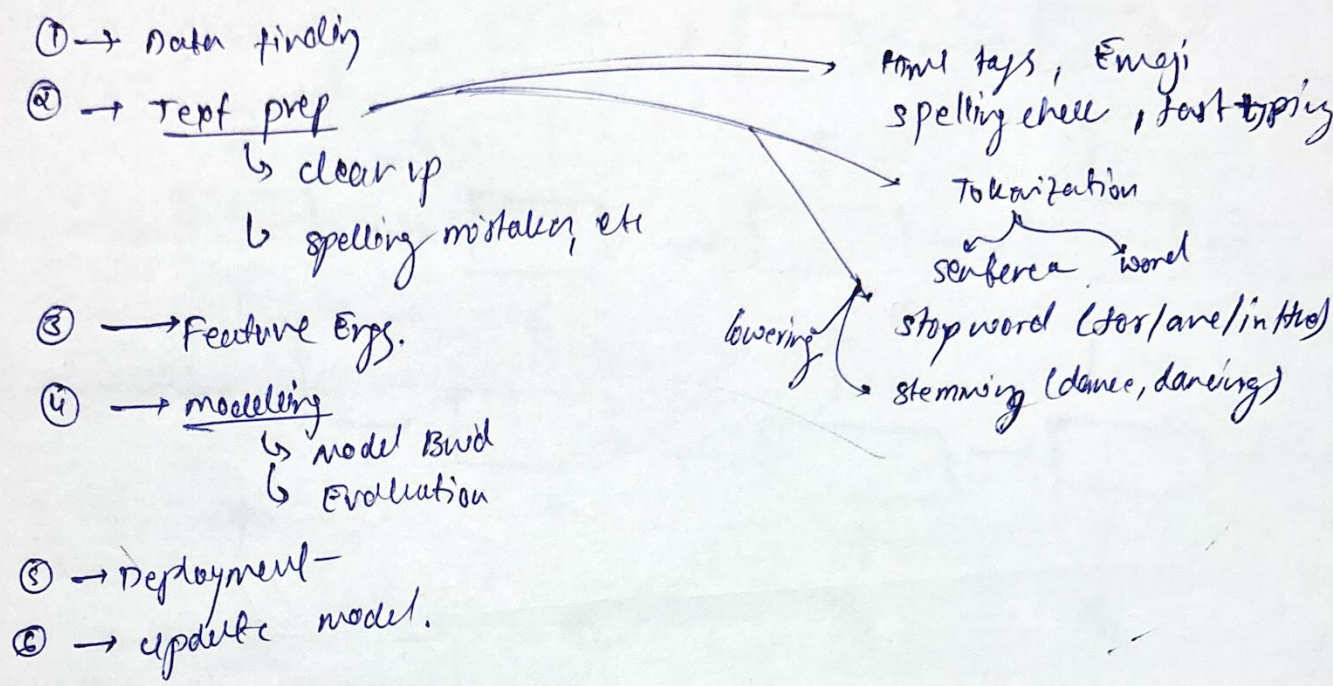


NLP Pipeline



Adv Text pre-processing

- Part of speech processing (noun, verb, subject, etc)
- Coreference resolution

ex: He my name is Tanishq and I'm 19

→ same person Tanishq

lect 3] # Text preprocessing (Basic)

① lower casing

A a → Both are different
 ↑ ↑
 start of sentence mid sentence
 tokenized, and interpreted different.

$$df['column'] = df['column'].str.lower()$$

② HTML Tag remove

→ use regular exp
import re

pattern = re.compile('<.*?>')
filtered = pattern.sub('', text)

③ Remove urls

→ using regular exp

④ punctuation

→ ! ' " # \$ % & ' () * + , - . / : ; < = > ? @ [\] ^ _ ` { | } ~

import string
↑ string.punctuation

for char in string.punctuation:
text = text.replace(char, '')

→ Too slow

text = text.translate(str.maketrans('', '', string.punctuation))

→ Fast

⑤ Char word treatment

① lmao / fyi / gn / ASAP (Examples)

↓

Have to do it manually (get list of these words and use wop)

⑥ Spelling check

from textblob import TextBlob

text = TextBlob(text).correct().string

use this method

⑦ Stop word

① a / the / I / we / my → sentiment

* ② pos tagging don't do this removal.

③ from nltk.corpus import stopwords
stopwords.words('english')

⑧ Handling Emojis

① use regex

pattern = re.compile("🔗")

🔗 emoji
🔗 symbols & photography
🔗 transport, map symbol
🔗 flags & lies
🔗 re.UNICODE

② use (emoji) package

import emoji

text = emoji.demojize('python is 🔗')

fire symbol

python is 🔗

⑨ Tokenization

\$20 → \$, 20

✓

Delhi! → Delhi, !

✓

New-York → New, York

✗

u.s → u, s

✗

ph.D → ph, D

✗

\$10.5 → \$, 10, 5

✗

import spacy

nlp = spacy.load("en-core-web-sm")

doc = nlp(text)

for word in doc:

print(word)

⑩ Stemming

(task)

process of inflected words to root word

(But here root word may be not a valid word
ex: story → stori)

root word

→ work
walking
walks

} this is called
inflection

inflected words

from nltk.stem.porter import PorterStemmer
ps = PorterStemmer()

text = " ".join([ps.stem(word) for word in text.split()])

① lemmatization

slow

Inflected word $\xrightarrow{\text{to}}$ root word (could root word)

from nltk.stem import WordNetLemmatizer

lect 4 # Text Representation

Common Tense

- ① corpus(c) \rightarrow count of total words in a dataset.
- ② vocab(v) \rightarrow count of unique words \rightarrow _____
- ③ Document(d) \rightarrow Individual row (sentence) in dataset.
- ④ word(w) \rightarrow word bucket
- ⑤ oov (out of vocab)

① one hot Encoding

- p1) people watch yt
 p2) yt watch yt
 p3) people write comment
 p4) yt write comment

V = 5

| | people | watch | yt | write | comment |
|-----|--------|-------|----|-------|---------|
| p1) | 1 | 1 | 1 | 0 | 0 |
| p2) | 0 | 0 | 1 | 0 | 0 |
| ... | | | | | |

② Bag of words

- p1) people watch yt
 p2)

\rightarrow [1 1 1 0]
 \rightarrow [0 0 2 0]

- p1)
 p2)
 p3)
 p4)

frequency of words

from sklearn.feature_extraction.text import CountVectorizer

cv = CountVec()

cv.fit_transform(text)

oov is ignored

③ N-gram
 → n-gram like Bag of words but contains N-words

Bag of words

docs

| | people | watch | yt | comment |
|-----|--------|-------|----|---------|
| d1) | 1 | 1 | 2 | 0 |
| d2) | 1 | 1 | 0 | 1 |

vocab of two words

2-gram (bi-gram)

saved doc

| | peoplewatch | watchyt | ytcomment | commentwatch |
|-----|-------------|---------|-----------|--------------|
| d1) | 1 | 1 | 1 | 0 |
| d2) | 0 | 0 | 0 | 1 |

④ TF-IDF (gives weightage to words)

term frequency

Inverse doc frequency

$$TF(t, d) = \frac{\text{Freq of term } (t) \text{ in doc } (d)}{\text{Total no. of terms in doc } (d)}$$

$$IDF(t) = \ln \left(\frac{\text{Total no. of docs in corpus}}{\text{Total no. of docs which has term } (t) \text{ in them}} \right)$$

- d1) people watch yt
- d2) yt watch yt
- d3) people write comment
- d4) yt write comment

from sklearn.feature_extraction.text
 import TfidfVectorizer

note

TF → gives prob of term (t) in that doc

high TF → term freq ↑ ⇒ more imp
 low TF → term freq ↓ ⇒ less imp

IDF → tell the value of that term in corpus

high IDF → term (t) freq ↓ ⇒ more valuable.
 low IDF → term (t) freq ↑ ⇒ less valuable.

$$TF(\text{people}, d_1) * IDF(\text{people})$$

| | people | watch | yt | write | comment |
|-----|--------------------------|--------------------------|----------------------------|-------|---------|
| d1) | $\frac{1}{3} \times 0.3$ | $\frac{1}{2} \times 0.3$ | $\frac{1}{3} \times 0.125$ | 0 | 0 |
| d2) | 0 | $\frac{1}{3} \times 0.3$ | $\frac{2}{3} \times 0.125$ | 0 | 0 |

lect 5] # word2vec

→ this can distinguish b/w Happy and Joy
(semantic meaning preserve)

→ low dimn vector
(unlike Bag of word / ETC)

→ Dense vec (unlike 2 which has sparse vector).

How to use it

use pre-trained model

train model in your dataset
then use it

↓
google news corpus
trained on 3 billion words
300 dimn vector

Import gensim

from gensim.models import word2vec, KeyedVectors

download google model (msab)

model = KeyedVectors.load_word2vec_format('googleNews - - ', binary=True)

model['cricket'] = [- - - -]
(1230)

we use deep learning to create these features

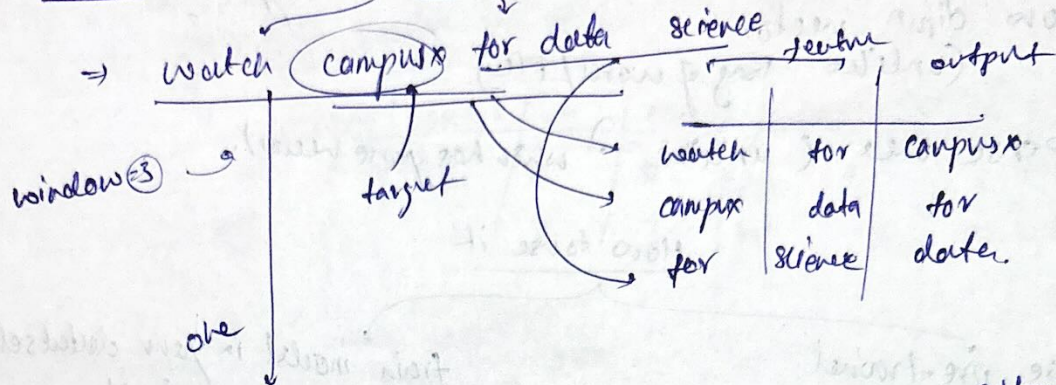
cy
these are features

| | money | power | Active | lender |
|--------|-------|-------|--------|--------|
| king | 1 | 1 | 1 | 1 |
| man | 0 | 0 | 1 | 1 |
| queen | 1 | 0 | 1 | 0 |
| women | 0 | 0 | 1 | 0 |
| Boy | 0 | 0 | 1 | 1 |
| girl | 0 | 0 | 1 | 0 |
| Bottle | 0 | 0 | 0 | 0 |

here
man/women
will have similar
vector.

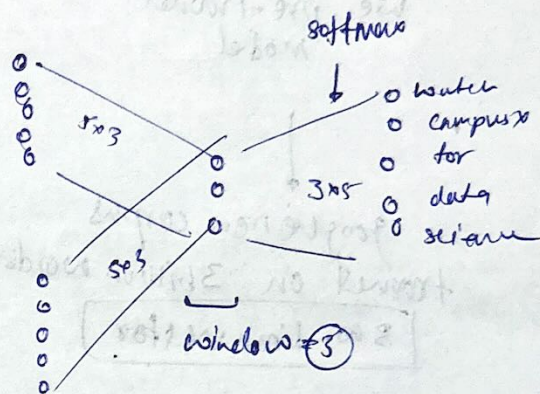
types of word2vec (training model)

① CBOW (for small dataset)



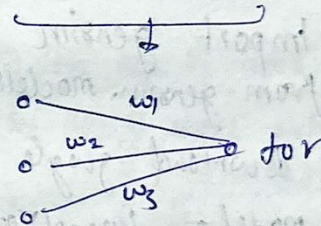
one

| | |
|---------|--------|
| watch | 1 0000 |
| campus | 010000 |
| for | 001000 |
| data | 000100 |
| science | 000010 |



vector representation of "for" is its weights

$[w_1, w_2, w_3]$

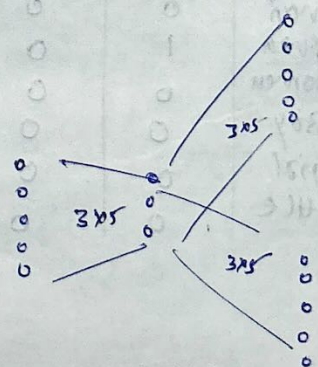


no. of windows = no. of N-vector

② skip-gram (for large dataset)

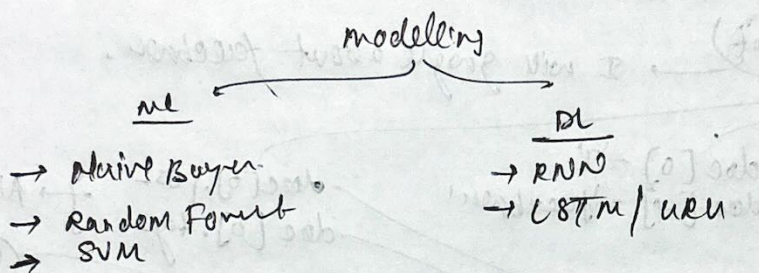
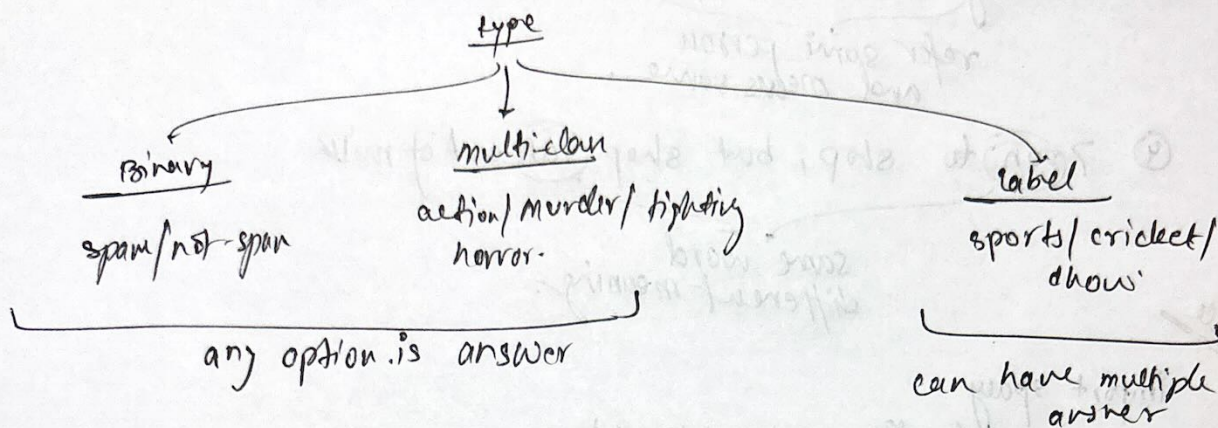
→ this is reverse of CBOW

| Input | output |
|--------|--------------|
| campus | watch, for |
| for | campus, data |
| data | for, science |



lect 6] # Text Classification

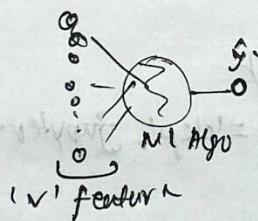
email content \rightarrow spam/not spam
 sentence \rightarrow (+ve)/(-ve) sentiment



methods

Bow/n-gram

each row is converted to vocabulary features.

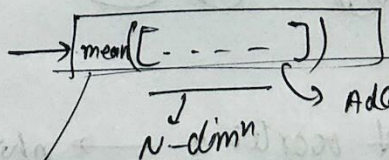


return n -dim vector for each word.

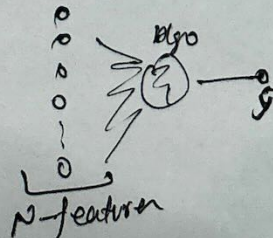
hsgrd 2 vec

train model / use pre-trained model.

P1) Hello Ross, I am good.



Add, axis=0



Text 7 # part of speech tagging (POS)

Applications of POS

- ① Named Entity recognition
- ② Question / Ans system
- ③ Hi my name is Tanish, I'm good.

refer same person
and means same.

- ④ 7 ran to shop, but shop ran out of milk

same word
different meaning.

user Hidden Markov models
Algorithm.
optimization: viterbi algo

code

```
import spacy
nlp = spacy.load('en_core_web_sm')
doc = nlp('text'), I will google about facebook.
```

doc.text
(will return text)

doc[0] = 'I'
doc[-1] = 'facebook'

doc[0].pos-
doc[0].tag-
→ Noun, verb,
etc
code

cats
grow
pos

five
grows
pos

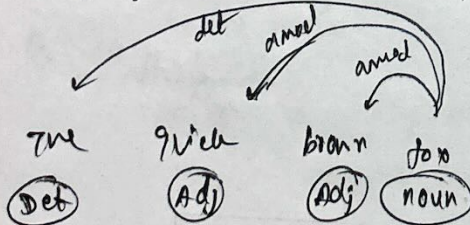
If you don't understand what this code is
spacy.explain('code')

= 'VB / PRP'
verb, base form

pronoun
person

from spacy import display

display.render(doc, style='dep', jupyter=True)



How it works?

→ nripadhar
(40-50 min ka video hai)