

NLP

Start with Example Use Case with define problem → then solve by NLP

→ 1st solve the
→ Solve NLP without any Specific Model like (RNN, LSTM)

① Solve NLP using basic ML

② then go NLP → Embedding
↳ Embed the words
↳ array like

• R.

③ How Represent the language in NLP ← Explore → by Embedding

eg: Translation

① Text generation

② text summarization

③ Speech to text or Text to speech

eg: Sentiment Analysis

Swaggy Review Analysis example code

① Clean text → Lower upper, special char Remove

② Tokenization → create list of words
↳ use NLTK
↳ word tokenize
↳ Stopword Removal (corpus from)
↳ Limitization (Root word)
↳ Stem

ML is not the feature vector sequence important nhi hota hai but in NLP word sequence matters

eg: pizza or pizzas (same na)

③ How ML Models work

① Intro ML

② Convert words / Text Representation

③ Bow (Bag of words) / TFIDF → convert sentences into number
↳ Count Vectorizer (Sklearn)

eg: John likes to watch movies to. → sent 1
→ too likes → sent 2

② TFIDF

→ $TF(t, d) = \frac{\text{Count of term } t \text{ in doc } d}{\text{total terms in doc } d}$

$IDF(t, D) = \log \left(\frac{\text{total docs in corpus } |D|}{\text{num of docs containing term } t} \right)$

TF-IDF Score: $TFIDF(t, d; D) = TF(t, d) \times IDF(t, D)$

ONE HOT REPRESENTATION OF words
• One Model
↳ Multinomial NB
Logistic, Random Sum

Sklearn feature-extract-text
import CountVectorizer
→ TfidfVectorizer

John: $\log \left(\frac{2}{1} \right)$ → 2 sent with docs
→ 1 time John

→ Now Start Predict Review is positive or Negative. ✓

① How Represent words

④ Use One Hot Representation of Words

eg cat $\rightarrow [0, 1]$
 Dog $\rightarrow [0, 0]$
 Giraffe $\rightarrow [1, 0, 0]$

What issue there

- ① millions of words
- ② dog or dogs same but here differs (no make similarity meaning)
- ③ impossible vector - of infinite size one hot
- ④ Also word meaning depends on whole sentence context.

→ Use Word Vectors: Representation of word \rightarrow A dense vector is built such that word vectors of words appearing in its context are similar to this word vector.

→ 2 major types \rightarrow ① Word2Vec & GloVe

① Word2Vec

- Meaning \leftarrow The words appearing around of words a word defines the word.
- we have extensive text corpus, serving as a collection of words (means \rightarrow जिसका word आप को जे उतना ही meaning आयेगा)
 - Each word within this text corpus needs to be denoted by vectors.
 - We identify a center word c and its surrounding context words. (window)
 - By leveraging the similarity b/w vectors represent c & o .

Similar words means \rightarrow probability occurring together is higher

Window size not so much large or small.

Loss func / loss objective function $\rightarrow J(\theta)$

T = total no. of words

$$= \frac{1}{T} \log(\text{Likelihood})$$

$$= -\frac{1}{T} \log \left(\prod_{t=1}^T \prod_{j=-m}^{+m} P(w_{t+j} | w_t, \theta) \right)$$

$$\text{obj. func} = -\frac{1}{T} \sum_{t=1}^T \sum_{j=-m}^{+m} [\log P(w_{t+j} | w_t, \theta)]$$

from gensim models import Word2Vec

$w = \text{Word2Vec}(\text{sentences} = \text{tokens}, \text{vector_size} = 100, \text{window} = 5, \text{min_count} = 1, \text{workers} = 4)$
 for Parallelization for process

$w \cdot wv['food']$

\rightarrow Get the vector repeats of a word (food)

$w2v$
 \rightarrow It is also model like other ML

100 Dimension

Likelihood
 \rightarrow इसका ये बनकर है then serve करने की Probability का है तो इस likelihood को देखते हैं

\rightarrow तो इसे Maximize करना है तो कैसे करेंगे
 \downarrow
 तो हम इसे Mini को -ve कर देते

ये find करना पड़ा है

② Glove Embedding

df = pd.read_csv('/content/glove.6B.txt', sep=" ", quoting=3, header=None, index_col=0)

• glove_vectors = {key: val.values for key, val in df.T.items()}

→ It is pretrained already → So load and use.

Reduce dimensionality using PCA or T-SNE

numpy, Matplotlib

from sklearn.decomposition import PCA
from sklearn.manifold import TSNE

↑
for Visual of 50 Dimension Glove into 2D.

• means इसे PCA के

through वो 2 Dimension
space में सबसे ज्यादा Variance
Cover करता है वो minimum
error देता है।

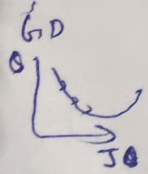
→ Cosine-Similarity

Word2Vec Continues :-

$$P(w_{t+j} | w_t)$$

- कभी word outside में होगा तो कभी वो center में होगा तो उसे time दोनो की probability different-different होगा इसलिए दोनो में P calculate कर रहे हैं

$J(\theta)$ - minimize this



$$P(o|c) = \frac{\exp(u_o \cdot v_c)}{\sum_i \exp(u_i \cdot v_c)}$$

$$P(o|c) = \frac{\exp(u_o \cdot v_c)}{\sum_i \exp(u_i \cdot v_c)}$$

Similarity b/w the 2 Vectors.

eg: Sigmoid $\sigma(x_i) = \frac{\exp(x_i)}{\sum_i \exp(x_i)}$

are occurring together such that 'o' is an outside word & 'c' is a centre word

normalization term:

Here
→ Train a binary logistic reg Model.

- Positive Pair (o & c which are occurring together)
- vs
- Noise pair (o & c randomly chosen).

- 2 kinds of Word2Vec
- ① Skip-gram $P(o|c)$ (more use)
- ② CBOW (Continuous Bag of words) $P(c|o)$ (less time, use for less frequent words in doc)
- which one should use

$$J(\theta) = \log \sigma(u_o^T v_c) + \sum_{i=1}^k E[\log \sigma(u_i^T v_c)]$$

noise pair

expectation → k times draw a sample from word prob distn

↓
So reduce -ve to maximize
↳ So called

→ Negative Sampling Method.

change here by using to reduce complexity (compute)

True Pairs

all pairs with c by k sample

$$\log [\sigma(u_o^T v_c)] + \sum_{i=1}^k [\log \sigma(-u_i^T v_c)]$$

$$\sum_{i=1}^k E[\log \sigma(-u_i^T v_c)]$$

k sample noise pairs

$V = 1000$ words

eg $k=5$

How to decide k-value by using Unigram

$P(w) = \frac{U(w)}{\sum_i U(i)}$

why 3/4 bec rare words have less prob of being drawn from distribution

Language Modeling:

if word **Back** ^{uses} ^{Object} ^{Verb} ^{Adj} ^{Whun.}

then

↳ vector =

→ Actual meaning of **Back** we are
get the context of **Back**
show below |

what is the word x_T
given $x_1, x_2, x_3, x_4, \dots, x_{T-1}, x_T, x_{T+1}$

write all previous words in a text
Predict next occurring words

eg: Students open their ^{books} ^{notebook} ^{pen} ^{Assignment}
eg: google keyboard, search bar prediction

in Probability terms

$$P(x_1, x_2, \dots, x_T) =$$

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

$$P(A \cdot B) = P(A|B) \cdot P(B)$$

$$P(\underbrace{x_1, x_2, \dots, x_T}_B) = P(x_T | x_{T-1} \dots x_1) \cdot P(x_{T-1} \dots x_1)$$

$$\downarrow$$

$$P(x_T | x_{T-1} \dots x_1) \cdot P(x_{T-1} | x_{T-2} \dots x_1) \cdot P(x_{T-2} \dots x_1)$$

$$\downarrow$$

$$\Rightarrow P(x_T | x_{T-1} \dots x_1) \cdot P(x_{T-1} | x_{T-2} \dots x_1) \cdot \dots$$

$$P(x_3 | x_2, x_1) \cdot P(x_2 | x_1) \cdot P(x_1)$$

\propto Probability

$$\prod_{t=1}^T P(x_t | x_{t-1} \dots x_1)$$

→ How to learn it.

① L.M (Pre Deep Learning)

→ Very Powerful for Scratch

use ① n-gram L.M :- word occurring at step 't' is

not dependent on All words before
it but only on 'n-1' word sequence
occurring before it.

eg: $x_1, x_2, x_3, x_4, x_5, (x_6)$
4-gram model

$$P(x_t | x_{t-1} \dots x_t) \sim P(x_t | \underbrace{x_{t-1} \dots x_{t-n+2}}_{n-1 \text{ words before } t})$$

$$P(x_6 | x_5, x_4, x_3, x_2)$$

$$\rightarrow P(x_6 | (x_5, x_4, x_3))$$

for 4-gram

→ be deep learning can use it

N-grams?

$n=1 \rightarrow$ unigram \rightarrow eg: 'student', 'opened', 'their', 'book'

$n=2 \rightarrow$ bigrams \rightarrow "Student opened", "opened their"

$n=3 \rightarrow$ trigrams \rightarrow "Student opened their"

$n=4 \rightarrow$ "Student opened their book"

$$P(x_t | x_{t-1} \dots x_{t-n+2}) \sim P(x_t | \underbrace{x_{t-1} \dots x_{t-n+2}}_{\substack{\uparrow \\ n-1 \text{ words before } t}})$$

↓
How it works

$$\frac{\text{Count}(x_t, x_{t-1} \dots x_{t-n+2})}{\text{Count}(x_{t-1} \dots x_{t-n+2})} \Rightarrow \frac{\text{Count}(\text{ngram})}{\text{Count}(\text{n-1 gram})}$$

eg: $n=4$

• Students opened their books on

$$P(\text{"on"} | \text{"opened their books"})$$

$$\Rightarrow \frac{\text{Count}(\text{"opened their books on"})}{\text{Count}(\text{"opened their books"})}$$

• What is problem in N-gram?

• Drawback

① Losing the context of words some times when depends on other part.

eg: Examiners asked the students to open their ?

book exam assignment

②

→ So its prob is 0 but rather should be
eg: $P(\text{"exams"} | \text{"students to open their"})$

"students to open their exams"

↓

So do Smoothing

ngram \rightarrow does not exist

Give some little probability to given every word - n-gram seq - is something

eg: Exams at correction exam

but - n-gram use most frequent word like \rightarrow book

\rightarrow to we are losing context

③ $n-1$ gram did not occur

eg: students to open their

~~start~~ ~~that~~ ~~port~~ ~~start~~ occurs ~~that~~ & ~~at~~

Prob of this is 0

↓

How to solve it

↓
using Backoff

eg: students to open their -

~~start~~ ~~5 gram~~ & ~~at~~ → 4 gram check ~~start~~ ~~start~~ ~~start~~
... - 3 gram

⑤ store all data

④ What if multiple n -gram have the same count

$$P(\text{"book"}) = P(\text{"exam"}) = 10$$

↓

Solve all drawback by using DL

by try to DL LM

① Neural NLM

② then RNN

7th Dec
2024
NLP