10th April 2023

# Text Preprocessing Using NLTK

Let's begin @ 9:05 PM

Images $\longrightarrow$ num. $\longrightarrow (0-255)$   RGB.

Text. $\longrightarrow$ string

Input x      Output y

1               10

2               20

3               30

4               40

5               50

.                .

.                .

.                .

$y = 10x$

$y \approx 10x$

numerical

text

mapping

- chat bots
  - QA
  - speech to text
  - Language models
  - text translation
  - text summarization

GPT - 4

* **High** value in stocks creats **high** return.

high

* **Cricket**
  'kohli'
  'dhoni'
  'catch'
  'win'
  'bowled'

**Politics**
  'Modi'
  'Putin'
  'UK'
  '
  '

freq mapping

bank / bank

$2

riverside

bank: 2

mouse / mouse

computer

animal

* King Kohli *is* *a* good batsman. Kohli *is* great.

* Modi *is* PM *of* India.

Stopword ✓ ⇔ Keep Keywords

Stopword ↓ remove

* Grammerly * → Stopwords may be drop.

$$\{w_1 : n_1 , w_2 : n_2 , \text{-------} \}$$

## Tokenization

Strings $\longrightarrow$ list of words.

$\hookrightarrow$ Aplit on "spaces"

tokenizer.

**Code mixed**

- En + Bengali
- Eng + Telgu
- Eng + hin → **hinglish**

**Romanized text**

( mai ) ( market ) ( jaa ) ( raha ) ( hu )

(En+Bn)

**Pure English**

* king Kohli good batsman. Kohli great.

* Modi pm india.

lower $\longrightarrow$ r SW    keywords $\longrightarrow$ unique.

Bag of Words *    Vocabulary $V$

| | india | pm | Kohli | good | king | batsman | great | modi |
|---|---|---|---|---|---|---|---|---|
| $d_1$ | 0 | 0 | 2 | 1 | 1 | 1 | 1 | 0 |
| $d_2$ | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |

Size = 10k X $V$

Only 3 type of articles.

Cricket

politics

Sci-fi



Cricket

Politics

Sci-fi

|  | 0th. | 1 | 2 | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | india | pm | kohli | good | king | batsman | great | modi |
| $d_0$ | 0 | 0 | 2 | 1 | 1 | 1 | 1 | 0 |
| $d_1$ | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |

$\vdots$

$d_n$
$\downarrow$
1M
$=$

$V = 10,000$

$(0, 2) \rightarrow 2$

$(0, 3) \rightarrow 1$

$(0, 4) \rightarrow 1$

$\vdots$

very sparse data $*$

# Stemming & Lemmatization : → (Lemma) .

warm / warmer / warmest → 1 col^m

play / playing / played . → 1 col^n .

- Stemming → Cutting . → English . ✓

  ↓
  Algo. → rules

  · Porter
  (diff. rules) ✓

  · Snowball
  ↓
  more robust.
  International ✓

\*    caresses $\longrightarrow$ Caress \*

rule. 'sses' $\longrightarrow$ 'ss'

By stemming
=.

rule. 'ies' $\longrightarrow$ 'i'

. ties $\longrightarrow$ 'ti'

not taking grammer into account

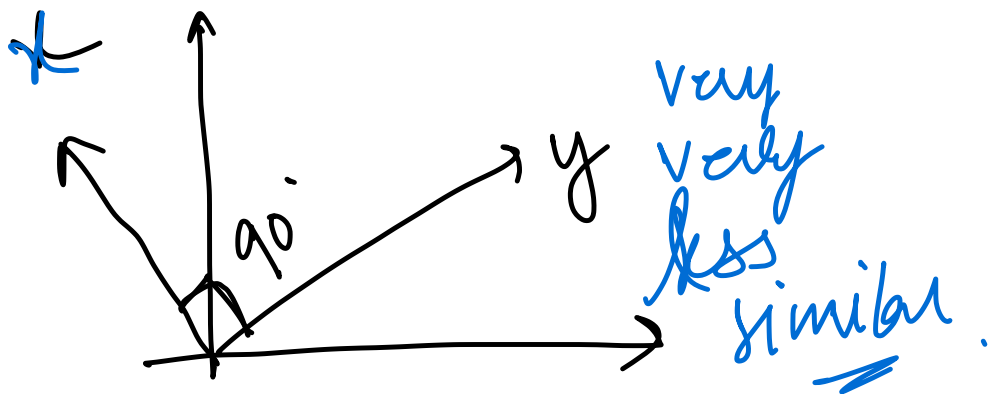② **Lemmatization** : ( grammer into account)

"lemma"

feet
feetest  } 'foot'
feetor —

very similar.

less similar

$\theta_1$

$\theta_2$

very very less similar.

$90°$

$180°$

worst

sim·score

$+1$

$0$

$90°$

$-1$

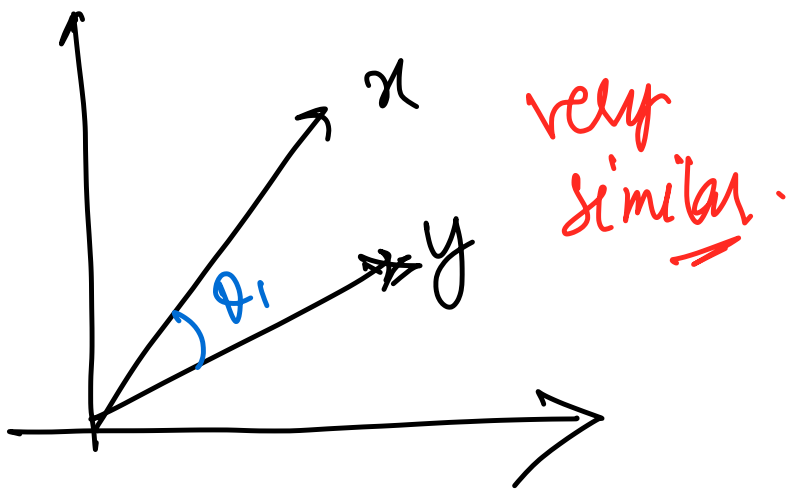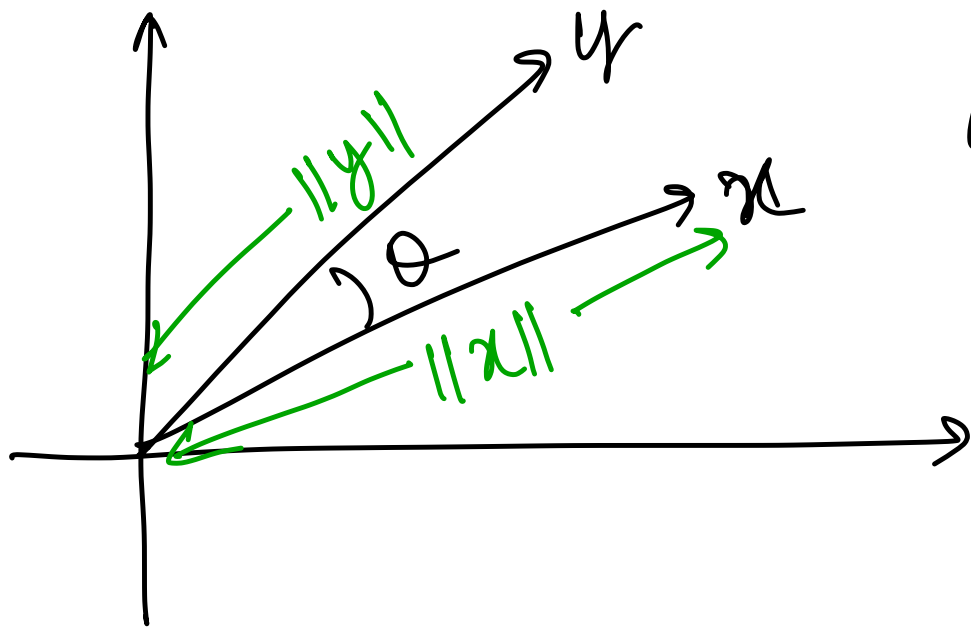$\theta$

$\cos\theta$

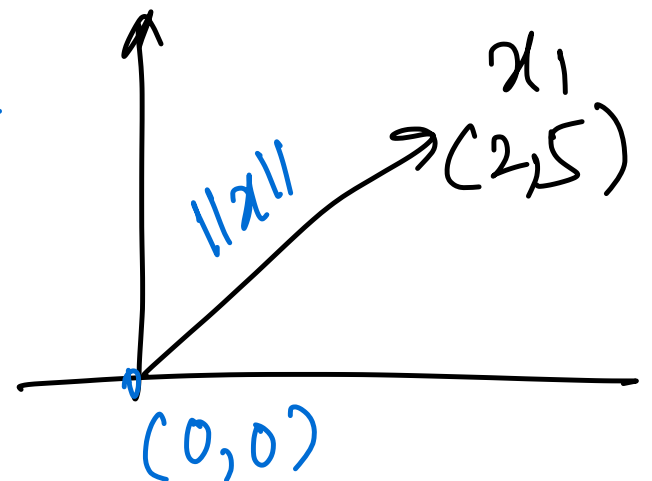$[-1 \text{ to } 1]$

$$\cos \theta = \frac{x \cdot y}{\|x\| \|y\|}$$

dot product

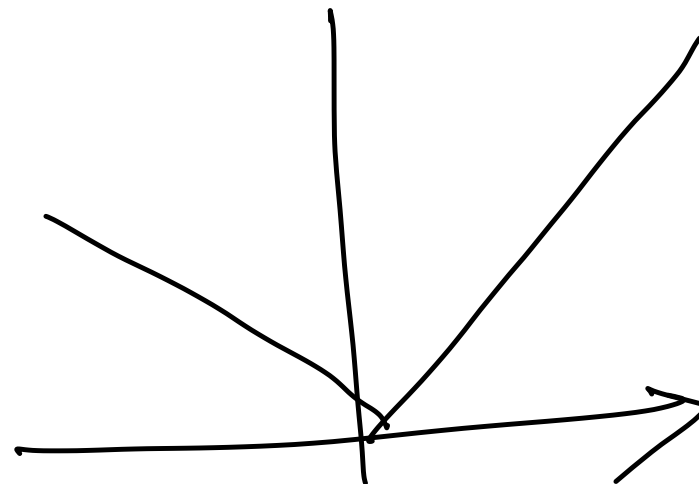· elementwise multiplication
↳ addition.

$x_1 \to (2, 5)$

$x_2 \to (-1, 3)$

$x_1 \cdot x_2 = 2 \times (-1) + 5 \times (3)$

$\|x\| = \sqrt{(2-0)^2 + (5-0)^2}$

$x_1$
$(2, 5)$
$\|x\|$
$(0, 0)$

$|x| \longrightarrow$ mod $f^n$

$|-5| \rightarrow +5$

$|5| \longrightarrow 5$

vector length.

$\|x\| \longrightarrow$

dist. formula

norm $\longrightarrow$ (L2) norm