

# **BILTEGIRATUTAKO DATUEN PROZESATZEA**

**Manex Sorarrain eta Aner Garatea**

## AURKIBIDEA

- 1.Sarrera
  - ◆ Proiektuaren deskribapena
- 2.Proiektuaren garapena
- 3.Emailzak
- 4.Gehigarriak
- 5.Ondorioak

## SARRERA

Lan honetan, Docker Compose bidez Kudu+Impala ingurune bat hedatuko dugu eta biltegitratutako datuak Bussiness Intelligence tresna baten bidez bistaratuko ditugu

Kudu abiadura handiko datu-biltegitratze sistema bat da, datu multzo handiekin denbora errealean lan egiteko diseinatua.

Bestalde, Impala errendimendu handiko datuak aztertze tresna bat da, Hadoop-en gordetako datuetan SQL kontsulta interaktiboak egiteko aukera ematen duena.

Kudu eta Impala batera erabil daitezke Hadoop-en gordetako datu multzo handietan denbora errealean datuen analisisa egiteko. Kudu-k datuen biltegi azkar eta eraginkor gisa jotzen du, Impala-k, berriz, SQL interfaze interaktibo bat eskaintzen du datuak kontsultatzeko eta aztertze. Elkarrekin, Kudu eta Impala-k datu handien analisirako irtenbide osoa eta eskalagarria eskaintzen dute.

Bestalde, Bussiness Intelligence-ren helburu nagusia datuak informazio baliagarri bihurtzea da, enpresek gertakari eta datu errealean oinarrituta erabakiak har ditzaten. Lan honetan, BI tresna bezala, Superset plataforma erabili dugu.

Datubasea, New York-eko taxiek egindako bidaiak biltzen dituen Parquet datubase bat da., <https://www.nyc.gov/site/tlc/about/tlc-trip-record-data.page> hain zuzen ere.

.

## PROIEKTUAREN GARAPENA

Lanaren oinarritzko azpiegitura honako elementuez osatuta dago:

- Kudu azpiegitura bat, docker bidez hedatutakoa.
- Impala irudi bat, kudura konektatutakoa.
- Superset irudi bat.

Azpiegiturak ongi funtziona dezan hurrengo pausuak jarraitu ditugu:

1. Parquet fitxategia deskargatu eta data karpeta batean gorde.
  - a. `wget`  
  
`https://s3.amazonaws.com/nyc-tlc/trip+data/yellow_tripdata_2022-01.parquet`
2. Enuntziatuko git clone aginduarekin errepositorioa jaitsi. Honen barruan **quickstart.yml** fitxategia aldatu dugu, honetan impala kontainerren konfigurazioa gehitu dugu, eta egitura guztia sare berean egon dadin, sare bat sortu dugu bridge motakoa. Behin hori eginda, sare hau esleitu diogu irudi bakoitzari. Impalaren konfigurazioan data bolumena irudi barruko data karpetara ere mapatu dugu.
3. Kudu eta impala kontainerrak jaurti
  - a. **`docker compose -f quickstart.yml up -d`** komandoa erabiliz
4. Ondoren, Impala kanpo taula bat sortu dugu Parquet datuetara seinalatzen duena, gainera, kudu taula bat ere sortu dugu Impalan aurrekoaren azpi-aula bezala. Hau dena impala-shell ean exekutatu dugu.
  - a. **`create external table taxi like parquet 'file:///data/taxi.parquet' stored as parquet location 'file:///data/';`**
    - "taxi" izeneko Impala-n kanpo taula bat sortzen du, 'file:///data/taxi/taxi.parquet' helbidean dagoen Parquet fitxategiaren egituraren itxura duena. Taulako datuak 'file:///data/taxi/' direktorioan gordeko dira, Parquet formatuan
  - b. **`create table taxik (pickup timestamp, dropoff timestamp, distance double, primary key(pickup,dropoff)) stored as kudu;`**
    - Impalan "taxik" izeneko taula bat sortzen du, hiru zutabe dituen: timestamp motako "pickup" eta "dropoff" eta double motako "distance". Gainera, taulak gako nagusi bat du "pickup" eta "dropoff" zutabeek osatutakoa. Hau, Kudu taula gisa gordeko da.
  - c. **`insert into taxik select tpep_pickup_datetime, tpep_dropoff_datetime, trip_distance from taxi;`**

- Azkenik, Parquet taulako datuak Kudu taulan txertatu ditugu goiko komandoa erabiliz.

- Behin Impalara konektatuta, proba-kontsulta bat exekutatu dezakegu honela:
  - SELECT COUNT(\*) FROM taxik;**
  - Dena behar bezala konfiguratuta badago, erantzun bat jaso beharko genuke taxik taulako errenkada kopuruarekin.
- Superseten konfiguraziorako, klaseko tutoriala jarraitu dugu, honako bi salbuespenekin: requirements-local.txt fitxategian impyla, thrift eta thrift\_sasl dependentziak jarri ditugu, eta lehen aipatutako sare bera konfiguratuta dugu, dena sare berean egon dadin.
- Ondoren, Superset-en datuen jatorria sortu beharko da. Hau egiteko, Superset-en interfazetik "Origen de datos"-en datu base mota bezala Impala jarri beharko dugu eta Datu baserako konexio informazioa sartu.

Gure entregagarria jaurtitzeko,

**export KUDU\_QUICKSTART\_IP=\$(ifconfig | grep "inet " | grep -Fv 127.0.0.1 | awk '{print \$2}' | tail -1)** komandoa exekutatu dugu KUDU\_QUICKSTART\_IP ingurune aldagaia definitzeko. Ondoren, beheko komandoarekin arkitektura guztia martxan jartzen dugu..  
**'docker compose -f quickstart -f docker-compose.yml up'**

Behin guztia martxan jarrita, kudun taxik taula agertzen da, impala eta superset arteko konexioa ondo dago, impalan ikusten dira superseteko queriak, eta gainera, superseten taxi eta taxik taulen gainean kontsultak egin ditzakegu. Hurrengo argazkietan ditugu hau baieztatzen duten kapturak:

Tables

There are 1 tables.

Table Name	Table Id	Owner
impala::default.taxik	29b05ee89687466e9b4e7057f2b7d623	impala

Superset

Databases

Database	Backend	AQE
Apache Impala	impala	x

1 select \* from taxi limit 10

RUN LIMIT: 1 000 00:00:02.20

RESULTS QUERY HISTORY

10 rows returned The number of rows displayed is limited to 10 by the query

vendorid	tpep_pickup_datetime	tpep_dropoff_datetime	passenger_count	trip_distance
1	2022-02-01T00:06:58	2022-02-01T00:19:24	1	5.4
1	2022-02-01T00:38:22	2022-02-01T00:55:55	1	6.4
1	2022-02-01T00:03:20	2022-02-01T00:26:59	1	12.5

1 select \* from taxik limit 10

RUN LIMIT: 1 000 00:00:00.22

RESULTS QUERY HISTORY

10 rows returned The number of rows displayed is limited to 10 by the query

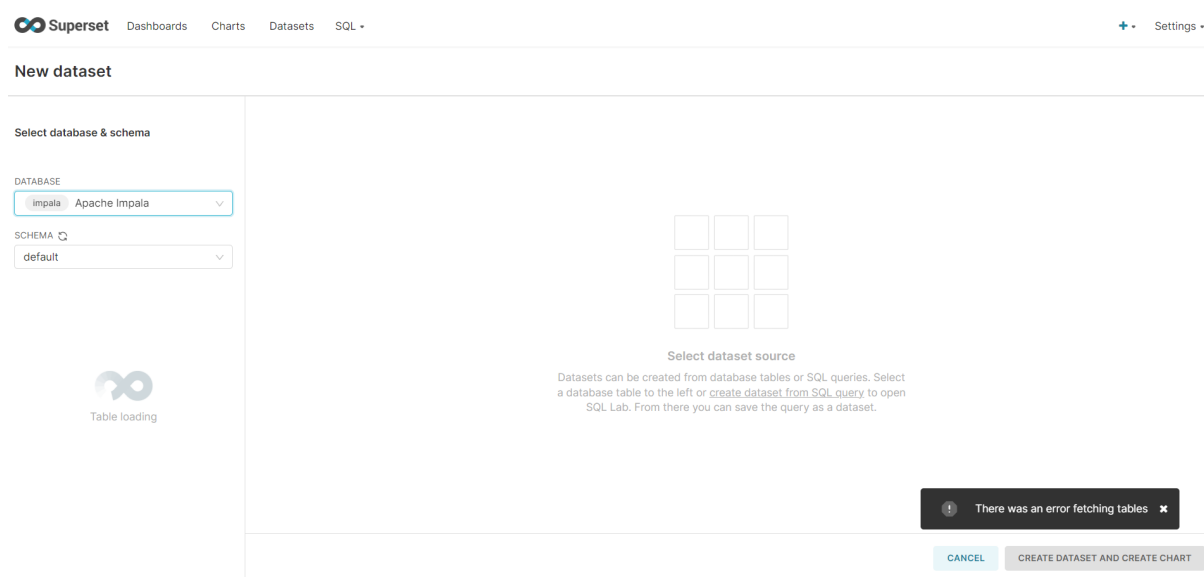
pickup	dropoff	distance
2003-01-01T00:10:06	2003-01-01T12:38:59	0
2003-01-05T04:05:18	2003-01-05T18:41:11	0
2008-12-31T23:02:01	2009-01-01T17:24:13	0.95

## Last 25 Completed Queries

User	Default Db	Statement	Query Type	Start Time	End Time	Duration	Scan Progress	State	# rows fetched	Resource Pool	Details
root	default	select * from taxik limit 11	QUERY	2023-05-01 21:28:01.206515000	2023-05-01 21:28:01.271550000	65.035ms	1 / 1 (100%)	FINISHED	11	default-pool	<a href="#">Details</a>
root	default	select * from taxi limit 11	QUERY	2023-05-01 21:24:58.111617000	2023-05-01 21:25:00.070570000	1s958ms	1 / 1 (100%)	FINISHED	11	default-pool	<a href="#">Details</a>

Baina grafikak sortzerakoan, honako arazoa sortu zaigu: Grafikak sortzeko taxik taulatik abiatuta dataset bat sortu behar genuen, baina hau sortzerakoan, honako mezua lortu eta ezinezko egiten zen hau sortzea.

Hemen dataseta sortu ezin dela erakusten duen kaptura:



Hainbat saiakera egin ondoren, ez dugu lortu hau konpontzea, beraz .parquet fitxategia eskuz igotzea erabaki dugu impalara 'upload columnar file to database' aginduarekin eta hemendik grafikak sortzeko asmoarekin.

Uste genuen horrela lortuko genuela baina, ezin izan dugu probatu izan ere, datuak Impalara eskuz igotzerakoan ondorengo errorea ematen zigun terminalean.

```
superset_app      | File "/app/superset/models/core.py", line 937, in _has_view
superset_app      | view_names = dialect.get_view_names(connection=conn,
superset_app      | schema=schema)
superset_app      | File
superset_app      | "/usr/local/lib/python3.8/site-packages/sqlalchemy/engine/interfaces.py", line 332, in
superset_app      | get_view_names
superset_app      | raise NotImplementedError()
superset_app      | NotImplementedError
superset_app      | 2023-05-01 21:13:28,358:WARNING:superset.models.core:Has view
superset_app      | failed
superset_app      | Traceback (most recent call last):
```

```

superset_app      | File "/app/superset/models/core.py", line 937, in _has_view
superset_app      | view_names = dialect.get_view_names(connection=conn,
schema=schema)
superset_app      | File
"/usr/local/lib/python3.8/site-packages/sqlalchemy/engine/interfaces.py", line 332, in
get_view_names
superset_app      | raise NotImplementedError()
superset_app      | NotImplementedError

```

Ez gara gai izan hau konpontzeko, hainbat aldiz probatu dugu datuak igotzea baina denbora askoz igotzen egon eta gero ,behin eta berriz errore berdinarekin topatu gara..

Dena dela, jadanik pentsaturik geneuzkan sortu beharreko grafikoak zeintzuk izango ziren eta nahiz eta ezin erakutsi, behintzat aipatu egingo ditugu.

- **Barra-diagrama:** barra-diagrama bat eraikitzea asteko egunaren edo eguneko orduaren arabera bidaia kopurua erakusteko.
- **Lerro-diagrama:** lerro-diagrama bidez denboran zehar egindako bidaia kopurua erakustea.
- **Bero-mapa:** bero-mapa bat sortzea hirian bidaia gehien egiten diren lekuak erakusteko.
- **Diagrama zirkularra:** diagrama zirkularra erabiltzea ibilgailu motaren edo ordainketa motaren arabera bidaien banaketa erakusteko.
- **Sakabanatze grafikoa:** sakabanaketa grafikoa erabilgarria izan daiteke bidaia-denboraren eta egindako distantziaren arteko erlazioa erakusteko.
- **Kutxa eta bibote grafikoa:** Kutxa eta bibote grafikoa erabilgarria izan daiteke bidaia-denboraren edo egindako distantziaren banaketa erakusteko.
- **Histograma:** Bidaien iraupenen banaketa adierazten duen histograma.

## GEHIGARRIAK

Behin oinarrizko atazak burututa, gehigarri batzuk egiten ahalegindu gara:

### Azureko hedapena

Azureko hedapena egiteko, honako pausoak eman ditugu: Lehenik eta behin, 'grupo de recursos' bat sortu dugu honako komandoarekin:

```
'az group create --name lana2 --location northeurope '
```

Ondoren, kontainer biltegi bat sortu dugu, aurreko errekurtsio taldean, komando honekin:

```
'az acr create --resource-group lana2 --name dmpabigarrenlana --sku basic'
```

Hona igoko ditugu erabilitako irudi guztiak, lehenik eta behin lehenik eta behin izena aldatuz:

```
'docker tag apache/kudu:latest dmpabigarrenlana.azurecr.io/kudu:latest'
```

eta ondoren push eginez:

```
'docker push dmpabigarrenlana.azurecr.io/kudu'
```

Hau, erabilitako irudi bakoitzarentzat egin dugu, eta hemen ez ezik, '.yaml' fitxategietan ere images atalean izenak aldatu dizkiegu, entregako azure karpetaiko '.yaml' fitxategietan ikus daitekeen bezela. Behin hau dena eginda, eta azuren login eginda, proiektua abiarazteko agindua eman dugu. Baina dirudienez, ez dugu baliabide nahikoa hau aurrera emateko, honako mezua lortu baitugu:

```
'containerinstance.ContainerGroupsClient#CreateOrUpdate: Failure sending request:
StatusCode=400 -- Original Error: Code="ContainerGroupQuotaReached"
Message="Resource type 'Microsoft.ContainerInstance/containerGroups' container group
quota 'StandardCores' exceeded in region 'northeurope'. Limit: '6', Usage: '0' Requested:
'6.01'.'"
```

Baliabide muga dela eta, ezin izan dugu gure proiektua azuren hedatu.



## ONDORIOAK

Nahiz eta pauso guztiak klasean edota tutorialetan ikusi bezala eman ditugun, oso zaila egin zaigu aurreratzea, hainbat errore desberdinekin topatu gara proiektuan zehar, gehienak soluzionatu ditugu baina superseteko bi errore horiek dirudienez honen implementazioko akatsak dira eta ez dugu ia dokumentaziorik bilatu. Hortaz, ezin izan dugu proiektua osotasunean garatu, panelaren sorkuntza egin gabe utziz. Hala ere, datu basearen sorkuntza, eta kontainerren arteko konexio guztiak ondo doaztela frogatu dugu.

Kode guztia github errepositorio honetan dago, ezin bainuen igo egelara.  
<https://github.com/manexsora/dmpa>