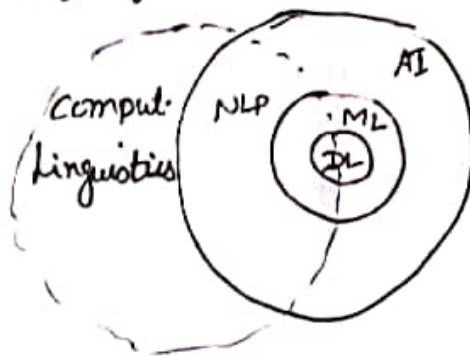


NLP (NATURAL LANGUAGE PROCESSING)

part of AI where computers analyse, understand and derive meaning from meaningful info from Human language



Libraries → NTLK (Natural Lang Toolkit)
 → ~~Scipy~~
 → Spacy

Application → Automatic Summarisation
 → Sentiment Analysis
 → Speech Recognition
 → Topic Segmentation
 → Named Entity Recognition

NLP → NLU (Natural Lang Understanding)
 → NLG (Natural Lang Generation)
 → Lexical Ambiguity
 → Syntactical Ambiguity
 → Referential Ambiguity

Lexical Ambiguity

Ambiguity of a single word

She bagged two silver medals.
 She made a silver speech.
 His worried had silvered his hair.

silver → noun
 → adj
 → verb

Syntactical Ambiguity

Two or more possible meanings

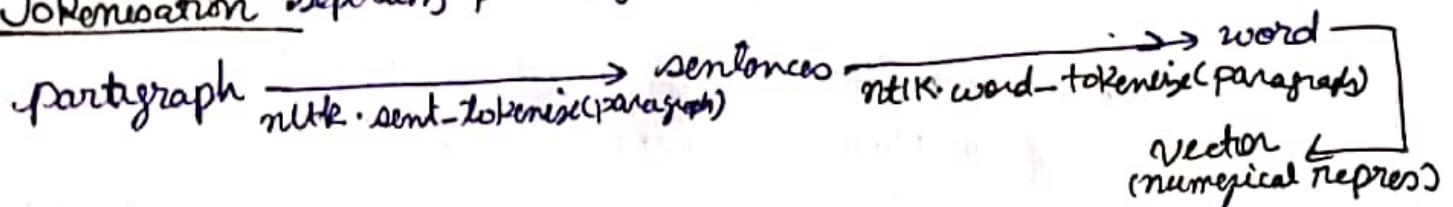
Eg. The chicken is ready to eat.

chicken ~~animal~~ bird is ready to eat
 chicken dish! prepared/cooked?

Referential Ambiguity also have many ref.
 John met Tom & Harry. They went to garden.
 John went with Tom/Harry or both to garden?

Corpus (collection of text documents) → Docs → Para → Sentences → Tokens

Tokenisation separating piece of text into smaller units called tokens



Stemming & Lemmatization used to analyse the meaning of word behind it

Stemming \rightarrow produce intermediate representation of word
 \rightarrow reduce words to their word stem

history \rightarrow histori
historical \rightarrow histori

finally \rightarrow fina
final \rightarrow fina
finalized \rightarrow

going \rightarrow go
goes \rightarrow go
gone \rightarrow

Lemmatization

history \rightarrow history
historical \rightarrow history

finally \rightarrow final
final \rightarrow final
finalised \rightarrow

Converts words to word that a is understood by human

Stopwords stopwords. words ('english')

Eg i, me, myself, in, out, you, you, wouldn't etc.

Bag of Words

Sent1: He is a good boy

Sent2: She is a good girl

Sent3: Boys & girls are good

stop keywords
↓
Lower Sentences
(except count words)
name

Sent1: good boy

Sent2: good girl

Sent3: Boy girl good

Words	Freq
Good	3
Boy	3
girl	2

Bag of Words
converts
this table
to vector.

	good f_1	boy f_2	girl f_3	f_4	f_5
S1	1	1	0	—	—
S2	1	0	1	—	—
S3	1	1	1	—	1

↓
Binary Bag of Word
(Values 1, 0, 1)

f_1, f_2, f_3 are independent variables
& o/p data will be used to train ML/DL
model

Disadv

(i) "good" & "Boy" have equal representation. Semantics
of these words are almost, we aren't able to derive
which word is more imp. (good >> boy)

Solⁿ → TFIDF (term freq & inverse document freq)
↓
Word2vec (huge dataset)

$$\underline{TF - IDF = TF \times IDF}$$

$$TF = \frac{\text{No of representation of words in a sentence}}{\text{No of words in sentence}}$$

$$IDF = \log \left(\frac{\text{No of sentences}}{\text{No of sentences containing words}} \right)$$

sent 1 → good boy
sent 2 → good girl
sent 3 → boy girl good

Histogram
or Freq Table

Words	freq
good	3
boy	2
girl	2

↓
tfidf
vectors

	TF		
	Sent 1	Sent 2	Sent 3
good	1/2	1/2	1/3
boy	1/2	0	1/3
girl	0	1/2	1/3

Words	IDF
	IDF
good	$\log 3/2 = 0$
boy	$\log(3/2)$
girl	$\log(3/2)$

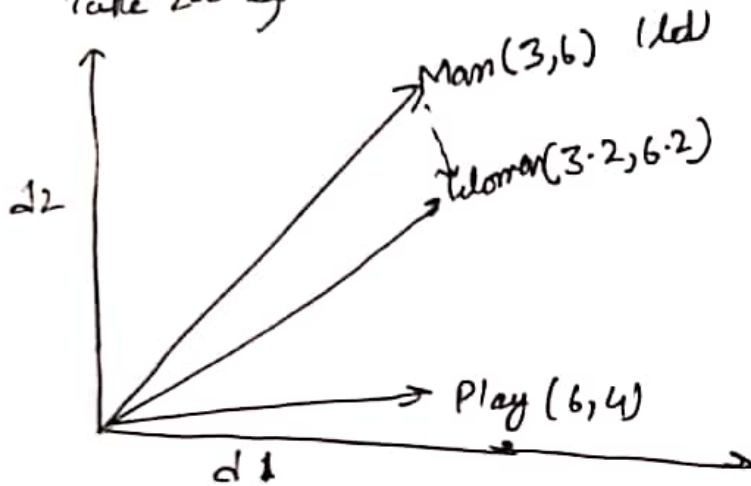
	f ₁	f ₂	f ₃	o/p
	good	boy	girl	
Sent 1	0	$(\log 3/2) \frac{1}{2}$	0	1
Sent 2	0	0	$\frac{1}{2} (\log 3/2)$	1
Sent 3	0	$\frac{1}{2} (\log 3/2)$	$\frac{1}{2} (\log 3/2)$	1

defn: ~~map~~ empty
given to good, boy,
girl, trying to bring
in semantic meaning

Word & Vec

- Both BoW & Tf-Idf semantic info is not stored properly (order of words)
- Tf-Idf gives input to uncommon words
- Chance of overfitting

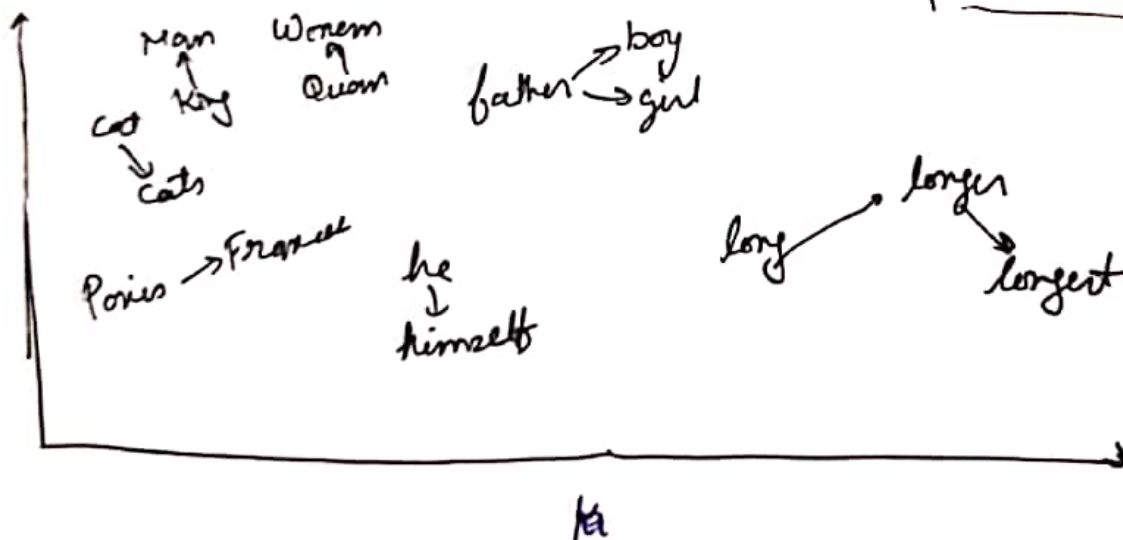
- Each word is represented as a vector 32 or 32+ dimension & not a single word.
- Hence semantic info & relation b/w diff^t word is preserved



- Men & Women are similar words, related with semantics
 \therefore distance b/w them is small

opt can be done

$$\text{King} - \text{Man} + \text{Woman} = \text{Queen}$$



tokenize \rightarrow histogram \rightarrow intake most freq words
 \downarrow
 matrix of all unique words