

④ validating the [K] value

(i) kneelocation

(ii) silhouette scoring

~~# NLP~~

- What about features that are completely text?

→ We can convert those text ~~int~~ to continuous values by using [OneHotEncoding]

→ But letter wise OHE is not viable,

so, we've to represent the data using vectors

- If we start OHE for letter wise we'll have

$[0, 0, 0, 1, 0, \dots, 0]$ very big.

⑤ Roadmap

- Token Text Preprocessing : Tokenization, Lemmatization, Stemming, Stopwords
- Text Preprocessing II (Vectors repr) : Bag of words, TFIDF, Unigrams, bigrams

Input → vectors [Better Advanced]

- Text Preprocessing : word2vec, AvgWord2Vec
- RNN, LSTM, GRU,
- Text Preprocessing : Word embedding {uses word2vec}
- Transformer, Best

→ Accuracy increases from bottom to top.



Deep Learning : RNN, LSTM, GRU
↳ TensorFlow, Pytorch

BERT, Transformer : [NLTK, Spacy] libs.



Tokenization in NLP

(i) Corpus : Paragraph

(ii) Documents : Sentences

(iii) Vocabulary : Unique words [Dictionary]

(iv) Words : all the words in the Corpus

Tokenization is process to convert Corpus, documents to tokens i.e. sentences.

[sent1 • sent2]

so,

[sent1]
[sent2] into token

- Split on the basis of full stops.

Then again, apply Tokenization to convert to [words] tokens.

- Once we've all the words we'll find the unique words i.e vocabulary for the corpus [prompt]

Text Preprocessing

- NLTK library [Advanced RegEx, Pattern Matching]
- Sent_tokenize (corpus) @ document

Uses "Punkt" model for tokenization.

- word_tokenize (document)

Converts every character, words to separate token.

[word, ;, !, ...]

Then we clean to make actual words.

- wordpunct_tokenize

• Punctuations are also considered tokens.

• TreeBankWordTokenizers()

→ considers E.g full stops with the word itself.

E.g I am good.

So, ["I",
"am",
"good."]

④ Stemming

Process of reducing a word to its word stem that affixes to suffixes and prefixes or to the roots of words unknown as a lemma.

In data, we can have same words with different verbs such as

eat, eaten, eat, eating, eats

whose root word is Eat

So, we can only use Eat instead of every single word.

• PorterStemmer

Stem. PorterStemmer

• Stem(word) ⇒ stemword

E.g. "eating" \Rightarrow eat

"History" \Rightarrow histori (disadv, for some words, not accurate)

- `RegexpStemmer()`

\rightarrow can use Regex for each word to generate the Ustem.

RegexpStemmer ('ing\$') checks if end is ing and give Stem.

- `Snowball Stemmer`

- Better than Postex Stemmer

- Stem English

Snowball Stemmer

SnowballStemmer ('english').stem(word)

\rightarrow still not efficient. for some words

Use "Lemmatization"

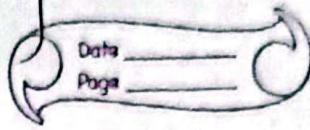
~~X~~ Lemmatization

- Similar to Stemming. The output after lemmatization is lemma which is root word, rather than stem.

`nltk.WordNetLemmatizer()`

\rightarrow wrapped around Wordnet Dataset.

→ Lemmatization preserves the meaning



- WordNet Lemma() lemmatizes ("work", Pos

noun verb adverb

Pos = n, v, a, s

We can find for noun, adverb, verbs etc.

→ WordNetLemma takes a lot of time.

can be used for Chatbots, text summarization.

④ StopWords

In a data, words like I, the, have, and, to does not play a big role. ④ we need to remove such words.

Stopwords

→ Every language has Stop words

e.g Nepali → [ए, ए, उँ]

We need to remove such words.

→ We can download or have our own set of Stopwords.

Download and remove those words.

⑤ Steps

- Tokenization
- StopWord [if not in Stopword apply Lemmatization]
- Lemmatization, stemming, Snowball
- Converts words to lowercase and then apply the Stopword, Lemma, Stemming

⑦ Parts of Speech Tagging

- Each parts of speech has a representation in short (Abbreviation)

`nltk.pos_tag(word) \Rightarrow Reps of Part of Speech Tuple.`

For ex,

`[('I', 'PRP'), ('three', 'CD'), ('vision', 'NNS')]`

- PRP \Rightarrow Personal Pronoun
- NNP \Rightarrow Proper Noun.

⑧ Name Entity Recognition

Earlier we've got Parts of Speech but we can also find the category or entity of each word. Perfect example of category.

For e.g: "The Eiffel tower was built by Bigut"

Person : Bigut

Place : India

Date :

Money:

- ne_chunk(nltk.pos_tag(word))

We can draw and visualize as well.

How to Approach a Dataset / Problem?

- Identify the Problem i.e. Summarization, sentiment, Classification, Translation.

(i) Start with Cleaning (Tokenization to words)

- Use RegEx for data cleaning [Rows]

- Lower CASE

(ii) ~~Tokenization~~ Stemming, Lemmatization, Stop words, Parts of speech, Name Entity.

(iii) Join and prepare the data for Embedding [vectors . xps]

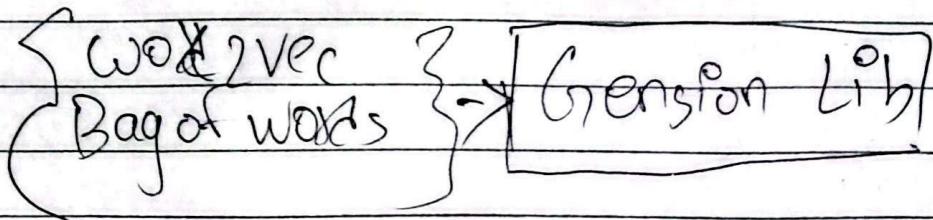
- Transformers, BERT [Word Embeddings]

- Word2Vec

- Bag of words, TF-IDF

- One Hot Encoding [very less category] [Not eff]

(iv) One word vector representation, we train with ML, Deep Learning Algos.



~~Text to Vector~~

- OneHotEncoding

→ Unique words in the dataset → Sparse matrix

→ Length will be length of unique words.

e.g if, Total unique words are 3
then,

To represent a sentence with 4 words.

$\begin{bmatrix} [0, 1, 0], [1, 0, 0], [0, 0, 1], [0, 1, 0] \end{bmatrix}$

↳ sent 1.

$3 \times 4 \Rightarrow 12$ i.e 12 values for sent of n words.

→ Not efficient.

→ Sparse leads to overfitting.

→ No semantic meaning.

- Bag of Words

→ Lower all words

→ Stop words [Reduce words]

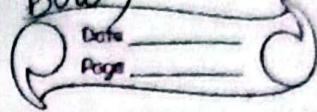
→ Calculate vocabulary (unique)

→ Calculate the frequency of vocabulary.

→ Sort on Descending Order

→

- Increase the count (Non Binary 'Bow')
- Boolean for (Binary)



eg

He is a good boy \rightarrow good boy good
She is a good girl \rightarrow good girl

* vocab

good \rightarrow 3
boy \rightarrow 2
girl \rightarrow 2

* Vector & Pk

	good	boy	girl	
s_1	1	1	0	\rightarrow vectors.
s_2	1	0	1	

cont...

→ This approach has less values, represents more with less

→ ~~Tokens~~

→ Does not have fixed size inputs.

- Ordering of the word is changing.
- Does not consider any new words while testing i.e out of vocab problem
- No capture Semantic meaning
[which is the most imp word]

- N-grams (How to build vocabulary)

e.g) The food is good

The food is not good

Vocab: [food not good]

S_1 | o |] vectors should be
 S_2 | | | - very diff, but they aren't
 [solve with ngrams]

so, what if while building vocabulary
 we take combination of vocab.

i.e Bigrams, Trigrams

Two words Three words \Rightarrow Combination

Then vocab becomes

food	not	good	food	good	food	not	not	good
S_1		o			o		o	
S_2				o				

Now, if we compare this vectors
 they are very different. This shows semantic
 meaning.

Sklearn \rightarrow n-grams (1,1) \rightarrow unigrams
 \rightarrow (1,2) \rightarrow Unigram, bigram
 \rightarrow (1,3) \rightarrow Unigram, bigram, trigram.

- feature_extraction.text import CountVectorizer.
 $\text{CountVectorizer(n_gram_range} = (1,2)\text{)} \text{ bigram}$

• TF-IDF [Term Frequency - Inverse Document Frequency]

(i) $TF = \frac{\text{No. of repetition of words in sentence}}{\text{total words in sentence.}}$

(ii) $IDF = \log_e \left(\frac{\text{No. of sentence}}{\text{No. of sentence containing the words}} \right)$

* TF

	S_1	S_2
good	1/2	1/2
boy	1/2	0/2
girl	0	1/2

* IDF

$$\text{good} \quad \log_e(3/2) = 0$$

$$\text{boy} \quad \log_e(2/1)$$

$$\text{girl} \quad \log_e(2/1)$$

We multiply TF * IDF

	good	boy	girl
S_1	$1/2 \times 0 = 0$	$1/2 \times \log_e(2/1)$	0
S_2	0	0	$1/2 \log_e(3/2)$

} vector

→ Train this vector

Adv's

Date _____
Page _____

- Fixed size
- Word importance is being captured [Repeated words have less value]

Dis

- Out of vocabulary problem

• text import & TfIdf vectorization.

Word Embedding

- Representation of text in vector.
- Text is encoded into vector.

- Words with close meaning are closer in the vector space
- Words with diff meaning are far in the vector space.

e.g

Happy ≈ excited [so closer vector space]
Angry ≠ Happy [far and opposite space]

→ Distance is calculated, cosine similarity.

Word Embedding

↓
Count or Freq
• OHE, Bow, TF-IDF

↓
Deep Learning
Trained models
• CNN, RNN, VAE

→ Deep Learning has better accuracy.

Word2Vec (Gensim library)

Two types (Based on Training)

- CBOW (continuous Bag of words)
- Skipgram

- Text to Vector
- Uses Neural Networks to learn word associations from a large corpus of text.
- Once trained model can detect synonyms, words or suggest additional words for a particular sentence.
- Represents each distinct words with a particular list of numbers called a vector.



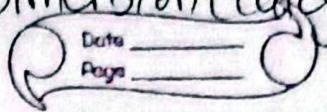
Intuition

→ First we've a set of dimensions or features [category] that word belongs to.

Say 300 dimension, such as Gender, Food, Animal, Location etc

→ Then our vocabulary (Input's vocab) would be of 300 dimension i.e one word would be represented by 300 dim vector.

- Relation of word with the dimension (category)
→ how strong the reln is.



e.g.

(dim)	Vocab	Boy	Girl	Mi	You	Mango
Category						
Gender	-1	1				
Royal	0.01					
Age	0.03					
Food						
Animal						

- These all values comes from the trained Word2Vec model.

Now,

After we've vector representation,

$$[\overrightarrow{\text{KING}} - \overrightarrow{\text{MAN}} + \overrightarrow{\text{QUEEN}}] = \cancel{\text{KING}}$$

CTRL ↪ closer to QUEEN.

A Cosine Similarity

- Distance based on point i.e Euclidean Dist

- Distance based on Angle

$$\text{i.e } d = 1 - \cos(\theta)$$

$$= 1 - \cos(45)$$

$$= 0.29$$

$$[0.29 < 1]$$

But,

$$d = 1 - \cos(0)$$

$$= @ 1$$

So, the vector with angle 45 is closer to another vector.

- ① The goal is to first find the total feature representation i.e Dimension
- ② Then create vector for each word in that Dimensions to understand the where each value is the relation in the vectors Shows it's relationship with n^{th} Dimension as Feature.

For e.g if any word has

$$[0.25, -1, 0, 1, -0.25] \Rightarrow \text{cat}$$

$\swarrow \quad \downarrow \quad \uparrow \quad \searrow$

Dog	Plant	Animal
-----	-------	--------

③ Ways to Train Word2Vec

- Train from scratch
- Train Pretrained model

(i) (Bow) [Continuous Bag of Words]

- Dataset [Large] [Input, Output] [Supervised]

We decide Window-Size = k [Feature Representation (Dimension)]

We take k word token, and find the middle token.

Middle token becomes output (target)
Other token becomes input.

e.g Dataset

Input: Fineuron Company is related to AI
Output: IS related.

[Fineuron company, related to]
[company is related to AI]

Then, we [OneHot] encode the inputs
and outputs.

Data 1 = $\begin{bmatrix} [1, 0, 0, 0, 0, 0], [0, 1, 0, 0, 0, 0], \\ [0, 0, 0, 1, 0, 0], [0, 0, 0, 0, 1, 0] \end{bmatrix}$

Output = $[0, 0, 1, 0, 0, 0]$

Then we fed, into the ANN

So, Now sending 1st. word of 1st sentence
is

$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

• We've fixed input i.e - window - 1

So, (Bow) looks like below.



[words]

For, $W-1$ input we need $\text{len}(\text{voc})$ for each word.

so, If $W=5$ then and $\text{len}(\text{voc})=6$
then,

Input for CBOW model would require

~~4x6~~ neurons $\Rightarrow 24$

i.e INPUT

S_1, E_0, I

0
0
0
0
0

hidden

E_3, E_1, I

0
0
0
0
0

S_2, E_2, I

0
0
0
0
0

$y - \hat{y}$

0	0	0
0	0	0
0	0	0
0	0	0
0	0	0

weights = 4×6

S_3, E_3, I

0
0
0
0
0

weight = $4 \times 6 \times 4$

- Hidden would be of size $[W-1]$
- Output would be of $\text{len}(\text{voc})$

- Then, $y-y'$, Loss, Backprop
- Update weight and Bias

Then our model is ready.

(*) Note:

- Feature Representation = Window-Size
- More window, better the model.

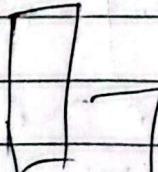
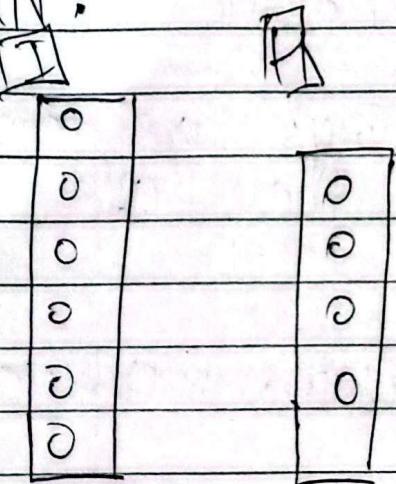
(ii) Skip Gram

~~Everything~~ While data prep in CBOW we take middle as output and rest as input.

Here, we flip ft.

OUTPUT

NN:



6×4 dimension

$$w = 4 \times 6$$

$$w = 6 \times 4$$



- Output \Rightarrow Softmax \Rightarrow Activation
- Hidden \Rightarrow ReLU \Rightarrow Activation

- Training is sum.

(*) When to use CBOW and SkipGram?

- Small dataset \rightarrow CBOW
- Large Dataset \rightarrow SkipGram

Google's Word2Vec

\rightarrow Trained with 3B dataset
 \rightarrow Feature Reps of 300 dim.

(*) Advantage of Word2Vec

- \rightarrow Not sparse matrix \rightarrow Dense Matrix
- \rightarrow Semantic is captured
- \rightarrow Inference is easier [Cosine]
- \rightarrow Fixed set of dimension vector.
- \rightarrow As it is trained with large dataset
out of vocabulary is not possible.

Now, yes the embeddings are generated for each word. But they are very big dimension.

So, if we use Pretrained Word2Vec for say, "The food is good"

[300d] [300d] [300d] [300d]

To represent a sentence with 4 words we need 4×300 values i.e. 1200 inputs required for our own training.

We need to represent a complete sentence or document in less dimension.
i.e say 300 only.

One way is to take average of all the vectors word wise and take average.

So, 300 dim will rep the sentence

This is known as AvgWordVec.

We can build Neural Net with input of only [300] neurons.

④ Pretrained Word2Vec

gensim.models import Word2Vec

.most_similar("word") \Rightarrow gives the top similax.

Once we've all the embeddings, we choose the machine learning or deep learning model to train the data.

Based on Classification or Regression choose model.