

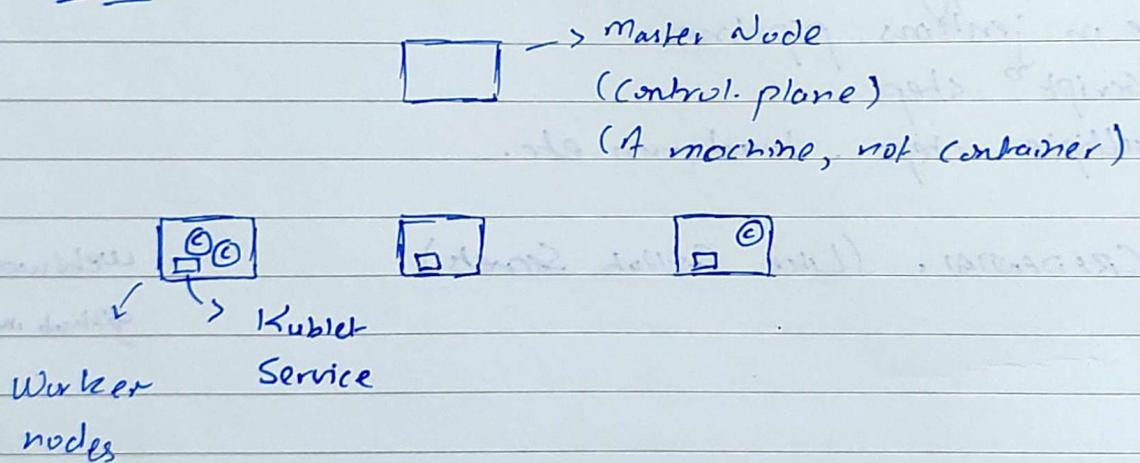
Kubernetes (128S)

- ⇒ Google's container orchestration tool.
- ⇒ Increasingly micro-service oriented applications nowadays
- ⇒ Kubernetes used to manage stuff between them

Expectations

- Availability (of servers) should be high / no-downtime
- ⇒ Scalability or high performance
- ⇒ Disaster Recovery (Backup & Restore)

Architecture:



- Everything in a Kubernetes cluster is a container.
- Every node has a service running known as Kublet service.
 - i) keeps track of machine stats and communicates to master.
- None of our work/containers is done/are run on master

STUFF RUNNING ON MASTER:

1) API Server:

- Entry point into the K8S cluster.

2) Controller Manager (CM)

- Keeps track of what is happening in the cluster

3) Scheduler

- Ensures pod placement

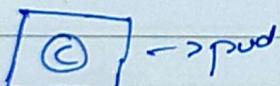
4) etcd

- K8S backup store (All others log their activity here)

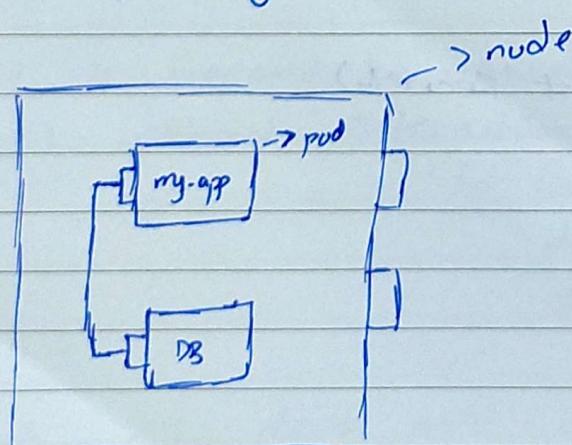
Accessing K8S:

- i) UI
- ii) CLI
- iii) API

Pod:



- Smallest unit in K8S cluster (for what its worth)
- An abstraction over container.
- Usually 1 container per pod (very rare multiple)
- Each pod gets its own IP address



• Pods are ephemeral (temporary). Can be created/destroyed. Part of a normal process of scaling & churn scaling.

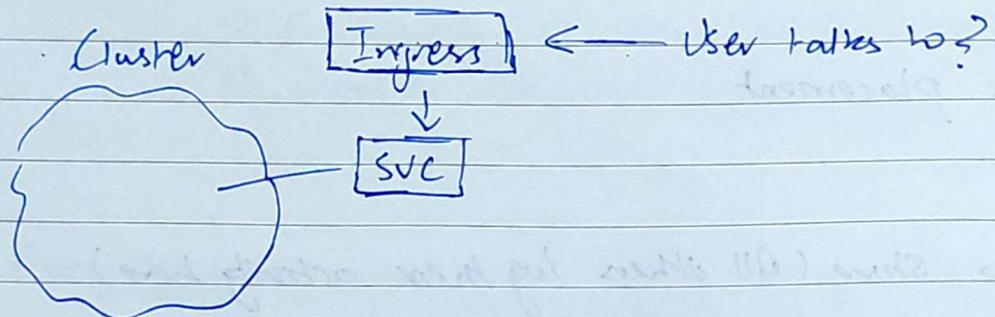
minikube

PAPERWORK

- Each pod gets a new IP on creation.
- This creates a problem, any other pod talking with it would be confused (when scaled / de-scaled)

Lifetime of SVC and Pods are not linked

2-9-2025



Ingress:

Config-map:

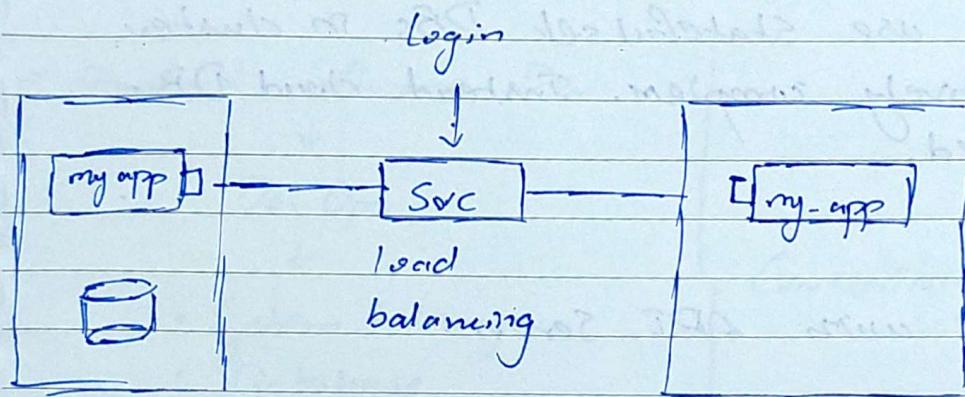
- Contains external configuration to your app
- Counterpart of environment variable / URLs of docker Compose
- Readable isn't a problem

SECRET: (base 64, encoding ≠ secret)

- Isn't human readable.

VOLUMES:

- Can be within or outside the node
- K8S doesn't maintain data persistence by default.
Has to be maintained by user.



• Distribution of load isn't a problem with a stateless login page

- How do we ensure we now have an issue
- What if we have a DB that needs state maintenance

K8s Keywords:

1) Deployment : When you need stateless deployment.

2) Statefulset : When you need stateful deployment.

⇒ SVC is needed to interact with pods, so if a pod runs without interaction, we don't need svc.

↳ for eg a crawler that gets data from website and put in DB.

⇒ Each SVC deals with one type of pods PAPERWORK
(multiple instances of that type)

Statefulset /

⇒ Deployment done through configuration.yaml

Abstraction of pods.

← a blueprint for working

(you can't directly give commands to create pods, you write in files).

⇒ we usually don't use Statefulsets in cluster since that is relatively complex. Instead cloud DBs like S3 are used.

configuration.yaml:

- communicates with API-Server.

Format: (rough)

1 - metadata,

name: _____

2 - spec

replicas:

image:

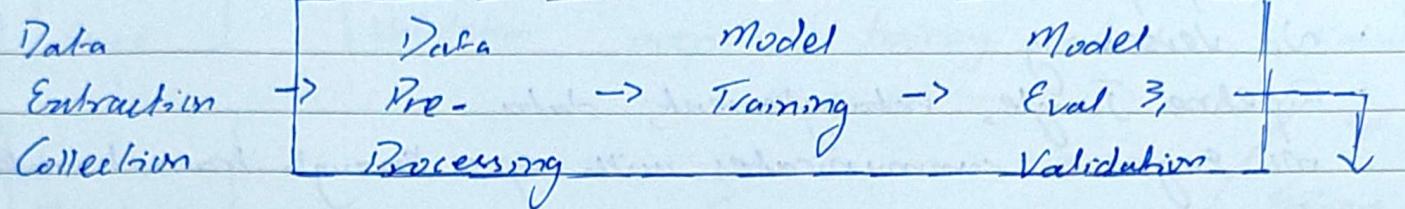
selection:

template:

3 - status

Desired + Status

MLOPS



⇒ Standard Software Dev



DevOps



- System latency
- Database
- Other conventional monitoring stuff
- etc

CONVENTIONAL

Trained Model
↓
(Deploy)

⇒ What's more in ML?

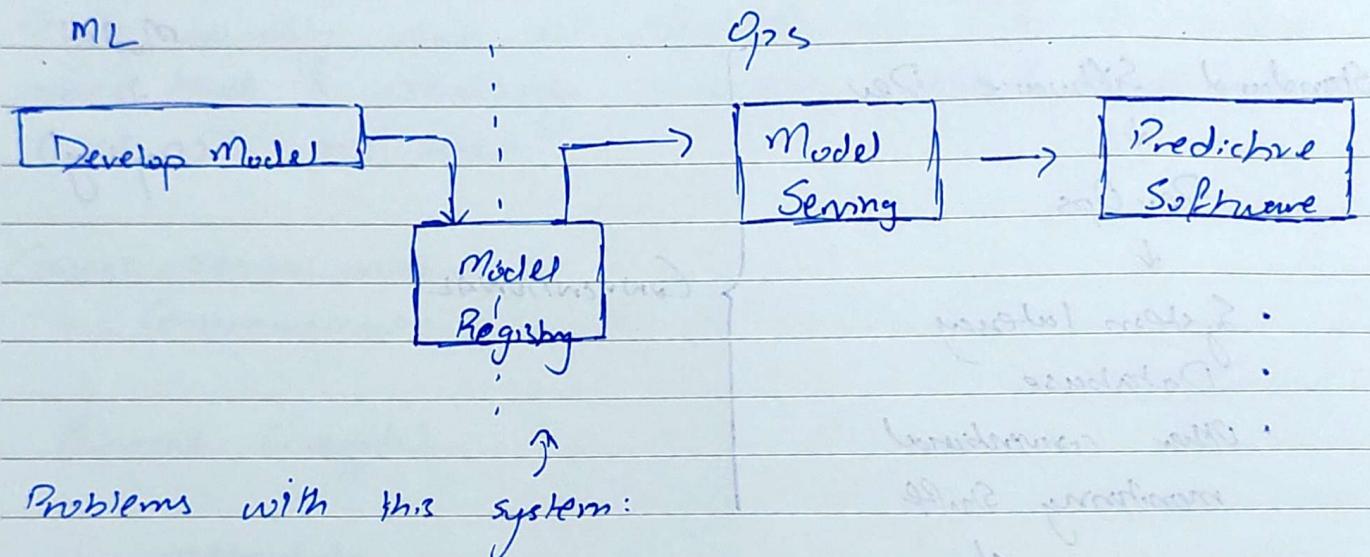
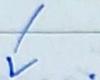


Performance Monitoring

→ Model versioning problems
→ Model serving problems

Level 0: (Manual)

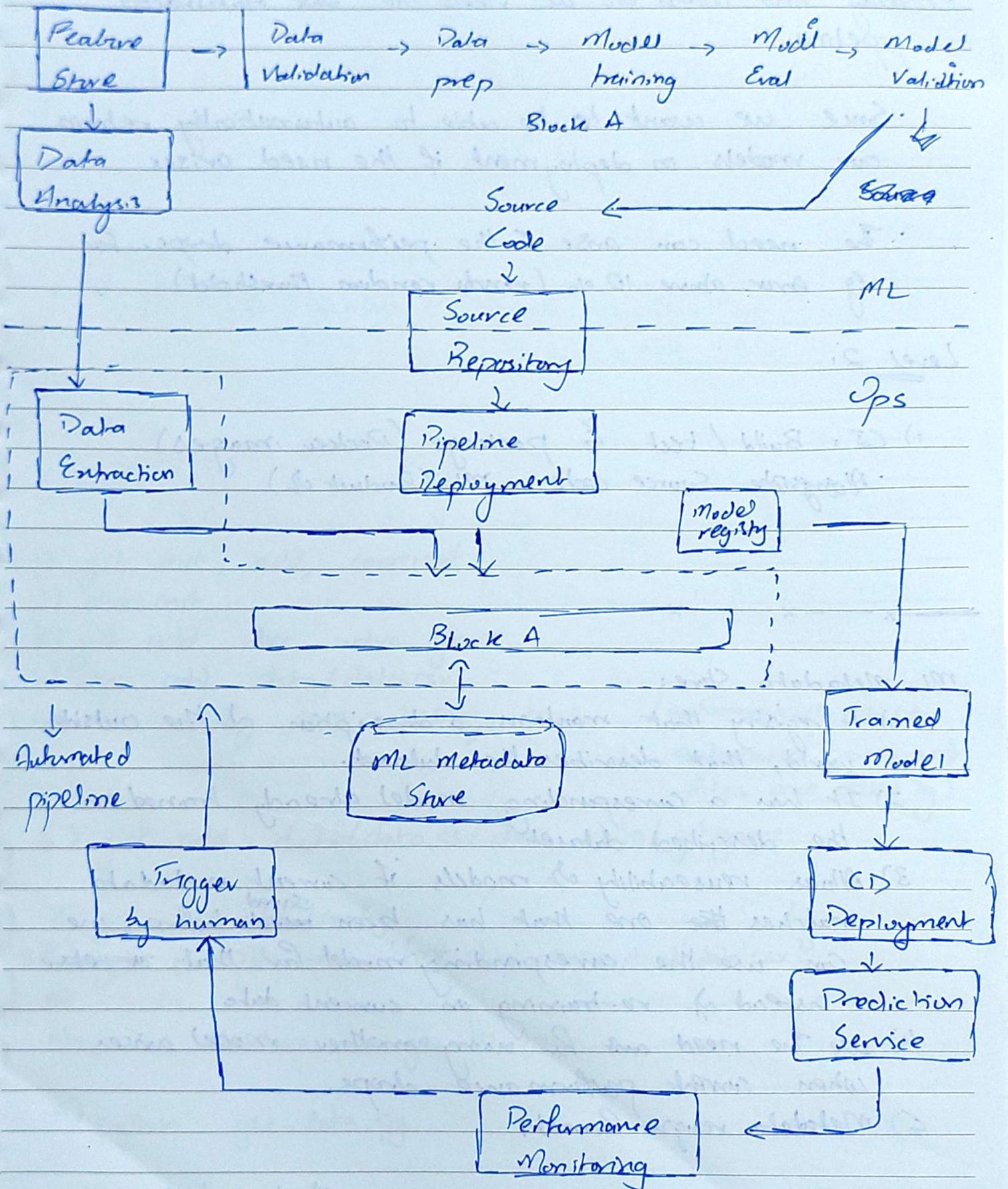
- Data preparation
- No versioning
- Pipeline Jungle, entanglement, data
- ML Engineer communicates with ops through trained models



Problems with this system:

- Ops just gets a file, no mechanism for transfer or return information.
- No Rollback etc. ↴ (Personal addition point)

Level 1:



=> Why and when do we need the ~~Big~~ Automated pipeline

- Since we want to be able to automatically retrain our models in deployment if the need arises
- The need can arise if the performance drops. for eg error above 10% (~~random~~ random threshold)

Level 2:

- 1) CI: Build / test & package (Docker images)
 - Alongside Source code in ML (Product dir)

ML Metadata Store:

- 1) A registry that maintains a description of the outside world, that describes the dataset.
- 2) It has a corresponding model already trained on the described dataset
- 3) Allows reuseability of models if current metadata matches ~~the~~ one that has been ^{saved} ~~used~~ before, we can use the corresponding model for that ~~or else~~ instead of re-training on current data
- 4) If the need ~~arise~~ for using another model arises when current performance drops.
- 5) Metadata rough format:

Date : Metadata: Model

Monitoring tool:

Graphana, prometheus

Data Version Control (DVC)

14-10-25

Why?

Why not git?

Commands:

- 1) git init, add, commit
 - 2) dvc init
 - 3) git add .dvc .dvcignore
 - 4) dvc add data/data.csv
 - 5) dvc config core.autostage true
 - 6) git add data/data.csv.dvc data/.gitignore
 - 7) git add .dvc/config
 - 8) git commit -m "first data commit of data.csv"
- Now we introduce change to track
- 9) python get-data.py

- ⇒ Git hasn't recognized change, although dvc has noticed change due to due to auto-staging
- ⇒ But auto-staging we still have to ~~add~~ new data

dvc status

10) dvc add data/data.csv

⇒ Now git noticed a change in data.csv.dvc (since it is tracking that)

11) git add data / data.csv.dvc

12) git commit -m "data is doubled"

⇒ Now if we wanna go back to previous data

13) dvc logs → find id of what data you want

14) git checkout <id> // Changes data.csv.dvc

15) dvc checkout // Changes data.csv by
// looking at changed
data.csv.dvc

⇒ DVC isn't technically a version control system, the version controlling is done by git, it just reflects the changes since git isn't meant to handle large data. and other reasons

ls -lah data/data.csv

Now we also need a remote location to hold our data, so that others can collaborate (dvc push \exists :stuff)

1) pip install dvc-s3

2) dvc remote add gdrive-remote gdrive://<UUID>

UUID = Universal Unique Identifier
(for google drive)

\Rightarrow Open any google drive folder in the URL
folder/ \backslash <UUID>

\Rightarrow Config file has been updated

3) git add .dvc / config

4) git commit -m "remote storage added"

5) dvc push -r gdrive-remote.

\Rightarrow Failed to push since we didn't authenticate dvc to push on google drive

\Rightarrow Let's authenticate.

1) Google Cloud (go to)

2) Create a project

3) APIs \exists Services

4) Enable API \exists Services

5) Enable google drive API

6) Now you have the dashboard

OAuth
7) Create ~~OAuth~~ client ID

8) Desktop app, name

9) Now you have a client ID

10) Now you have to download

json which has client ID
 \exists secret

6) dvc remote modify gdvire-remote gdvire-client-id <client-id>

8) dvc push -r gddrive-remote

⇒ The other collaborator will now download step one and start creating his/her own.

7) git pull

8) dvc pull // after sorting credentials.

Possible Errors

- 1) first user did git push but not dvc push.
 - 2) Second user does git pull, dvc pull - Will get error related to m3.

DVC Scripting YAML

- We want to automate the retraining process on new data.
- You don't need to manually push every gray DVC file.

1) DVC dvc repro

2)

dvc params show
dvc metrics show
dvc metrics diff

Issues:

- 1) Dependency didn't change, but the output of that stage hasn't changed
- 2) We add an "externally-changed" file in depend dependencies which triggers the whole pipeline.

yaml

get-data: python get-data.py

preprocess:

train: cmd: python train.py, deps:, params: train.params

outs:

metrics:

- metrics.json

PAPERWORK

- Q: Why is train.py stage running on every "dvc repro"?
- => It only updates md5 values when stage is completed.
 - => If md5 hashes are different, the stage is run again when "dvc repro"
 - => but when Stage was run it wouldn't complete since it didn't find all the files it needed.
 - => So the stage didn't complete, the hash didn't synchronize
↳ the train stage is triggered every time because of the hash difference

GitHub Action yaml explanation:

run: | train.py

make install

dvc repro

git fetch --prune

dvc metrics diff --md master > report.md

gh pr comment \$((github.event.pull-request.number))
-F report.md

Explanation:

1) Install dependencies

2) Re-run training (whole yaml of dvc). We now have new statistics.

3) git fetch --prune Fetches all branches except the one it is running on (to compare old metrics with new one).

4)

- You can have a specific file which is purely for triggering GitHub actions.
- ~~like~~ when you need to repro, which is done in GitHub actions, we can push to that file which will trigger GitHub Actions.
- Later we'll automate the trigger based on monitoring. (Graphene etc.)

EXPERIMENT TRACKING

23-10-25

- ⇒ We have hundreds of parameters (hyperparameter level)
- ⇒ So infinite Combinations of hyperparameters
- ⇒ How to find the best ones? (all of them are not good)

(Concepts: we have to track all the metrics. Now what?)

- 1) ML Experiment
- 2) Experiment Run
- 3) Run Artifact
- 4) Experiment metadata

WHY IS IT IMPORTANT:

- 1) Reproducibility
- 2) Organization
- 3) Optimization

WHAT'S THE PROBLEM WITH SPREADSHEETS?

- 1) Error prone
- 2) No standard format
- 3) Visibility & Collaboration

SOLUTION 1: ML Flow (Python package)

- 1) Tracking
- 2) Models
- 3) Model Registry : Model versioning is
- 4) Projects : A complete unit which is pushed on GitHub and can be run by just one command. No manual repo cloning and stuff.

about mlflow

Storing the stuff:

1) Local Repositories on your device. Different repositories for different projects and a dashboard tying them together (problem: not accessible for other people).

2) Server based. Can use through: (live ip needed)

- 1) ~~live ip~~
- 2)

mlflow ui

mlflow server --host=127.0.0.1 --port:

Demo of mlflow: (Python file)

- 1) ~~live~~ pip install mlflow
- 2) Import
- 3) Initialize (where you can also set server ips or local)
- 4) Normal train
- 5) Use mlflow syntax to log what you want
- 6) Also, watch

AUTO LOG:

~~Tempo~~ => Don't have to manually set ~~which~~ which parameters
to log and any code
=> Just call autolog, it will detect model & do it
for you

Tracking:

Hyper parameter tuning using GridSearch in MLFlow:

⇒ Example

- 3 hyper parameters
- 3 values for each
- Run one
- Compare 3, pick best-combo (automatically save track as well).

Plots:

- 1) Without steps for measures
- 2) With steps for showing improvement or reduction in metrics.

Projects:

- 1) MLproject
 - 2) conda.yaml
 - 3) python-env.yaml
 - 4) requirements.txt
 - 5) train.py.
- # project config file } Standardize environment.
 }
 - Run from anywhere
 - Experiment in diff envs environments

|| mlflow run . -e projects -P n_estimators=200 -P max_depth=10

|| Can also run a github repo (and its commit version) which are compliant with MLFlow

|| mlflow run (can run through python as well).

|| Can give any environment on runtime: local / conda,

Models

⇒ Saving models (Pytorch, sklearn...)