

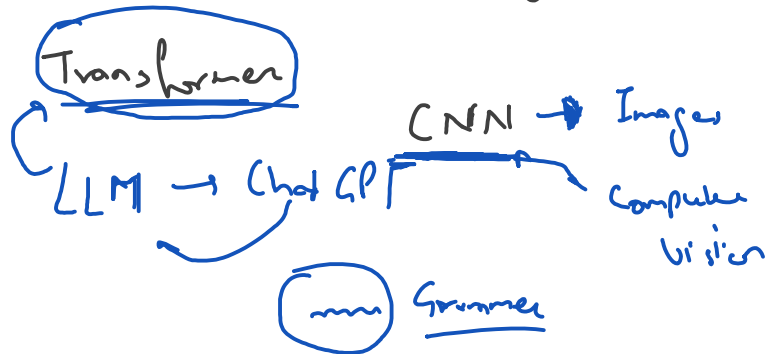
NLP

Tuesday, July 08, 2025 1:15 PM

Natural Language Processing :

Text → Numbers

RNN

Google - IRTransformerNLP : Translator

news Summarization

spam

autoCorrect

Text: I like the movie I like the movie but acting was subpar

↑ → P

→ N

Pipeline

Text → num.

① text label

① Get the Entire Data/Corpus ✓

② Tokenization :

② Tokenization ✓

③ stopword removal ✓

④ stemming/lemmatization ✓

⑤ Building a vocab ✓

⑥ Vectorization

My name is Shubham. I have
an interest in teaching ML.

① Sent tokenization

② word tokenization

① My name is shubham

I have - - - - -

② My I :

nltk

sklearn

have
is
shubham
have
an
interest
.

Stop words : I like the movie
↑

$l = []$

for i in range(10)
l.append(i)

for i in range(10)
 i

Stemming

run
running → run

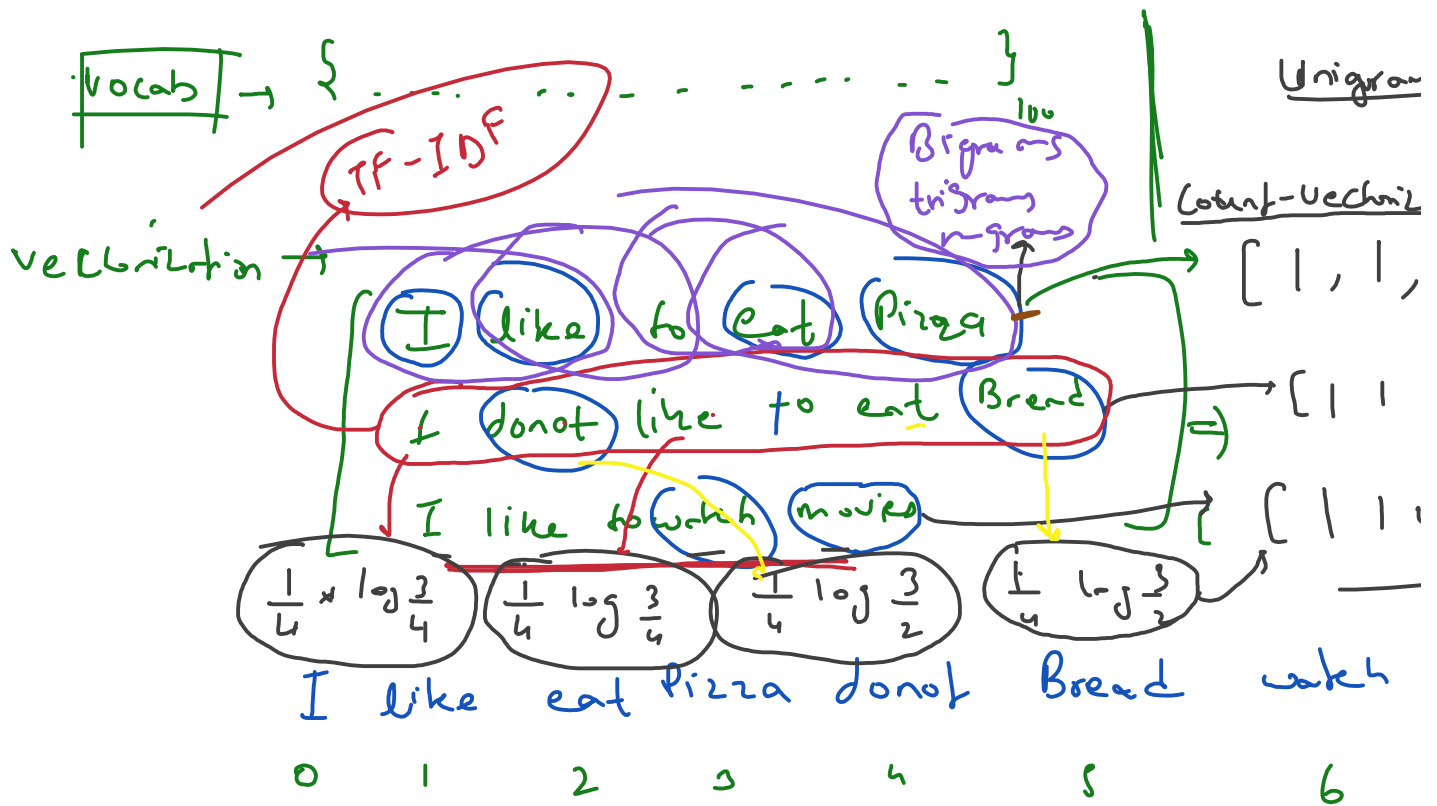
give

gave

swim → swim
swimming

Lemmatization

car/c



I like to eat pizza and I like to watch

[2, 2, 1, 1, 0, 0, 1, 1]

1000 → 800 (like) ↓ TF-IDF

TF * Term Frequency

$tf(t, d) \rightarrow \frac{\text{No. of times } t \text{ occurs}}{\text{Total no. of terms/word}}$

$IDF(t, c) \rightarrow \log \frac{\text{Total number of documents}}{\text{Number of documents containing } t}$

TF-IDF

new York

Bigrams

Word 2 Vec

Word → Vector

RNN, Transformers

Trying to
capture
meaning ↓

the → [.]₃₀₀

1000

king → [.]₅₀

1000

male 0.8
Female 0.1
Food 0.2
kingdom 0.9
city 0.5
country 0.7
:
:

Apple

fruit
Red
taste
mobile
sweet

① Dim

②

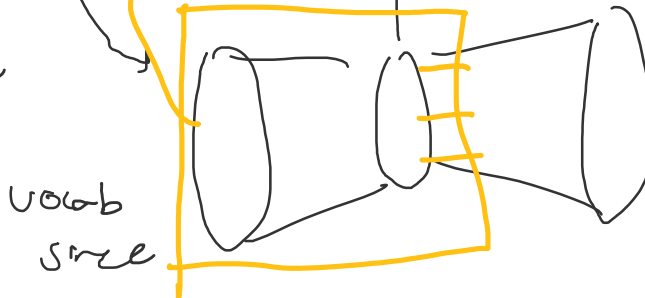
C Bow
Skip gram

(0 0 1 0 0 0)

hidden
layer
of size d

(0 1 0 0)

Google Word2Vec
Wikipedia



vocab
size

is o/f

The girl is dancing

girl

dancing

I/P

dim

[0 0 1 0 0 1 0 0 0 . . .]

