

Natural Language Processing



Textual data

↳ Linguistics that enable the Computers to understand, interpret & generate the human language.

① Google Translation

I love data science



मुझे डेटा साइंस पसंद है

② Sentiment Analysis

Today's movie is quite boring. I have just wasted two hours watching this movie.



Negative

③ Language Model → Predict next word

↳ Autocompletion task (Email)

I hope you are doing well.

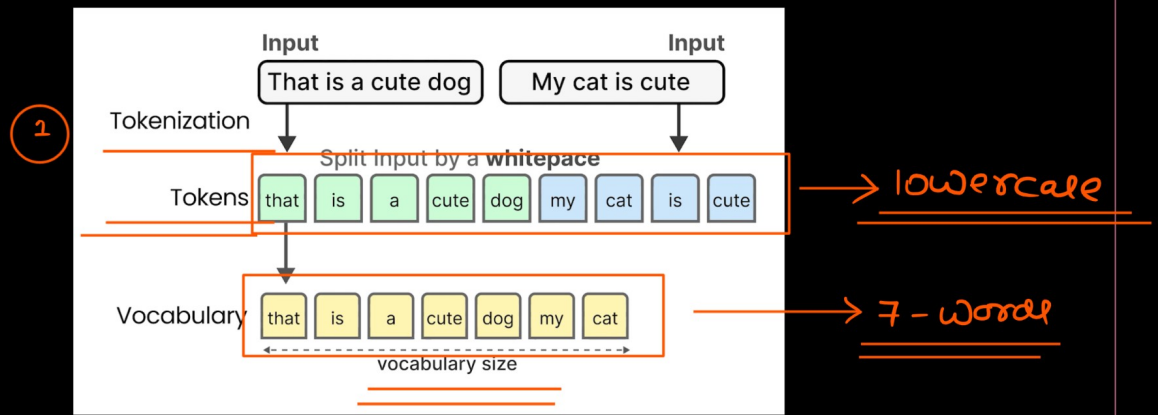
Code → further depending upon function name it gives all the further code.

④ Text-to-image

Step 1 → Computers/machines will be
able to understand the
textual data.

Embeddings → textual data
↓
numeric vector

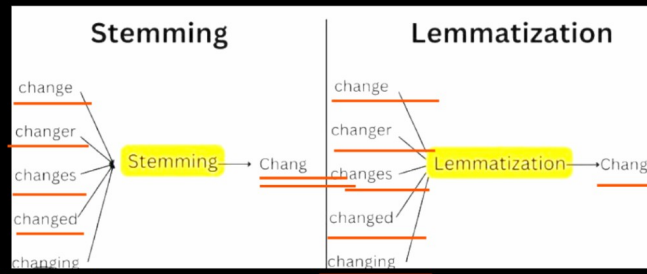
Textual Data Preprocessing



2 Stopword removal (a, an, the, is)

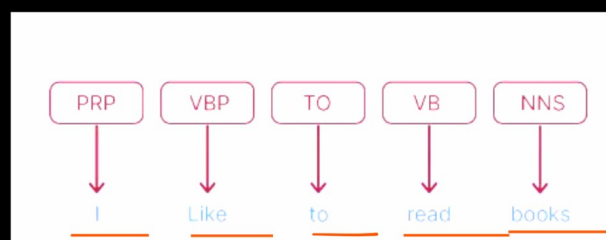
3 Stemming & Lemmatization
(Root word)

[run, ran, running] → run



→ has a literal meaning

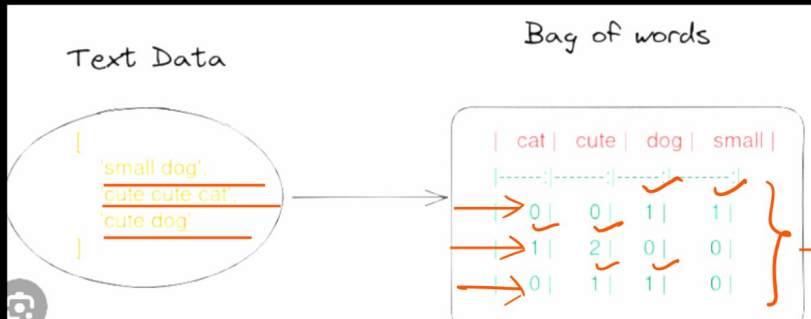
4 POS Tagging



⑤ Embedding → to the numeric
vectors

Bag of words → no context { frequency
based }

✓ ✓ ✓ ✓
small, dog, cute, cat → unique words



→ frequency of
the individual
words

word2vec → static embeddings

The table displays word embeddings for five categories: cats, puppy, houses, apple, and baby. Each category has a vertical column of values for different attributes: animal, newborn, human, plural, and fruit. The values are color-coded: red for animal, blue for newborn, green for human, orange for plural, and purple for fruit. A vertical double-headed arrow on the right indicates the number of dimensions (properties) is 1024 values!

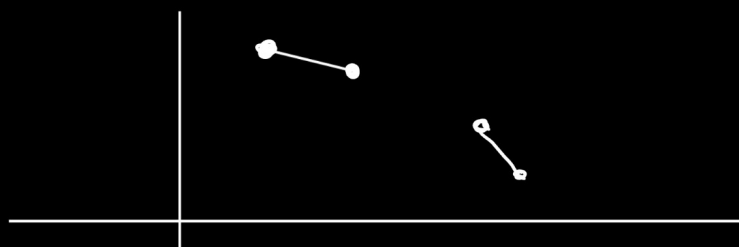
	cats	puppy	houses	apple	baby
animal	.91	.93	-.56	-.67	.01
newborn	-.11	.71	-.32	-.1	.90
human	.19	.36	.31	.29	.87
⋮	⋮	⋮	⋮	⋮	⋮
plural	.94	-.82	.94	-.51	-.11
fruit	-.51	-.91	-.5	.89	-.51

Pretrained
model

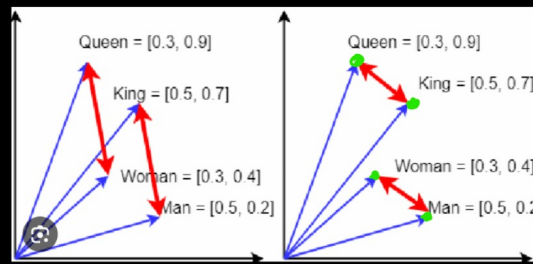
Number of
dimensions
(properties)

1024 values!

Dimensional



$$\text{king} + \text{Queen} - \text{Man}$$



vector

Language

ambiguous

apple (I am eating
company apple)

apple (I am

currently working
in apple)

single
vector

word2vec

Dynamic
embeddings

RNNs

Attention based models

Transformers

GPT → Generative Pretrained

Transformers

Attention is all
you need

Pretrained Models

transformers

1000 Records

finetuning

→ Customize the

above pretrained
model for your
dataset.

I love Data Science.

[I, love, Data, Science, .]

[40, 3047, 4833, 13993, 13]

Task + ML Project Ready

↳ credits → \$5

↳ secret key → create

↳ .env → define

↳ tokenization, embedding → execute

↳ deep analysis → dimensional &

What each dimension

represent
