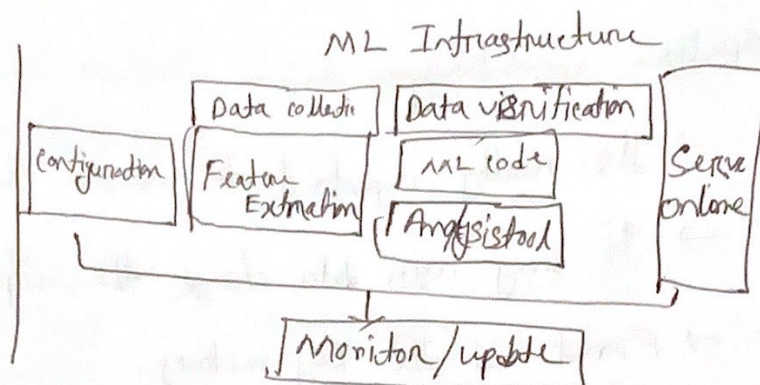
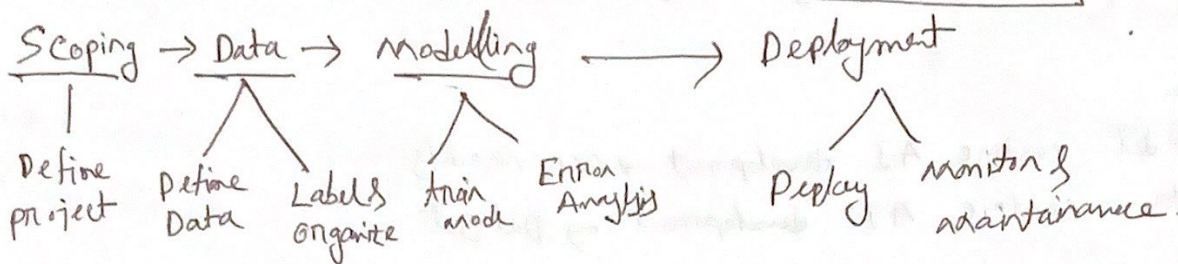


pipeline → W1

→ Data drift problem

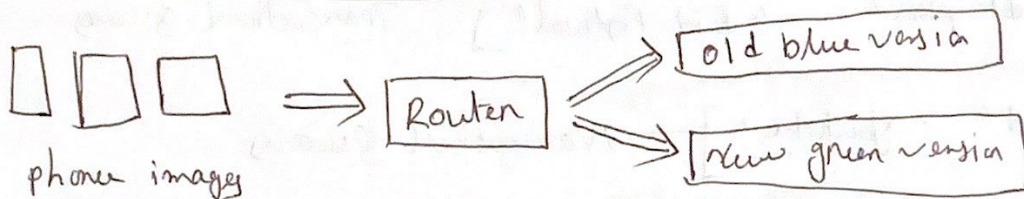
ML ops Lifecycle



Data Drift → your data change after your model deployed.

Canary development → deploy for a subset of users.

Blue green development



Monitor metrics | monitor dashboard

→ set thresholds for alarms

→ Adapt metrics and threshold over time.

~~Manual~~

Manual retraining

Automatic retraining



→ pipeline monitoring

→ Its really important to monitor the metrics.

→ If any user data change the output will change.

→ Find out the key metrics.

W2

→ Model - centric AI development → NN models

→ Data - centric AI development → Dataset

⇒ AI system ⇒ code/model + Data

→ web search example

"Apple pie recipe"

"Latest movie"

"Wireless data plan"

"Eid Festival"

⇒ Informational & Transactional Queries

"Stanford" "Youtube" ] — Navigational Queries

⇒ Sanity check for code & Algorithm

↳ try to use small amount of data at first

⇒ Error Analysis →

→ Add/improve data for specific category

→ Skewed dataset (imbalanced data)

→ Try to Brain storm what could go wrong?



Data Augmentation → create realistic example ~~that~~ that human can check easily.

Experiment tracking →  
Algorithm / code  
Dataset used  
Hyperparameters  
Results

Tracking tools →  
Test files  
Spreadsheet  
Experiment tracking system  
Tools  
Weights and biases  
→ Comet  
→ MLFlow  
→ Sage make studio

F1 Score recommended  
for skewed dataset

W3 → Data

Data →  
Unstructured data  
Structured data

⇒ Data and Label consistency

⇒ Human Level performance (HLP) → Estimate Bayes Error / irreducible error

→ Beat the HLP is good for research purpose  
but for production it's better to improve HLP so that model  
outperformance increase

# Data → Don't increase data by more than 10x at time

→ PoC (proof of concept)

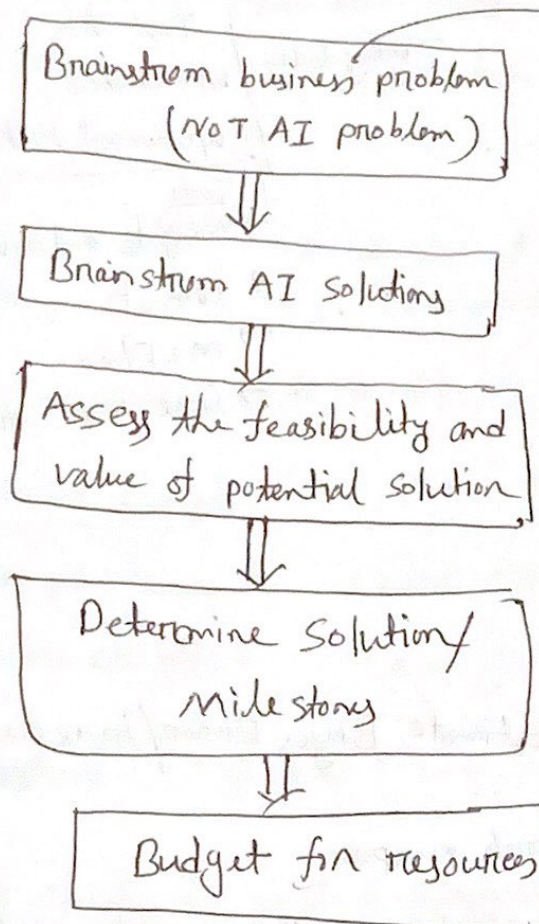
→ TensorFlow Transform, Apache Beam, Airflow

# Meta-data, data, provenance and lineage ⇒ Data pipeline  
(File types, origin points) where data comes from Sequence of steps.



- A high HLP metric implies good label consistency.

## # Scoping process



Some questions to ask

→ what thing you wish were working better

- Increase conversion
- Reduce inventory
- Increase margin (profit)

→ make sure your feature map to prediction

↪ give past purchase, predict future purchase  $x \rightarrow y$

$x \rightarrow$  Given DNA info, predict heart disease  $(x \rightarrow y)$

→ Query-level accuracy



## Specialization #2

→ Machine Learning Data Lifecycle

### TF Metadata (w1)

→ MLMD is a library to track of full lineage of entire ML like data ingestion, preprocessing, validation, training, evaluation, deployment and so on.

PoC (proof of concept):

→ Goal is to decide if the application is workable and worth deploying

→ Focus on getting the prototype to work

→ pre-process data manually if needed. But take extensive notes/comments

### # Machine Learning pipeline (ML pipeline)

→ Pipeline Automation → Airflow, Argo, Celery, Luigi, KubeFlow

### # TensorFlow Extended (TFX)

E2E platform for ML pipeline

Data ingestion → Data validation → Feature Engineering → Train model → validate model → push it → serve model.

⊛ TensorFlow Data validation (TFDV) ⊛ Chebyshev Distance

$$D_{\text{chebyshev}}(x, y) = \max_i (|x_i - y_i|)$$