

16/01

① Prerequisites: Python, SQL, command line

② Dataset: NY Taxi trips dataset

③ Syllabus

→ Week 1 : Intro & Prerequisites

→ Week 2 : Workflow Orchestration

→ Week 3 : Data Warehouse

→ Week 4 : Analytics Engineering

→ Week 5 : Batch Processing

→ Week 6 : Stream processing

→ Project: Putting everything we learned in practice

④ Tech Stack

① GCP ← Google Cloud Storage : Data Lake
Bigquery: Data Warehouse

② Terraform: Infrastructure as Code

③ Docker: Containerization

④ SQL : Data analysis & exploration

⑤ prefect : Workflow orchestration

⑥ dbt : Data transformation

⑦ Spark : distributed processing

⑧ Kafka : streaming

22/01

Data Engineering Zoomcamp

Week 1 : GCP, Docker, SQL, terraform.

→ 1.4.1 : Setting env on GCP

- Compute engine → VM instances
 - Add ssh key to CE metadata (create)
 - Region, 16gb ram, 30gb space, Ubuntu
 - Ssh -i private_key user@External ip
 - curl & bash Ansible.sh (yes to all)
 - SSH Config (File in ~/.ssh)
- Host de-zoomcamp
- HostName <External ip of VM>
- User <user name>
- Identityfile <private key Path>
- ssh de-zoomcamp

• Docker installation

- snap install docker sudo apt-get update &
→ remote ssh extension in VS code
→ google "Docker without sudo" & follow VM
→ docker run hello-world
→ if external ip not displayed, check columns selected.

Install Docker Compose

- curl <Docker Compose link> -O docker-compose
→ chmod +x docker-compose
→ Add bin folder to Path in .bashrc

- P stop vs apt - get
- ④ apt vs apt - get
- ③ pip vs cargo
- ④ default sch identityfile is id_rsa
- ⑤ Users added to docker group to use docker without sudo
- ⑥ Ubuntu VM APM, AMD, X86_64
- ⑦ SFTP vs SCP
- ⑧ When VM instance is stopped after it is not charged but resources like disks & IP's are charged. They get charged as they get destroyed.

Containerization

Containerization is a technology that allows multiple applications to run on a single host machine. It achieves this by creating isolated environments (containers) that share the host's operating system and hardware resources. This results in higher efficiency, faster deployment, and better resource utilization compared to traditional virtual machines.

Types and benefits of containerization:

- **Portability:** Containers are portable across different hosts and environments, making it easier to move workloads between cloud providers or on-premises infrastructure.
- **Resource Efficiency:** By sharing the host's operating system, containers require less memory and disk space than virtual machines, leading to cost savings and improved performance.
- **Deployment Speed:** Containerized applications can be deployed much faster than traditional monolithic applications, allowing for more frequent updates and releases.
- **Scalability:** Containers can be easily scaled up or down based on demand, providing a more flexible way to manage resources.
- **Isolation:** Each container runs in its own isolated environment, which helps prevent one application from affecting another.

- ① Clone dec git repo
 - ② cd week1 basic n setup / 2 docker sql
 - ③ docker-compose up -d detached mode
Run multiple containers apps at once.
 - ④ Run containers in background
 - ⑤ This is to install Pg database & Pgadmin Services as Containers
 - ⑥ Use Cmd to install Pgcli & Pipfdart for myclis
 - ⑦ Use pip to install Pgcli
 - ⑧ forward 8080(Pgadmin port) in VS Code and go to localhost:8080 in chrome
 - ⑨ In jupyter Create Schema & insert 100 rows of yellow taxi data (download CSV) into Pg. Check using Pgcli

Terraform

Terraform

① Dockerfile

```
FROM python:3.9
RUN pip install pandas
WORKDIR /app
COPY pipeline.py pipeline.py
ENTRYPOINT ["python", "pipeline.py"]
```

② Build Image from Dockerfile

```
docker build -t test:pandas
```

If "Can't stat errors"; add folder to .dockerignore file (due to permission of folder)

③ Docker run Container

```
docker run -it test:pandas 2023
```

O/P: ['Pipeline.py', '2023']

Job finished successfully for day = 2023

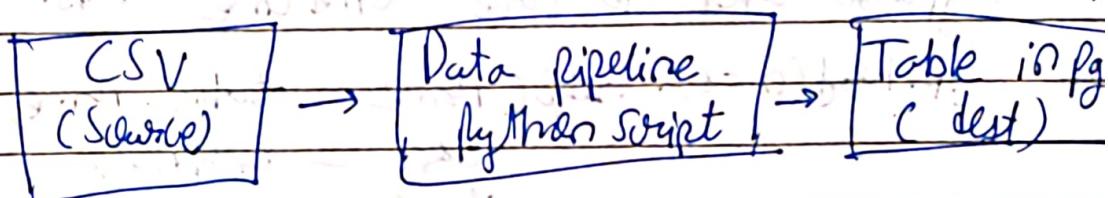
Image: • Snapshot of filesystem includes starting cmd
• just disk space no CPU/memory
• Created using "docker build Dockerfile"
• Includes everything needed to run on app

Container: • Created from images
• takes up CPU/memory
• always runs the same, regardless of where

1.2.1 \Rightarrow Introduction to Docker

Jan 23

Data pipeline



why Docker?

- ① Reproducibility
- ② Local experiments
- ③ Integration tests (CI / CD)
- ④ Running pipelines for the cloud
- ⑤ Spark & Serverless (AWS Lambda)

Docker

- provides ability to package & run an app in a loosely isolated env called a Container
- Dockerfile : text doc with instructions to build an image.
- FROM : set Base image
- RUN : execute Cmds
- WORKDIR : set working dir
- COPY : copy files / dirs
- ENTRYPOINT : configure container to

In26 1.2. 4 \Rightarrow Dockerizing the Ingestion Script

① jupyter notebook to python script

jupyter nbconvert -- to=script

② argparse used to take Cmdline args.

③ df.head(n=0).toSQL(name = table_name,
con = engine,
if_exists = "replace")

drops table before inserting new values.

append : insert new values to existing table

So this command is used to create a fresh new table before adding records to it.

④ Pass passwords using env vars/password store instead of passing it in Cmdline arg

⑤ Dockerize ingestion Script

- add wget & sqlalchemy & psycopg2 to Dockerfile apt-get -> pip install

- Build image

- Run Container

⑥ Local Server (Using python -m http.server) can be used for testing

⑦ localhost for Container would be Container's localhost. So to connect to host machine's network use --network=host

⑧ If --it flag is absent, interactivity is not possible

1.2.2 \Rightarrow Ingesting My taxi data to Postgres

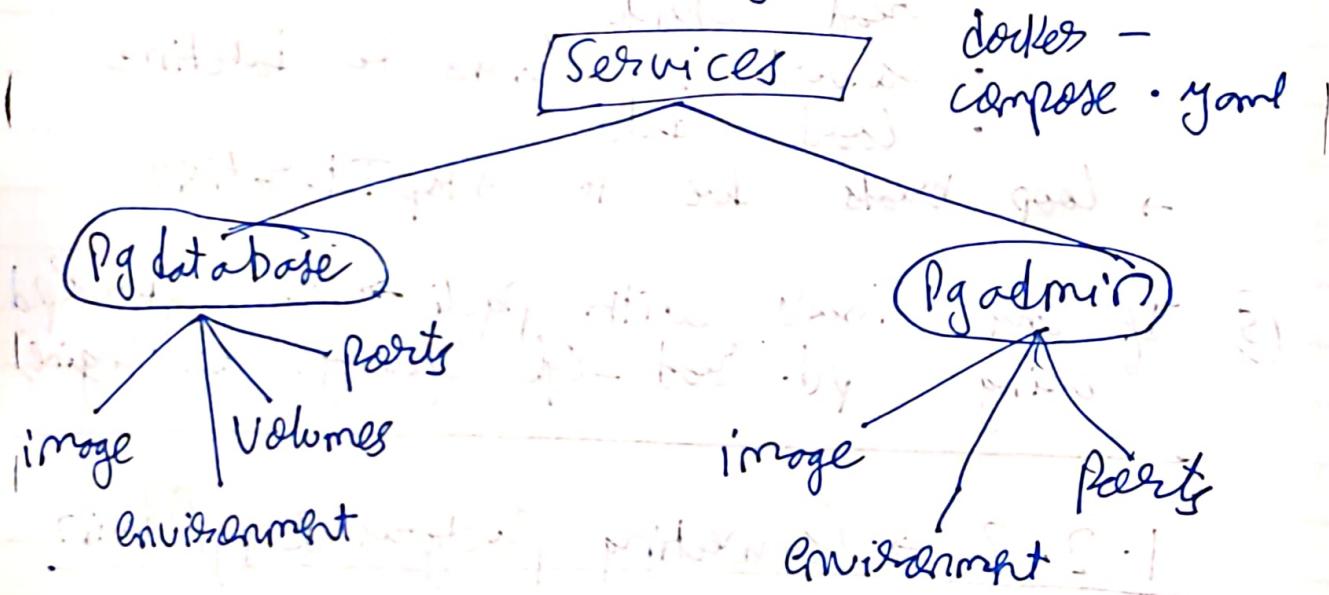
- ① postgres run as Container
- ② pgcli used to interact with db
- ③ wget CSV data
- ④ → SQLAlchemy engine created
→ CSV read in chunks by specifying
 iterator = True , chunksize = 100000
→ In a infinite while loop
 - read chunk
 - Convert 2 columns to datetime
 - load data
- loop breaks due to StopIteration
- ⑤ If any issues with Pgcli \rightarrow Use jupyter & pd
using pd.read_sql(query, con=engine)

1.2.3 \Rightarrow Connecting Postgres & PgAdmin

- ~~Forward~~ forward 8080 (Pgadmin Port) in vsc
- Create Servers on Pgadmin
- If postgres & pgadmin are individual
 Containers started using Command line, then
 a network has to be created
- If docker - Compose is being used, network
 is automatically created
- Run queries.

Ques 2.5 \Rightarrow Running PostgreSQL & PgAdmin with Docker Compose

- ① Docker-Compose is a tool for defining & running multi-container Docker applications.
- ② a YAML file is used to configure the app
- ③ Then, with a single command, all services from config can be created & started
- ④ network is automatically created



- ⑤ docker-compose up -d to start
- ⑥ docker-compose down to stop
- ⑦ docker network ls to list networks
- ⑧ by default, "<dir-name>-default" network is created for each docker Compose file

1.4.2 \Rightarrow Port mapping & Networks in Docker

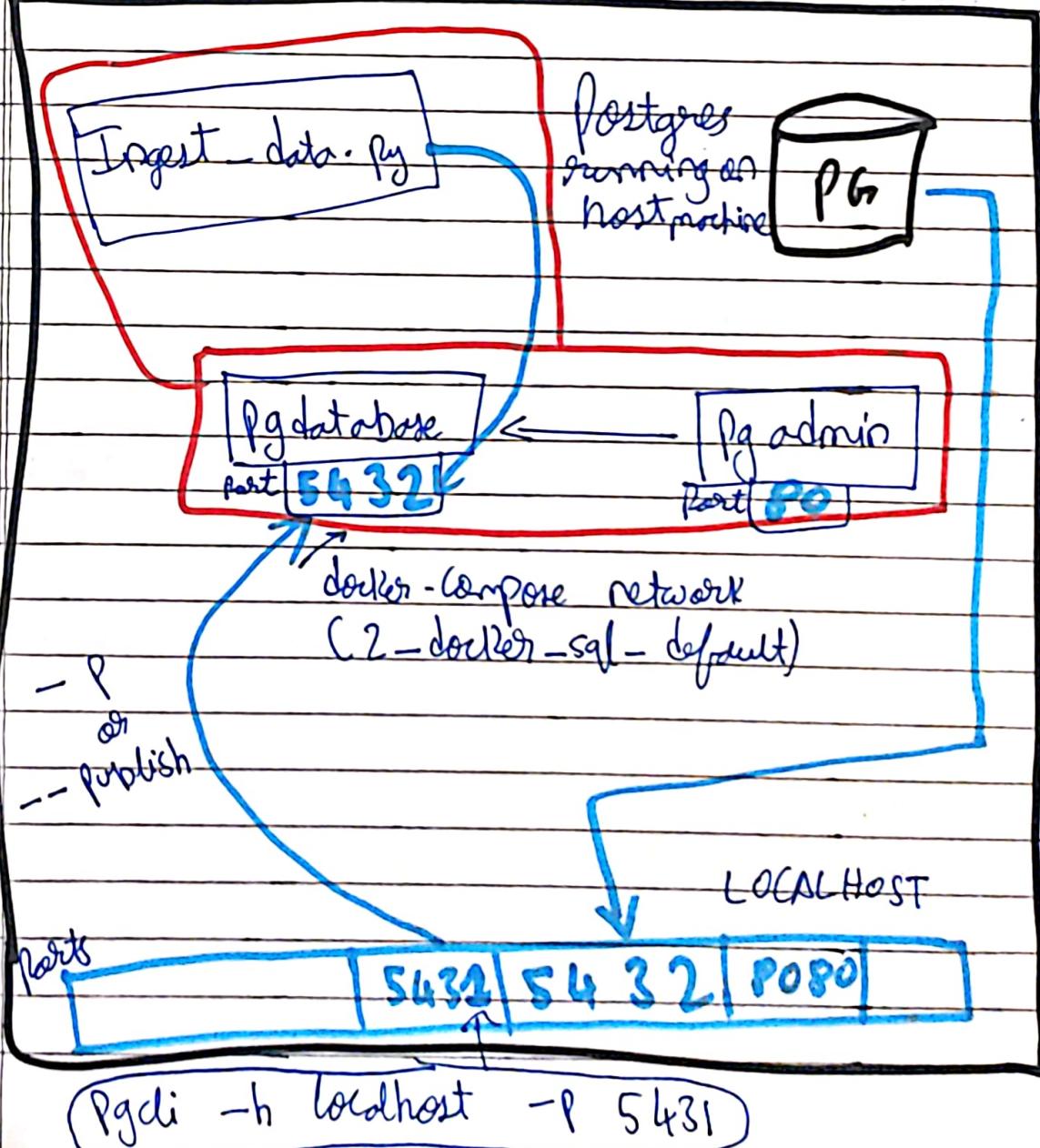
① errors checking context: Can't stat < dir-path>

This is caused as permissions are different
dir is owned by root.

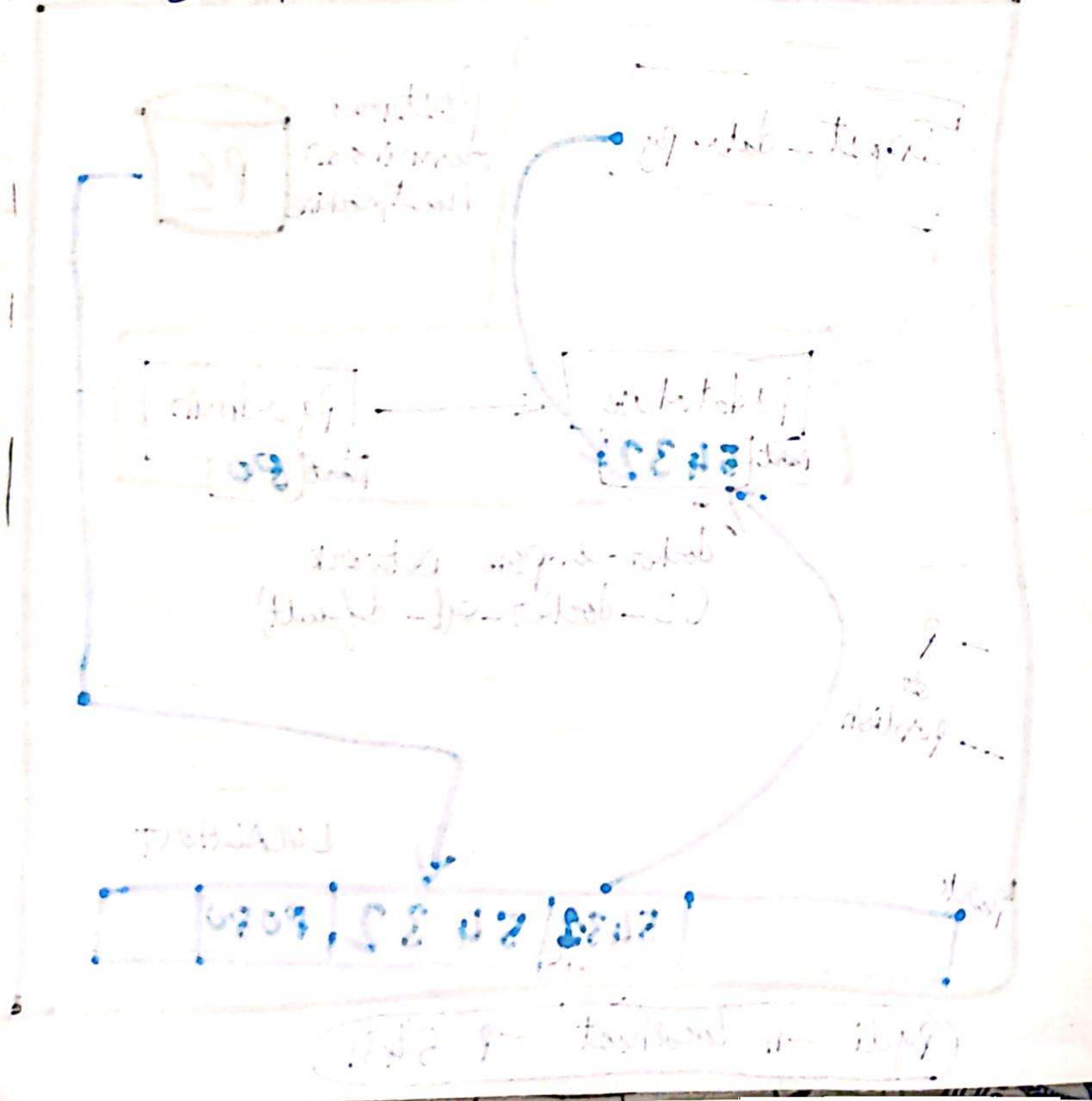
To solve it create a data folder,
move Pg folder inside data and
add data folder to .dockerignore.

② Docker networks diagram

Host Computer (VM instance w/ Ubuntu on GCP)



- On the host computer 5431 is used (Assuming there is a Pg instance running on 5432 on host) by mapping "5431 on host" → "5432 on container" (This is explicitly written in PgDockerfile)
 - So PgDocker would use 5431 to communicate with Pg inside Container which will run on 5432.
- ingestion Script goes via docker - Compose network (as it's network is specified)
- Hence it uses 5432



Jan 29

1.2.6 → SQL Refresher

- ① Zones Data in taxi_zone_lookup.csv is downloaded and added to db
- ② Implicit inner join b/w trips & zones to view zone name instead of zone id
- ③ Explicit inner join to do the same. This has the JOIN Keyword and is preferred over implicit join
- ④ JOIN clause is used to combine rows from two or more tables based on a related column between them.
- ⑤ Types of joins
 - (INNER) JOIN: Matching Values in both tables
 - LEFT JOIN: All from left, matched from right.
 - RIGHT JOIN: All from right and matched from left
 - FULL JOIN: All records when there is a match in either left or right table
- ⑥ CAST → Convert data from one type to another
- ⑦ Group By → Group rows with same values.
- ⑧ ORDER BY → Sorting data in ASC / DESC

1.1.1 \Rightarrow Intro to GCP

- ① Google Cloud Platform includes a range of hosted services for Compute, Storage and application development that run on Google hardware.
- ② Navigate to Services using left navbar or top search bar.

1.1.3.1 \Rightarrow Intro to terraform

- ① Open Source tool by Hashicorp
- ② Used for provisioning infra resources
- ③ IaC : Infrastructure as code \Rightarrow

Build
change
Manage \rightarrow your infra in a $\begin{cases} \text{Safe} \\ \text{Consistent} \\ \text{Repeatable} \end{cases}$ way

by defining resource config that you

can $\begin{cases} \text{Version} \\ \text{Reuse} \\ \text{Share} \end{cases}$

④ Adv:

- Infra lifecycle management \rightarrow state based approach to
- Version Control Commit \rightarrow track resource changes

- Service account : Identity that an instance / app can use to run API requests on your behalf
- IAM & Admin → Service Accounts → Create Service Account
 - Add name → Grant viewer → Done access
- Action(?) → Manage Keys → Add → Create, JSON Key → new Key
- terraform binary & gcloud SDK (required)
 - terrforn -v
 - gcloud -v } to view
 - verif } verify
- Add JSON file path to env var & run auth command in gcloud & enter token in terminal
- IAM → edit permission → Add
 - Usually predefined rules not used in prod
Custom rules are created
- enable IAM & IAM credentials API's

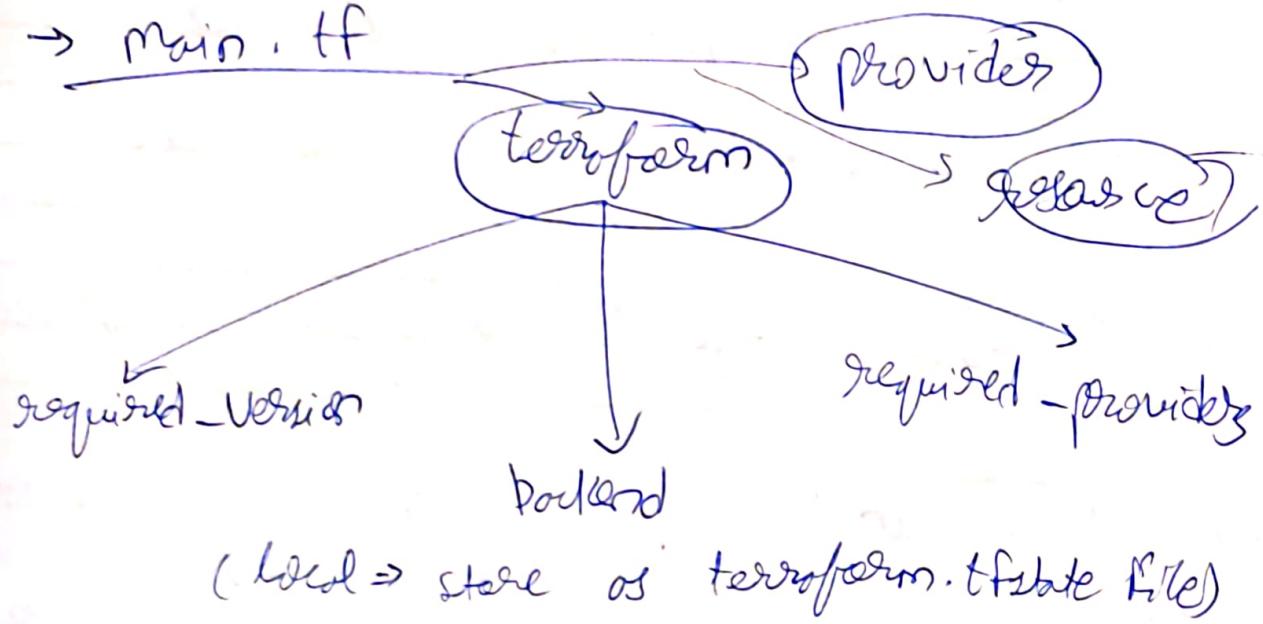
1.3.2 \Rightarrow Creating GCP infra w/ TF

\rightarrow Change .terraform - version to Version being used

\rightarrow Providers Plugins: binaries written in go lang, present in terraform registry

\rightarrow Docs for providers available for reference

\rightarrow main.tf



\rightarrow variables.tf $\xleftrightarrow{\text{local}}$ Variable

\rightarrow change region in Variables.tf

\rightarrow terraform init : $\begin{cases} \text{Initialize Backend} \\ \text{Install Plugins/Providers} \end{cases}$

\rightarrow terraform plan : preview ^{local changes} against remote stat

\rightarrow terraform apply : Ask for approval to proposed plan and applies changes to cloud

\rightarrow terraform destroy : Remove stack from cloud

Jan 30

Week 2 : Workflow Orchestration

2.1.1 \Rightarrow Data Lake

- Data lake holds data from many places.
- can be structured, semi or unstructured
- metadata is usually associated for faster access.

Ex: GCS, S3, Azure Blob

Data Lake	Data Warehouse
* Contains data with minimal processing	Contains cleaned and pre-processed data
* Large amount of data	Small amount of data
* ML, Streaming analytics	Batch processing, Reporting
* ELT	ETL
* Schema on read	Schema on write

2.2.1 \Rightarrow Intro to WO

From docs If you move data you may need these:

- Scheduling
- Retrials
- Logging
- Cocking
- notifications
- Observability.

prefect offers this functionality

2.2.2 \Rightarrow Introduction to Prefect Concept

- Create Conda env \Rightarrow conda create -n zoomcamp
- || → Activate " \Rightarrow conda activate zoomcamp
Python=3.9
- Clone Code from discrivers repo
- pip install -r requirements.txt
- || Rebs → Before: ingest_data that takes params for a postgres connection, pull the data & load to db.
- || → Its to be manually triggered, no visibility, retries, caching etc.

After:

- ① A prefect flow is a function with @flow,
has workflow logic
- ② Flow run is a single execution of a flow
- ③ Task represents unit of work, function
with @ task. Not required i.e. prefect
workflow can contain only flows - tasks
are called from a flow
- ④ extract_data downloads file from github, hence
retries & caching is added
- ⑤ transform_data removes null passenger_count
and logs before & after values
- ⑥ Blocks enable storage of configurations.
- ⑦ SQLAlchemy Connector block is used to store db
config.

Feb
5

2.2.3 \Rightarrow ETL w/ GCP & prefect

① prefect orch start : Start prefect UI

② etl-web-to-gcs() : ETL flow

③ fetch : gets data from web & returns df.
retries & waits (Exponential backoff etc) on failed

④ clean : Takes df, fix datatype issues (Convert
obj to datetime), write locally to Parquet
log_prints = True , logs the print statements

Parquet is zipped , using pyarrow for
best compression

⑤ write-gcs : Upload parquet file to GCS.

Create a bucket in GCS
register pref-gcp block in prefect

Create Service account with BigAdmin &
Storage Admin roles

Create GCS Bucket block , add \rightarrow bucket name.
run 'python etl-web-to-gcs.py'

⑥ local Parquet file created

⑦ This file is uploaded to bucket in GCS

2.2.4 \Rightarrow from GCS \rightarrow BigQuery.

① Flow : etl_gcs_to_bq

② task 1 : extract_from_gcs

Download parquet file from gcs to local

③ task 2 : transform

Read parquet file using Pandas, fill NA passenger count with 0, log pre & post count and return df

④ task 3 : write_bq

get creds from prefect block

use df.to_gbq() to load df to BQ in chunks.

⑤ Dataset has to be created in BQ UI
- & table

⑥ Change Bucket block name, table & project id

Feb
6

2.2.5 \Rightarrow Parameterizing flow &

deployments with ETL, into GCS Flow

- ① Pass parameters instead of hardcoding (Ex: Colour, year & month to construct URL for NY taxi dataset).
- ② ETL - web - to gcs Called inside etl = parent - flow, with different formats.
- ③ Cache key for added to avoid fetching pre-fetched data
- ④ Flow is run manually & observed via UI / API
- ⑤ Prefect deployment: triggers & schedules flows from UI instead of terminal.
Can be done through CLI / Python
- ⑥ prefect deployment build <filename.py>: < deployment >
- A "`<name of deployment>`"
- ⑦ prefect deployment apply <name of generated YAML file>
- ⑧ Work queues organize work that agents pick up
Agent processes are lightweight polling services that get scheduled work from a WQ
- ⑨ Notifications can be sent to different destinations.

offline notes -

① Code formatting

→ Block + VS Code extension

→ MLops ZenCamp → Best practices.

→ typehints

→ docstrings

→ Also midjour blog.

② Prefect Recipes Catalog

③ Why Parquet?

compression w/ Snappy
Columnar format

metadata about columns
partitioning

Apache Delta Lake

Schemas definitions has to be strict

④ "prefect register" is used to tell the server about available prefect blocks.

⑤ Prefect vs Airflow

founder was main contributor

read why not Airflow blog

prefect cloud vs open source

→ updates

→ collaborators

→ Actions: trigger something on event

Feb

8

2.2.6 \Rightarrow Schedules & Docker Storage w/ Infrastructure.

- ① Schedules are used to run flows automatically
- ② Can be added via CLI during deployment build or via UI
- ③ Cron, Interval, or Grade can be used
- ④ Use prefect deployment build --help
- ⑤ Dockerfile with prefect base image, requirements.txt, flows code and empty data folder.
- ⑥ docker image build -t heisenhub/image-name. (Should be logged in to dockerhub)
- ⑦ docker image push heisenhub/image-name
- ⑧ Docker block with image created.
- ⑨ Deployment of dockerized flow done through python script (can also be done via CLI)
- ⑩ PREFECT-API-URL set to local endpoint
- ⑪ To avoid Caching issue in dockerized deployment, remove caching code
- ⑫ Add GCP Credentials block containing Service account data, instead of path to JSON file

12

Week 2 Homework Solutions

① Hardcode Params or parameterize Script ideal

- Column names in green & yellow are different
(lpep) (tpep)
- if else based on colors to clean these columns
- Need to create data/colors dict

② '0 5 1 * #' \Rightarrow 5 AM on day 1 of the month

③ count sum can be added in code itself
instead of adding 2 manually.

④ Before: code locally or code files in docker
after: code in Github
stored

GIT block Created through UI / Python

⑤ Add webhook URL, notified when flow
reaches given state

Data Warehouse & Big Query.

3.1.1 OLAP vs OLTP

↳ Online Transaction Processing

OLTP	OLAP
Purpose Control & run essential business operations in real time	Plan, solve problems support decisions, discover hidden insights
Data updates Short, fast, initiated by user	Data periodically refreshed w/ scheduled batch jobs.
db design Normalized for efficiency	Denormalized for Analysis
Space req Generally small if history of data is archived	Generally large due to agg large datasets
User role Clerks, Online Shopper	data analysts, BA, Creatives.

- DWH is a OLAP Solution for reporting & data analysis.
- BD : Serverless DWH - No servers to manage
No db software to install
- Software & infra - Scalability
High availability.

Built-in features like ML, geospatial, BI Analysis

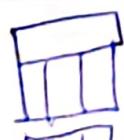
- Maximizes flexibility by separating compute engine analyzing data from storage.

- Hierarchy : Project
 - Schema (dataset)
 - Table

- Cost < On demand flatrate

Partitions (Segment)

- Separate data into smaller chunks
- If most of our queries filter based on a column value (say date), the dataset can be partitioned on this column.
- That way BQ would go through only the required data. So ↑ perf & ↓ cost
- Partitioned table query scans through less data
- Partitions can be viewed by "dataset.INFORMATION_SCHEMA.PARTITIONS where table_name = table_name."
- Partitioned table displays table size, no of rows, type & field
- No Partition icon Partition icon



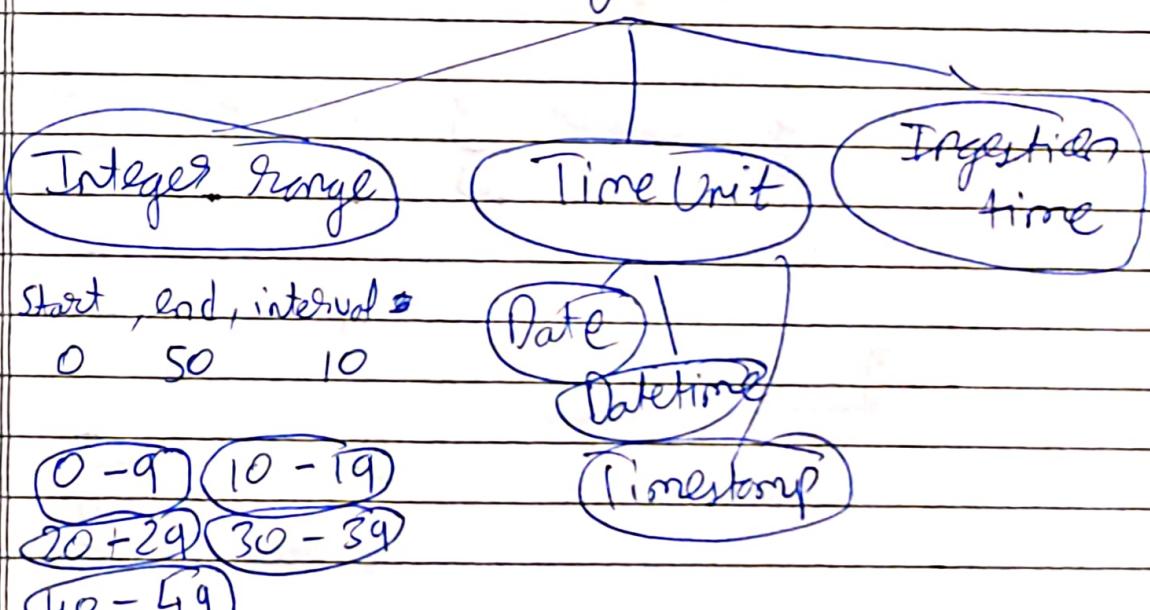
Clustering (Optimize)

- Sort data using 1 or more columns ^{Store} together
- User defined Column Sort Order
- improves query performance & reduces cost

In P : Query Cost Known before
In C it is after

3.1.2

Types of Partitioning



time Unit

or

ingestion time

hourly

daily

monthly

yearly

- No of Partitions can be 4000 max

For Clustering:

- Columns specified are used to co-locate related data
- Order of columns is important - ~~order~~
order of specified columns determining sort order of data
- Clustering improves filter & Agg queries
- Table with $< 1\text{GB}$ data size do not show significant improvement with P&C
- Up to 4 Clustering can be specified

Clustering over partitioning

- If P results in small amount of data per P less than 1GB
- If P results in large number of P (> 4000)
- If P results in modifying majority of partitions frequently

Automatic Reclustering

- When data is added to a C table
 - Newly added data can be written to blocks that contain key ranges that overlap with key ranges in previously written blocks
 - This violates Sort property of table
- To maintain performance
 - BD performs Automatic RC in background to restore Sort property
 - For P tables, C is maintained for data within scope of last partition

Feb 19

3.2.1 ⇒ BD Best Practices

Cost reduction

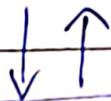
- ① Avoid Select *. Use Select <column names>
- ② Price queries before running them
- ③ Clustering & partitioning
- ④ Use Streaming inserts w/ caution
- ⑤ Materialize query results in stages

Query Performance

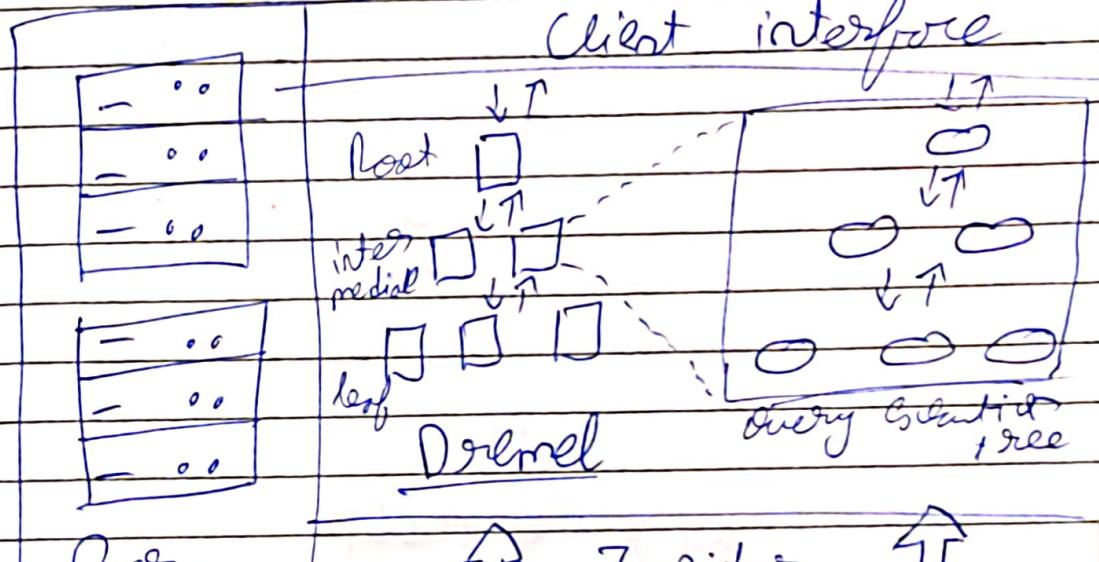
- ① Filter on partitioned columns
- ② Denormalizing data
- ③ Use nested or repeated columns
- ④ Use external data sources appropriately
- ⑤ Reduce data before using a join
- ⑥ Do not treat WITH clauses as prepared statements
- ⑦ Avoid oversharding tables
- ⑧ Avoid JS User defined functions
- ⑨ Use Approx Agg fns (HyperLogLog+)
- ⑩ Order last in query
- ⑪ Optimize join patterns
- ⑫ Place table with most rows first then table with the fewest rows then remaining tables by decreasing size

3.2.2 \Rightarrow Internals of BD

Clients - Rest API, web UI, bq cmdline.

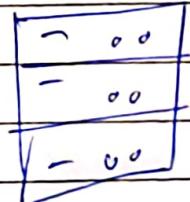


Client interface

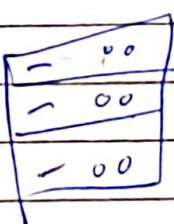
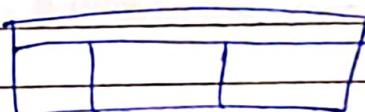


Borg

Jupiter



colossus



Dremel Execution engine : SQL query to Selection tree

Colossus Distributed storage : Distributed file system

Borg Compute : Cluster management System

Jupiter Network : Networking Infra

BD : These infra + Databases of high level tech from BD

Feb
25

3.3.1 \Rightarrow BigQuery Machine Learning

- Can build ML model in DWH itself, avoiding extra step of exporting the data.

Steps

- Create a table with appropriate columns data types. Categorical features are considered as string
- Create model with default setting
- Check features
- Evaluate the model
- Predict using the model
- Predict & Explain model
- Hyper parameter tuning

3.3.2 \Rightarrow PO ML Deployment

- Get cloud auth login
- Export project to GCS
- Create dir & copy model to local dir
- Create Serving dir & copy model here
- Pull tensorflow Serving docker image & run it
- Call API using POST request w/ parameters & get predicted tip amount

Week 3 HW Solution

① Modified web to GCS preflight Script to work for F1V data

Set ② Created deployment & run with 1-12 months 8 2019 parameter. Data now uploaded to GCS bucket

③ In BD, Create external table from 12 GZ files stored in GCS

④ Create BQ table from external table.

Q1 Count of records in table.

Sol: 4,324,696

BQ table took 26 Sec

BD table took 366 ms (faster)

Q2 distinct count of columns in table

Sol: BD table processes 0 bytes

BQ table processes 317.94 MB

As data is not in BD, estimation is not possible

Q3 Count where both column values are null

Sol 717,748

④ Strategy for optimize query.

filter hints partition

order hints cluster

Sol: P by pdt, C by abn

⑤ Performance of optimized table.

non optimized \Rightarrow 647.87 MB

optimized \Rightarrow 23.05 MB

⑥ Sol: Data in External table is stored in GCS bucket

⑦ Sol: Not a Best Practice to always cluster data, ex: when data is less than 1GB, creates overhead instead of optimized perf

2.1 \Rightarrow office hours \Rightarrow prefect, Jeff Hale.

① Better approach than download to local & upload to BD.

sd Ingestion tool like Airbyte, Fivetron
prefect - gcp

② K8's like Argo

Sol: Collections Catalog \rightarrow Ray, Dask, Kubernetes.

③ prefect & spark

sol: Collections Catalog \rightarrow Databricks, Fugue

④ prefect opns, deployment, profile

⑤ prefect helm-chart

office hours #3 \Rightarrow Alsey.

① Data cleaning \leftarrow Soda Core
great expectations

② dft Spark

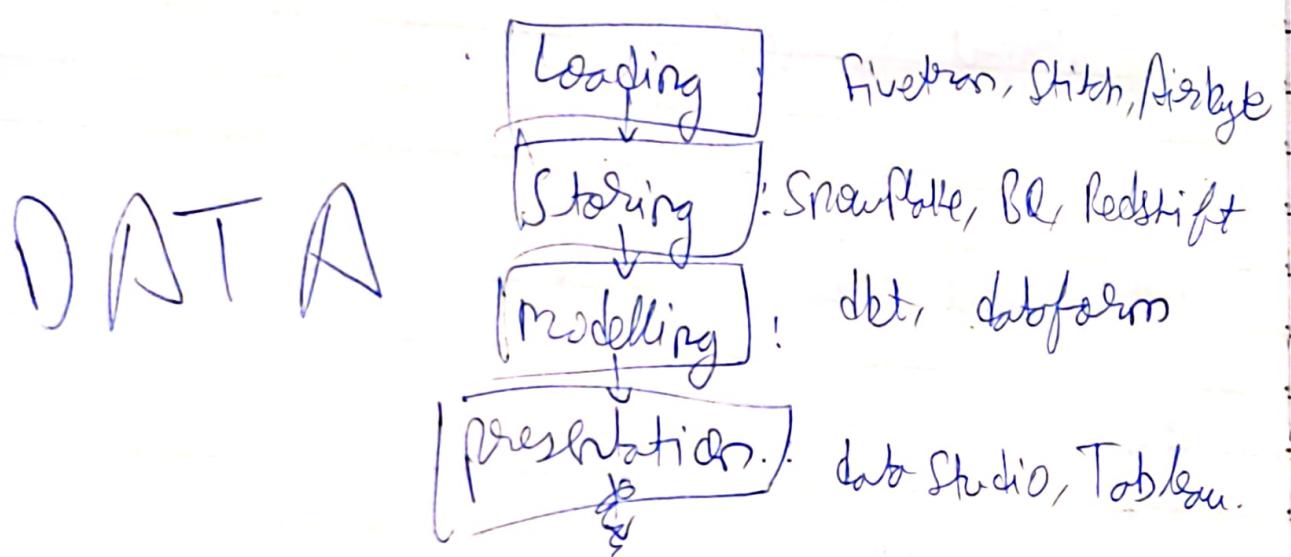
- Compute in DW/H needs infer
- Data never leaves DW/H $DW/H \rightarrow$ clean \rightarrow DW/H

Spark

Use prefect for orchestration, delegate compute
to K8s, AWS batch, spark etc.

④ for project

Bluie dataset → bluie streaming → Think
or
Batch
Criteria
Points



Data modelling Concepts.

- ETL : Extract → Transform → Load → Reporting
 - Slightly more stable & Compliant data analysis
 - Higher Storage & Compute Costs
- ELT : Extract → Load → Transform → Reporting
 - Faster and more flexible data analysis
 - Less maintenance

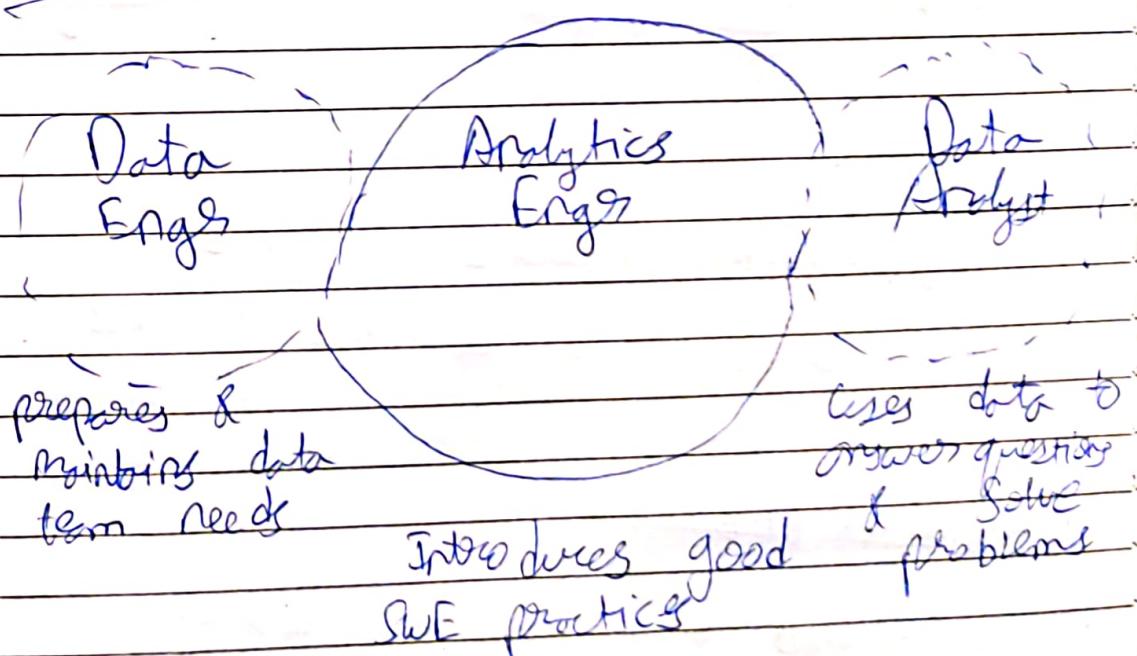
Mar 7

Week 4 : Analytics Engineering.

Setup dbt for BQ

- ① Create dbt cloud account
- ② Create BQ Service account with BQ Admin role , add & download JSON Key
- ③ Create project in dbt → Select BQ Data → Add SA Key JSON file.
- ④ Copy SSIT git URL of your repo → paste ~~the~~ generated key in git-hub settings Under Deploy Keys, with write access.

4.1.1 ⇒ Analytics Engineering Basic



Kimball's Dimensional Modelling

- Objective :
- Delivers Data understandable to the business users.
 - Delivers fast query performance

Approach : User Understandability & Query performance

→ Non redundant priority data (3NF)

Other approaches -

→ Bill Inmon

→ Data Vault

Elements of Dimensional modelling

Fact table

- Measurements, metrics or facts
- corresponds to a business process
- Verbs



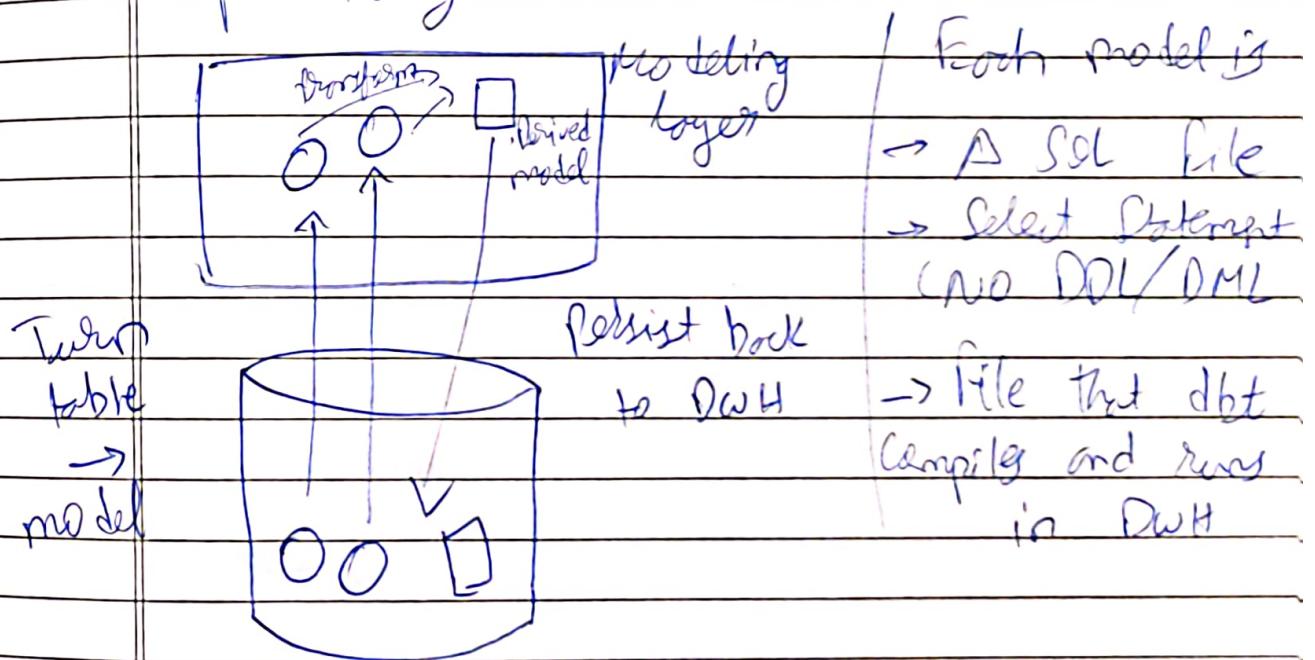
Dimension table

- corresponds to a business entity
- provides context to a business process
- Nouns

Architecture: [Staging] → [Processing] → [Presentation]

4.1.2 → What is dbt?

- dbt: data build tool.
- transformation tool, uses SQL to deploy analytics code following SfE best practices like docs, CI/CD, modularity, portability



dbt Core: Open Source, allows data transformation

dbt Cloud: Send app to develop & manage dbt projects

MoS8] Week 5 : Batch Processing.

5.1.1 ⇒ Intro to BP

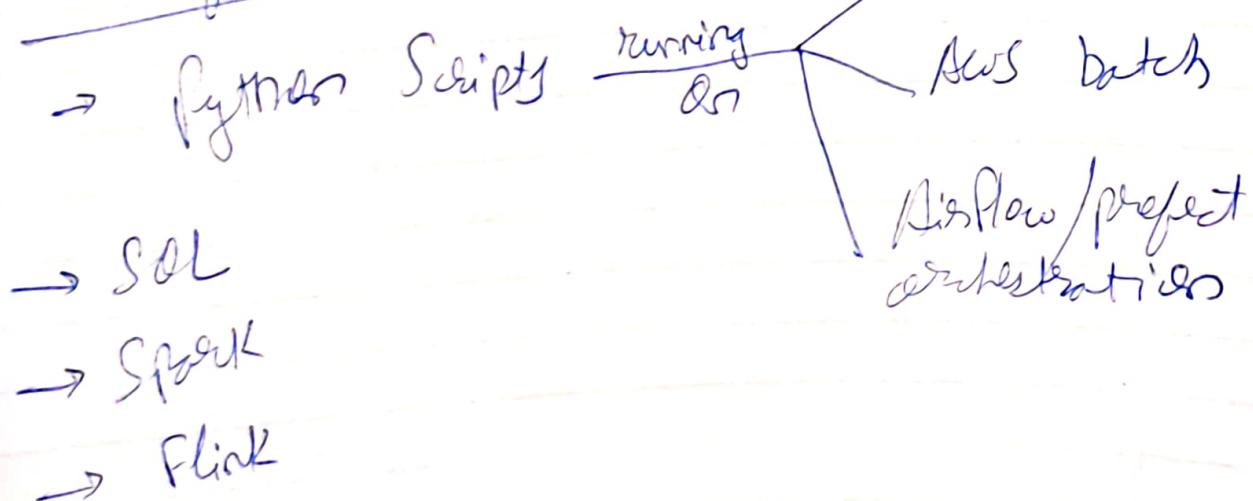
Batch processing : Processing & Analysis happens on data already stored over a period of time

Ex: Payroll & Billing Systems processed weekly or monthly.

Stream processing : Happens as data flows through a system. Analysis & reporting of events as it happens.

Ex: Fraud detection.

Tech for Batch Processing.



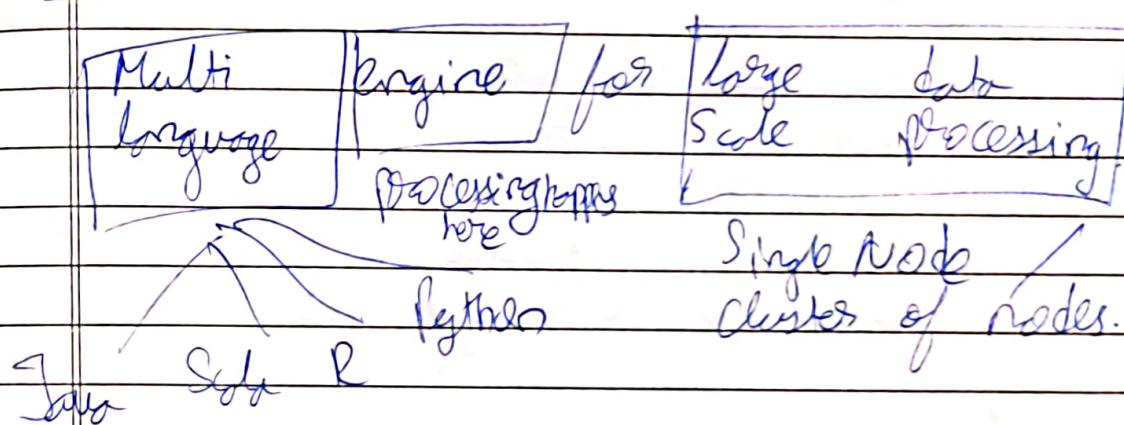
Adv

- ① Easy to manage
- ② Better
- ③ Scaling is easy.

Rise

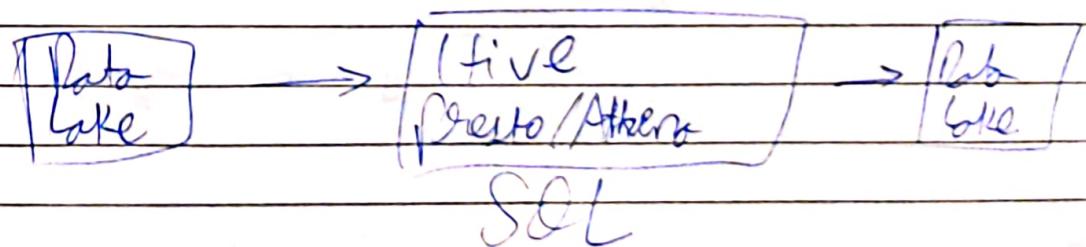
Waiting for data to be processed, ex:
to sort data from previous hour we
need for it to get processed

5.1.2 \Rightarrow Intro to Spark.

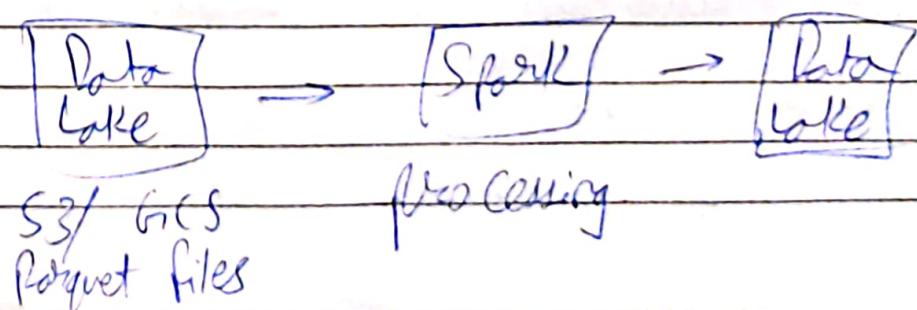


When to Use Spark?

If Batch job can be expressed as SQL then
use SQL



Otherwise use Spark.



- Spark can be used for ML work load

5.2.1 ⇒ Installing Spark on Linux (GCP VM)

① Install Java (OpenJDK)

Wget file → tar xzvf → Set path Var → java --version

② Install Spark

③ Add path vars to .bashrc so that when bash is run these vars get set

④ PySpark

Set pathways & add to .bashrc

Download CSV file and run python code in

Jupyter →

- ① Create SparkSession
- ② Read CSV file to df
- ③ Show df
- ④ Save df to Parquet

⑤ Spark UI : When Spark Session is created a UI is hosted on port 4040. Expose port 4040 in VS Code & open localhost:4040 in chrome to view Spark Web UI

5.3.1 \Rightarrow First look at Spark & PySpark

- Create Spark Session
- Download & Extract FHV HV 2021-01 data
- Read this CSV file using Spark. 11 million rows
- Spark reads all columns as `String` type.
- To fix this:
 - ① Read Sample data (1000 rows) in pandas and print df. types
 - ② Create df from pd df
 - ③ Create Schema with required types
 - ④ Pass Schema while reading CSV in Spark

Partitioning

- To achieve high parallelism, spark splits data into smaller chunks called partitions.
- These are distributed across different nodes in the Spark cluster.
- Use `df.rdd.getNumPartitions()` to see no. of partitions
- Use `repartition()` to change & `coalesce()` to decrease no. of partitions
- repartition involves shuffling, which is expensive, Coalesce does not up shuffle
- df is repartitioned to 24 & stored as parquet

5.3.2 ⇒ Spark Dataframes

- Parquet knows schema so each column has a type.

Actions	Transformations
→ lazy: not executed immediately	→ eager: executed immediately
→ select	→ show
→ filter	→ take
→ join	→ head
→ groupBy	→ write

- PySpark comes with many in-built functions.
from pyspark.sql import functions as F
- UDF: User defined functions can be used to

5.3.3 ⇒ Preparing Yellow & Green taxi data.

- download-data.sh: to download green $\frac{1}{2020}$ yellow $\frac{1}{2021}$ data
 - set -e used to stop execution if non 0 states code
 - Takes arg as colour & year
 - loops through all months
 - format number as 2 digits $\frac{1 \rightarrow 01}{10 \rightarrow 10}$ using "%02d"
 - Create dir using mkdir -p
 - wget csv.gz file to created path

- Read CSV and store as Parquet.
- get Schema from CSV file using pandas.
Create Spark DF using pd DataFrame
- Schema created, edited to add types.
- Convert Long to Integer, string to datetime
- Loop over each month 1-12
Read CSV with Schema
Repartition it into 4 files
Write to Parquet

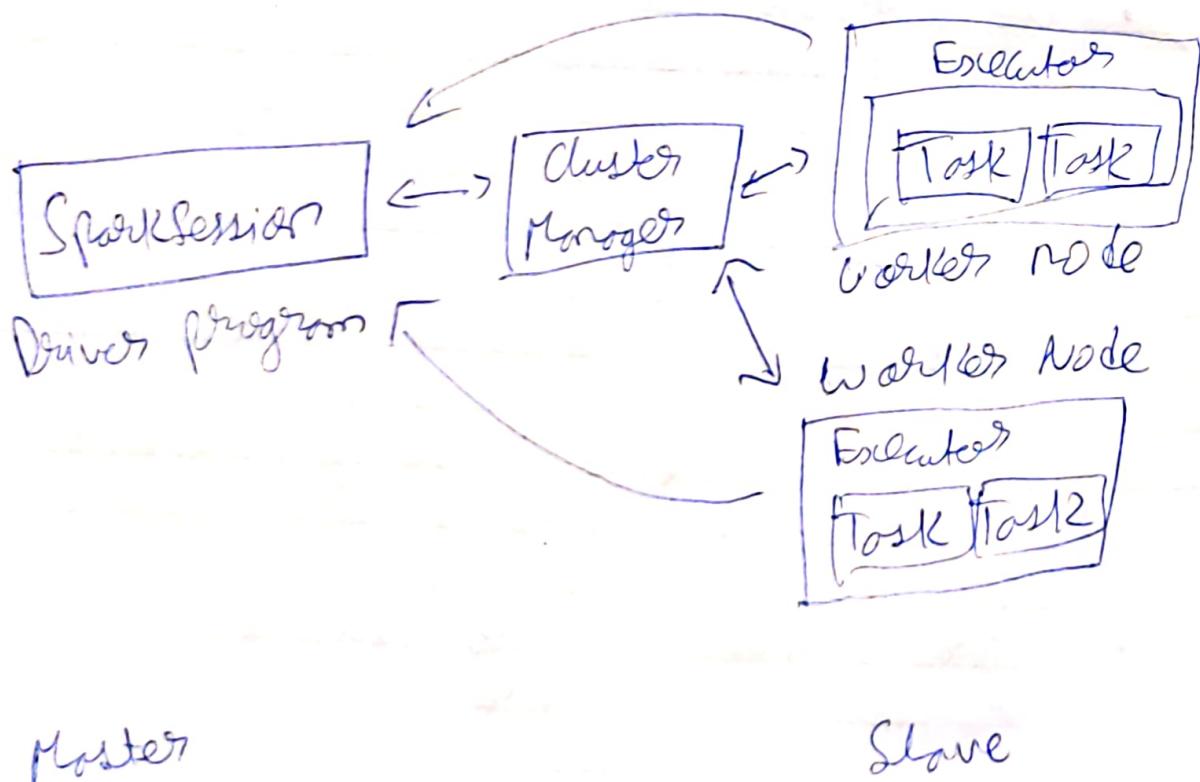
5.3.4 ⇒ SQL with Spark

- We will use the files created in prev video
Using download data.sh & taxi schema.ipynb
- The files are 2020 & 2021, yellow & green
taxi data in Parquet format
- Create Spark Session, read above data
- Rename Datetime columns so green &
yellow column names match (lprep vs tprep)
- Create Common-Columns list. Set intersection
not used to preserve order
- Select only Common Columns & combine both
to Single DF. Service-type used to know where
columns come from

- Register it as table using `registerTempTable` to run SQL queries on table
- Run the grouping query (from code 4) is run on Spark
- result is saved to parquet file.
- This creates hundreds of small files.
- Coalesce is used to decrease it to 1
- So the result is stored in a single parquet file

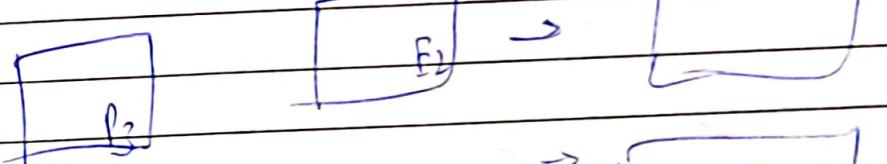
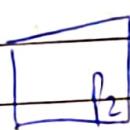
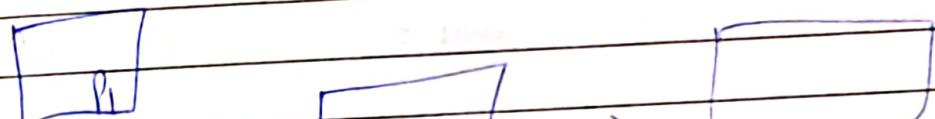
5.4.1 \Rightarrow Anatomy of a Spark Cluster

`local[*]` : Run Spark locally with as many worker threads as logical cores on your machine



5.4.2 \Rightarrow GroupBy in Spark

Stage 1

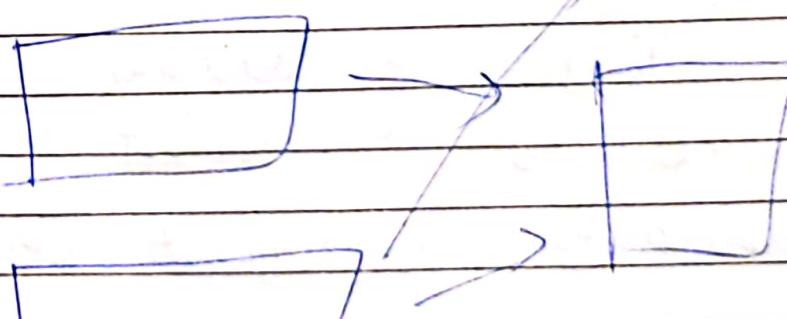
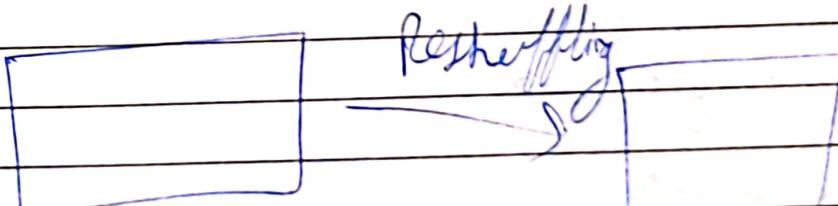


partitions

executors

Intermediate
results.

Stage 2



External
merge
sort

Intermediate
results

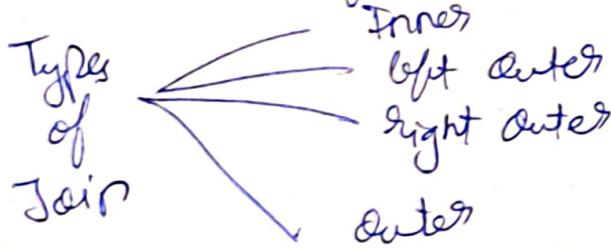
Apr
12
2023

5.4.3 \Rightarrow Joins in Spark

- Joins are used to combine 2 or more datasets based on a common key.
- The idea is to match records from multiple datasets and return a new dataset that contains information from both sources.
- Join uses reshuffling which uses External merge sort Algo.

Materialize

- Convert RDD/Dataframe \rightarrow Concrete data structure
- can be stored in memory/disk and reused for faster processing



(Join strategy)

Based on how data is distributed across nodes

Partitioned Broadcast

Partitioned & Sort merge Join	Broadcast Join
For large tables	<ul style="list-style-type: none">For 1 large & 1 small tableSmaller table sent to all nodes
Default strategy	

5.5.1

Operations on Spark RDD's

→ RDD: Resilient Distributed Dataset

→ Way of representing data on spark that allows you to distribute it across multiple nodes in a cluster so that you can work w/ it in parallel

→ Collection of objects that are partitioned & distributed across different nodes in a cluster

→ DataFrames have become the preferred API for working with structured data.

However RDD's are still important in certain scenarios.

take(n)

returns list of n Rows. Row is a spark object representing a row of data

map()

Applies function to each row & returns new RDD w/ the results.

Ways by
SQL query is expressed as set of steps
in Spark RDD query.

- ① Spark Session built
 - ② Green taxi data read from parquet files
 - ③ Only required ~~taxi~~ columns are selected as RDD
 - ④ A chain of operations are done:
 - a) filter-outliers: Consider only rows with pickup datetime after 01/01/2020
 - b) map: prepare for-grouping:
 - return (key, Value)
 - ↳ (hour, zone)
 - ↳ (amount, count)

⑥ deduce By Key : Elevate - Freesame

Group elements of a pair RDD by key and apply a reduction function to the values associated with each key.

Ex: $\left[\begin{array}{l} ('a', 1), \\ ('b', 2), \\ ('c', 3), \\ ('a', 4), \\ ('b', 5) \end{array} \right] \xrightarrow{\text{reduceByKey}} \text{Sum operation} \left[\begin{array}{l} ('a', 5), \\ ('b', 7), \\ ('c', 3) \end{array} \right]$

it is used to sum amount
and calculate orders grouped by
keys

① ~~map~~: Single row created
from nested structure using tuple
unpacking

② ~~toDF (Schema)~~: Converts result to
DF w/ given Schema

Result is also written to Parquet

5.5.2 \Rightarrow Spark RDD mapPartitions

\Rightarrow map() applies a function to each individual
element of an RDD. The function is called
once for each element in the RDD.

\rightarrow mapPartitions() applies a function to each
partition of an RDD, allowing the function
to process multiple elements of the partition
at once

Eg:
To run an ML model on partition of
a dataset.

5.6.1 \Rightarrow Connecting to GCS

- gcloud auth login
- gsutil -m cp -r "pq/" "gs://dez-bucket/pq"
- All parquet files gets uploaded to GCS bucket
- GCS Connectors: Open Source Java library that can run Spark jobs directly on data in Cloud Storage.
This is downloaded
- SparkConf Modified to include Connectors and Credentials file
- Hadoop Config done
- now data can be read directly from GCS bucket using "gs://<path>"
- Sample rows shown & total count is shown

5.6.2 \Rightarrow Creating a Local Spark Cluster

Before: A local Spark cluster is created. We connect to this.

master is specified as local

Now: Creating Spark cluster outside
of jupyter notebook.

AKA: Running Spark in Standalone mode

① Start master & workers

- cd to SPARK_HOME/sbin dir
- ./start-master.sh
- UI is on localhost:8080
- Instead of master as local[*], the URL of Spark master is specified (spark://)

Jobs won't be run yet, as there are no running workers.

So ⇒ ./start-worker.sh <master-spark-URL>

② Convert notebook to python Script

jupyter nbconvert --to=script notebookname

③ Running Script

- Make sure no other script is running as it might take up all resources

- Use args instead of hardcoding green, yellow & output path

- Master URL is hardcoded & no of executors is not specified in script

④ Using Spark - Submit

- Used to launch applications on a cluster
- master URL can be specified here instead of in code
- Run command where .py file is present, Spark-Submit is present in PATH
- Results are stored after job is run

Stop masters & workers as done

- ./ stop-worker.sh
- ./ stop-master.sh

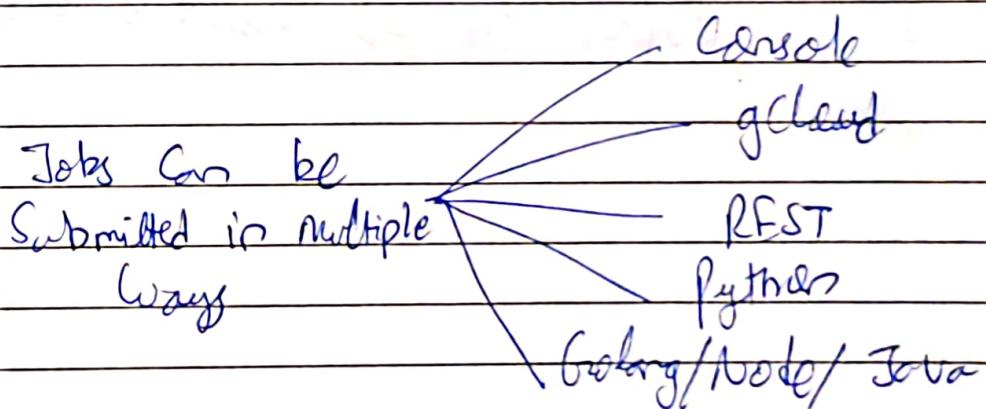
5.6.3 ⇒ Setting up a Dataproc Cluster.

Dataproc is a fully managed & highly scalable service for running Spark jobs (among many others)

- enable the API
- Set cluster name & region
- ideally we use Standard (1 master, N workers) but for practice/testing we go with single node (1 master, 0 workers)
- enable Jupyter notebook & docker
- Create Cluster

Submitting Jobs

- Click on Created Cluster → Submit job
- Job type is pyspark, Main python file is the gs:// path of python script
- Add arguments (change paths from local file system to GCS buckets)
- Once the job is run & complete, the output report parquet is stored in GCS bucket
- In case of permission errors : Add IAM roles for dataproc.



5.6.4 → Connecting Spark to Big Data

- we saw how to get data from GCS process using spark & put result back to GCS
- We might need to put results data to DWH like BD, instead of GCS
- One way : Create external table from files in GCS and then create Normal table.

A better way would be to use Spark
to write to Bd

- Set temporary GCS bucket (bucket is automatically created during cluster creation) in spark config
- Coalesce not required, as data is written to Bd
- write format is bigquery, options is "table" along with table name

Once the command is run, a Spark job is sent to Dataproc, the job result is stored in Bd.