11th April 2023

# Text Preprocessing Using NLTK

Let's begin @ 9:05 PM

**Image:** ↱ numerical ⟶ (0 – 255)  R.G.B

**Text:** string.

Input $x$      Output $y$.

1        10

2        20

3        30

.        .

.        .

.        .

$$y = 10x$$

$$y \approx 9.98x$$

num.

num

text

map    text ⟶ numerical
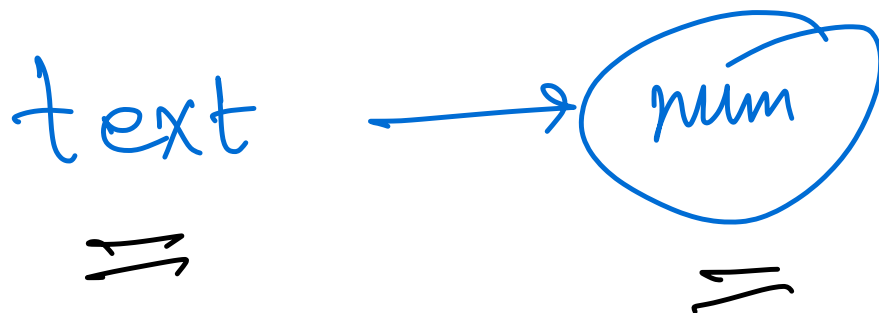
- chat bots
- SA

:)
::|
::(

- recommendations
- text to speech / vice versa.
- language models.
- text summarization.

Chat GPT - 4

text ⟶ (num)

☆ flow. text ⟶ numerical
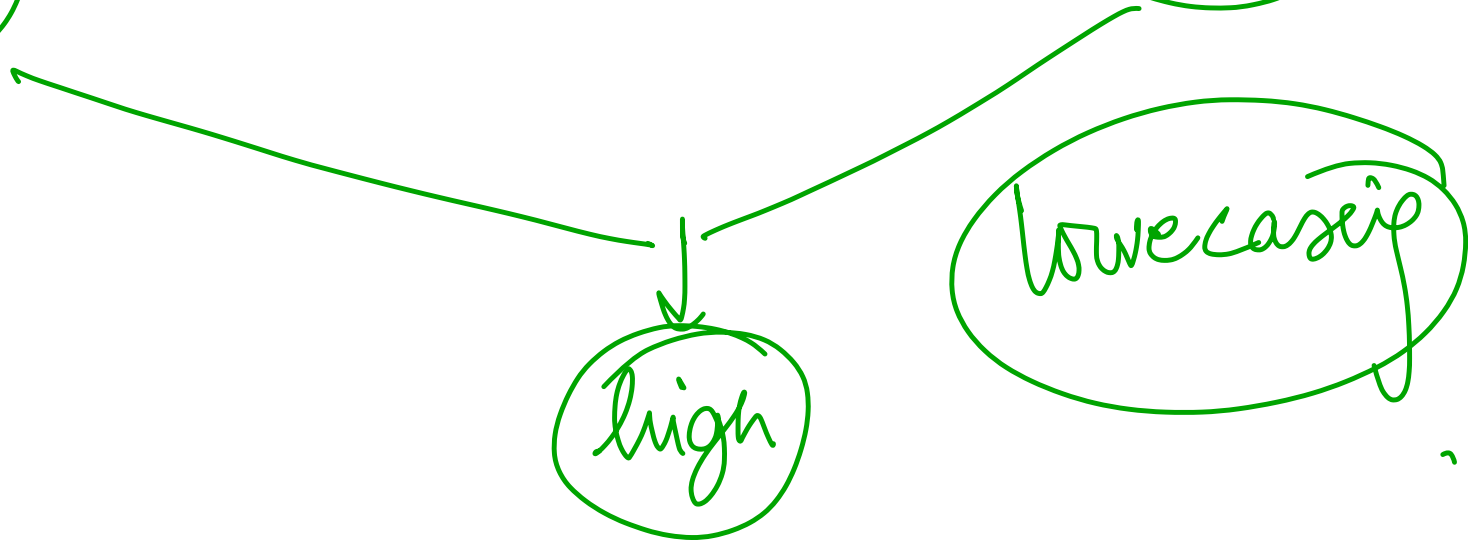
↓

raw — lowercase the data.

= stopwords.
= tokenization
— Bag of words (BOW)
— Stemming / Lematization
— Case study.
  ↳ Similarity score.

eg.) High value in stocks creats high return.

high

lowecaseip

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

⟹ Modi is PM of India.

is / a / the / of ⟶ stopwords. ⟶ remove ✗.

(Keep keywords)

Grammerly : → stopword may be drup.

$d_1:$ King Kohli good batsman Kohli great.

$d_2:$ Modi PM India.

$d_i:$

$d_n:$

lower → Removed SW → Keywords → unique

Bag of words *

Vocabulary V

| | King | kohli | good | batsman | great | modi | pm | india |
|----|------|-------|------|---------|-------|------|-----|-------|
| $d_1$ | 1 | 2 | 1 | 1 | 1 | 0 | 0 | 0 |
| $d_2$ | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |

$:d_n$

8 # unique Keywords : ✓

$n$ :

$(10000 \times 8)$ → numerical.

2 types of documents → Cricket ✓

Politics. ✓

|        | 0: King | 1: Kohli | 2: good | 3: batsman | great | Modi | pm | india |
|--------|---------|----------|---------|------------|-------|------|-----|-------|
| 0: $d_1$ | 1 | 2 | 1 | 1 | 1 | 0 | 0 | 0 |
| 1: $d_2$ | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |

$(0,0) \rightarrow 1$

$(0,1) \rightarrow 2$

$(0,2) \rightarrow 1$

$(0,3) \rightarrow 1$

$(0,4) \rightarrow 1$

$(0,5)$

$(0,6)$ $(0,7)$

$1000 \times 500$

↪ 90% are zero )

only stores non zero positions in matrix.

① Stemming : ↪ cut down.

warm / warmer / warmest → only one word =.

play / playig / played

→ Root Word

Algorithm    Cuttig with some rules.

① PORTER
↓
(english) old

② SNOW BALL
↓
(chinease, japanese,
- - - - - .)

* Caresses ⟶ 'Caress'

rule: 'sses' ⟶ 'ss'

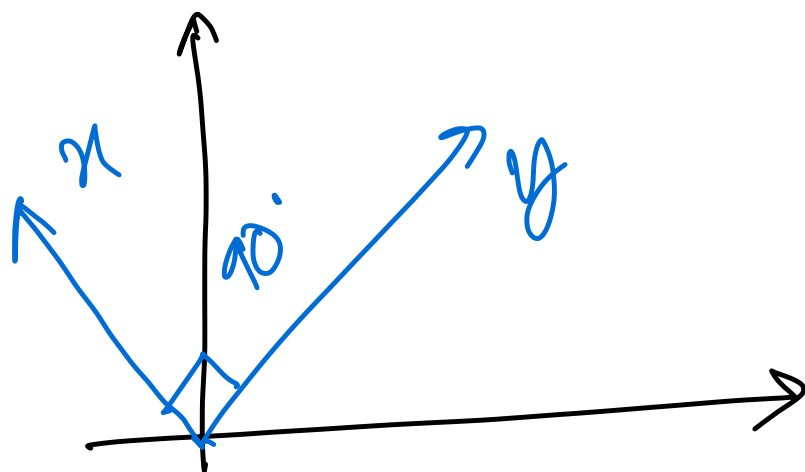rule: 'ies' ⟶ 'i'      ⟶ ties ⟶ 'ti'
                                   ⇃
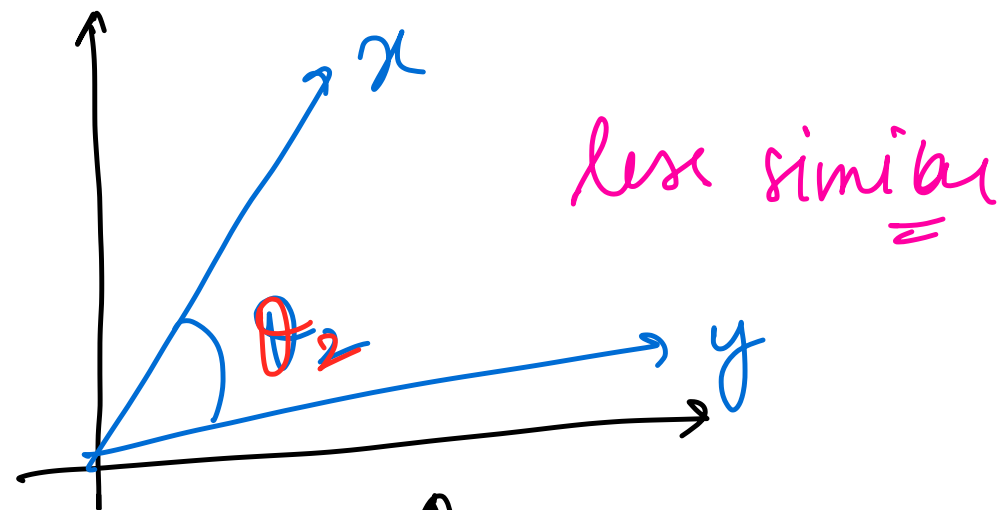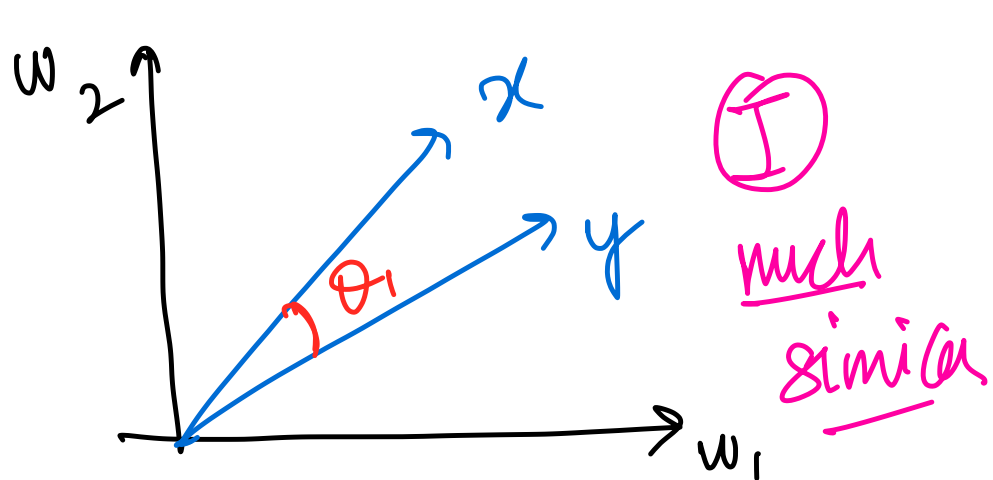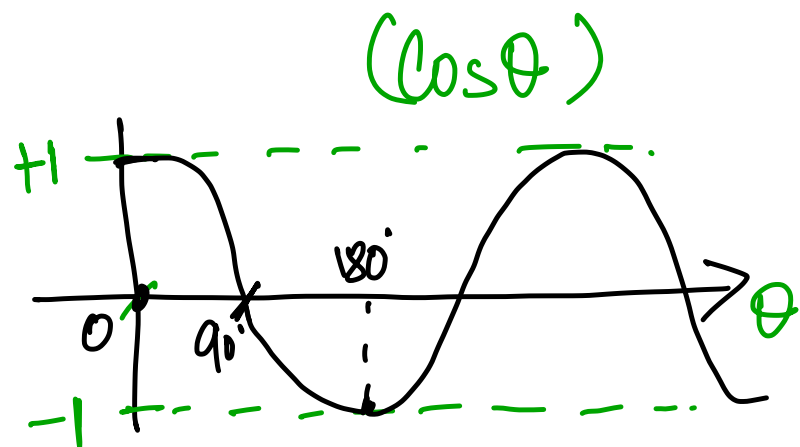                              maki sense?

Not taking GRAMMER into account

② **Lemmatization** → lemma Base word

taking grammer into account.

following rule of Grammer.

$w_2$ $x$ $\theta_1$ $y$ $w_1$

① much similar

$x$ $\theta_2$ $y$

less similar

$x$ $90°$ $y$

Very very very similar.

$180°$ $x$ $y$

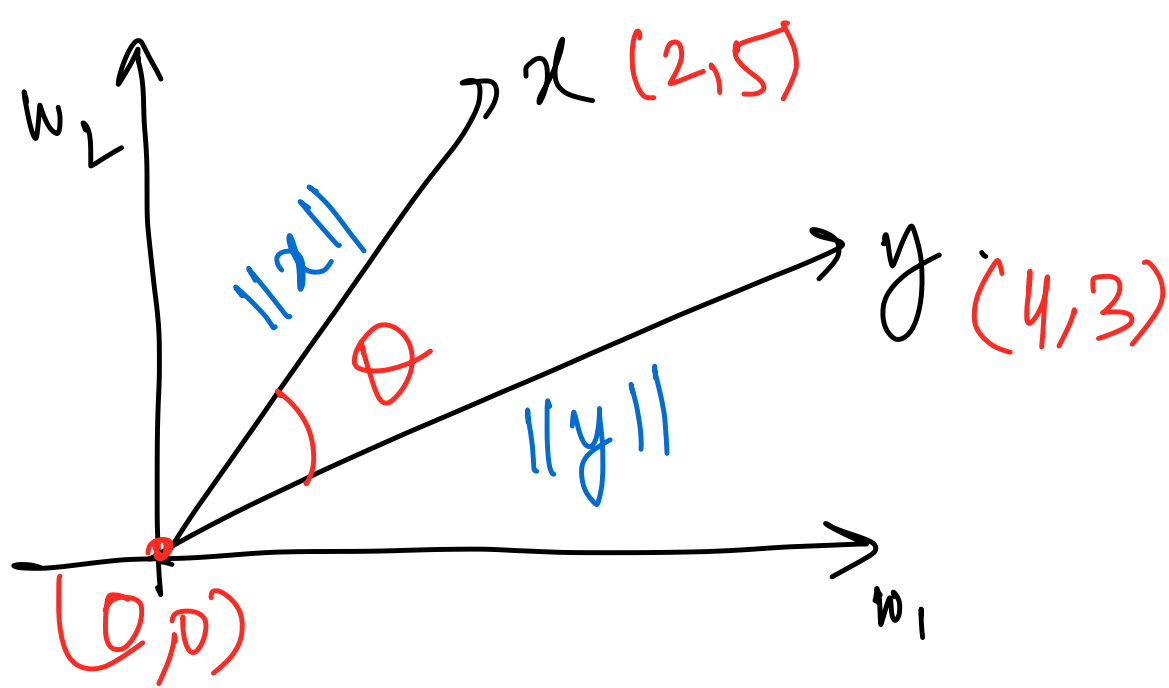dissimilar worst.

$(\cos\theta)$

$+1$ $0$ $90°$ $180°$ $\theta$ $-1$

@ $\theta = 0 \Rightarrow +1$

@ $\theta = 90° \Rightarrow 0$

@ $\theta = 180° \Rightarrow -1$

Similarity score $= [-1 \text{ to } +1]$

$\cos \theta$

$w_2$

$x$ (2,5)

$\|x\|$

$\theta$

$y$ (4,3)

$\|y\|$

(0,0)

$w_1$

$x \rightarrow (2,5)$

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

norm/magnitude

$x \cdot y = 2 \times 4 + 5 \times 3$

$= 23$

$x \cdot y$    dot product

$\cos\theta$ → Similarity score    $(-1 \text{ to } +1)$

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