

28/03/24

Feature engineering & NLP Algo.

gathered corpus

Pre processing our data

divide the classes & finds the imp. pt.

challenges:

- detect good features
- Feature selection : selected no. of features
- Manual feature engineering is time consuming
- should have basic knowledge on domain.

7 techniques used in FE for NLP

dependency parsing

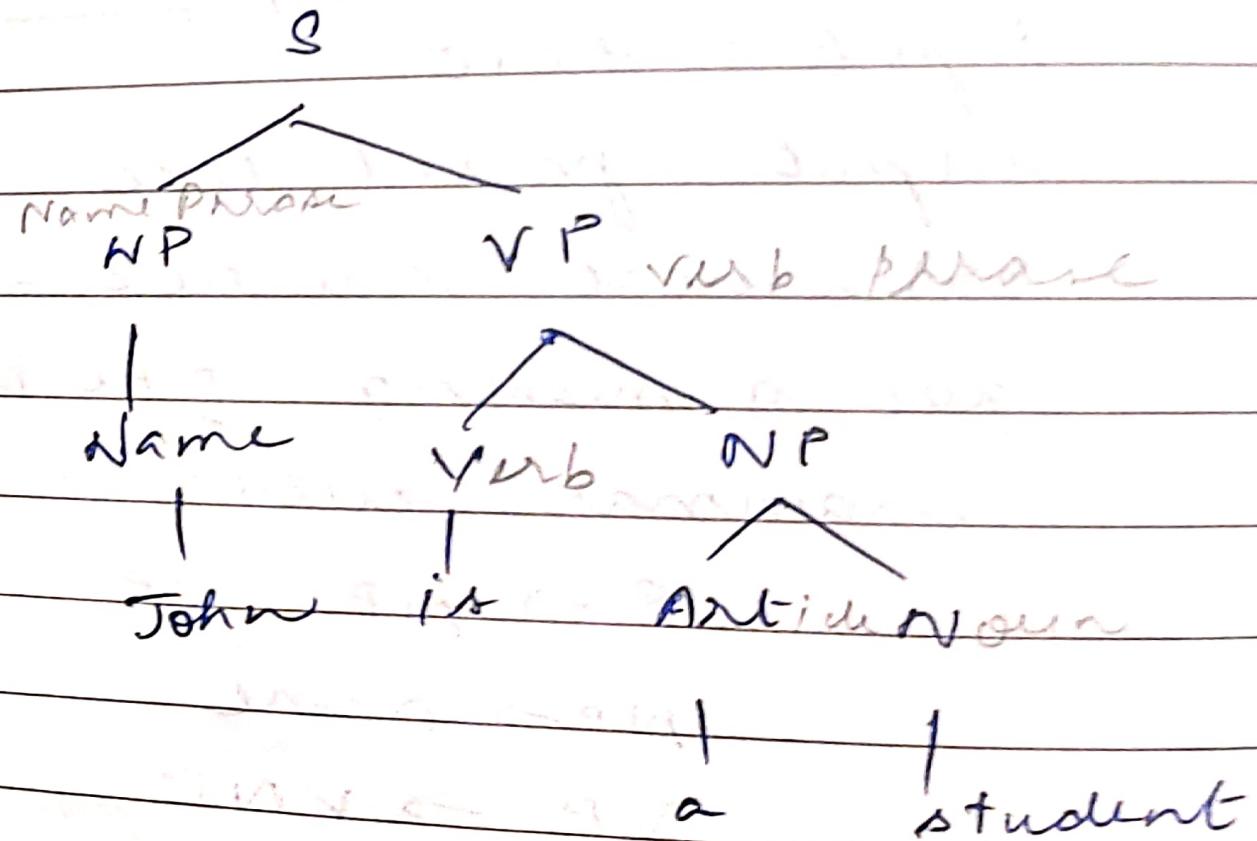
input: text

output: parsed tree

concurrent structure - combining words to give a meaning [He is ^{Tom} ~~falling~~]

grammar rules:

 $S \rightarrow NP VP$ $S \rightarrow \text{sentence}$ $NP \rightarrow \text{Name}$ $NP \rightarrow \text{Noun phrase}$ $VP \rightarrow VNP$ $NP \rightarrow \text{A/S C N/Object}$
nou



This happens which is followed by dependency parsing
 Lexical analysis is part of this

Name → John

Verb → is

mapping in table

2. Bag-of-Words (BoW)

Aim: To verify the syntactical structure of the sentence (need not be grammatically correct.)

Detects ambiguity in the sentence

Used in: regression based algo.

For 2 text document

d_1 : John likes to watch cricket. chris likes cricket too.

d_2 : John also likes to watch football.

BoW: ["John", "likes", "to", "watch", "cricket",
 "chris", "too", "also", "football"]

Repetition is avoided.

Picking out the unique words.

"John", "likin", "to", "watch", "cricket", "Chris", "also", "too",
Frid, 1, 1, 1, 1, 0, 07

"John", "likes", "to",
Eng. g di [1, 1, 1, 1, 1, 1, 1, 0, 0]
d2 [1, 1, 1, 1, 0, 0, 0, 1, 1]

3) Term Frequency - Inverse Document Frequency

$$TF_{(t)} = \frac{\text{no. of times } t \text{ appeared in d}}{\text{total no. of terms in d}}$$

$$IDF_{(t)} = \log \left\{ \frac{\text{Total no. of documents } (D)}{\frac{\text{no. of documents with}}{\text{term t in it}}} \right\}$$

$$TFIDF = TF * IDF$$

eg: d1: This is a car

d₂: This is a bike.

It = this

$$TF(t) = \frac{2}{8} + \frac{1}{4}$$

$$\text{IDF}_{(E)} = \log\left(\frac{2}{2}\right) = \log 1 = 0.$$

$$\text{TF-IDF} = \frac{1}{4} \times 0 = 0$$

$t = \text{car}$

$$\text{TF}_{(E)} = \frac{1}{4}$$

$$\text{IDF}_{(+)} = \log\left(\frac{2}{1}\right) = \log 2 = 0.30$$

$$\text{TF-IDF} = \frac{1}{4} \times 0.3 = 0.25 \times 0.3 = 0.075$$

$\times \quad \times$

Text / Word rep.

- ① dependency Parsing
- ② BOW [0, 1, ...]
- ③ TF-IDF []
- ④ word embedding

④ Word embedding
(converting words to vector)

Word 2Vec

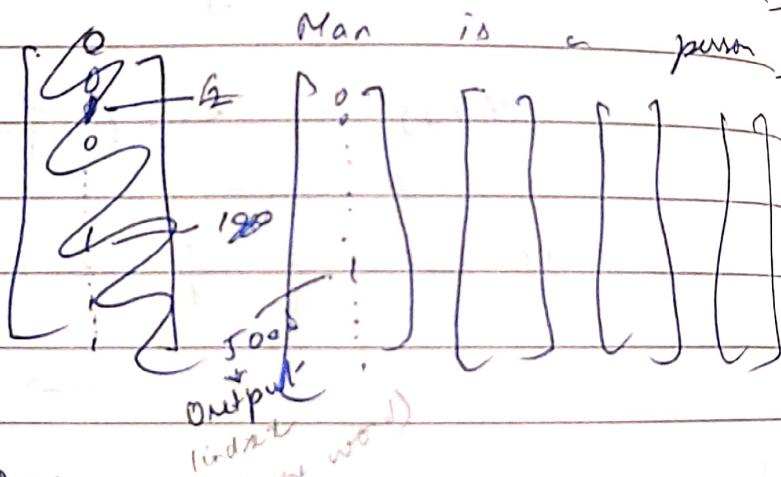
glove

one hot rep.

$$|V| = 10,000$$

Man is a person
index: [5000] [100] [=] [200]

1 hot rep:-



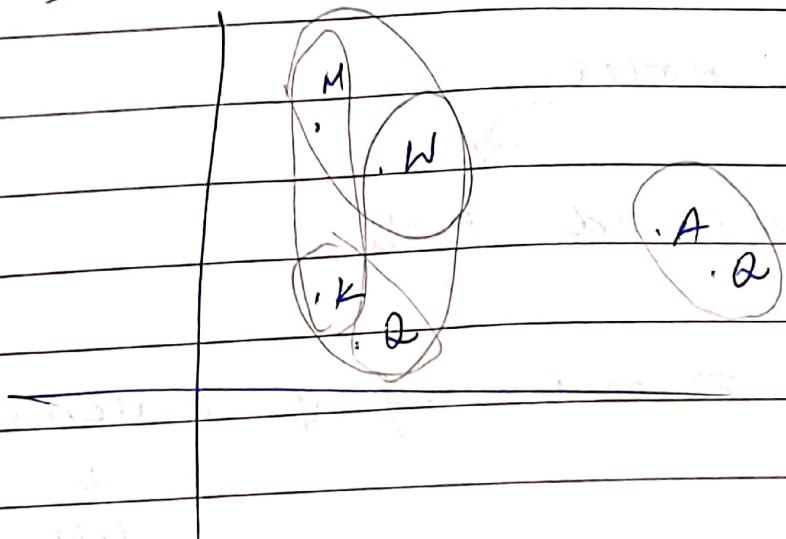
Disadvantage:

- Vector [high dimension]

Now comes Word embedding

Features	Words	Man	Woman	King	Queen	Apple	Orange
Gender							
Human		-1	1	0.93	0.92	0.01	0.01
Fruit		0.98	0.97	-1	0.01	0.0	0
;		0.01	0.0	0.00		-1	1
;							
300	dimension						

Backend -



steps:

- vector dimension mention ($10,000$)
- One hot rep.
- Out of hot rep.
- no. of dimension (no. of features)

Word embedding using keras

① sentences

② 1 Hot Rep. \Rightarrow Index from dictionary③ 1 Hot Rep. \Rightarrow embedding layer to keras
to get Embedding Matrix④ emb. matrix \Rightarrow $\sqrt{\text{size}} = 10,000$

points = 10

technique



n -gram model



no. letter / word Prob.

n -gram \Rightarrow conti. seq. of n items



letter / word

syllable

Types

$n=1 \dots$ unigram

$n=2 \dots$ bigram

$n=3 \dots$ trigram

:

:

1P = The boy ate a chocolate

uni = The boy ate a chocolate

Bi = the boy ate a chocolate

Tri = the boy ate boy ate a ate a

Application of ngram

→ plagiarism tool
(copying others work)

→ computational biology

Language Model: To predict next words in our suggestion

$$P(w) = \frac{\text{count}(w)}{\text{count}(N)}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$\text{Bigram: } P(w_i | w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$$

$$\text{Trigram: } P(w_i | w_{i-2}, w_{i-1})$$

$$= \frac{\text{count}(w_{i-2}, w_{i-1}, w_i)}{\text{count}(w_{i-2}, w_{i-1})}$$

Eg:

corpus:

The girl bought a choc.

The boy ate a choc.

The girl bought a toy.

The girl played with the toy.

step 1: build vocabulary.

B.O.W

{ the, girl, bought, a, toy,
boy, ate, toy, played, with }

Step 2: predict next words, for all
words in vocabulary

I/P : the girl

$$P(\text{the} / \text{the, girl}) = \frac{\text{count}(\text{the, girl, the})}{\text{count}(\text{the, girl})}$$

$$= \frac{0}{3} = 0.00$$

$$P(\text{girl} / \text{the, girl}) = \frac{\text{count}(\text{the, girl, girl})}{\text{count}(\text{the, girl})}$$

$$= \frac{0}{3} = 0.00$$

$$P(\text{bought} / \text{the, girl}) = \frac{\text{count}(\text{the, girl, bought})}{\text{count}(\text{the, girl})}$$

$$= \frac{2}{3} = 0.67 = 67\%$$

$$P(\text{a / the, girl}) = \frac{0}{3} = 0.$$

$$P(\text{croca. / the, girl}) = \frac{0}{3}, 0.$$

$$P(\text{boy / the, girl}) = 0$$

$$P(\text{atv / the, girl}) = 0$$

$$P(\text{toy / the, girl}) = 0$$

$$P(\text{played / the, girl}) = \frac{1}{3} = 0.33 = 33\%$$

$$P(\text{with / the, girl}) = 0$$

-04-24

(2)

⑥

POS Tagging



Parts of speech

POS tagger \Rightarrow netk

Assignment 1

word embedding:

Word ~~vec~~ rec:

~~pre-train word~~ Pre-trained word vec:

- continuous bag of words

'in', 'cat', 'over'

predict → 'jumping'

- skip-gram

How its Trained

- Initialization

- Sliding window

- context prediction

- Optimization

- Iterate

Application:

Advantages of word vec:

- semantic understanding

- Reduced dimensionality

Word2Vec Implementation

Prerequisite : generative models

9.05.24

Applications

Rule-based system

Components

- rules

- parsers

- linguistic resources (dictionary)

Diff b/w Machine learning + deep learning

Rule-based vs. Machine learning approach
self driving cars

Project:

what is ?

Domain ?

data set ?

Next action ?

rule based architecture

01/05/2019

general
archi.

practical
archi.

classmate

Date

Page

general architecture of the rule-based system as an expert system

