Problem statement

$\downarrow$

Problems $\longleftarrow$ Data $\longrightarrow$ Data formats

- json
- images
- .txt
- .csv

$\downarrow$

Initial preprocessing
- Language detection
- Language translation

$\downarrow$

Ngrams $\longleftarrow$ E D A

Ⓐ Ngrams
- Unigram
- Bigram
- Trigram

$\longrightarrow$ Domain Specific stopwords

$\longrightarrow$ Key phrases - Root cause analysis

Ⓑ Word cloud

Ⓒ Key phrase extraction

RAKE        YAKE

Preprocessing

↳① Remove spaces,

→ ① Remove spaces, newlines, blanklines

→ ② Contraction mapping

→ ③ Handling accented characters

→ ④ Cleaning

ⓐ Tokenization
- Sentence
- Word
- Whitespace
- Regex tokenizer

ⓑ Remove Punctuation

ⓒ Remove stopwords
- language specific
- domain specific

ⓓ len (word) > 2

ⓔ Normalization

→ ⑤ Autocorrection
- autocorrect
- text blob
- Symspell

→ ⑥ Stemmer
↳ Lemmatizer

↳ (6) ~~Stemmer~~
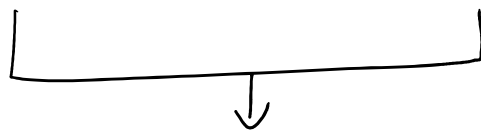Lemmatizer
Porter stemmer
Snowball stemmer

↓

Clean data

Target — yes ↙

No Target ↘

evaluation :
- Silhouette score
- Silhouette visualizer
- Dunn index
- Kappa index

① Clustering
② K means — Count Vectorizer
        — T F I D F
        — Word 2 Vec
(b) Hierarchical clustering
(c) DB Scan
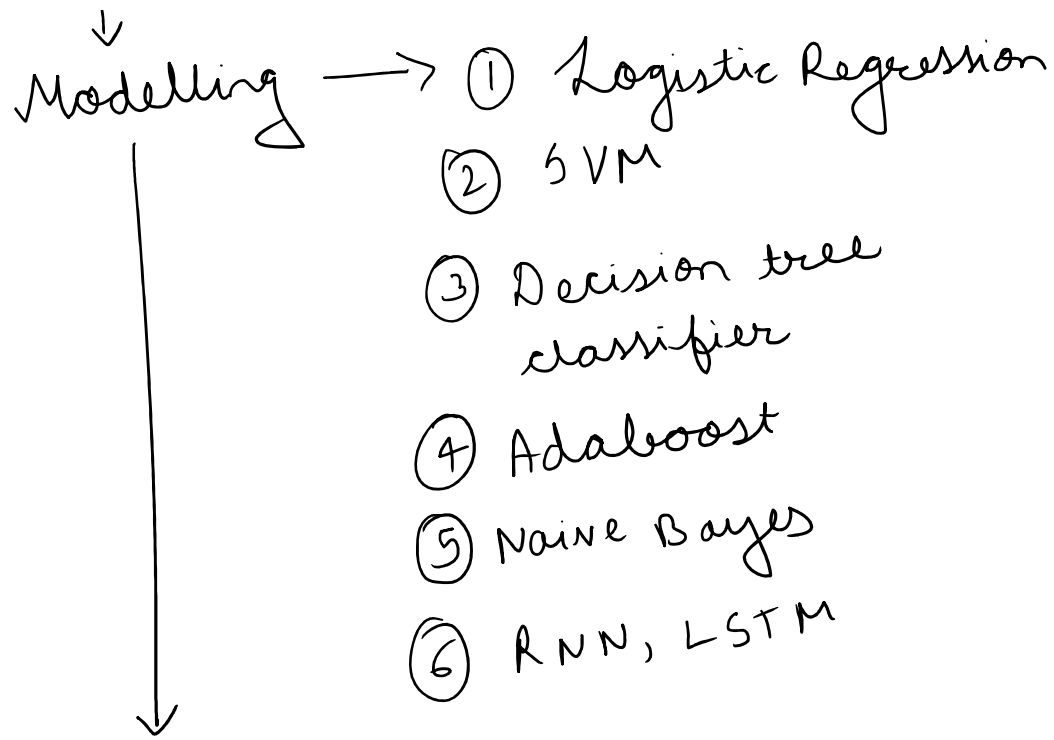(d) SOM → Self Organising maps

② Supportive column
    — Star / Rating column

↓

Training data

↓

Modelling ⟶ ① Logistic Regression

Modelling ⟶ ① Logistic Regression

② SVM

③ Decision tree classifier

④ Adaboost

⑤ Naive Bayes

⑥ RNN, LSTM

↓

evaluation

↓

create pickle file

↓

Deployment

↓

CICD

Model Build ⟶ Deployment

3 - 6 Months

CICD ⟶ Continuous integration, Continuous deployment

Production model $\longrightarrow$ Customers use

Non- Production model $\rightarrow$ updates/changes

Whatsapp $\longrightarrow$ Playstore

Non - Production $\rightarrow$ Developer
model                         updates/ changes