

27 Mar 2023

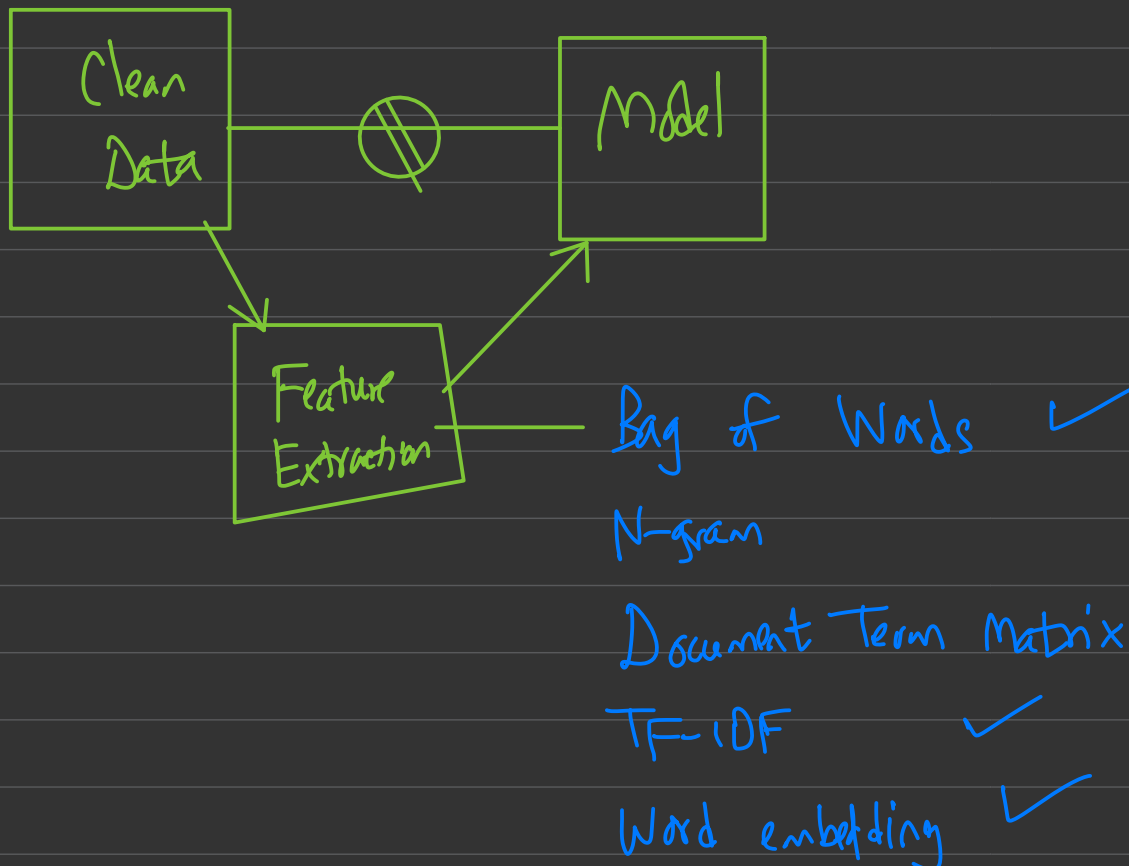
Goal: compute how close two pieces of text are

- (semantic similarity) meaning
- surface closeness
- (lexical similarity)

The cat ate the mouse
 1 2 3
The mouse ate the cat food.
 1 2 3 4

(Clustering, Redundancy removal, Information Retrieval)

What is Feature Extraction?



Cosine Similarity: Recommendation Systems

Similarity
↓
↑

Distance
↑
↓

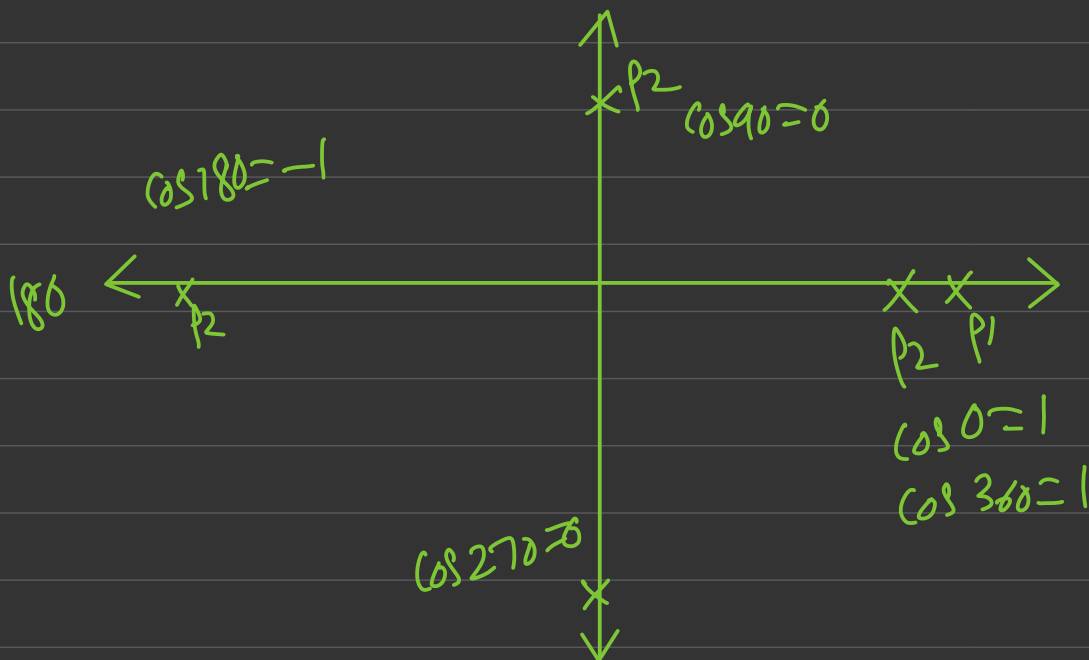
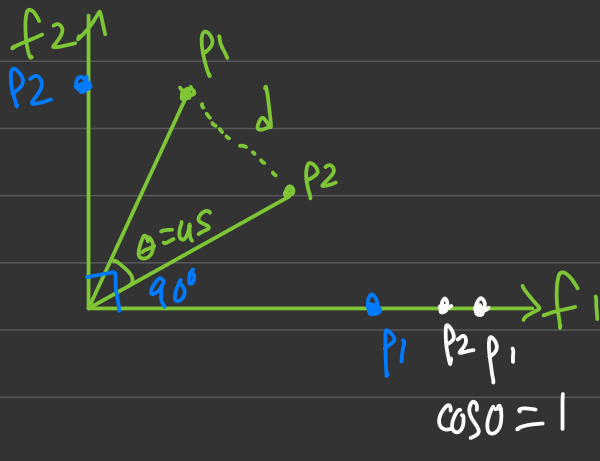
$[-1 \text{ to } +1]$

$$\text{cos_sim} = \cos \theta$$

angle b/w b/w p_1 & p_2

$$= \cos 45$$

$$= 0.53$$



Count Vectorizer: (Bag of Words):

S1: I love playing
S2: We love to play cricket
S3: Playing is fun

Text Preprocessing
→
(lower, rem. punct, rem. stopwords, stemming/lemmatization)

S1: love play
S2: love play cricket
S3: play fun

Unique words \Rightarrow {love, play, cricket, fun}

	love	play	cricket	fun
1	1	1	0	0
2	1	1	1	0
3	0	1	0	1
Σ	2	3	1	1

Desc. ord.

	play	love	cricket	fun
1	1	1	0	0
2	1	1	1	0
3	1	0	0	1
Σ	3	2	1	1

S1: I love playing
S2: We love to play cricket
S3: Playing is fun

Cvect \rightarrow

S1: [1 1 0 0]
S2: [1 1 1 0]
S3: [1 0 0 1]

Comedy ↑ Angry birds
* [0, 1]

$$\begin{aligned}\cos_{sim} &= \cos \theta \\ &= \cos 90 \\ &= 0\end{aligned}$$

90°

[1, 0] Avenger
X X → Action
Ironman [1, 0]

$$\cos \theta = 1$$

TF - IDF Term Frequency - Inverse Document Frequency

TF = # of rep of words in sentence

S1: good boy
S2: good girl
S3: boy girl good

$$\text{IDF} = \log \left(\frac{\text{\# of words in sentence}}{\text{\# of sentences containing words}} \right)$$

TF	S1	S2	S3
good	1/2	1/2	1/3
boy	1/2	0	1/3
girl	0	1/2	1/3

IDF	words	IDF
	good	$\log(3/3) = 0$
	boy	$\log(3/2)$
	girl	$\log(3/2)$

words Freq
good 3

boy 2

girl 2

	f_1 good	f_2 boy	f_3 girl
S1	0	$\frac{1}{2} \log(\frac{3}{2})$	0
S2	0	0	$\frac{1}{2} \log(\frac{3}{2})$
S3	0	$\frac{1}{3} \log(\frac{3}{2})$	$\frac{1}{3} \log(\frac{3}{2})$

X