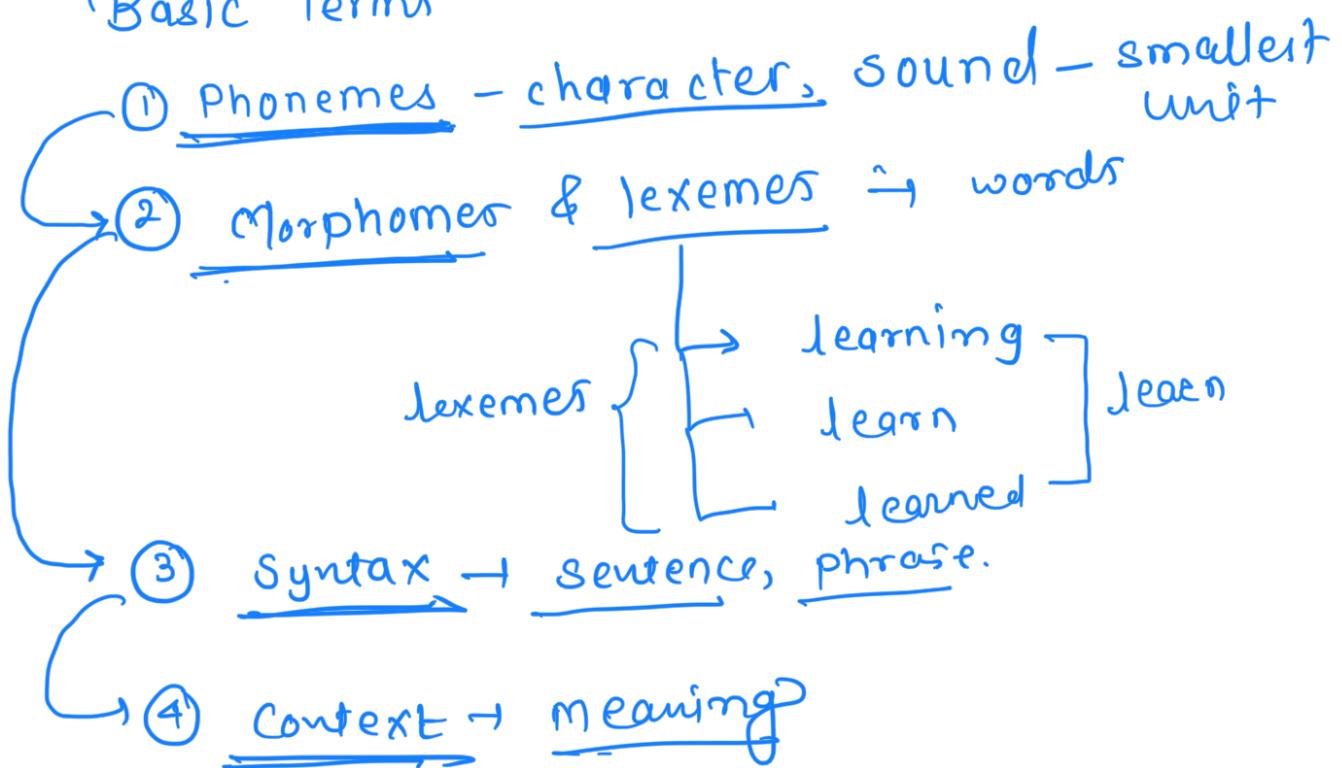


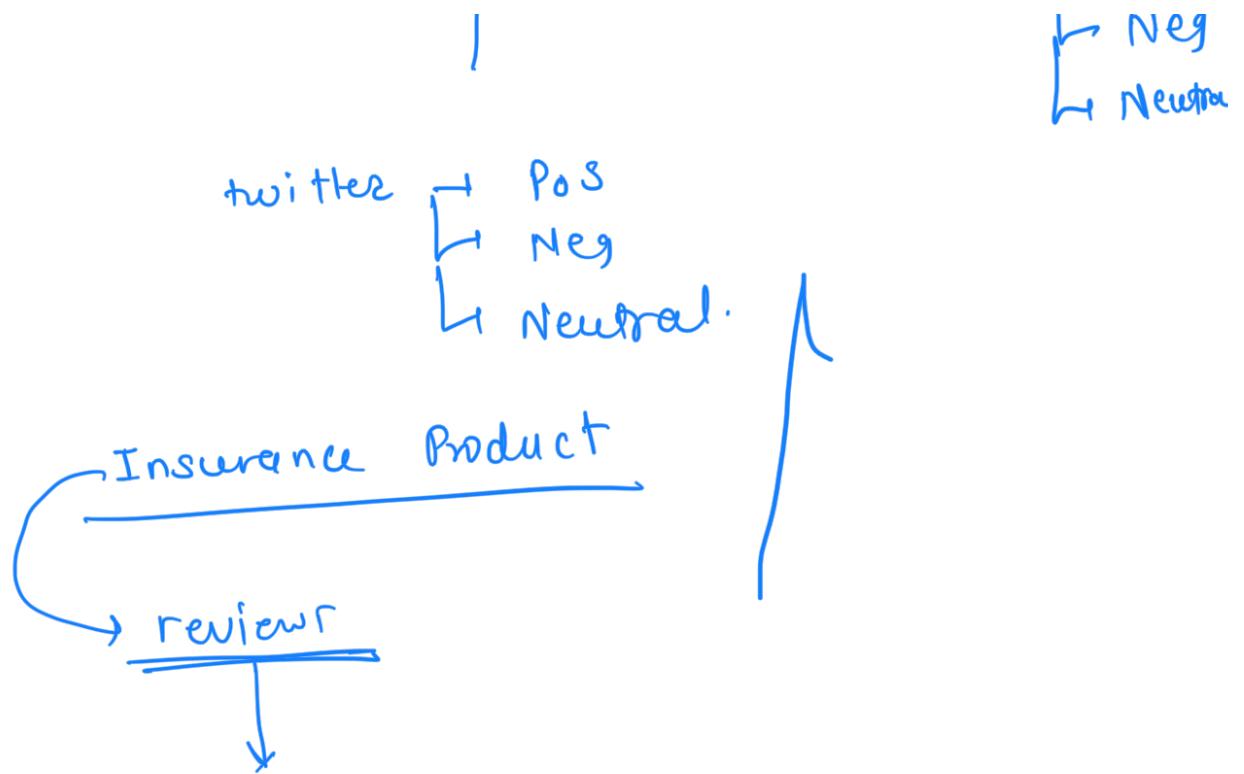
Basic Terms



Applications

① Sentiment Analysis → classification

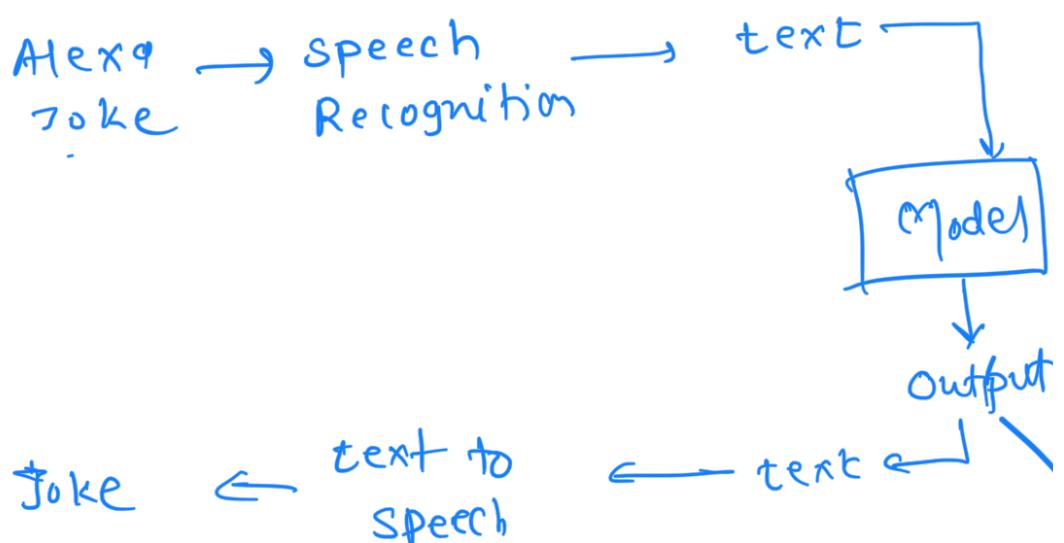
text	target	
review1	0	0,1,2
review2	1	review → pos



② Document classification

text	target
Doc 1	Adhar
Doc 2	D.L.
Doc 3	I.P.

③ Google Assistant.



④ chatbot → NLU

- Greeting Hi → Hello
- Order 1 2 3 4 → ≡

⑤ Text Summarization

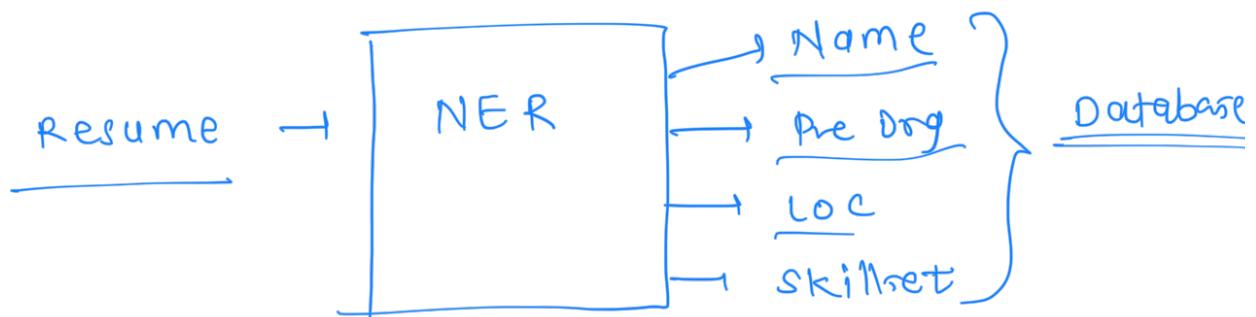
1000 doc → 10, 15

→ ① Extractive → NLU

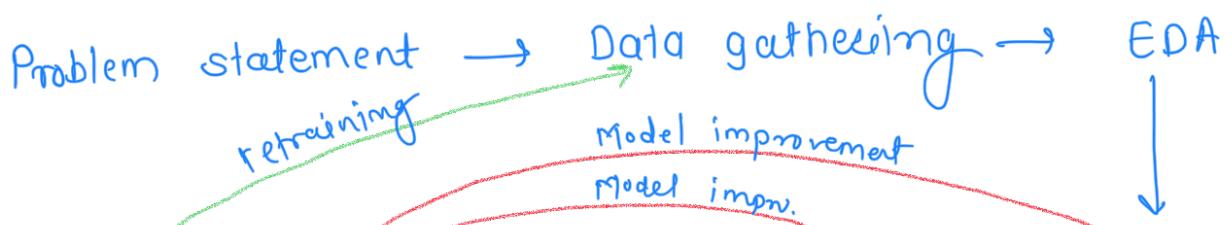
→ ② Abstractive → NLG

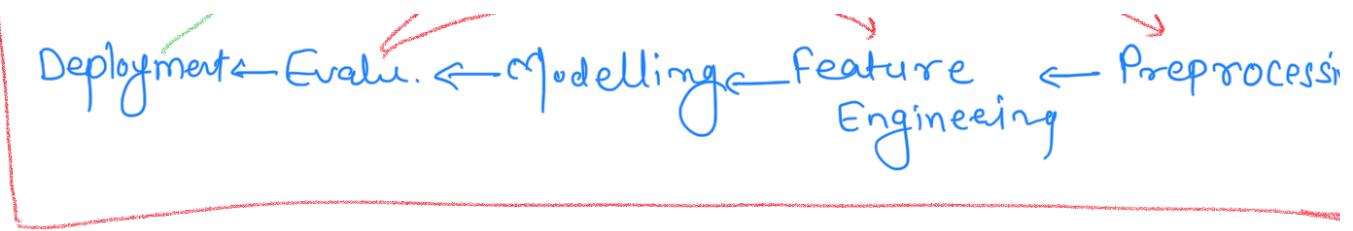
⑥ Named Entity Recognition

- Resume Parsing



NLP Pipeline





NLP Roadmap

① Preprocessing

- Tokenization
- stemming/lemmatization

② Text to numeric

- Count vectorizer
- TFIDF
- word2vec
- doc2vec

③ m.l. → classification

④ Deep learning

- RNN
- LSTM
- GRU

⑤ model Building with deep learning

⑥ Transformer

- Attention

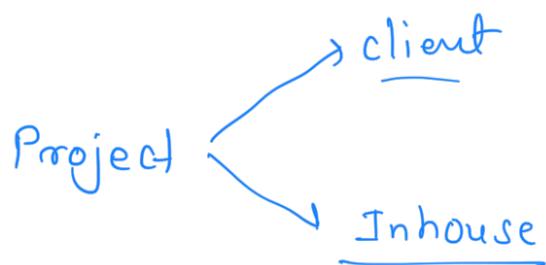
- ...
- self Attention
- multi head Attention
- encoder & decoder

1

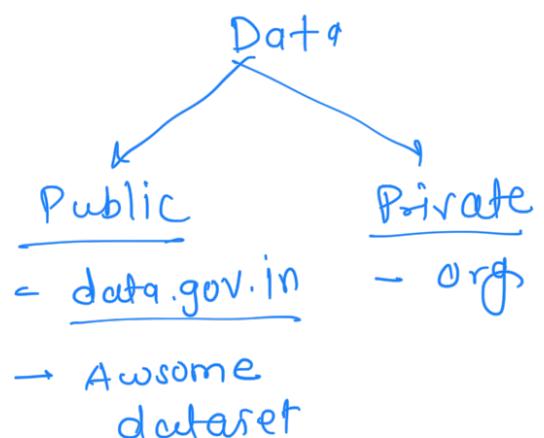
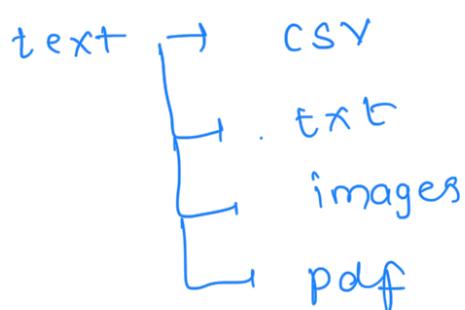
⑦ BERT

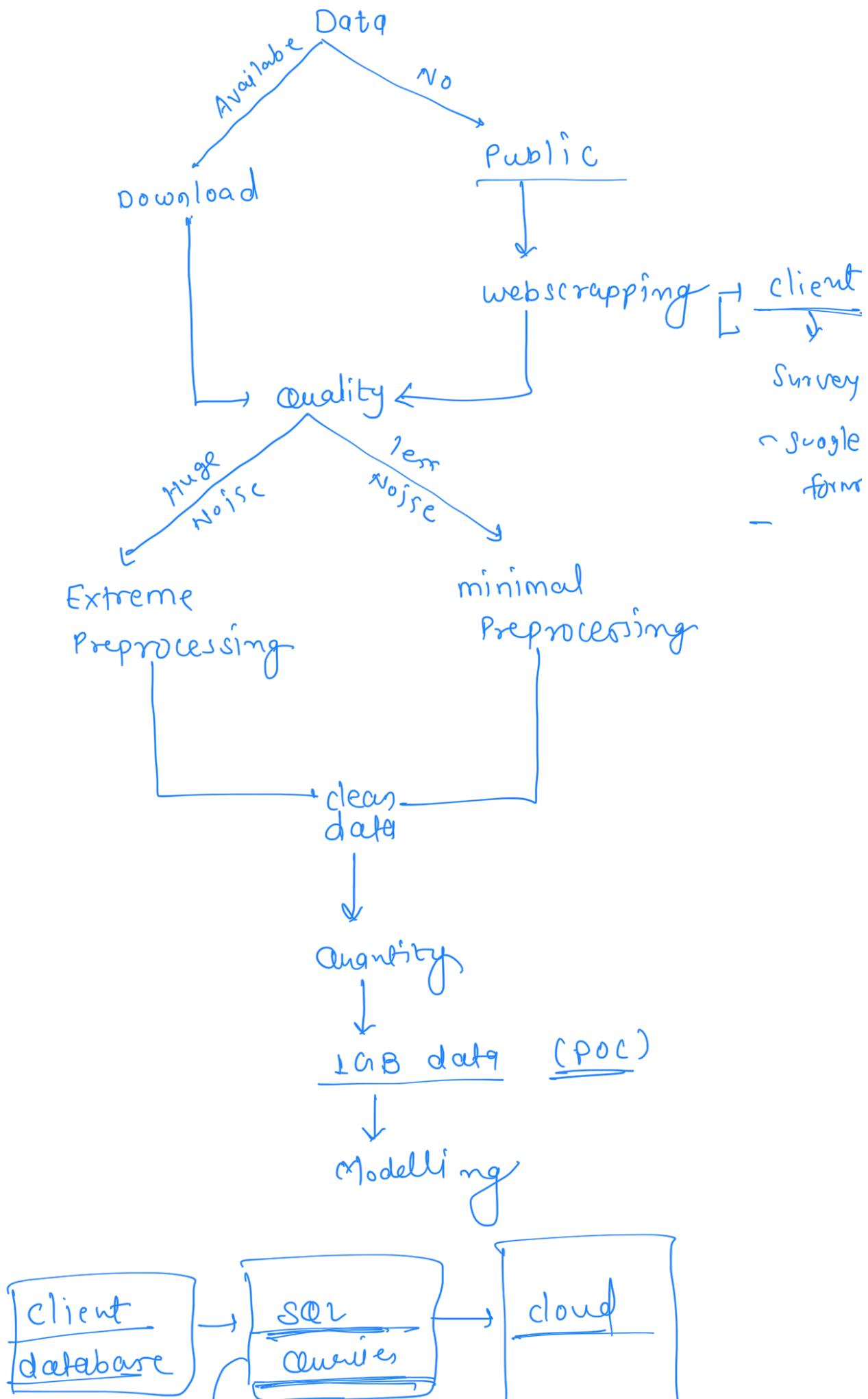
- GPT 3
- T5

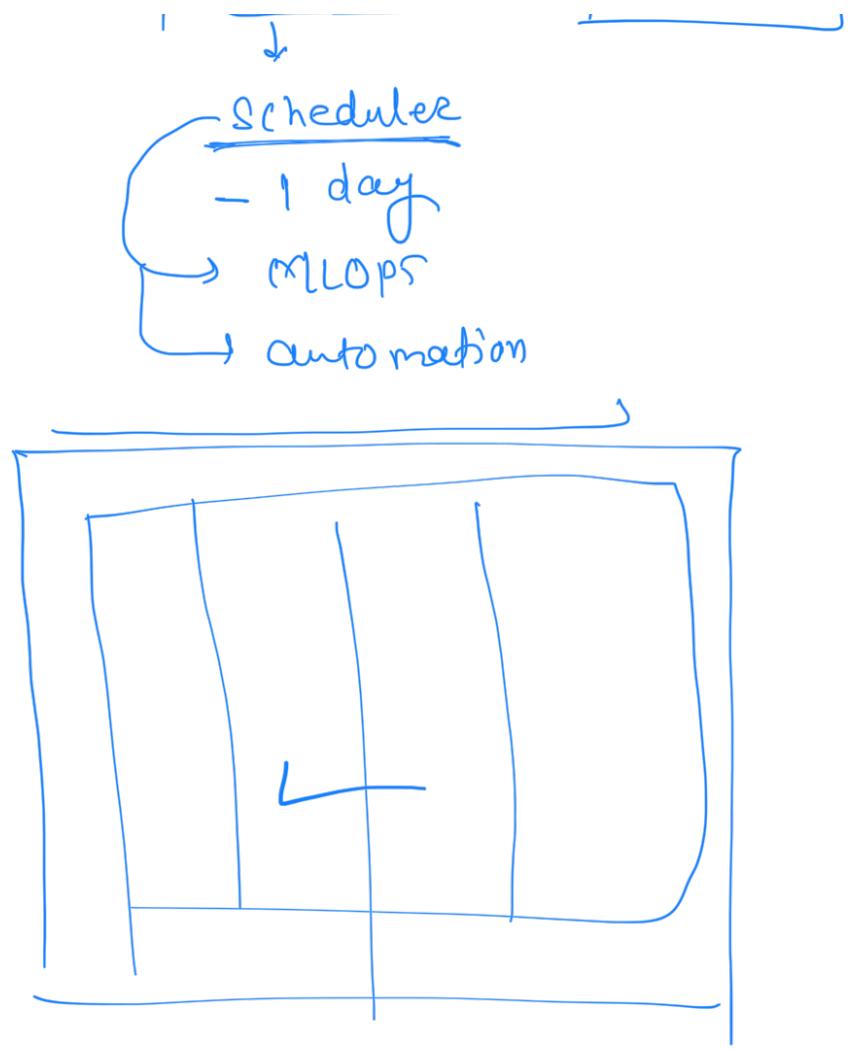
① Data Gathering



- website →
- database →
- api







- Beautiful.

② EDA

- ① Ngram
- ② Wordcloud
- ③ Keyphrase Analysis

① Ngram

"Minakshi is working in google"

(i) Unigram

- ["minakshi", "is", "working", "in", "google"]

(ii) Bigram

- ["minakshi is", "is working", "working in", "in google"]

(iii) Trigram

-

① POS → words analysis

Neg →

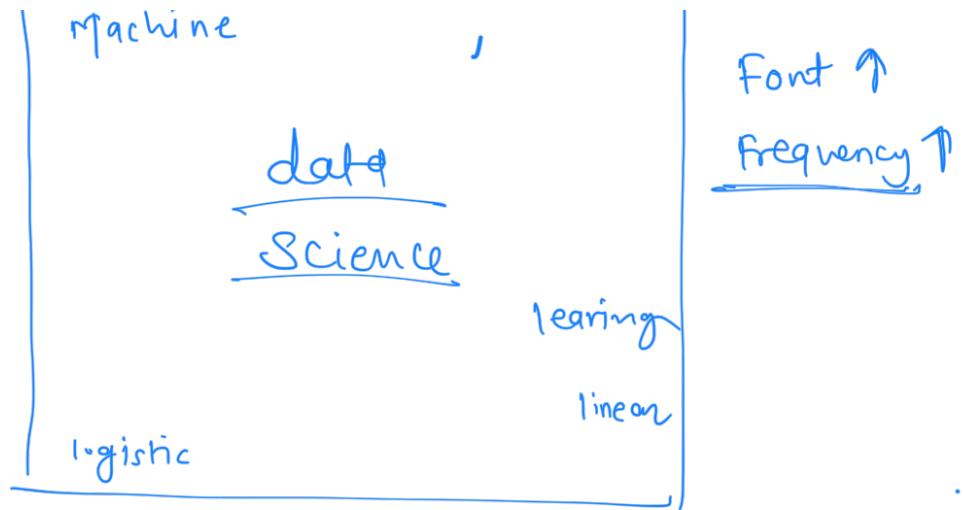
② extract stopwords.

review① - We like the product

review② - Awsome product, I will recomm
to everyone

review③ - Great customer service, staff
is quite polite.

② Word cloud



④ keyphrase analysis

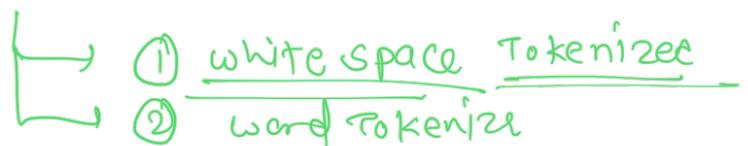
- We are learning nlp

(*) Preprocessing

① Tokenization

① Sentence Tokenizer → Punctuation, Syntax,

② Word Tokenizer. Conjunction



③ Regex Tokenizer

② Normalization → lower case,

Great → num1

great → num2

- Stop word Removal

⑧ Stopwords

Language Specific Domain specific

" Arun is suffering from a cancer.
doctor manish is treating w/m. Nurse
Rutika just gave capsule XYZ"

Healthcare

- Nurse, doctor, patient, bed, appointment,

Corpus → Collection of words, sentences

Corpora → collection of corpus