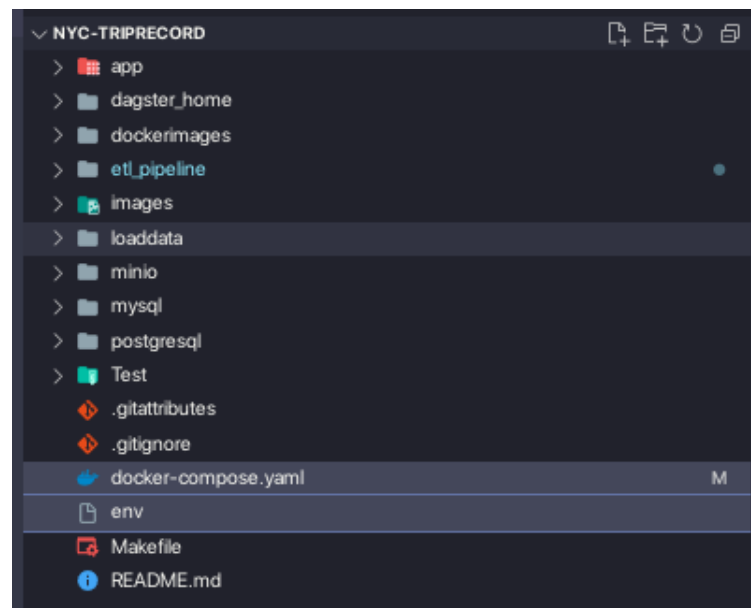
warehouse_fhv_info PickUpID DropOffID dispatching_base_num affiliated_base_number sr_flag	warehouse_pickup pickup_datetime PULocationID PickUpID longitude latitude	warehouse_dropoff Dropoff_datetime DOLocationID DropOffID LocationID longitude latitude	warehouse_info VendorID PickUpID DropOffID PaymentID passenger_count trip_distance store_and_fwd_flag trip_type	warehouse_payment fare_amount mta_tax improvement_surcharge payment_type RatecodeID extra tip_amount tolls_amount total_amount congestion_surcharge airport_fee PaymentID
---	---	--	--	--

- Data Lineage

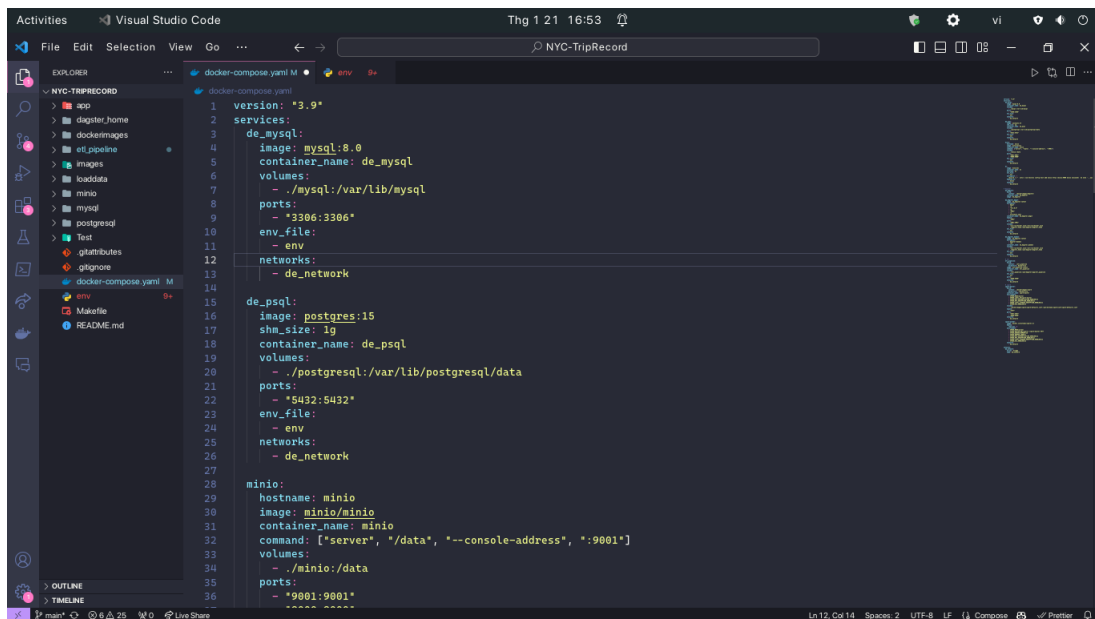


3. CÀI ĐẶT

3.1 Ngoài thư mục lớn NYC-TRIPRECORD

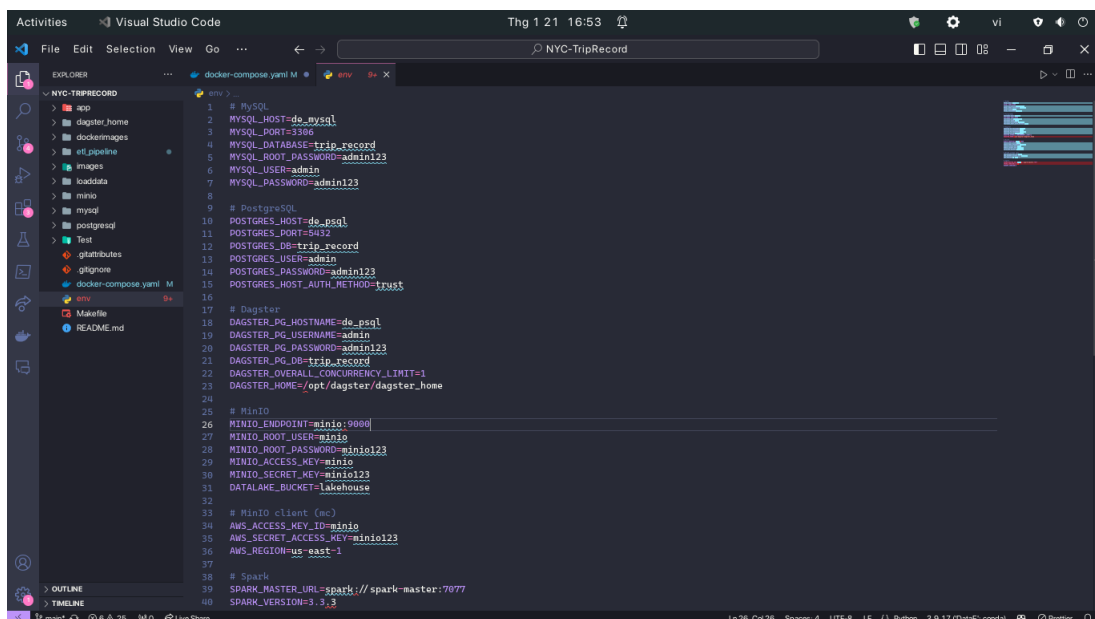


Tại 2 file env và docker-compose như hình trên



Picture 3.6 <https://github.com/TrungNotHot/NYC-TripRecord/blob/main/docker-compose.yml>

Sử dụng docker để tạo các docker container sử dụng cho project



Picture 3.7 <https://github.com/TrungNotHot/NYC-TripRecord/blob/main/env>

Tạo file env để chứa các biến môi trường như tên database, username, password ...

3.2 Trong thư mục etl_pipeline



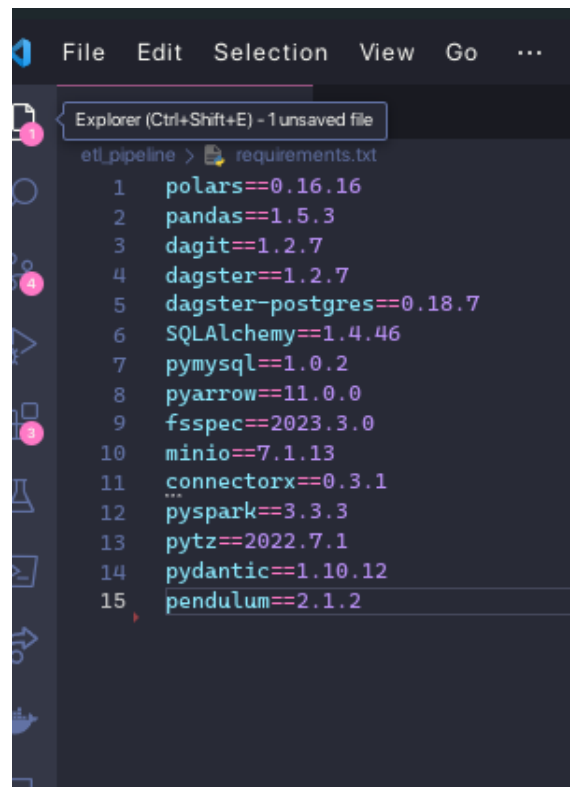
Tạo file requirements.txt để cài các thư viện python cần thiết cho etl_pipeline, và 1 Dockerfile

```

1 FROM python:3.9.16-slim
2
3 # Install spark and java
4 ARG openjdk_version="17"
5
6 RUN apt-get update --yes && \
7     apt-get install --yes curl "openjdk-${openjdk_version}-jre-headless" ca-certificates-java procps && \
8     apt-get clean && rm -rf /var/lib/apt/lists/*
9
10 # Download spark necessary jars
11 RUN curl -O https://d1cdn.apache.org/spark/spark-3.3.3/spark-3.3.3-bin-hadoop3.tgz \
12     && tar xzvf spark-3.3.3-bin-hadoop3.tgz \
13     && rm -rf spark-3.3.3-bin-hadoop3.tgz \
14     && mv spark-3.3.3-bin-hadoop3 /usr/local/ \
15     && rm -rf /usr/local/spark \
16     && ln -s /usr/local/spark-3.3.3-bin-hadoop3 /usr/local/spark
17
18 # Download PostgreSQL JDBC driver and add it to the Spark classpath
19 RUN curl -O https://jdbc.postgresql.org/download/postgresql-42.7.1.jar \
20     && mv postgresql-42.7.1.jar /usr/local/spark/jars/
21
22 RUN curl -O https://repo1.maven.org/maven2/software/amazon/awssdk/s3/2.18.41/s3-2.18.41.jar \
23     && curl -O https://repo1.maven.org/maven2/com/amazonaws/aws-java-sdk-1.12.367/aws-java-sdk-1.12.367.jar \
24     && curl -O https://repo1.maven.org/maven2/com/amazonaws/aws-java-sdk-bundle-1.11.1026/aws-java-sdk-bundle-1.11.1026.jar \
25     && curl -O https://repo1.maven.org/maven2/org/apache/hadoop/hadoop-aws/3.3.2/hadoop-aws-3.3.2.jar \
26     && curl -O https://repo1.maven.org/maven2/org/apache/hadoop/hadoop-aws/3.3.2/hadoop-aws-3.3.2.jar \
27     && mv s3-2.18.41.jar /usr/local/spark/jars \
28     && mv aws-java-sdk-1.12.367.jar /usr/local/spark/jars \
29     && mv aws-java-sdk-bundle-1.11.1026.jar /usr/local/spark/jars \
30     && mv delta-core-2.12-2.2.0.jar /usr/local/spark/jars \
31     && mv delta-storage-2.2.0.jar /usr/local/spark/jars \
32     && mv hadoop-aws-3.3.2.jar /usr/local/spark/jars
33
34 # Add required environment variables
35 WORKDIR /opt/dagster/app/etl_pipeline
36 COPY requirements.txt /opt/dagster/app/etl_pipeline
37 RUN pip install --upgrade pip && pip install --default-timeout=300
38 COPY . /opt/dagster/app/etl_pipeline

```

Picture 3.8 https://github.com/TrungNotHot/NYC-TripRecord/blob/main/etl_pipeline/Dockerfile



```
File Edit Selection View Go ...
Explorer (Ctrl+Shift+E) - 1 unsaved file
etl_pipeline > requirements.txt
1 polars==0.16.16
2 pandas==1.5.3
3 dagit==1.2.7
4 dagster==1.2.7
5 dagster-postgres==0.18.7
6 SQLAlchemy==1.4.46
7 pymysql==1.0.2
8 pyarrow==11.0.0
9 fsspec==2023.3.0
10 minio==7.1.13
11 connectorx==0.3.1
12 pyspark==3.3.3
13 pytz==2022.7.1
14 pydantic==1.10.12
15 pendulum==2.1.2
```

Picture 3.9 Các thư viện python được cài trong etl_pipeline

4. DEMO KẾT QUẢ VÀ VISUALIZE

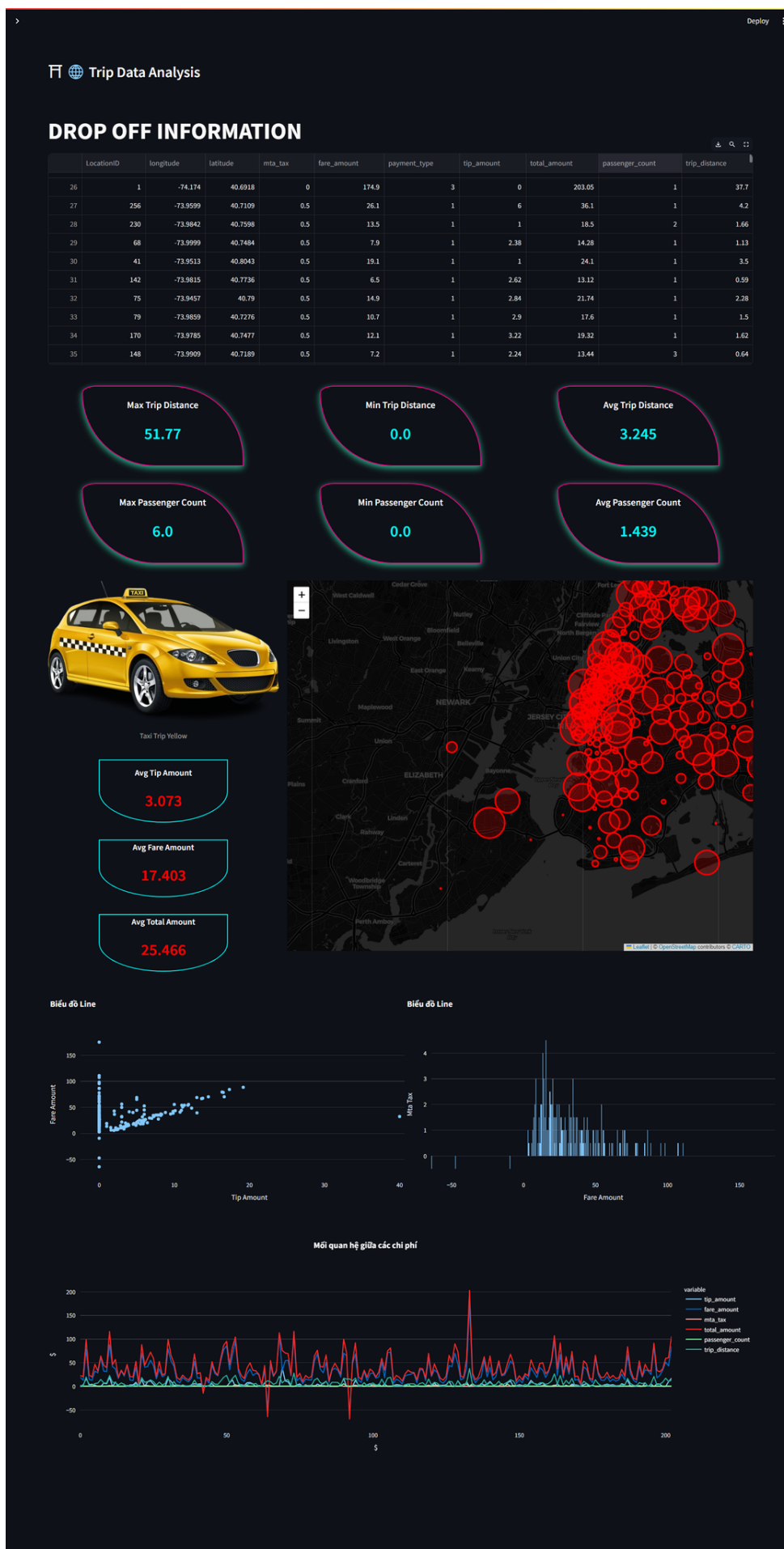
4.1 DEMO

Demo link: https://studenthcmusedu-my.sharepoint.com/personal/21280012_student_hcmus_edu_vn/_layouts/15/stream.aspx?id=%2Fpersonal%2F21280012%5Fstudent%5Fhcmus%5Fedu%5Fvn%2FDocuments%2FDE%2Fdemo%2Emp4&referrer=StreamWeb2Web&referrerScenario=AddressBarCopied%2Eview

4.2 VISUALIZE

Sử dụng **Streamlit** Library để visualize các dữ liệu

Dữ liệu được visualize được connect đến Postgres để lấy dữ liệu và chuyển đổi xử lý theo như mong muốn



TÀI LIỆU THAM KHẢO

- [1] Tài liệu Python cho KHD thầy Hà Văn Thảo
- [2] PySpark Overview
- [3] Dagster doc
- [4] <https://github.com/vsonwork/yellow-taxi-trip-ETL>
- [5] <https://github.com/lelouvincx/goodreads-elt-pipeline>
- [6] Docker doc