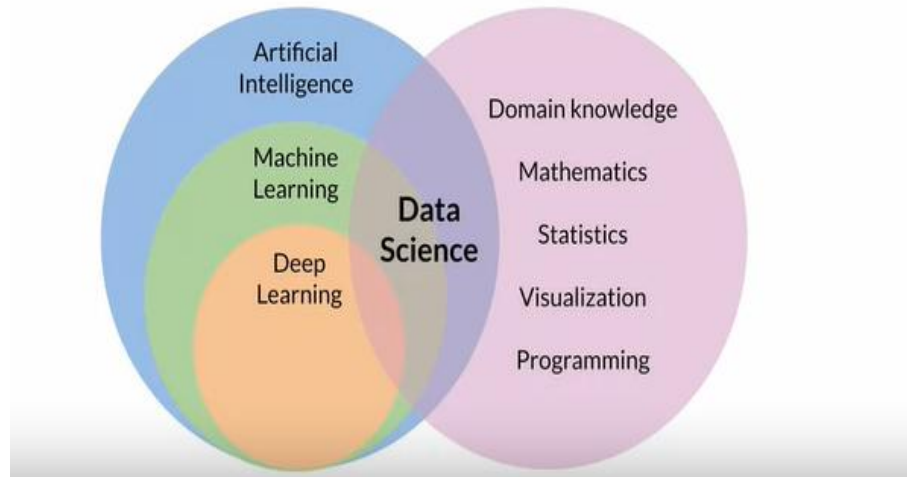# Practical Data Science with AWS Cloud

*(Hand-written notes and Course Gist)*

**Divyanshu Vyas**
**Energy Data Scientist, Petroleum From Scratch**

# 1. Baiscs about AWS cloud : Scaling Up & Scaling Out abilities
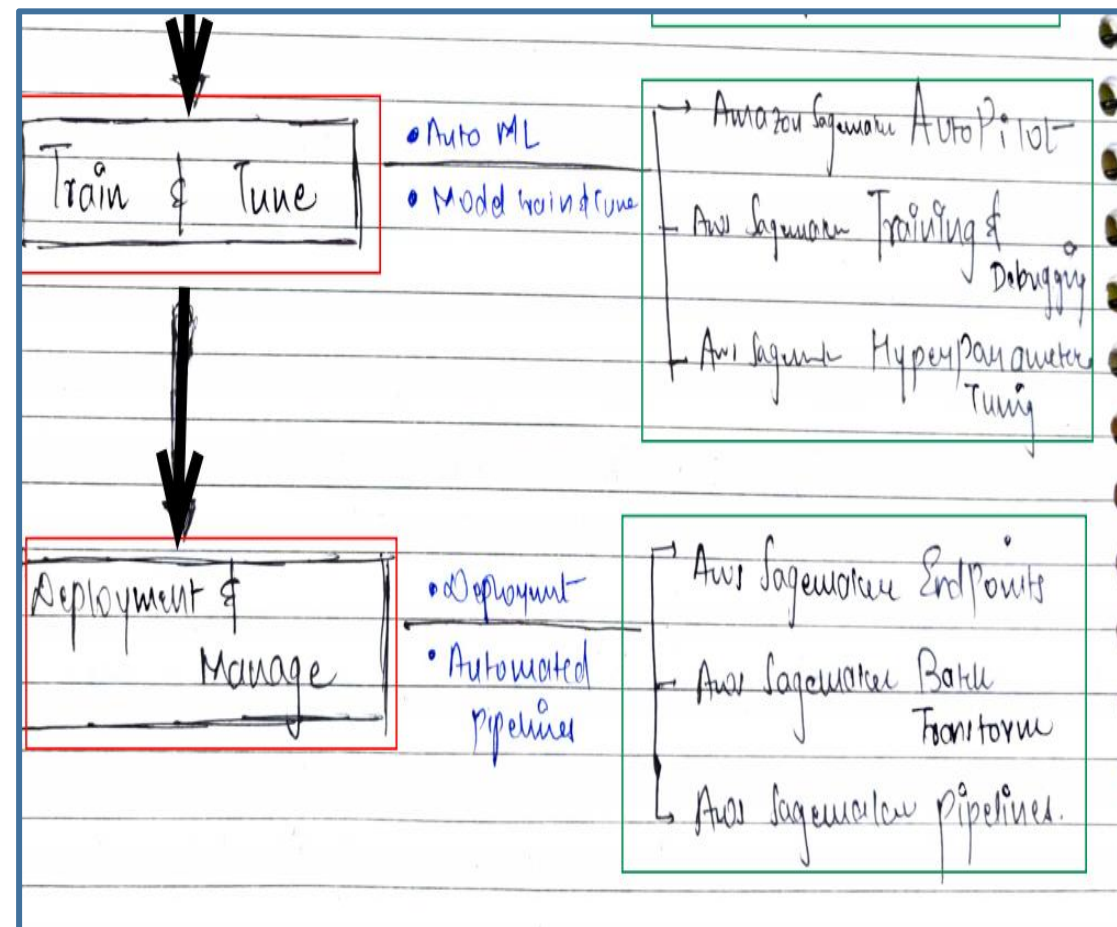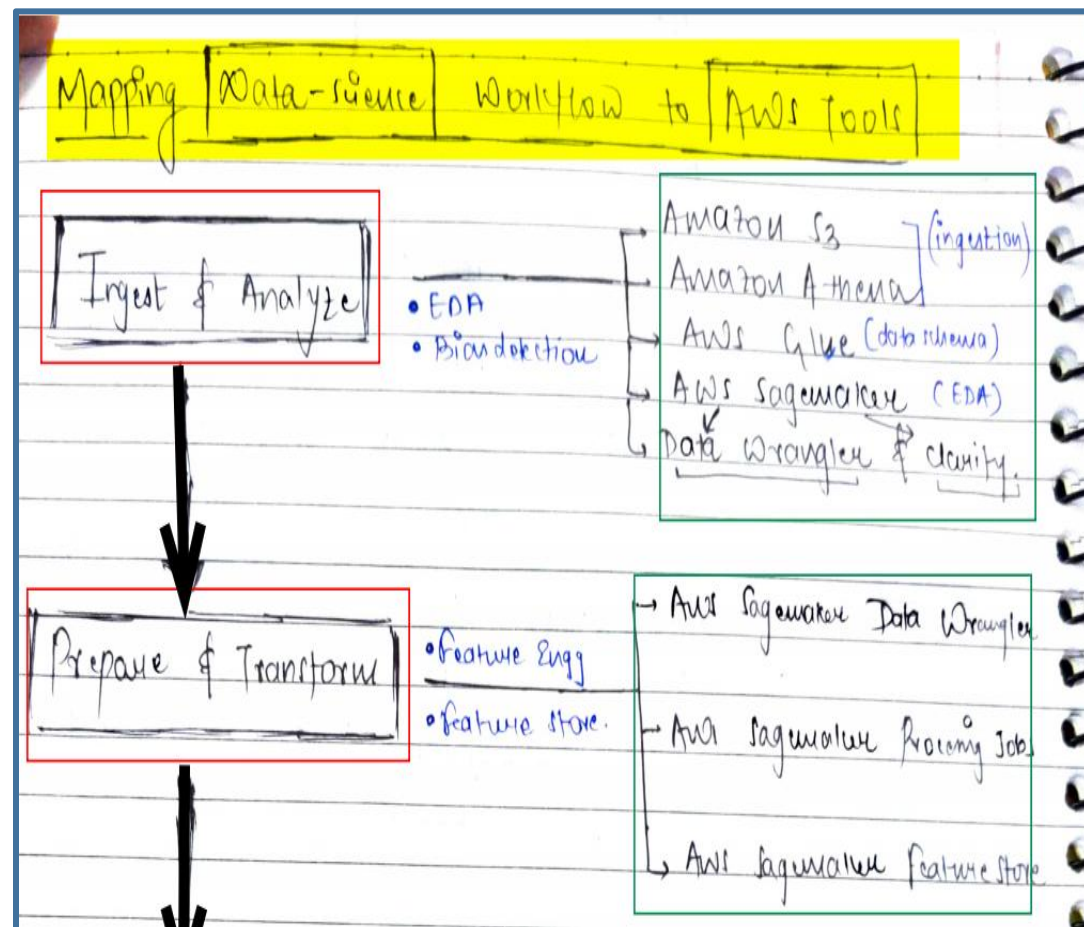


PRACTICAL DATASCIENCE ON AWS - CLOUD!

- Local Datascience → limited by Hardware & Memory

- Data science → Agility | Elasticity
  on cloud → scale up | scale out

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

↑ se size of compute instance (single CPU)

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



- 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 Toolbox that saves time & makes it Quick & easy.

# 2. Mapping Data Science workflow with AWS Tools



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 & Clarity.

**Prepare & Transform**
- Feature Engg
- Feature Store.

→ AWS Sagemaker Data Wrangler
→ AWS Sagemaker Procesing Job
→ AWS Sagemaker Feature Store

**Train & Tune**
- Auto ML
- Model train & tune

→ 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.

# 3. Mapping Data Science workflow with AWS Tools
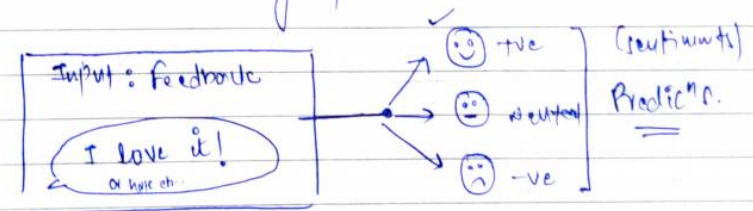


Typical Data ingestion & Query / Explore Workflow.

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-com company like Myntra / flipkart.
- Your customers leave feedbacks via all online channels.

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

Example

Input : feedback

I love it!
or hate etc.

→ 😊 +ve  ⎤ (sentiments)
→ 😐 neutral ⎬ Predic"n.
→ ☹ -ve  ⎦

Observation 1 → The feedbacks generated / reviewed can be millions every minute. So, we might need a data (ingestion) Repository that's ELASTIC enough to expand as per data size & formats.

# 4. Example Use case



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
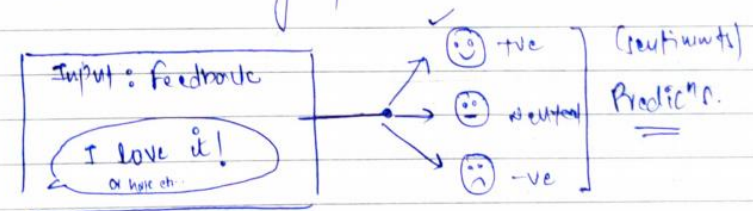
↓

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-com company like Myntra / flipkart.
• Your customers leave feedbacks via all online channels.

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

Example

Input: feedback

I love it!
or hate etc.

→ ☺ +ve
→ ☺ neutral    (sentiments)
→ ☹ -ve         Prediction.

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

ELASTIC enough to expand as per data size & formats.

# 5. Data Lake (AWS S3) , Data Wrangler, AWS Athena & AWS Glue

**Data Lakes** : data is ingested into data lakes.

↳ 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 ie 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
- Converts pandas + Aws powerel.

→ 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 & )
→ info that you need to register before importing

# 5. Data Lake (AWS S3) , Data Wrangler, AWS Athena & AWS Glue



**Left page:**

(*) **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_wv_ table (.....)

Catalog
Name = ~
Database = ~
Classifica = wv
loc^n = s3: || < > | ...

**Right page:**

(*) **SQL Queries using Aws ATHENA** → (SQL Queries in Python)

→ import aws wrangler 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.