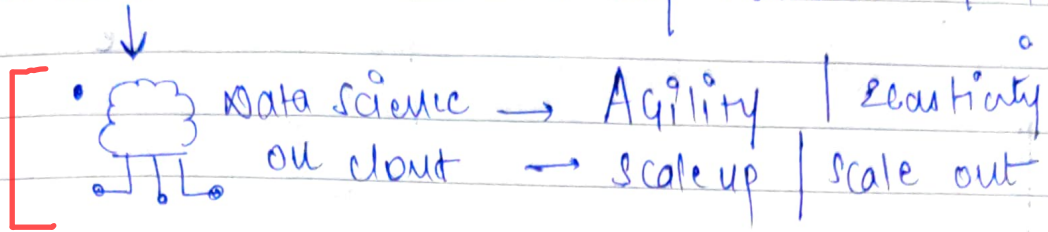


PRACTICAL DATA SCIENCE ON AWS - CLOUD

- Local Data Science → limited by Hardware & Memory



- Scaling up → If Model training takes too long → means u've exhausted your CPU limit on the current compute instance, so you can
↓
↑ the size of compute instances (single CPU)
↳ pick a compute instance with higher CPU resources.
↳ pick up a GPU based compute instance

- Scaling out → Move from single CPU to PARALLEL computing CPU instances.
↓
(parallel X CPUs)

"Scaling Up & Scaling Out Are possible within seconds in cloud"

- Once model training is finished, the instance is killed.
↓
No extra payment. { Pay As you Use }

- Cloud has (AWS has) its own ML Toolkit that saves time & makes it quick & easy.

Mapping Data-science Workflow to AWS Tools

Ingest & Analyze

- EDA
- Bias detection

Amazon S3 (Ingestion)
Amazon Athena
AWS Glue (data schema)
AWS SageMaker (EDA)
Data Wrangler & Clarify

Prepare & Transform

- Feature Engg
- Feature store.

→ AWS SageMaker Data Wrangler
→ AWS SageMaker Processing Jobs
→ AWS SageMaker Feature Store

Train & Tune

- Auto ML
- Model hyperparameter tuning

→ Amazon SageMaker AutoPilot
→ AWS SageMaker Training & Debugging
→ AWS SageMaker Hyperparameter Tuning

Deployment & Manage

- Deployment
- Automated pipelines

→ AWS SageMaker Endpoints
→ AWS SageMaker Batch Transform
→ AWS SageMaker Pipelines

Typical Data Ingestion & Query / Explore Workflow

① Ingest : Ingest the data with AWS S3



② CATALOG : Catalog the data into desired schema by using

AWS Glue



③ EXPLORE : Explore & run SQL queries on data using)

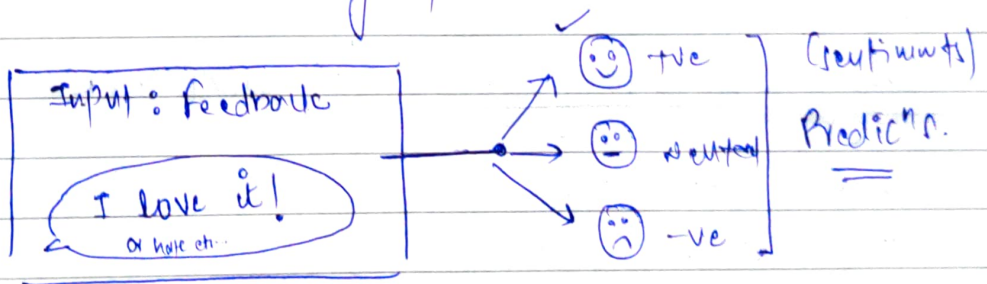
AWS Athena

Use-Case Description → Multi Class Classification for sentiment Analysis of Product reviews.

- You work at an e-commerce company like Myntra / Flipkart.
- Your customers leave feedback via all online channels.

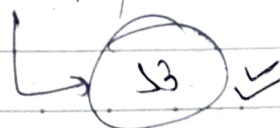
↓
Aim is to analyze these feedback quickly & alert if there are any product issues.

Example



Observation 1 → The feedbacks generated/received can be millions every minute. So, we might need a data (ingestion) Repository that's

FLAT enough to expand as per data size & formats.



Data LAKE : data is ingested into data lake.

↳ centralized & secure repository

↳ store / discover / share data @ ANY SCALE

- structured data (CSV...)
- semi-structured data (XML, JSON...)
- unstructured (Image, video...)
- streaming data (Live).

→ • Needs to be Governed i.e. managed to be used by further teams.

↳ Data lakes on AWS : S3

↳ stores data chunks as objects → each object with unique identifier & metadata for easy extraction.

* AWS Data Wrangler (pip install awswrangler)

- is a python library
- connects pandas + AWS power.

```
→ import awswrangler as wr
→ import pandas as pd
→ df = wr.s3.read_csv(path = '-/-')
```

* AWS Glue Data Wrangler catalogue

→ Catalogue is the metadata of your data (path, schema etc)

↓
info that you need to register before importing

→ import awswrangler as wr

→ wr.catalog.create_database(name = ---)

→ wr.catalog.create_mv_table(- - - -)

Catalog ↓
Name = ~

Database = ~

Classifier = mv

Location = s3:// < > / - -

(*) SQL Queries using AWS Athena → (SQL Queries in Python^{SW})

→ import awswrangler as wr

→ wr.athena.create_athena_bucket()

- df = wr.athena.read_sql_query(sql = 'select * - - -',
database = ' - - -')

↓
This queried df is first stored by Athena into S3
↓
then stored in variable df.