

# Natural Language Processing

1

# Regular expression

^ → starts with

\$ → ends with

[] → Find within range of characters

\ → Signals a special sequence (numbers, words etc)

\ can also be used to escape special characters.

. → Any character (except newline)

\* → zero or more occurrences

+ → one or more occurrences

? → zero or one occurrences

{ } → Exactly the specified number of occurrences

| → Either or

( ) → capture and group

use of several functions  
findall, search, split, sub, span, group, match, string ... etc

# collected from w3schools.com

official Documentation of Python.

Do follow

## # Why NLP is booming

- Freely available & Pretrained models
- Open source Ecosystem (spacy, Gensim, NLTK)
- Cheap Hardware, Cloud Resources
- Learning Resources
- Huge investment by Big Tech

## # Learn Regular expression (must)

Go to Link: [regex101.com](https://www.regex101.com) to use chatbot for Regex

#import re to import re module:

## Prerequisites

1. Python
2. Machine Learning
3. Deep Learning (RNN, CNN, LSTM, Word2Vec)

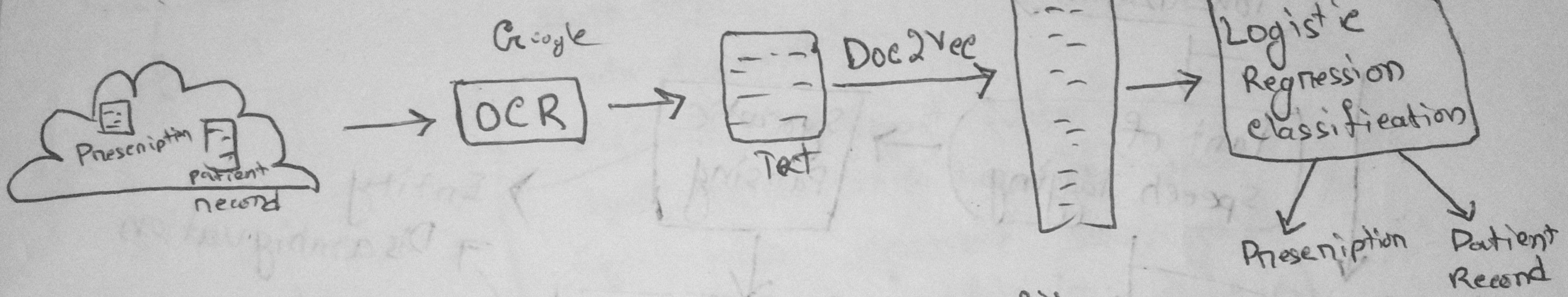
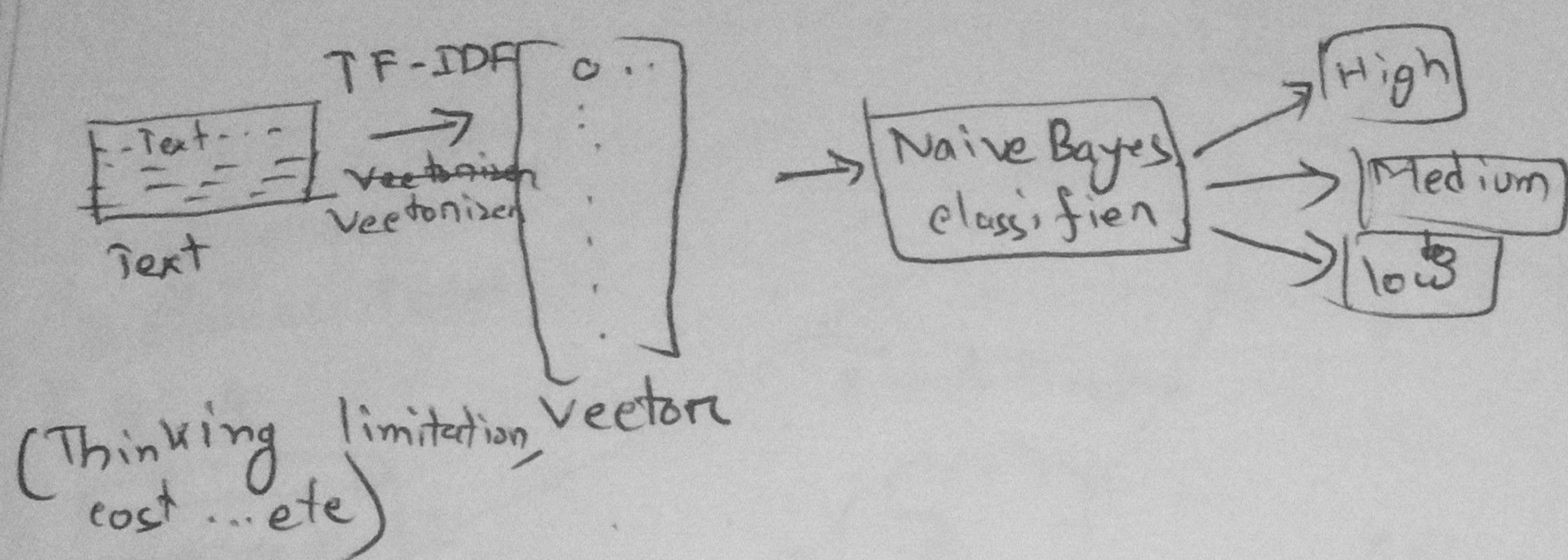
## \* Information Extraction using Regular Expression

## NLP Techniques

1. Rules & Heuristics (Use Regular Expression)
2. NLP Pipeline (Raw Text → Number Vector → Statistical Machine Learning)
3. Sentence Embedding / Word Embedding (Cosine Similarity based, Using count vectorizing approach)

## NLP Tasks

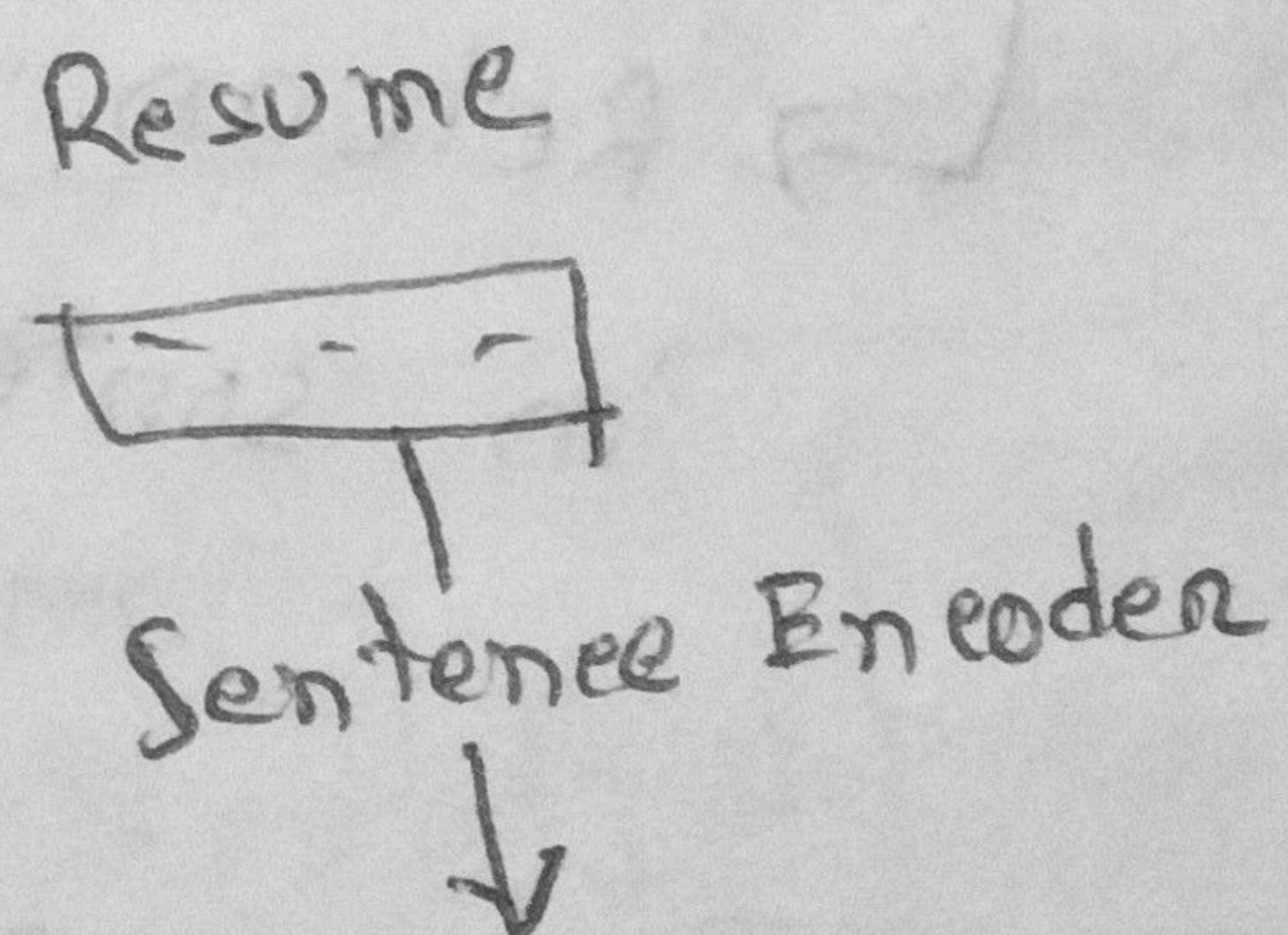
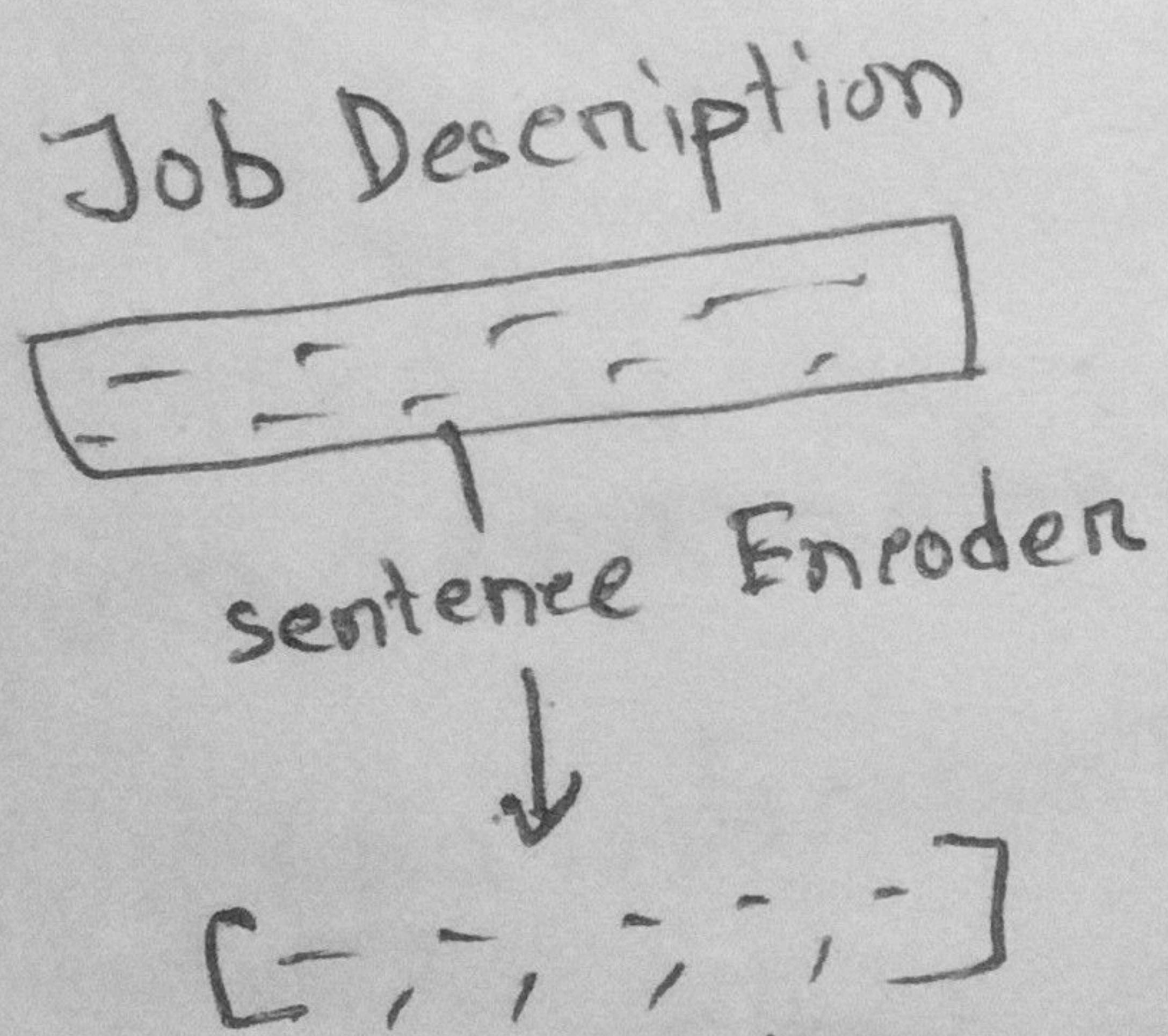
### Text classification



Application:

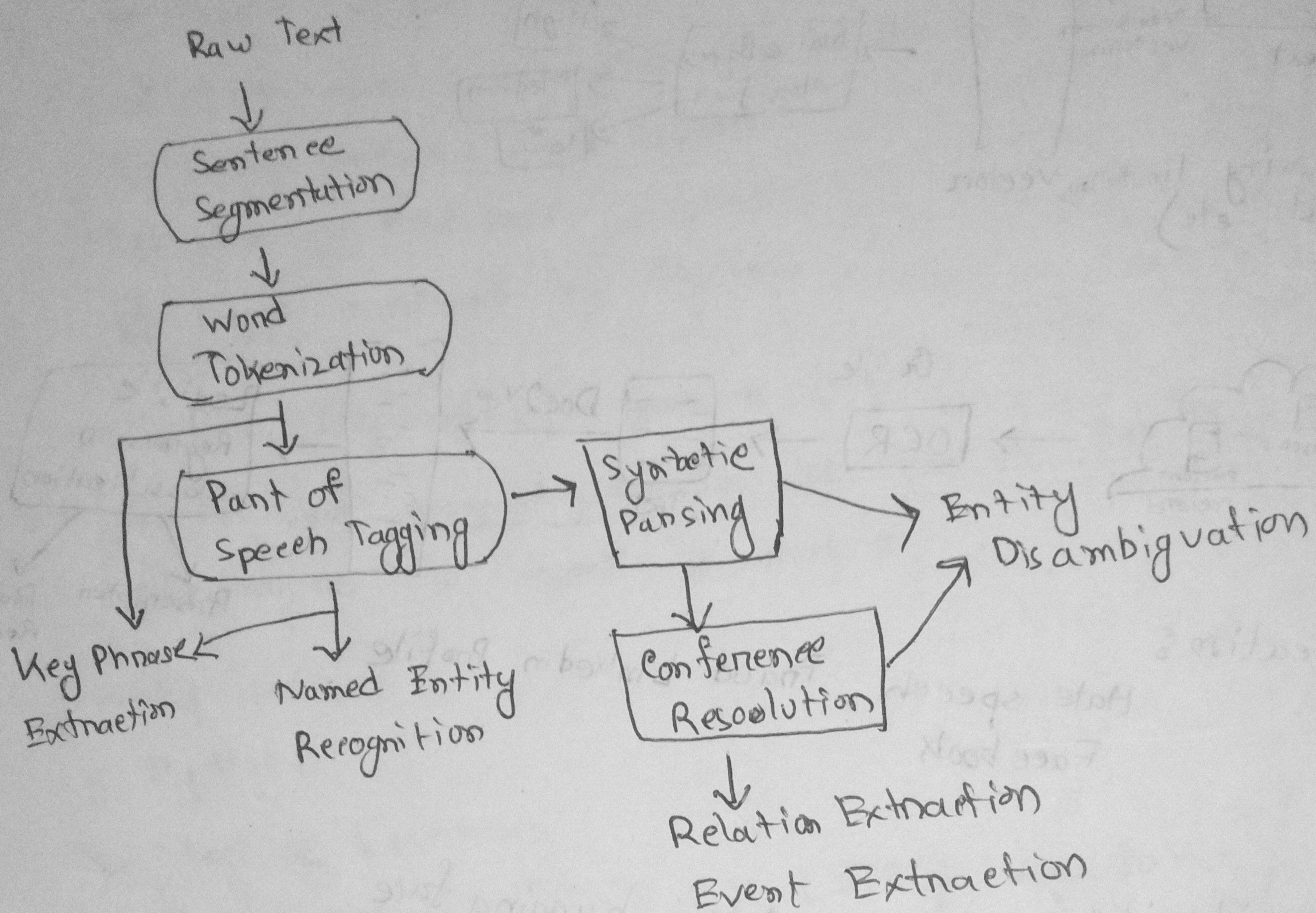
- Hate speech, Fraud LinkedIn Profile
- Face book

Text similarity: search in google: hugging face  
sentence encoder transformer



e cosine  
similarity

## Information Extraction



## Information Retrieval

↳ Returns the relevant website  
in sorted order

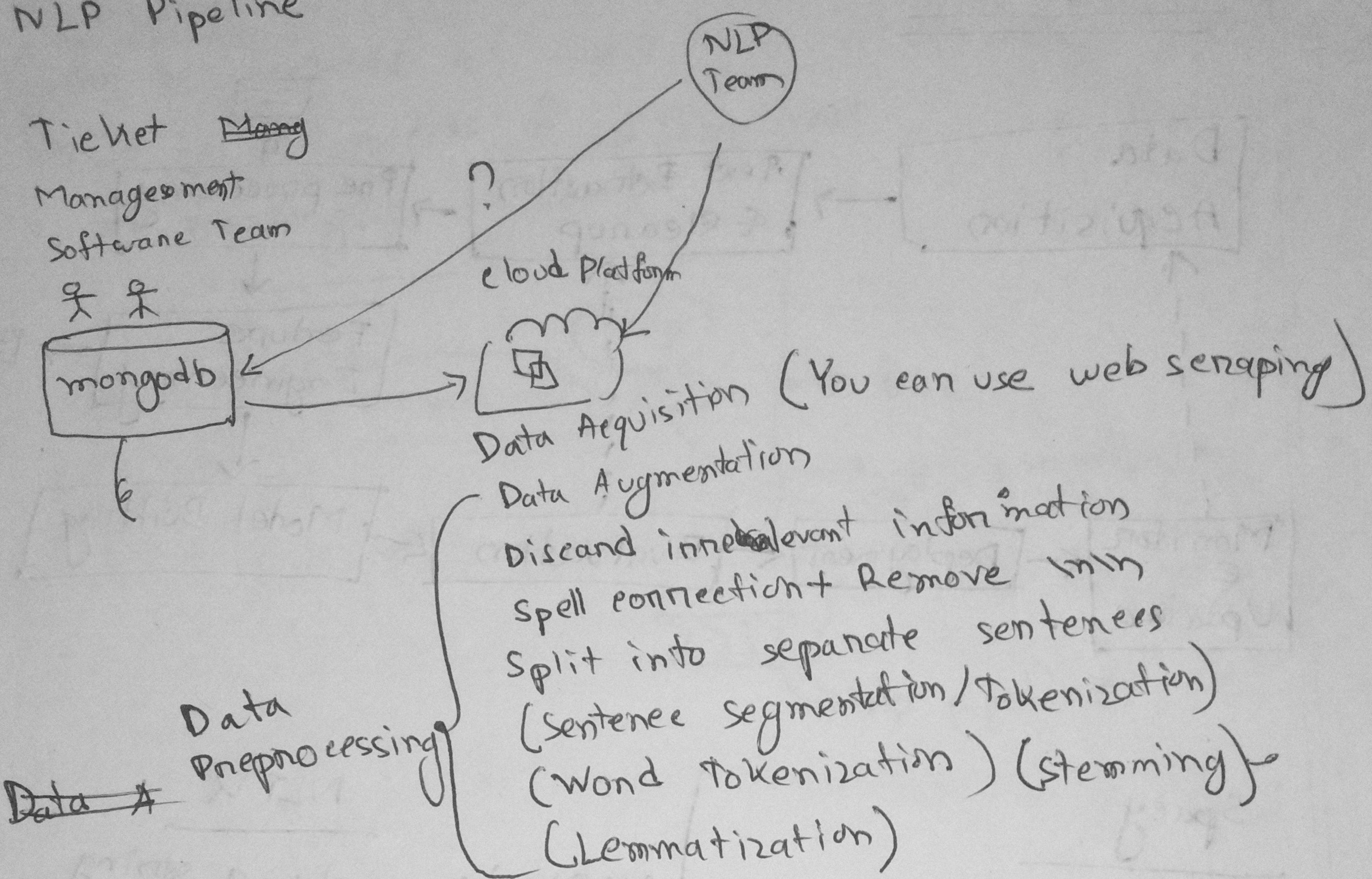
TF-IDF score

BERT

chatbots,  
Google Translator,  
Language Modelling,

Statistical Model | Neural Model  
Text summarization,  
Topic Modelling,  
Voice Assistants,

## NLP Pipeline

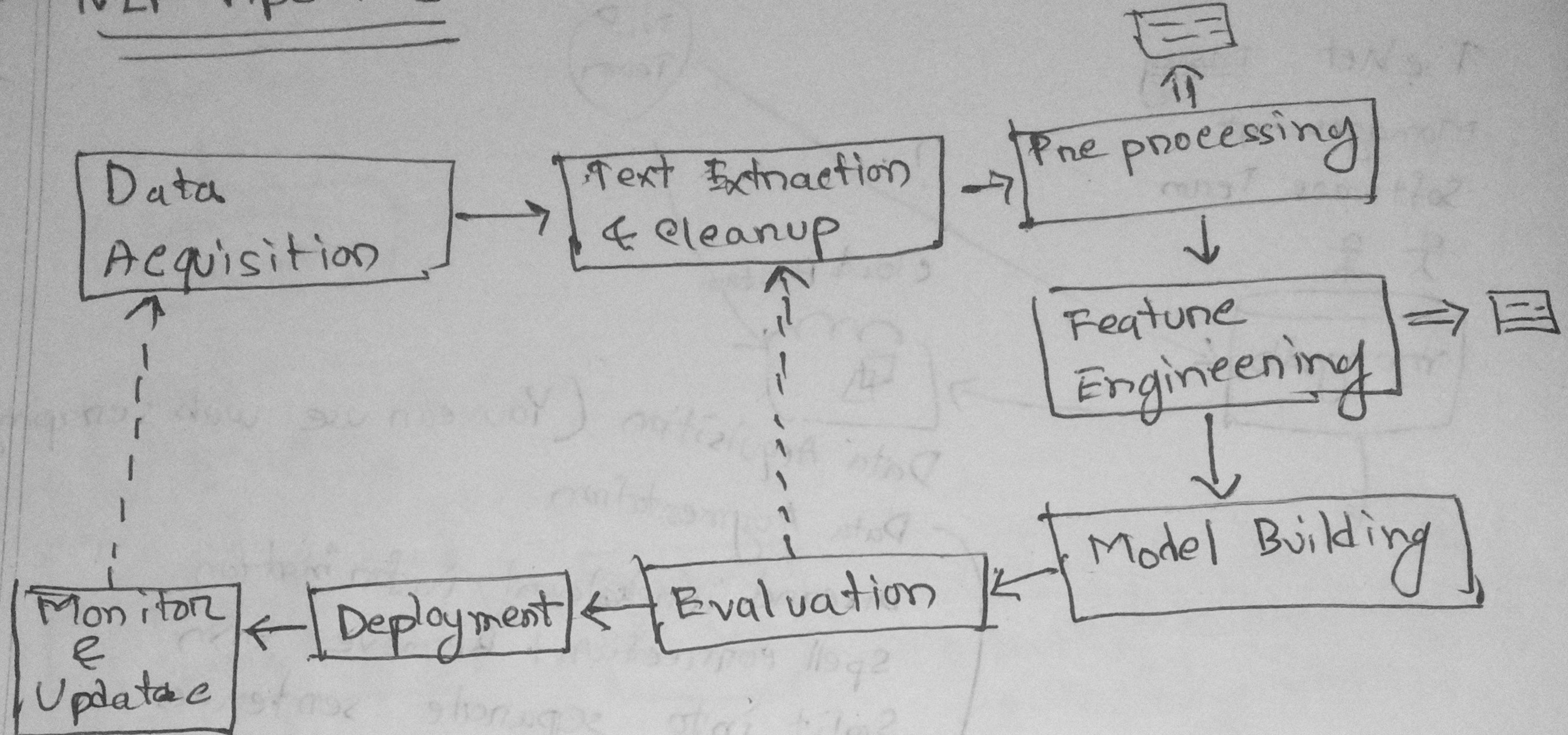


- Hyperparameters Tuning (Randomized search, Grid search)
- Evaluate the Model

Lastly, Deploy it on Google cloud, AWS  
Maintain and update the model

Stemming	Lemmatization
adjustable → adjust	was → (to) be
formality → formal	better → good
airline → airline	meeting → meeting

## NLP Pipeline



### Spacy

→ Spacy is Object Oriented

→ Provides most efficient NLP algorithms for a given tasks. Hence if you care about the end result, go with Spacy.

→ User friendly

→ Perfect for app developers

→ New library, very active user community

→ NLTK is mainly a string processing library.

→ Provides access to many algorithms. If you care about specific algo and customization go with NLTK.

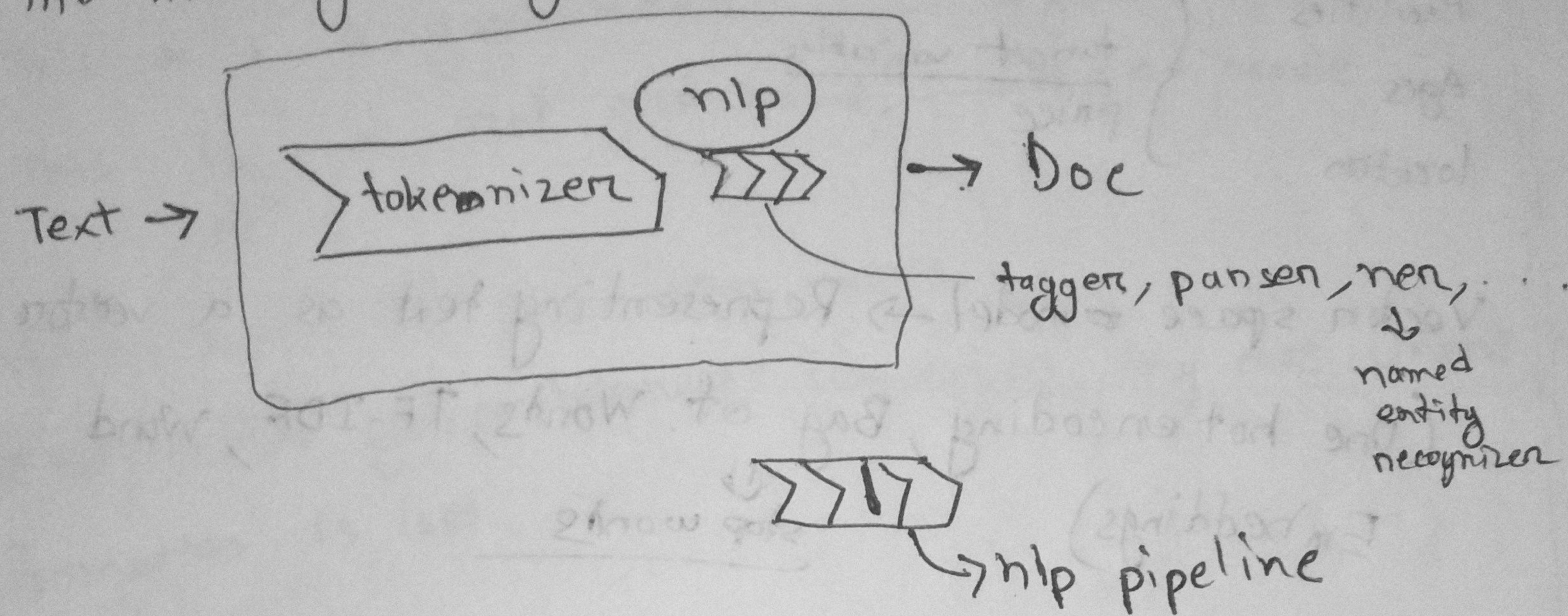
→ less user friendly

→ Perfect for researchers

→ Old library, not ~~more~~ much active as Spacy

## Tokenization

- ↳ can be sentence or word type → stemming, Lemmatization
- ↳ It is a process of splitting text into meaningful segments



## nlp.pipe\_names

- [ 'tok2vec', 'tagger', 'parser', 'attribute\_noun', 'lemmatizer', 'ner' ]
- ↳ named Entity recognition

fixed rules / heuristics to derive the base word.

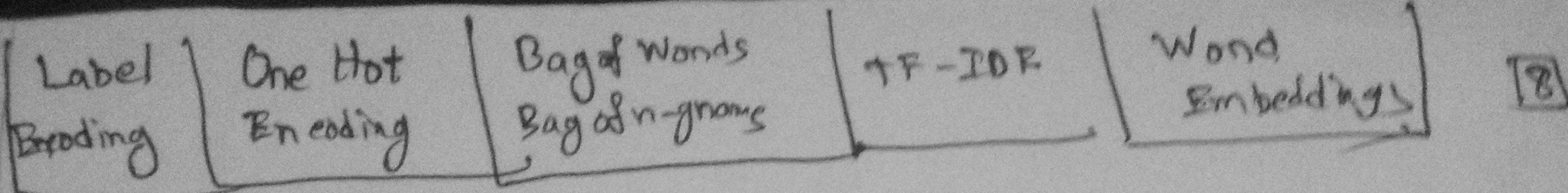
↳ Stemming

lemmatization → linguistic knowledge to derive the base word.

stemming and lemmatization → both are important

stemming → NLTK (Spacy does not support)

lemmatization → Both (Spacy, NLTK)



feature Engineering → Extracting features from raw data.

In NLP, Text Representation

### Features

Area  
Facilities  
Ages  
location

target variable  
Price

Vector space model → Representing text as a vector

(One hot encoding, Bag of Words, TF-IDF, Word Embeddings)

stop words

Label encoding One Hot encoding ⇒ In NLP, You can not use it.

### #Disadvantages

Similar words does not have similar representation

Consume too much memory and compute resources

Out of vocabulary (OOV) problem

No Fixed length representation.

#Bag of words → Sparse representation

(lot of memories,  
computer resources)

(Does not capture the meaning of the sentence)

↳ Count vectors

musks tesla iphone ipad

article1

14 0 1 5

article2

0 5 0 9

Stop words in BOW (To remove extra words)

When should I not remove stop words (losing important information)  
(Does not need to translate)

{ This ~~a~~ is a good movie → good movie  
This is a not good movie → good movie

where should not use stop words

{ chat bot, Q & A system, Language Translation, Any case where valuable information is lost.

Bag of n-grams → order of ~~not~~ words is important, capturing the relationship with words

(One) 1-grams	Dhaval	sat	on	a	couch	and	ate
(bi) 2-grams							
(Tri) 3-grams							

# several n grams combined

	thon eat	eat pizza	loki fall	loki eat
Doc1	1	1	0	0
Doc2	0	0	1	0
Doc3	0	1	0	1

limitation of bag of n-grams model

⇒ As n increases, dimensionality sparsity increases

⇒ Does not address out of vocabulary (OOV) problem

10

TF-IDF

	musik	that	price	market	investment	iphone	itunes
Apple	0	(32)	95	98	26	(7)	3
antideg	0	4	3	7	8	6	3
articles	15	(31)	44	43	25	0	0
Total	43	0	0	6	6	0	0

that  $\rightarrow 3$ iphone  $\rightarrow 2$ 

Document frequency (DF) = Number of times term  $t$  is present in all docs

\* log is used

to dampen

the importance

of term that

has high frequency

$$\text{Inverse Document frequency (IDF)} = \log \left( \frac{\text{Total Documents}}{\text{Document frequency}} \right)$$

$$\text{that} \rightarrow \frac{4(\text{total})}{3} \quad \text{iphon} \rightarrow \frac{4(\text{total})}{2}$$

To ~~non~~ normalize thing (to avoid biasness, for begin)

$$\text{To reduce } \text{TF}(t, d) = \left( \frac{\text{Total number of time term } t \text{ is present in doc A}}{\text{Total number of tokens in doc A}} \right)$$

$$\text{TF-IDF} = \text{TF}(t, d) \cdot \text{IDF}(t)$$

↓  
multiplication.