

Adv of using cloud :-

- any amount of data
- no CPU, memory limitation
- CPU → GPU , hardware expansion.

Ingest → Analyze → Explore

Manive datasets [social media channels, internal company sources]
web applications

Messy, poorly documented.

On Server/laptop
Limited by resources

hardware
CPU
memory
buy additional resources.

Use Cloud :-

Scale Up:- Switch to GPU or use instance of multiple CPUs

store and process

Scale Out:- Distributed training on multiple CPU instances.

any amount of data

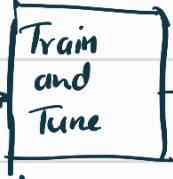
Machine Learning Workflow :-



- Data Exploration
- Bias detection



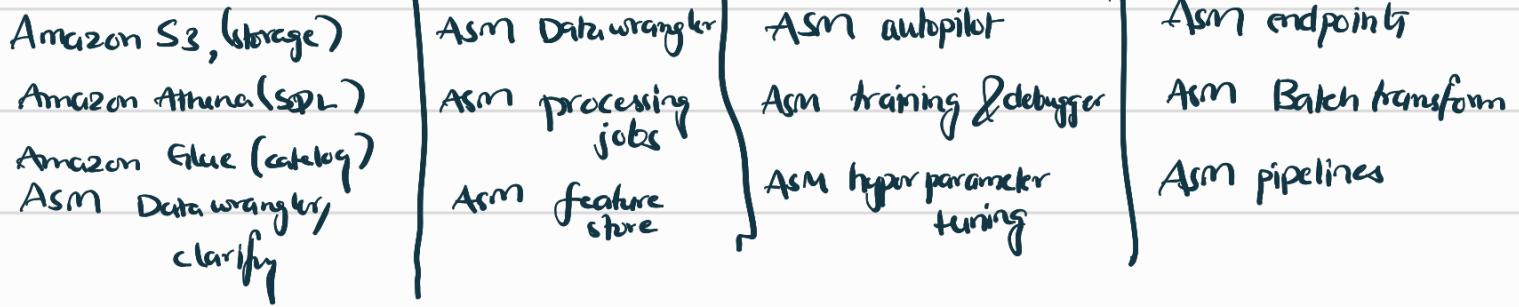
Feature Engr
Feature store



AutoML
Model train, tune



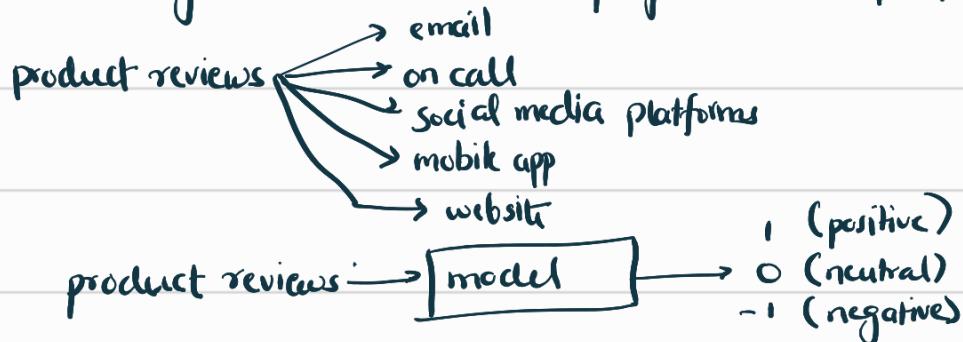
Model deployment
Automated pipelines



Use Case [Text Analysis]

Multi-class Classification.

Assume working on a ecommerce company with multiple products



Ingestion and Exploration:-

Ingest data into data lakes (centralized and secure repository that can store data, discover and share data at any scale and any type)

- structured (csv, excel)
- semi structured (xml, json)
- unstructured (img, audio)
- streaming data (app which gives continuous log files) /social media feed

Data lake to be governed.

- ways to discover and catalog new data
- secure and control access

Data lakes build on object storage (S3) Simple Storage Service.

File storage [Stores and manages data as individual files in a hierarchical order]

Block storage [Stores and manages data as chunks called blocks]
block → unique identifier but no metadata

object storage [data is stored and managed as objects which consists of
data, metadata (last modif etc.,), identifier]

Object storage is helpful for storing growing data of any type.

Durable, Available, exabyte scale

Secure, Compliant, Auditable

AWS Data Wrangler :-

- open source python library.
- connects pandas DF with AWS services.
- load/unload data from
 - data lakes
 - databases
 - data warehouses.

AWS Glue :-

Data Catalog	what data is stored in the S3 bucket / database
Name	reviews
Database	dsaws-deep-learning
Location	s3://<bucket>/
Classification / type	csv

Amazon Athena :-

- Query the data in S3 using Athena
- using SQL
- serverless (no infrastructure to set up)
- schema lookup in AWS Glue

Data Catalog

- Creating a reference to data (S3 to table mapping)
- metadata / schema stored in tables
- No data is moved. Only metadata is present in Glue database.
- Automatically setup using glue crawlers to
 - infer the schema
 - update the data catalog

Pros:-

Athena can deal with

- highly complex analytical queries
- GB > TB > PB
- need not worry about compute, memory resources as athena will scale out and split query into simpler queries and runs in parallel

- based on presto (open source distributed SQL engine to run queries on any data sizes)
- no infrastructure setup / data movement is required.

Lab:-

boto3 :- AWS SDK for python used to create, config, manage AWS services

sagemaker :- python SDK with high level abstractions for working with sagemaker

aws wrangler :- python open source that extends the power of pandas to AWS connecting dataframes to amazon related services.

athena :- can be used to interact the data in S3 bucket with standard SQL .

write SQL → Amazon Glue → S3

dataframe is returned.