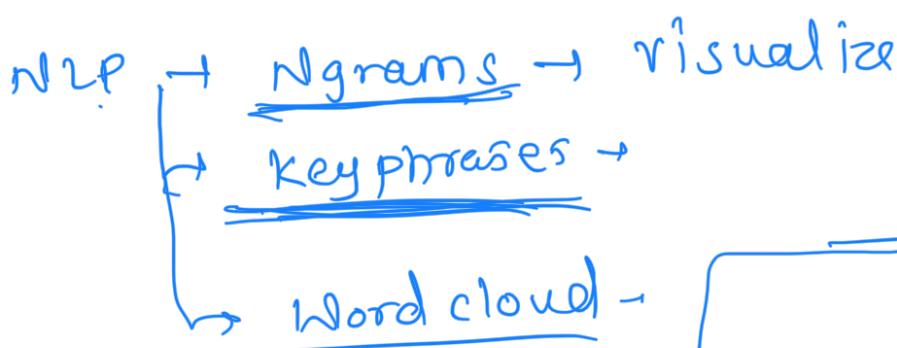
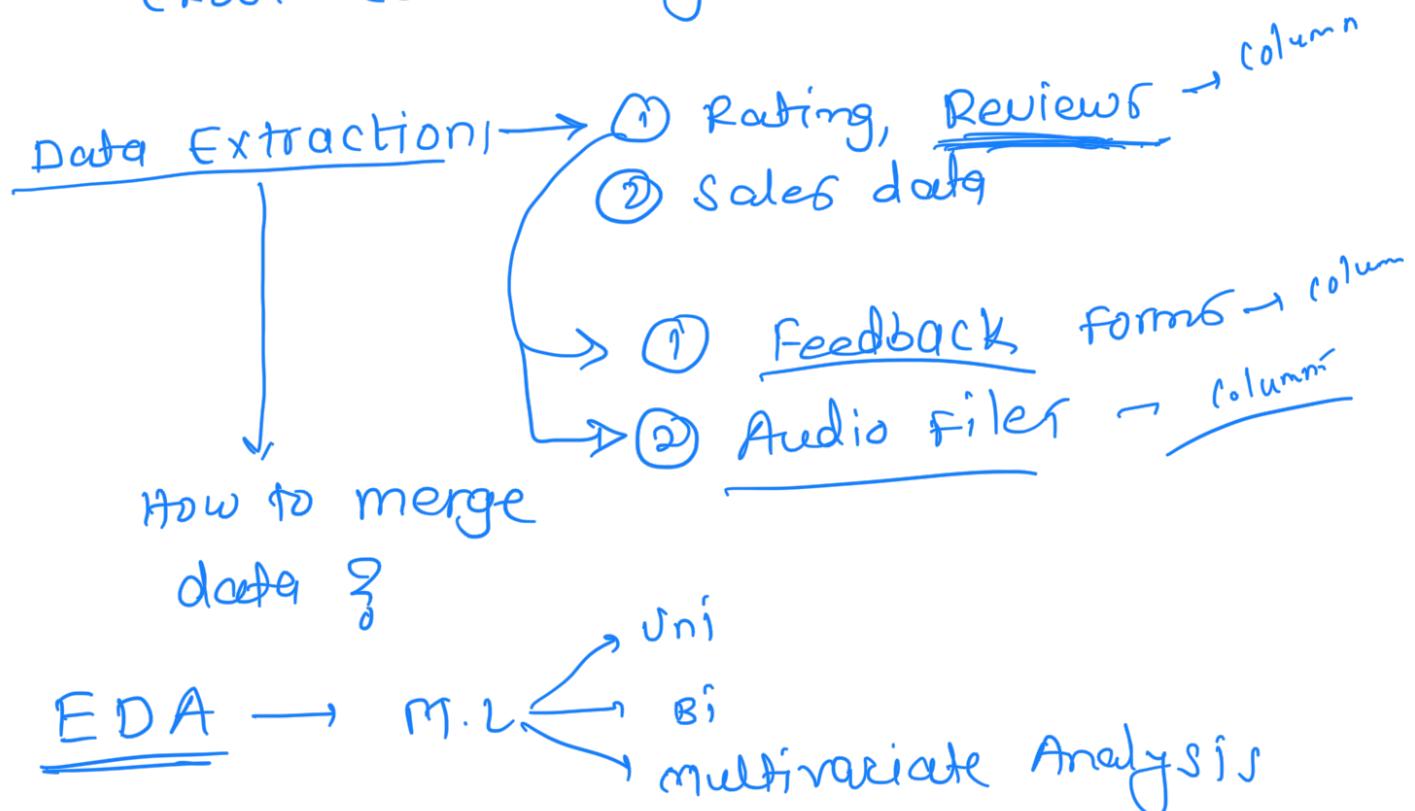


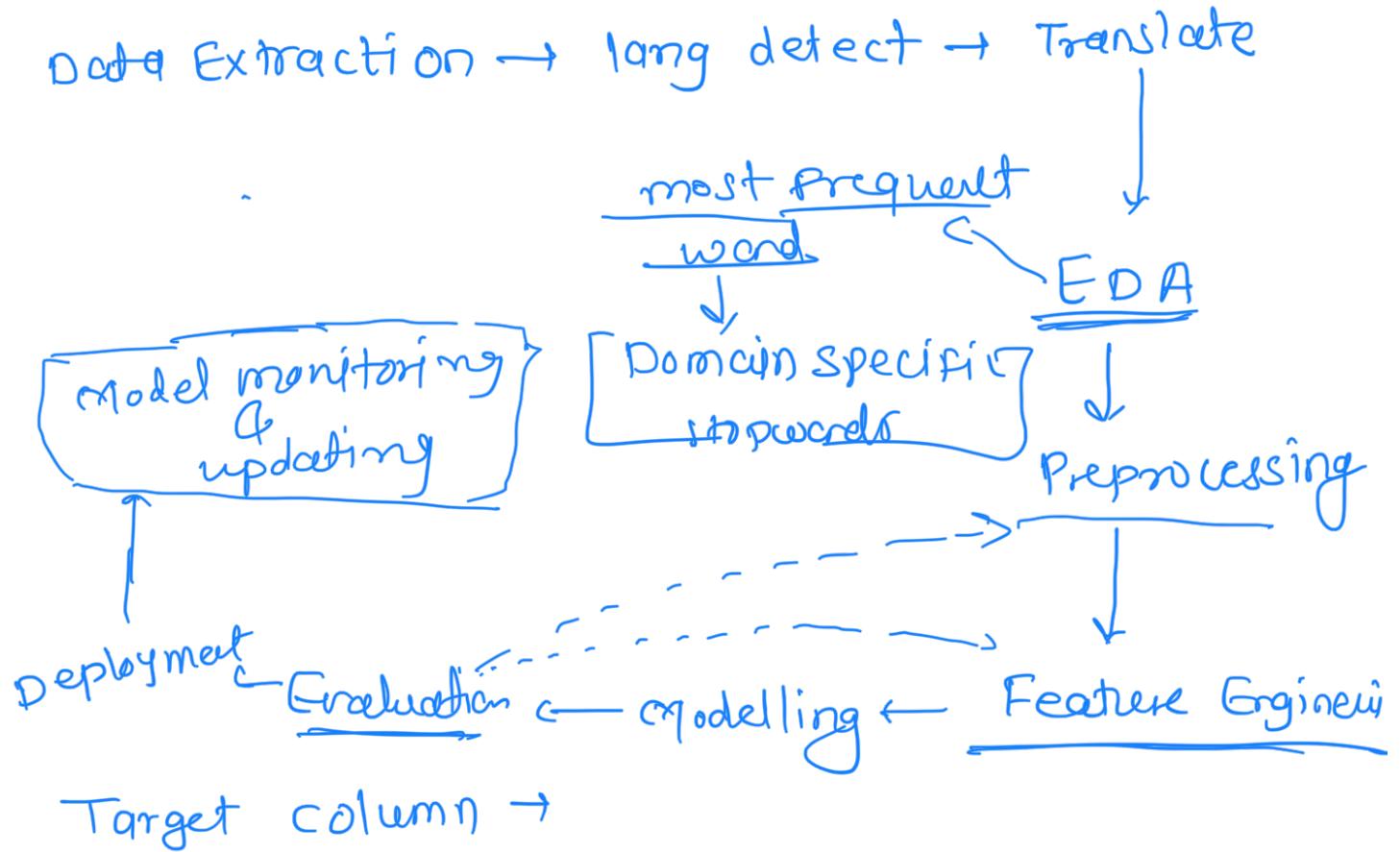
## Problem Statement :-

① client wants to check if customer are satisfied with product or not. If not then what is reason of that (Root cause Analysis).



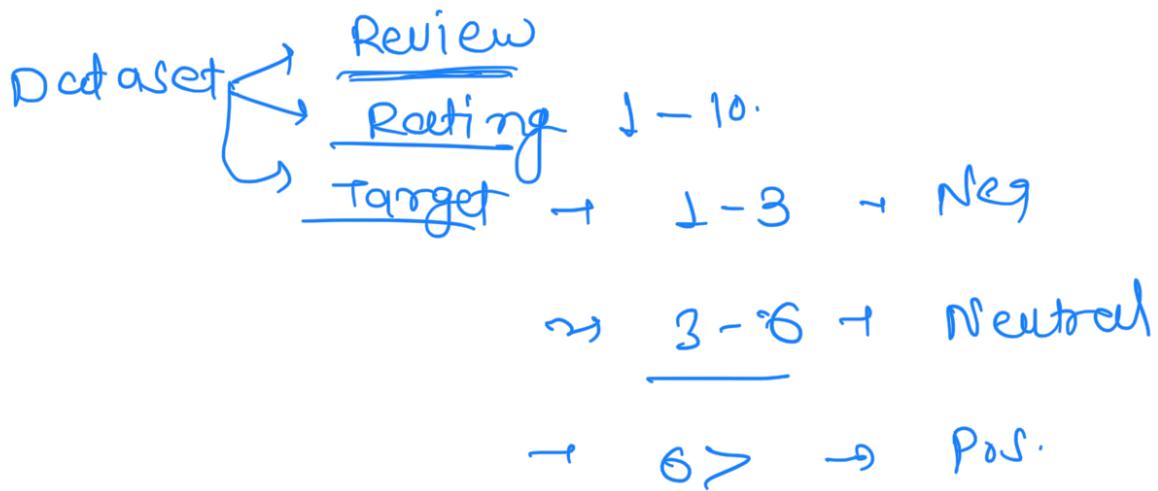
Reviews = multiple language





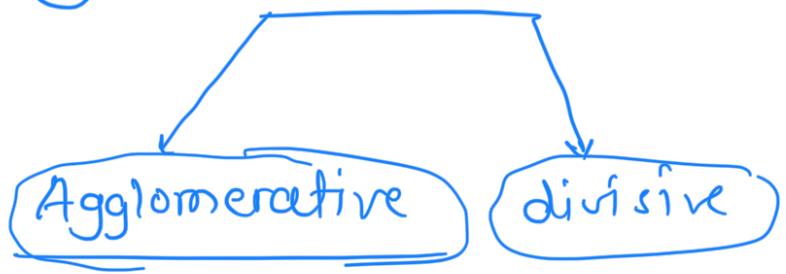
Target column →

- ① clustering
- ② creating target columns with help of other columns in dataset

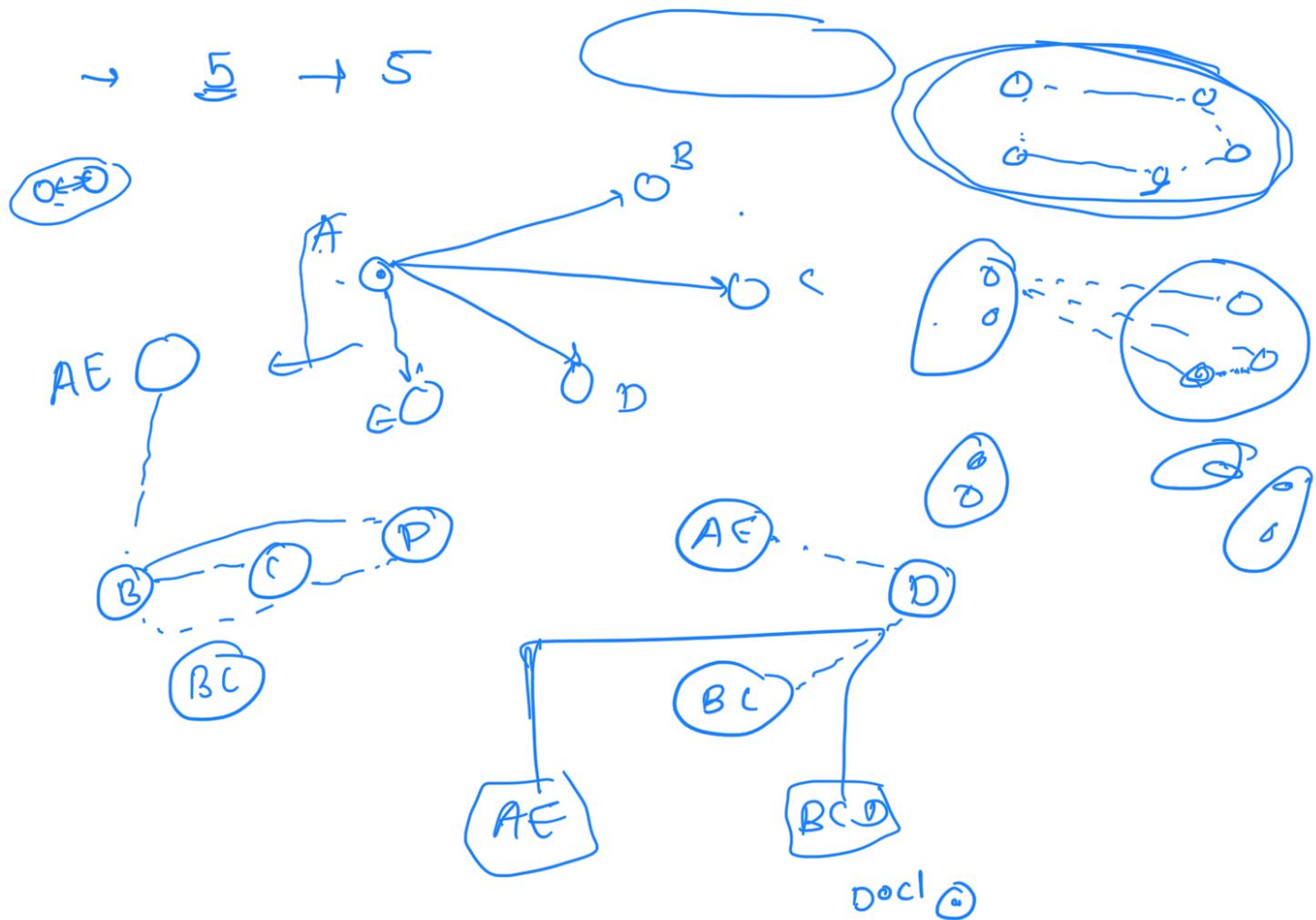


clustering | → ① kMeans

↳ ② Hierarchical



→ 5 → 5

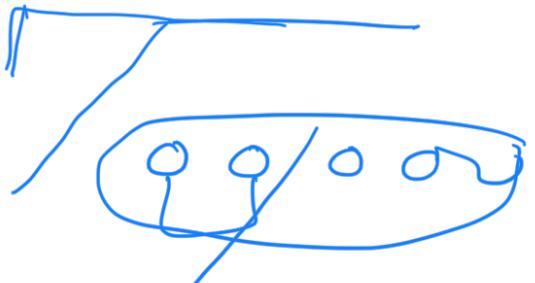


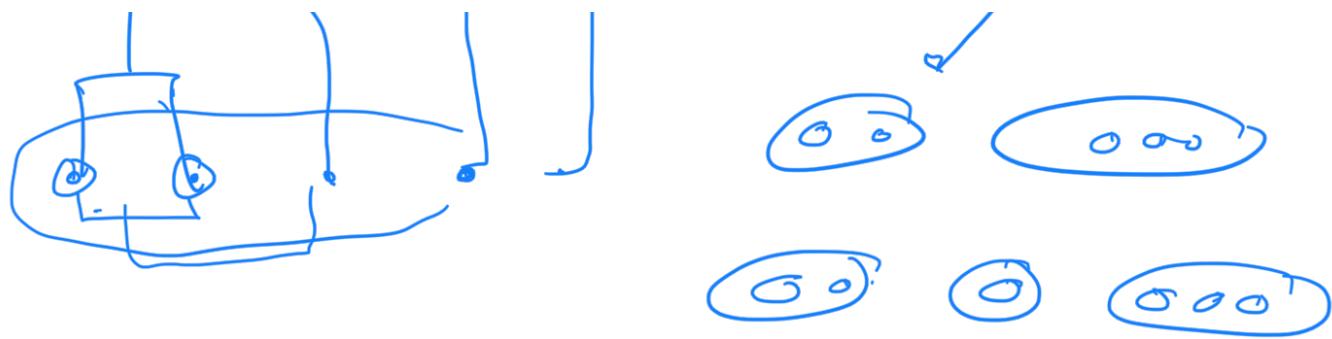
Text → Numerical

Agglomerative



Divisive

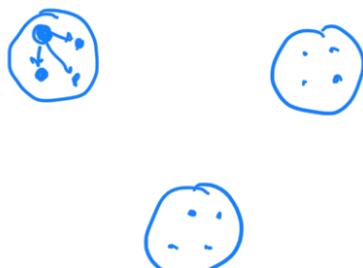




## Evaluation Metrics for clustering

' Silhouette Score' = -1 to +1

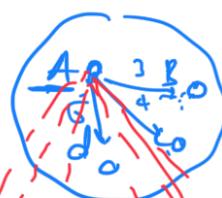
a)



$$= \frac{(b-a)}{\max(a,b)}$$

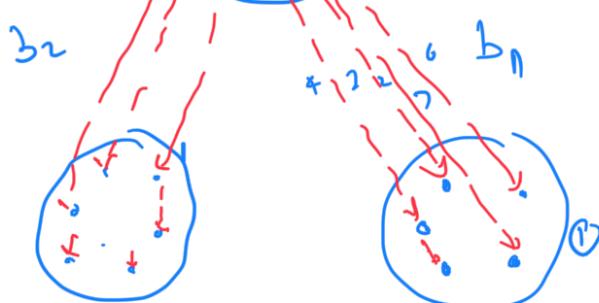
-

$$a =$$



$$a = \frac{3+4+5}{10} =$$

$$b =$$



$$b_1 = \frac{4+3+2+6}{4} =$$

$$b_2 = \frac{7+2+6+4}{5} =$$

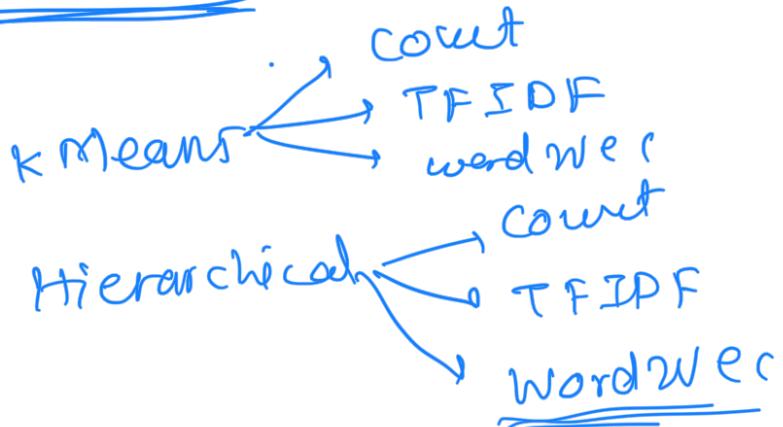
$$b = \min(b_1, b_2) =$$

$$\text{score} = \frac{b - g}{\max(a, b)}$$

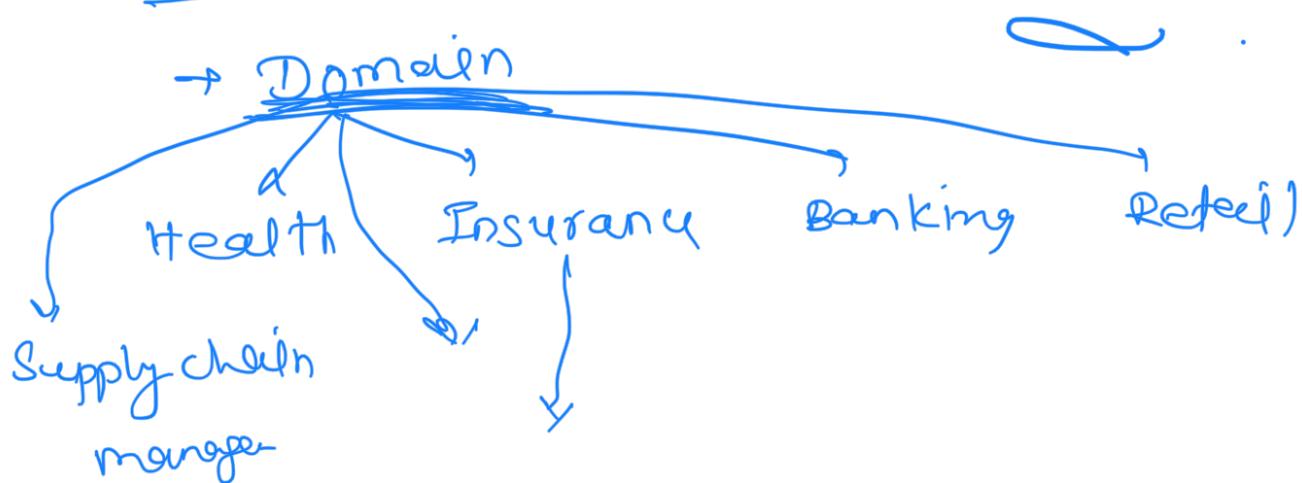
$g = \underline{\text{cohesion}} = \text{mi}$   
 $b = \underline{\text{separation}}$

$$= -1 \text{ to } +1$$

Yellowbricks → silhouette visualize



## ② Document classification



Insurance → Aadhar, PAN, Doc, reporting, POC,  
 Identification  
 Proof, RC Book

Column → Civil score,

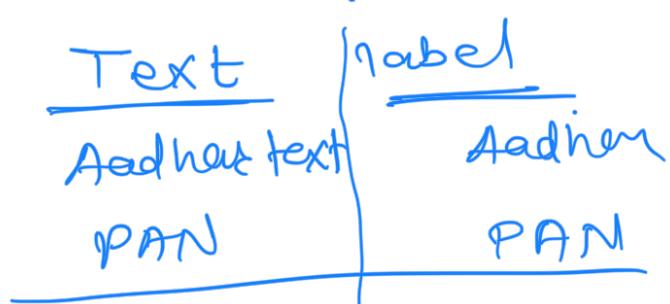
DOC → Training data

Data Extraction

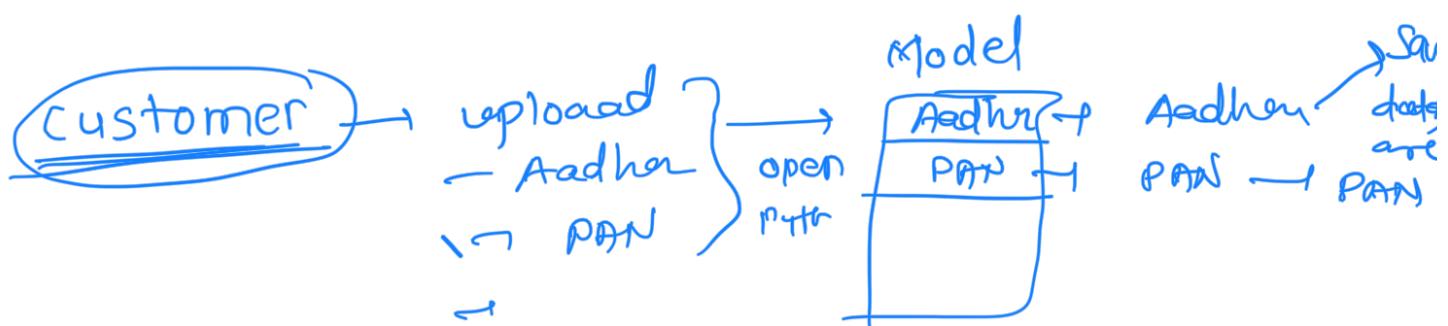
↳ Image + pytesseract, OpenCV



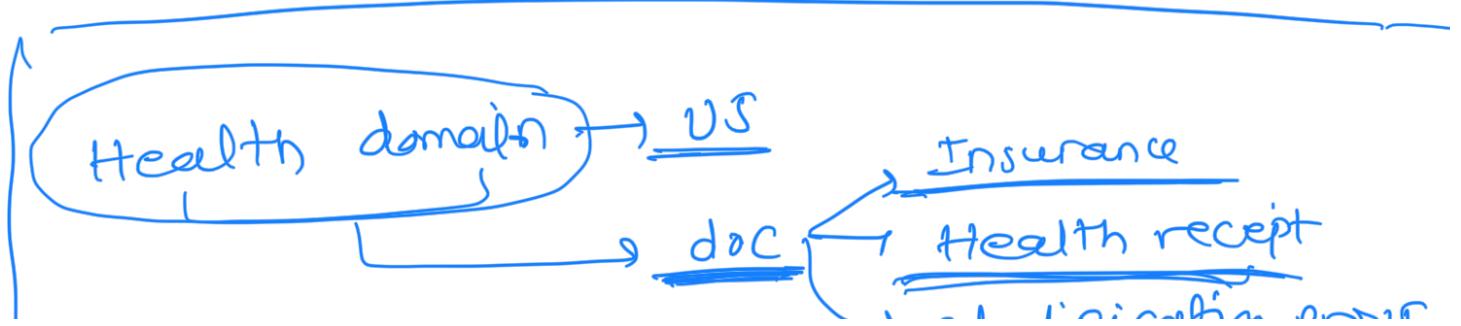
Regular Expression



client



Aadhar → model.predict(Text\_data) = PAN card



→ Identification point

Key phrases = ① Poor customer Service,  
Faulty product,  
② damaged product,  
Bad quality of product,  
Customer service is not good,

### Ontology creation

{ "Poor Customer Service": [ "customer service" <sup>keyphrases</sup>  
is not good; ] ,  
"damaged Product": [ "Bad quality of  
product",  
"Faulty Product" ] }  
}

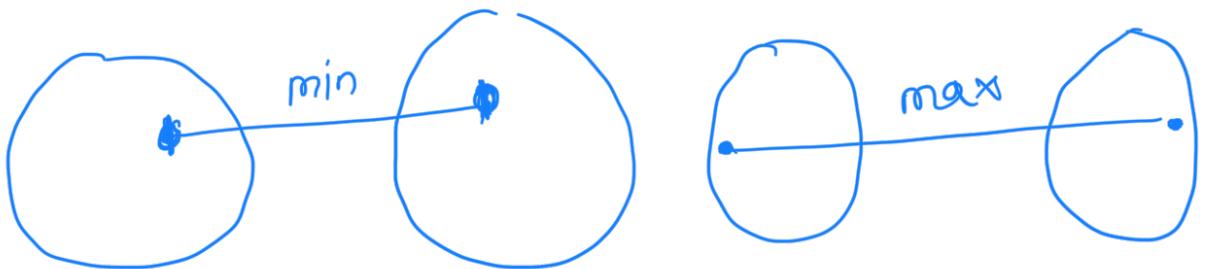
Mapping = Training data.

" Tent "                  " Root Cause "  
[ Faulty Prod. ] = damaged Product

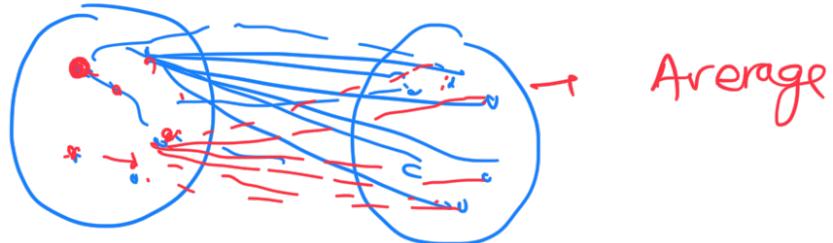
Pn

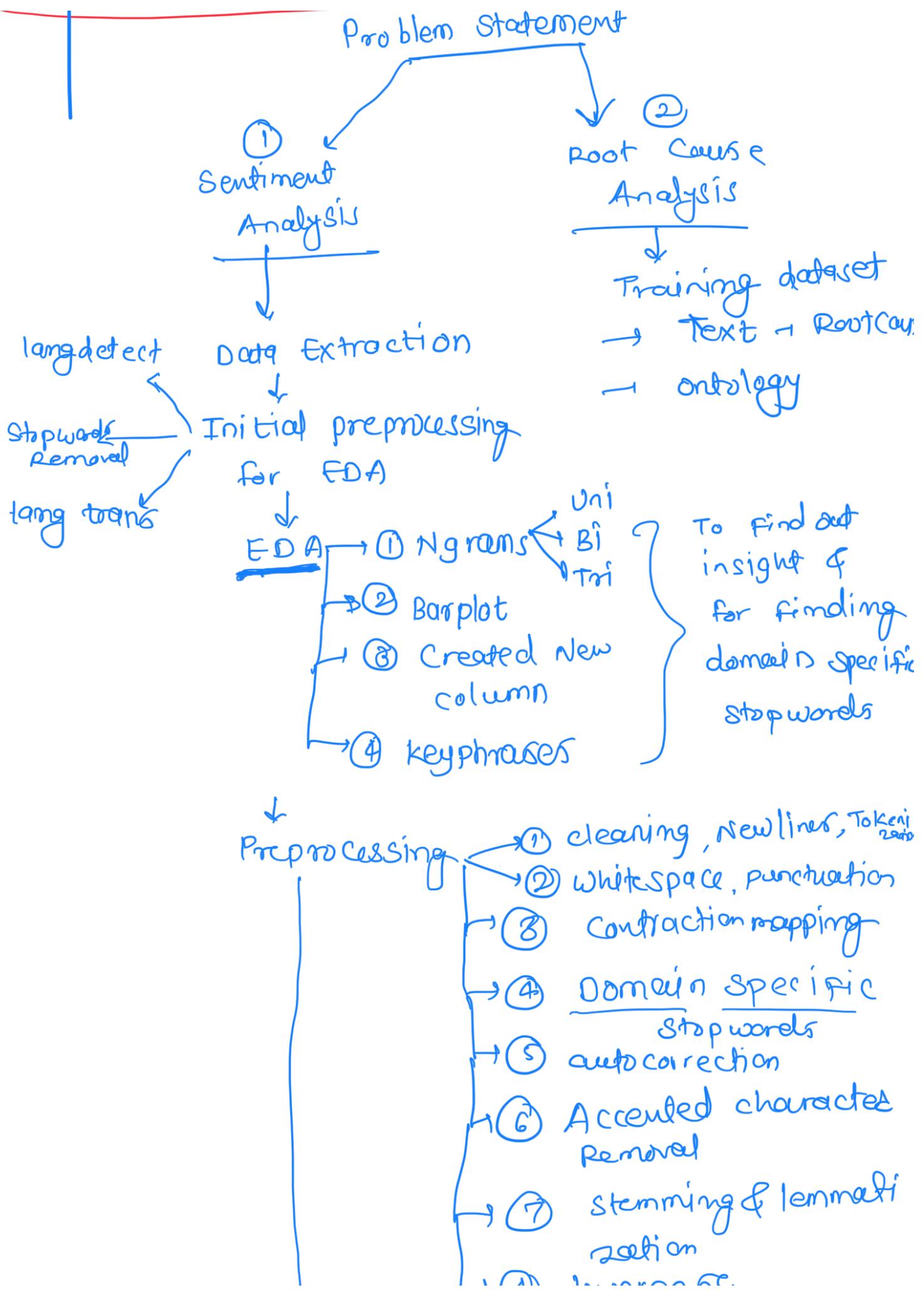
	YAKE	RAKE	Page	Kegph <sup>h</sup>
Text	-	-	-	

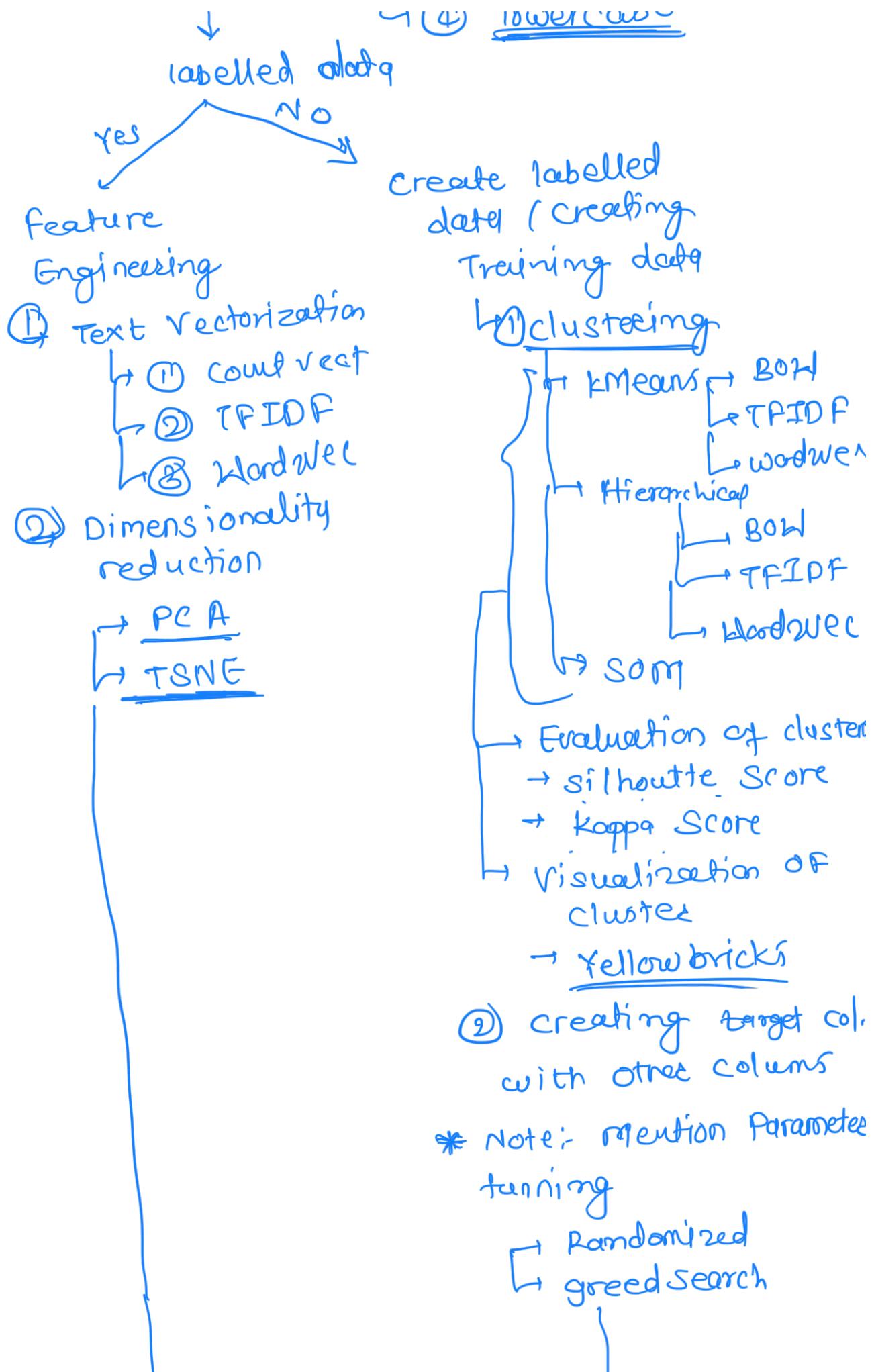
linkage = single, complete  
↓  
max

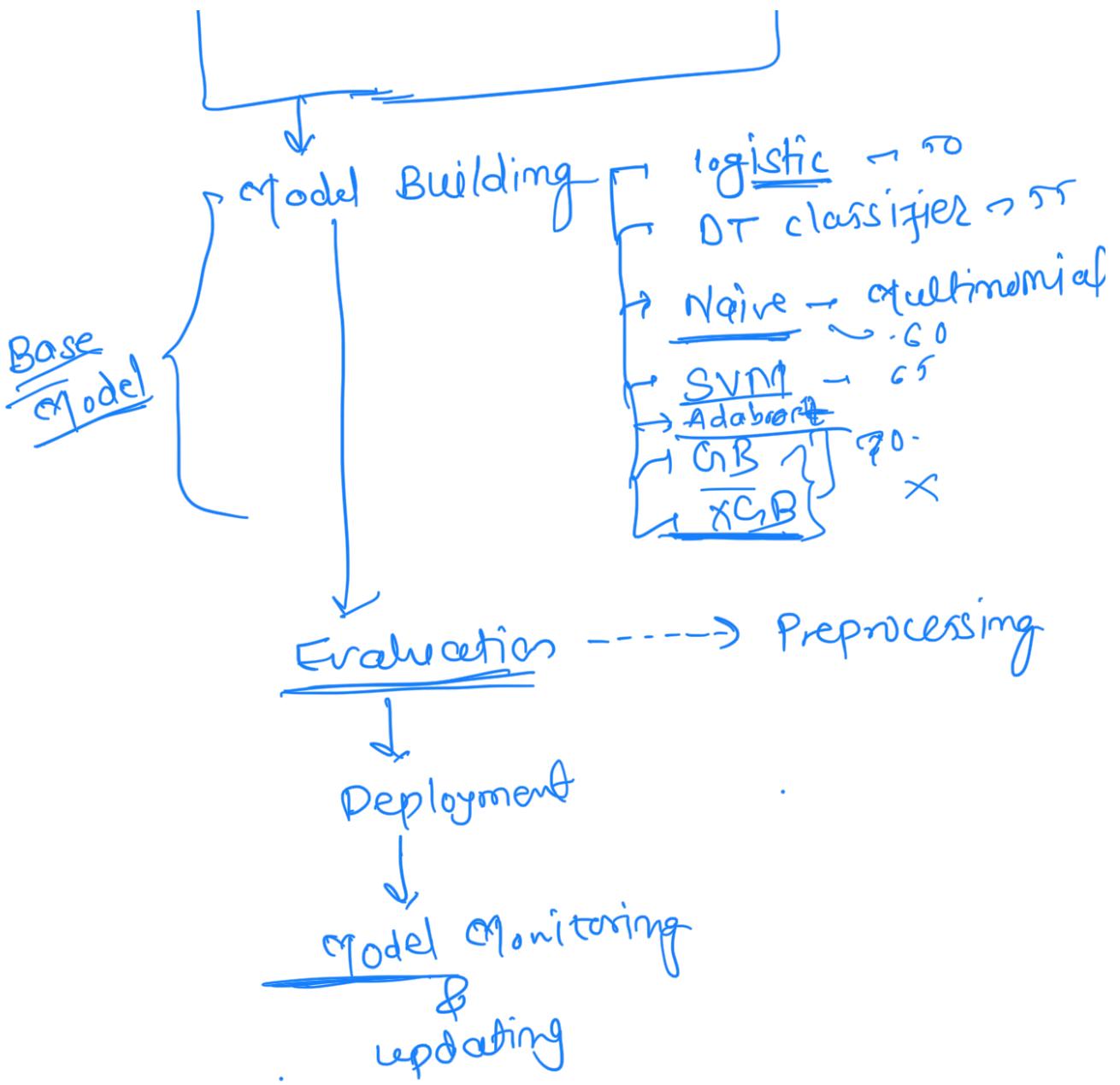


Average



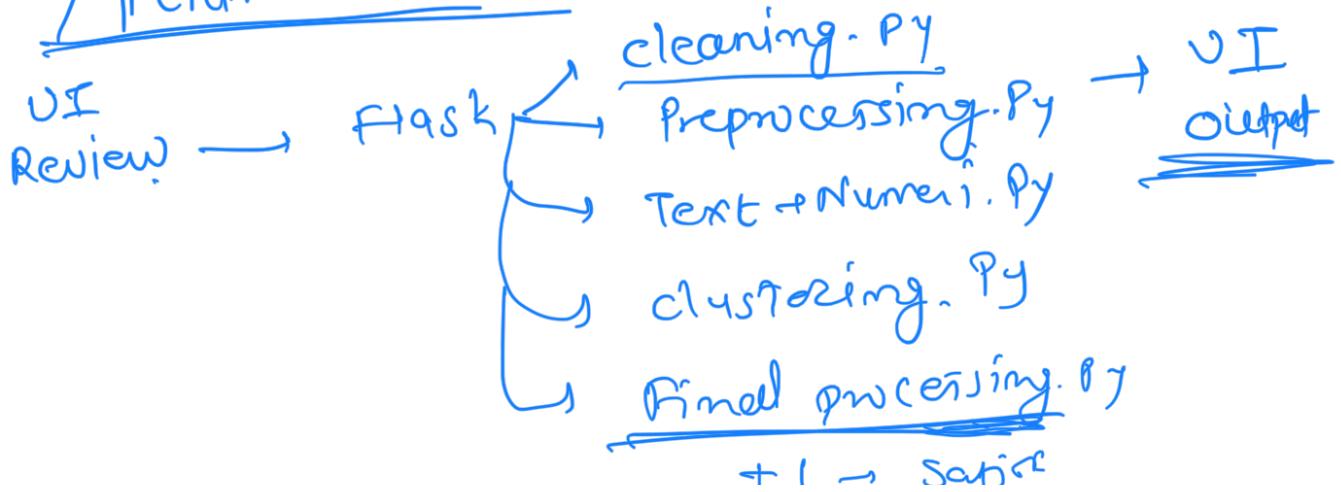






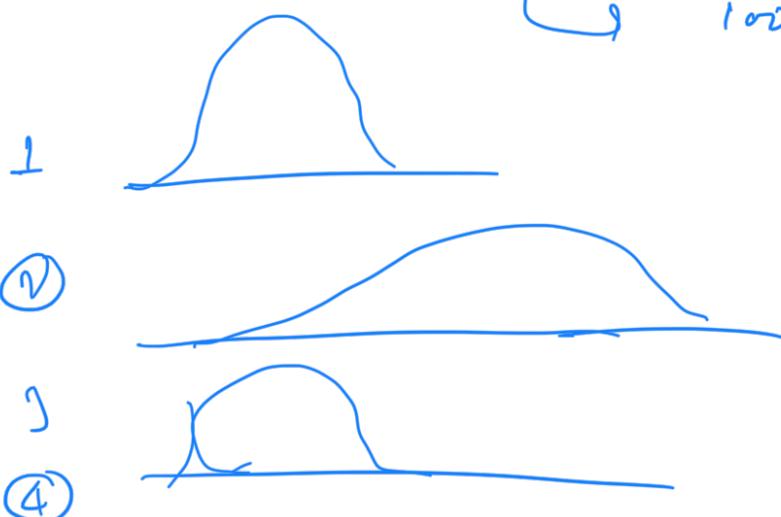
\* SymSpell → for auto correction

## Architecture



\* Monitoring & Model updating:

data drift  $\rightarrow$  column = Normally  
skewed  
column = 30-40  
100-150 } Model Impact



$\rightarrow$  Model Switching  $\rightarrow$  Adaboost  $\rightarrow$  85

SVM  $\rightarrow$  80  
Naive  $\rightarrow$  75

Adaboost  $\Rightarrow$  80  
SVM  $\Rightarrow$  85  
Naive  $= 70$

script  $\rightarrow$  max (SVM, Naive, Adaboost)

Jump

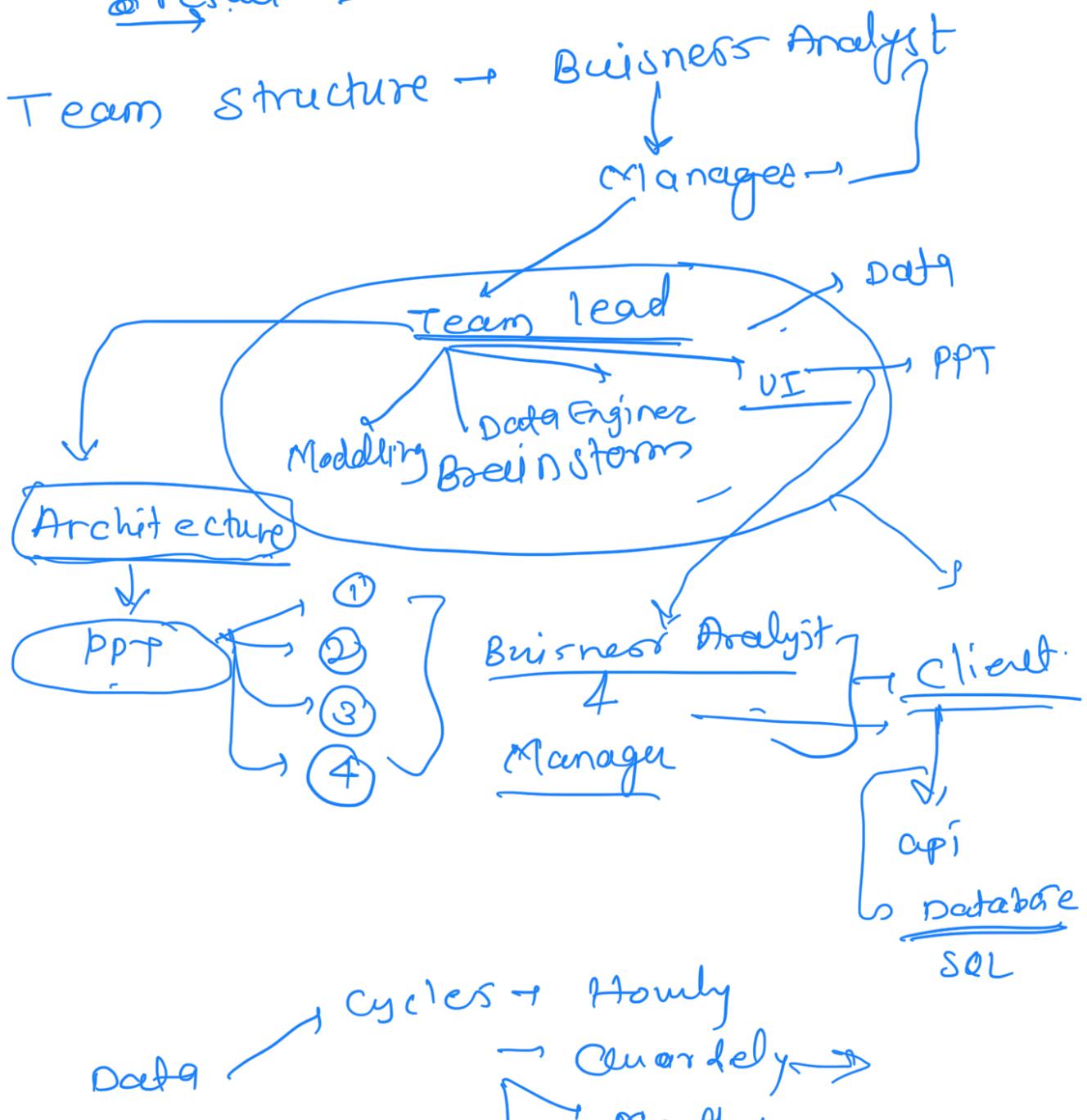
Jupyter = Ada  $\rightarrow$  Pickle, joblib

Frank

@app

import cleaning, clustering.. etc  
- cleaning (data)  
- cln

@result =



Scheduler → api.cloud → data saved

