

NLP

Natural Language Processing

NLP helps machine understand the human language

- both Text and Speech

Important; Machines do not understand the raw text so NLP converts words or sentences into Numerical Representations (vectors), which a machine can learn from.

NLP BUILDING BLOCKS

Task	Example
Tokenization	Breaking text into words: "This is NLP" → ["This", "is", "NLP"]
Stopwords removal	Removing common words like "is", "the", "a"
Stemming/Lemmatization	Reducing words to their root form: "running" → "run"
Part-of-Speech Tagging	Identifying word roles: "book" as a noun or verb
Named Entity Recognition (NER)	Detecting names, places, dates
Text classification (spam, ham)	Spam vs. Not Spam, Sentiment: Positive/Negative
Text generation GPT	Writing text automatically like ChatGPT does

VECTORIZATION PROCESS

How Text is Turned into Vectors (Simplified):

Method	Description	Example
One-hot encoding	1s and 0s representing presence/absence	"cat" → [0, 1, 0, 0]
TF-IDF <i>Weighted Frequency</i>	Weighted frequency of words	More weight to rare, important words
Word2Vec / GloVe	Words in vector space based on meaning	"king" - "man" + "woman" ≈ "queen"
BERT / Transformers	Contextual word embeddings	"bank" has different vector in "river bank" vs. "money bank"
Sentence embeddings	Whole sentence as a vector	Used for similarity, summarization, etc.

TFIDF = Term Frequency Inverse Doc. Frequency

NLP PIPELINE

Example NLP Pipeline (for text classification like sentiment analysis):

1. Raw text

→ "This is a fantastic product!"

2. Tokenization / Lemmatization

→ ["This", "is", "a", "fantastic", "product", "!"]

3. Lowercasing

→ ["this", "is", "a", "fantastic", "product", "!"]

4. Stopword removal

→ ["fantastic", "product"]

5. Stemming or Lemmatization

→ ["fantast", "product"] or ["fantastic", "product"]

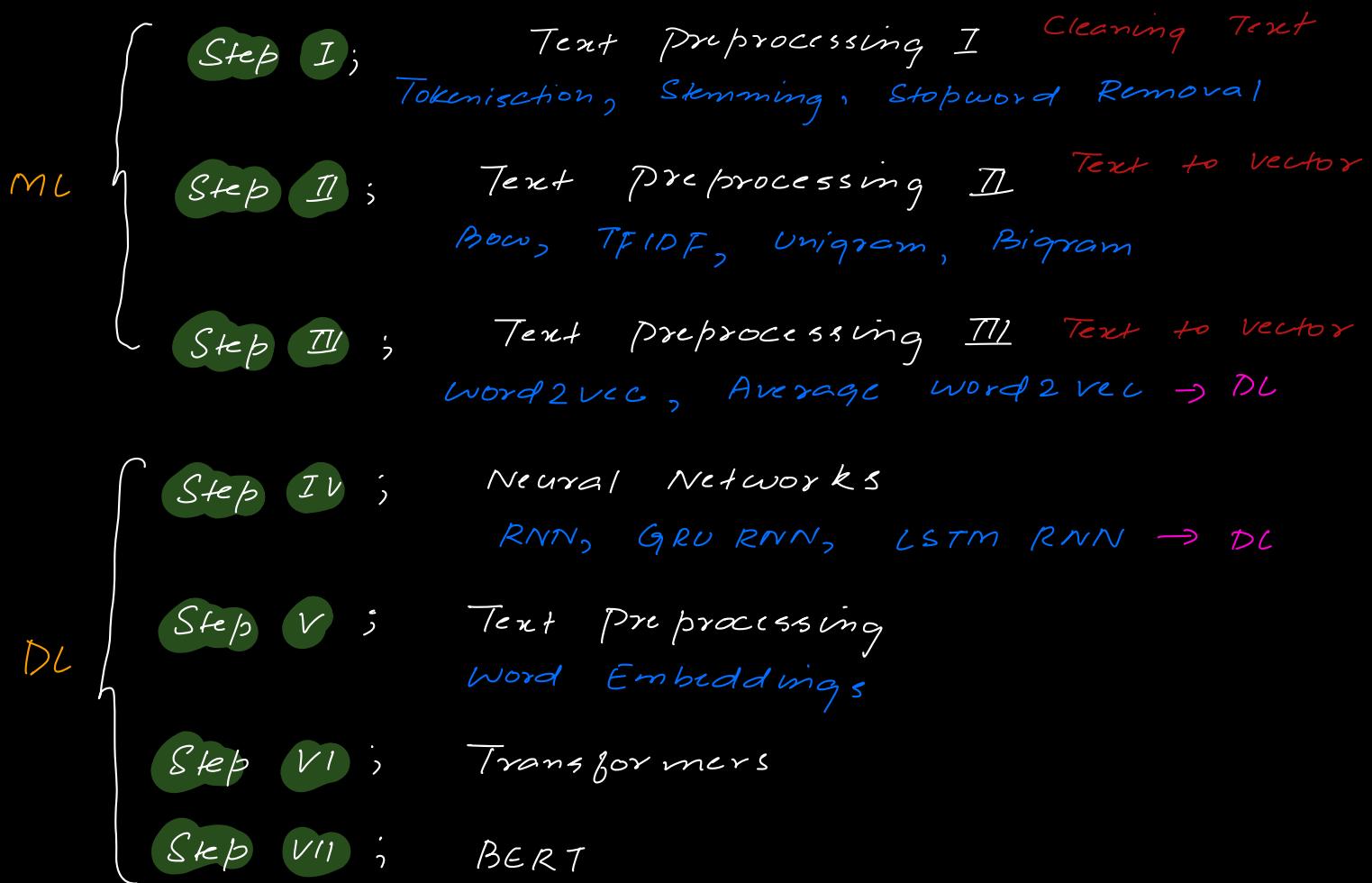
6. Vectorization (TF-IDF / Word2Vec / BERT)

→ [0.01, 0.34, 0.12, ...] (a vector)

7. Model input (e.g. Sentiment Classifier)

→ Predicts: Positive ✓

ROAD MAP



Libraries for Machine Learning

ML = NLTK, Spacy

DL = Tensorflow, Pytorch
↓ ↓
Google Facebook

USE CASES; Text to voice
Text translator
Aurra, GPT etc.

1. TOKENIZATION

breaks down text into smaller units (tokens)

Keywords

outputs in list

- Corpus - Paragraphs
- Documents - Sentences
- Vocabulary - Unique words
- words - Total words

Important we can have words or Sentences in tokenisations as tokens.



Example; "Tajamul is a Data Scientist . (sentence 1)
(Corpus) 1 2 3 4 5
He stays in Hyderabad " (sentence 2)
 6 7 8 9

Tokens = Sentence 1 and Sentence 2
(Sentences) 1 2

words = 9 total words

Vocabulary = 9 unique words (Vocabulary = 10)
(unique words)

TOKENISATION PRACTICALS

NLTK = Natural Language Toolkit

Two types of Tokenize = Sent - tokenize } sentence
word - tokenize } words

```
corpus=""Hello Welcome,to Krish Naik's NLP Tutorials.  
Please do watch the entire course! to become expert in NLP.  
"""
```

✓ 0.0s

```
## Tokenization  
## Sentence-->paragraphs  
from nltk.tokenize import sent_tokenize
```

✓ 0.0s sent_tokenize (corpus)

```
## Tokenization  
## Paragraph-->words  
## sentence--->words  
from nltk.tokenize import word_tokenize
```

✓ 0.0s word_tokenize (corpus)

Wordpunct - tokenize ; considers punctuations
as separate words as word - tokenize does not

wordpunct - tokenize = 's = '① , s ②
word - tokenize = 's = 's only one

```
from nltk.tokenize import wordpunct_tokenize
```

✓ 0.0s

Considers punctuation as separate word!

```
wordpunct_tokenize(corpus)
```

Imp; Most of the we will be using just

Sent - tokenize ()
word - tokenize ()

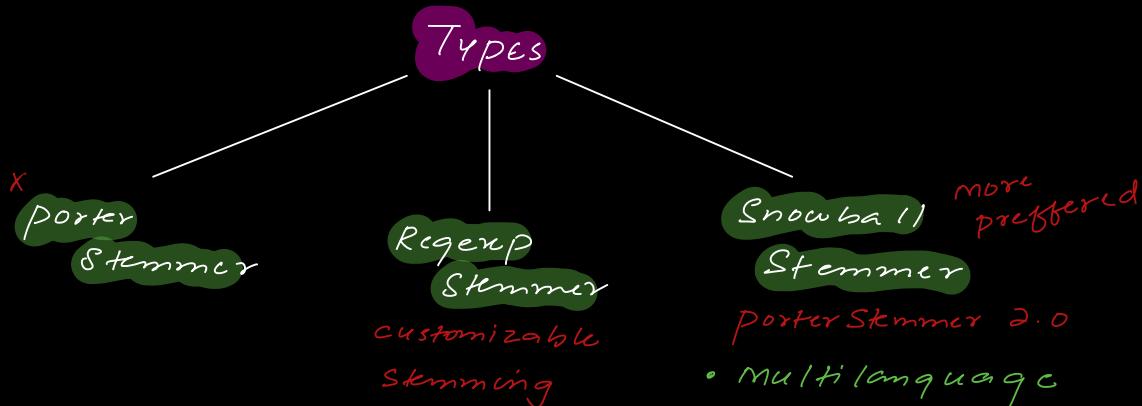
2. STEMMING

is the process of reducing word to its stem

example:

going }
 goes } Stem
 gone }

Stemming helps reduce redundant inputs



Imp. Porter Stemmer = some words may not be stemmed properly like

Congratulations → Stem → Congratul.

STEMMING PRACTICALS

Porter Stemmer

```

from nltk.stem import PorterStemmer
stemming=PorterStemmer()
stemming.stem('congratulations')

```

✓ 0.0s

'congratul'

Regexp Stemmer

```

from nltk.stem import RegexpStemmer
#ing$ means remove last ing
reg_stemmer=RegexpStemmer('ing$|s$|e$|able$', min=4)
reg_stemmer.stem('ingesting')
reg_stemmer.stem('eating')

```

✓ 0.0s

'eat'

Snowball Stemmer

```

from nltk.stem import SnowballStemmer
snowballs_stemmer=SnowballStemmer('english')

# Porter Stemmer
stemming.stem("fairly"),stemming.stem("sportingly")
# Snowball Stemmer
snowballs_stemmer.stem("fairly"),snowballs_stemmer.stem("sportingly")

```

✓ 0.0s

('fair', 'sport')

Problems with Snowball Stemmer

```

# Snowball Stemmer
print(snowballs_stemmer.stem("goes"))
print(snowballs_stemmer.stem("history"))

```

✓ 0.0s

goe
histori

→ Problems with snowball

Imp. Snowball Stemmer is better but it is still not reliable for spam classification projects, so we do Lemmatization, because it covers vast dict of words.

wordnet
dict.

3. LEMMATIZATION More Accurate

is more like advanced version of stemming

✓ In short:

- Stemming = Fast, crude, not always real words. → NO word meaning ^{Imp}
- Lemmatization = Slower, smarter, returns actual base forms. → word mean Aware

example

Better

Stemming

Better

Lemmatization

Good (more Accurate)

- The output we get after lemmatization is "lemma"

^{Imp} NLTK provides **WordNetLemmatizer()**

- It also keeps (POS) into consideration

```
lemmatizer.lemmatize("goes",pos='v')
```

Pos (Part of Speech) Tag	
Noun	n
Verb	v
Adjective	a
Adverb	r

...
Noun-n → Names
verb-v
adjective-a
adverb-r
...
lemmatizer.lemmatize("going",pos='v')

- Use nltk.download('wordnet')

LEMMATIZATION PRACTICALS

```
import nltk
#nltk.download('wordnet') → Download wordnet

## Q&A, chatbots, text summarization
from nltk.stem import WordNetLemmatizer Imp
lemmatizer = WordNetLemmatizer()

lemmatizer.lemmatize("fairly",pos='v'), lemmatizer.lemmatize("sportingly")
✓ 0.0s
('fairly', 'sportingly')
```

Q,, which one takes more time wordnet lemmatizer or stemming?

wordnet lemmatizer takes more time because it compares words to **wordnet Corpus**

^{Imp} Mostly used for Q&A, chatbots, text summarize

4. STOPWORD REMOVAL

where we remove common, insignificant words from the text

example; is, a, the, on, in, for, to

```
stopwords.words('english')  
✓ 0.0s  
['a',  
'about',  
'above',      etc.  
'after',  
'again',  
'against',  
'ain',
```

```
from nltk.corpus import stopwords  
nltk.download('stopwords')  
✓ 0.0s
```

All these words won't play huge impact on our spam classification model

```
# Tokenize  
sentences=nltk.sent_tokenize(paragraph)  
  
# Lemmatizer  
from nltk.stem import WordNetLemmatizer  
lemmatizer=WordNetLemmatizer()  
  
## Apply Stopwords And Filter And then Apply Snowball Stemming  
for i in range(len(sentences)):  
    sentences[i]=sentences[i].lower()  
    words=nltk.word_tokenize(sentences[i])  
    words=[lemmatizer.lemmatize(word.lower(),pos='v') for word in words if word not in set(stopwords.words('english'))]  
    sentences[i]=' '.join(words)# converting all the list of words into sentences
```

Imp. It is more preferred to have a customized list of stop words.

5. PARTS OF SPEECH TAGS

Grammatical Role of each word in sentence

e.g., NN = Noun Singular, NNS = Noun plural

```
import nltk  
nltk.download('punkt')  
nltk.download('averaged_perceptron_tagger')  
  
text = "The quick brown fox jumps."  
tokens = nltk.word_tokenize(text)  
pos_tags = nltk.pos_tag(tokens)  
  
print(pos_tags)  
✓ 0.0s  
[('The', 'DT'), ('quick', 'JJ'), ('brown', 'NN'), ('fox', 'NN'), ('jumps', 'NNS'), ('.', '.')]
```

Imp If you want to know pos-tag, avoid using lemma / stemming.

5. NAMED ENTITY DIAGRAM

Identify & classify Named Entities into Categories

Example

Entity Type	Examples
Person	"Elon Musk", "Sachin Tendulkar"
Organization	"Google", "UN", "NASA"
Location	"India", "New York", "Amazon River"
Date	"June 29, 2025", "last Friday"
Time	"5 PM", "noon"
Money	"\$100", "₹5000", "10 euros"
Percent	"40%", "half"
Miscellaneous	"iPhone", "Python", "FIFA"

REVISION

1. Tokenisation sent - tokenize word - tokenize

 ↳ corpus → sentences / words

2. Stemming snowball stemmer. stem()

 ↳ porter, Fast but less accurate

 Rqcxp,

 Snowball → More preferred

3. Lemmatization slower & more accurate

 wordnet lemmatizer

 word-net Lemmatizer

 ↳ based on wordnet_corpus

4. Stop word Removal import stop words.

LIFE CYCLE OF NLP PROJECT

1. Text preprocessing

2. Vectorization

3. Machine Learning models.

ONE HOT ENCODING; is a text preprocessing technique used to convert categorical data into numerical format that ML models can understand.

WORKING Food is Good D_1
Amazing Food D_2

STEP I; calculates vocabulary

Food is Good Amazing

1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1
0	0	0	0

$$V = 5$$

STEP II; consider sentences

$$D_1 \rightarrow \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \rightarrow 3 \times 4$$
$$D_2 \rightarrow \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{bmatrix} \rightarrow 2 \times 4$$

*I/P
No fixed
size hence
can't store.*

Q, why is it not preferred for NLP?

- Overfitting (Sparse matrix)
- No fixed size input
- How words are related is not captured
- OOV = Out of vocabulary.

Imp; For any ML Algo fixed I/p is required

BAG OF WORDS

V-Important Text classification

is a technique used for vectorization of text

Imp. It is frequency based, more frequently occurring has

Imp

more value

Binary BOW vs Normal BOW

Feature	Normal BOW (Count-based)	Binary BOW (Presence/Absence)
Meaning	Counts frequency of words	Only marks presence (1) or absence (0)
Values	Integers (0, 1, 2, ...)	Binary (0 or 1)

WORKING OF BOW

Text	O/P		
He is a good boy	1	lower all other words case	S1 → good boy
She is a good girl	1	→	S2 → good girl good
Boy and girl are good	1	Stopwords	S3 → Boy girl good [school] [Test]
Vocabulary	frequency	[good boy girl]	O/P
good	3	↓	S1 [1 1 0] 1
boy	2	↓	S2 [1 0 1] 1
girl	2	↓	S3 [1 1 1] 1
Binary Bow and Bow	{ 1 and 0 }	{ Count will get updated back on frequency }	

ADVANTAGES;

- O/P is fixed
- Simple & Intuitive

new word in
test data X

DISADVANTAGES;

- Sparse Matrix
↳ Overfitting
- Ordering of words
change means
the meaning changes
- OOV - Out of Vocab

Difference b/w OHE & BOW?

✓ Your corrected understanding:

- ✗ "OHE → vector for every word" → ✓ Correct
- ✗ "BOW → vector for each sentence" → ✓ Correct
- ✓ But they are not the same kind of vectors — OHE is word-level, BOW is sentence-level.

Bow PRACTICALS

Bow = Count vectorizer

```
## Create the Bag OF Words model
from sklearn.feature_extraction.text import CountVectorizer
## for Binary BOW enable binary=True
cv=CountVectorizer(max_features=100,binary=True)
```



Most frequently occurring words

```
X=cv.fit_transform(corpus).toarray()
```

Imp Max features means, I need 100 words that are more intensive

N Grams

continuous sequence of n items (usually words or characters)

💡 Types of N-grams:

N	Name	Example (from sentence: "I love NLP")
1	Unigram	["I", "love", "NLP"]
2	Bigram	["I love", "love NLP"]
3	Trigram	["I love NLP"]

You can go higher: 4-grams, 5-grams, etc., but they're less commonly used due to sparsity.

⚙️ Why Use N-grams?

- Capture context: Unigrams lose context, but bigrams/trigrams capture word relationships (e.g., "New York" vs. just "New").
- Better feature extraction for models like Naive Bayes or Logistic Regression in text classification.
- Useful in:
 - Language modeling
 - Text generation
 - Spelling correction
 - Plagiarism detection

INTUITION

Unigrams don't capture context

$S_1 = \text{Food } \text{is } \text{Good}$	$= \begin{matrix} 1 & 0 & 1 \end{matrix}$
$S_2 = \text{Food } \text{is } \text{not } \text{Good}$	$= \begin{matrix} 1 & 1 & 1 \end{matrix}$
Vocab Food Not Good	Not much difference

Bigram Vocab = Food Good
 Food Not
 Not Good

$S_1 = \begin{matrix} 1 & 0 & 0 \end{matrix}$	more interpretable
$S_2 = \begin{matrix} 0 & 1 & 1 \end{matrix}$	

N-gram Range	Meaning
(1, 1)	Unigrams only – single words like ["I", "love", "NLP"]
(2, 1)	✗ Invalid – the min value (2) > max (1); will raise an error
(2, 2)	Bigrams only – pairs of consecutive words like ["I love", "love NLP"]
(2, 3)	Bigrams and Trigrams – combinations of 2 and 3 words like: ["I love", "love NLP", "I love NLP"]
(1, 3)	Unigrams, Bigrams, and Trigrams – all n-grams from 1 to 3 words: ["I", "love", "NLP", "I love", "love NLP", "I love NLP"]

PRACTICALS

Hyperparameters = n-gram-range
 max_features
 binary = True

```
## Create the Bag OF Words model with ngram
from sklearn.feature_extraction.text import CountVectorizer
## for Binary BOW enable binary=True
cv=CountVectorizer(max_features=100,binary=True,ngram_range=(2,3))
X=cv.fit_transform(corpus).toarray()
```

1, 1 = Unigrams only
 1, 2 = Uni, bi
 1, 3 = Uni, bi, trigrams

TF-IDF

Term Frequency - Inverse Document Frequency

→ is a technique to convert text data into vectors

Bow → word frequency

TF-IDF → weighted frequency

Bow

TF-IDF

Focus	Frequency only	Frequency and uniqueness
Importance Handling	Treats all words equally	Penalizes common words, boosts rare but important words
Stop Words Effect	Common words dominate vector	Common words are down-weighted even if not removed
Captures Meaning?	No context, just counts	No deep meaning, but better than raw counts
Vector Output	Sparse vectors with word counts	Sparse vectors with weighted scores
Example Tool	CountVectorizer in sklearn	TfidfVectorizer in sklearn extraction

INTUITION

$$\text{Term Frequency} = \frac{\text{Frequency of words in sentence}}{\text{No. of words in sentence}}$$

$$\text{IDF} = \log_e \left(\frac{\text{No. of sentences}}{\text{No. of sentences cont. the word}} \right)$$

$$\begin{aligned} S_1 &= \text{good boy} \\ S_2 &= \text{good girl} \\ S_3 &= \text{boy good girl} \end{aligned}$$

TF.	S ₁	S ₂	S ₃
good	1/2	1/2	1/3
boy	1/2	0	1/3
girl	0	1/2	1/3

Two matrix products

- hadamard (element wise)
- Dot

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

FINAL $TF - IDF$ \longrightarrow $TF \times IDF$

	good	boy	girl
S_1	$1/2 \times 0$	$1/2 \times \log_e 3/2$	$1/3 \times \log_e 3/2$
S_2	$1/2 \times 0$	0	$1/3 \times \log_e 3/2$
S_3	0	$1/2 \times \log_e 3/2$	$1/3 \times \log_e 3/2$

VV Imp In $TF - IDF$, the word importance is captured which means the word which occurs in every sentence is given less value, and more value to uniqueness

🧠 How TF-IDF Still Helps:

Even after removing stopwords, your data will still have many words that:

- Appear frequently in one document (important)
- Appear across many documents (less informative)

TF-IDF gives more weight to the first kind.

Advantages

word importance is captured

Disadvantages

- Sparsity exists
- No Semantic Context

TFIDF is better than BOW

PRACTICALS

```
from sklearn.feature_extraction.text import TfidfVectorizer
tfidf=TfidfVectorizer(max_features=100)
X=tfidf.fit_transform(corpus).toarray()
```

WITH N GRAMS

N-Grams

```
tfidf=TfidfVectorizer(max_features=100,ngram_range=(2,2))
X=tfidf.fit_transform(corpus).toarray()
```

Dense

WORD

EMBEDDINGS

Semantic
Aware

word Embeddings are dense vector representations of words, similar words have similar vectors

💡 In Short:

Technique	Word Meaning Captured?	Vector Type	Example
BoW	✗ No	Sparse (1s/0s)	[0,1,0,0,1,0...]
TF-IDF	✗ No	Sparse (weighted)	[0.3, 0.0, 0.7, ...]
Embeddings	✓ Yes	Dense (float)	[0.12, -0.33, 0.88...]

💡 What's special about Word Embeddings?

They capture semantic meaning, so:

- "king" - "man" + "woman" ≈ "queen"
- "happy" and "joyful" are close in vector space
- "apple" is close to "fruit", not "car"

word meaning
is captured

V-V. Important

✓ Can you say BoW, TF-IDF, and OHE are types of word embeddings?

✗ Technically, no. They are not true word embeddings.

But...

✓ They are text vectorization techniques — they convert words into vectors.

Imp

💡 What is true word embedding?

✓ A word embedding is a dense, learned, and semantic-aware representation of a word — typically trained using large text corpora.

⌚ Let's break it down:

Technique	Converts Text to Vectors?	Captures Meaning (Semantics)?	True Embedding?
One-Hot Encoding (OHE)	✓ Yes	✗ No (all words equally distant)	✗ No
Bag of Words (BoW)	✓ Yes	✗ No	✗ No
TF-IDF	✓ Yes	✗ No (weights importance, not meaning)	✗ No
Word2Vec, GloVe, BERT	✓ Yes	✓ Yes (semantic relationships)	✓ Yes

Frequency

Semantic Awareness; means model understands the relationship b/w the words based on

- their meaning
- not their frequency / count

Google word2vec
↳ 3 billion words

WORD2VEC

- Dense vectors
- Semantic Aware

is a deep learning model — which helps to represent each word as a dense vector, where words with similar meaning have similar vectors.

Example:

After training Word2Vec, you might get:

- `vector("king") ≈ [0.25, 0.51, -0.33, ...]`
- `vector("queen") ≈ [0.26, 0.49, -0.30, ...]`

They're numerically close because their meanings and usage are similar.

Why Word2Vec is powerful

Unlike BoW or TF-IDF, Word2Vec:

- ✓ Understands context
- ✓ Captures meaning & similarity
- ✓ Produces dense, low-dimensional vectors

Types

CBOW

Skip Gram

How Word2Vec Works (Two Models)

Model	Idea
CBOW (Continuous Bag of Words)	Predict the current word from surrounding context
Skip-gram	Predict surrounding words from the current word
Both models are trained on a large corpus using neural networks.	

Imp In word2vec, each word is represented by vector of 300 Numbers

These numbers capture hidden features like:

- Gender
- Royalty
- Profession
- Power
- And many others (not named, just learned from usage)

V-imp

→ Vector Representation

Let's say; Vocab = Man, Woman, King, Queen

Feature Represent.	Gender	Man	Woman	King	Queen
	Royal	-0.01	+0.02	0.95	+0.96

A very popular example;

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

Removes Man
Left with Royal

Royal + Woman ≈ Queen

MATHEMATICAL INTUITION

King [0.95, 0.96]

Man [0.95, 0.98]

Queen [-0.96, -0.95]

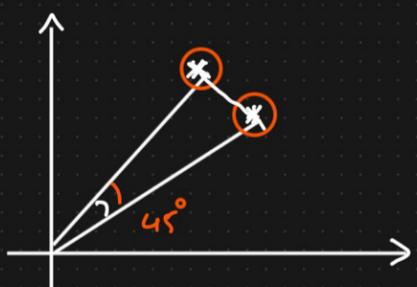
woman [-0.98, -0.95]

Imp To understand distance b/w 2 vectors, we need **Cosine Similarity**

Cosine Similarity

Distance = 1 - Cosine Similarity (Angle b/w 2 vectors)

Cosine Similarity



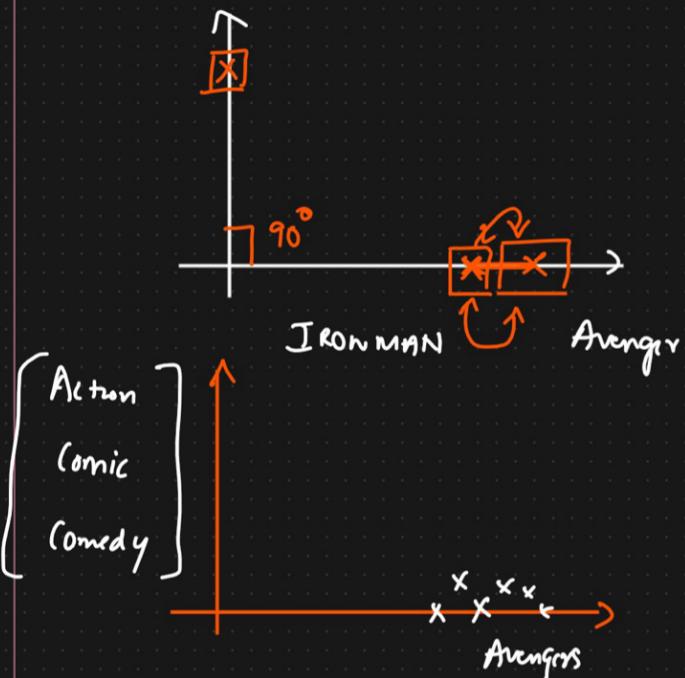
$$\text{Distance} = 1 - \text{Cosine Similarity}$$

$$\text{Cosine-Sim} = \cos 45^\circ = \frac{1}{\sqrt{2}} = 0.7071$$

$$\text{Distance} = 1 - 0.7071$$

$$\hookrightarrow = 0.29 \rightarrow \text{Almost Similarity}$$

$$|1-1|=0$$



$$\begin{aligned} \text{Distance} &= 1 - 0 \\ &= 1 \rightarrow \text{Different} \end{aligned}$$

$$\begin{aligned} \text{Distance} &= 1 - \cos 90^\circ \\ &= 1 - 1 \\ &= 0/1. \end{aligned}$$

Before getting into CBOW intuition, lets clear our heads with → DC concepts
Activation funcs
Optimizers

manual ML VS DL → black box

🧠 Basic Definition:

Concept	Machine Learning (ML)	Deep Learning (DL)
What is it?	A branch of AI that uses algorithms to learn from data	A subset of ML that uses neural networks
Inspired by	Statistics and math	Human brain (neurons and layers)
Feature Handling	Manual feature engineering often required	Learns features automatically

⌚ Key Differences:

Aspect	Machine Learning (ML)	Deep Learning (DL)
Data Need	Works well with small/medium data	Needs large data to perform well
Model Type	Decision Trees, SVM, KNN, etc.	Neural Networks (ANN, CNN, RNN)
Feature Engineering	✓ Manual (you design features)	✗ Not needed (learns features automatically)
Training Time	Fast	Slow (requires GPUs)
Hardware Needed	CPU is usually enough	GPU/TPU preferred for speed
Interpretability	Easier to understand	Hard to interpret (black box)

📦 Examples:

Problem	Use ML?	Use DL?
House price prediction	✓ ML (Linear Reg)	✗ DL not needed
Spam detection (emails)	✓ ML (Naive Bayes)	✓ DL (LSTM for text)
Image classification	✗ ML fails	✓ DL (CNN shines here)
Chatbot/Voice assistant	✗ ML limited	✓ DL (transformers, RNNs)

- "ML is like teaching a child with rules ("If weather = rain → carry umbrella")"
- "DL is like letting the child observe and learn everything by themselves — they just figure it out!"

Types of DL :

Recurrent Neural Net.

💻 Summary Table

Feature	ANN	Word2vec	CNN	RNN
Full Form	Artificial Neural Network		Convolutional Neural Network	Recurrent Neural Network
Best for	Tabular data		Images, videos	Time series, text, audio
Has memory?	✗ No		✗ No	✓ Yes
Sequence-aware?	✗ No		✗ No	✓ Yes
Key component	Dense (fully connected) layers		Convolutional + pooling layers	Loops & hidden states
Data requirement	Medium		High	Very high (needs sequences)

📌 Example Use Cases:

Task	ANN	CNN	RNN
House price prediction	✓	✗	✗
Image classification	✗	✓	✗
Stock price prediction	✗	✗	✓

ANN

- Basic form of Neural Network
- No Memory
- Tabular Data, simple classification problems

Basic Word2vec

3. Loss Function — What the model tries to minimize

? Why do we need a loss?

It tells the model:

"Hey, your prediction is off by this much — fix it!"

📦 In Word2Vec CBOW:

- The loss is high if the predicted word is far from the actual word
- It gets lower as the model learns better word associations

✓ Common Loss in Word2Vec:

- Cross-Entropy Loss (when predicting words)
- Negative Sampling Loss (when using efficiency tricks)

4. Optimizers — How the model learns better over time

? What do optimizers do?

They update the weights to reduce the loss — like giving the model feedback.

💡 Most common optimizer:

- SGD (Stochastic Gradient Descent) — basic optimizer used in CBOW
- Others: Adam, RMSProp (used in advanced models)

⚙️ Example in CBOW:

- It sees: ["the", "on", "the"]
- Predicts: "mat"
- Compares prediction with actual word
- Computes loss
- Optimizer updates weights to get better next time

⚡ What is an Optimizer?

An optimizer = smart version of gradient descent

It tells the model:

"Here's how to update your weights using gradient descent (and some extra tricks) to learn faster and better."

⌚ So:

- Gradient Descent = The basic idea
- Move weights opposite to the gradient of the loss.
- Optimizer = An algorithm that uses gradient descent (and often improves it)

Q1 Difference b/w RNN and Transformers ?

Feature	RNN	Transformer
Architecture	Sequential	Attention-based
Processing	Step-by-step	All at once (parallel)
Speed	Slow	Fast
Long Context	Struggles	Handles well
Use Cases	Time-series, speech	NLP, translation, chatbots
Examples	LSTM, GRU	BERT, GPT, T5

- ✓ RNN = good for short, sequential data
 ⚡ Transformer = best for long, complex text tasks

Q2 Difference b/w forward & back propagation ?

🧠 1. Forward Propagation (Prediction Phase)

This is when the input data flows through the network to produce an output (prediction).

Steps:

- Input data is passed into the model.
- Each layer applies weights and bias, then an activation function.
- The final output (prediction) is produced at the last layer.

📦 Think of it like a factory line where raw material (input) goes through various stations (layers) and gives a final product (output).

🤝 Relationship:

- Forward Propagation: Calculates predictions.
- Backward Propagation: Improves model by reducing error.

⌚ 2. Backward Propagation (Learning Phase)

This is when the model learns from the error by adjusting weights to reduce the loss.

Steps:

- Calculate loss (difference between prediction and actual value).
- Compute gradients of the loss w.r.t. each weight using the chain rule (calculus).
- Update weights using an optimizer like Gradient Descent.

⚡ Think of it like a feedback loop where you see the mistake and adjust the system to perform better next time.

One word
prediction

WORD 2 VEC

CBOW

INTUITION

softmax
window size

What CBOW Does

Given the context (surrounding words), predict the center (target) word.

word2vec can be pretrained or it can be built from scratch

Example; Let's assume window size = 2; this means take 2 words before & 2 words after

Left

Right

window_size = 2

2 words left + 2 words right = 4 context words

STEPS IN CBOW

Step-by-Step (with Intuition)

Let's say you have this sentence:

"The cat sits on the mat"

You choose a window size = 2, and want to predict the center word = "sits".

1. Context Selection

Context words (2 before, 2 after):

- ["The", "cat", "on", "the"]

Target word:

- "sits"

2. Convert Context Words → One-Hot Vectors

Each word in your vocabulary becomes a one-hot vector (e.g., a vector with one 1 and the rest 0's).

So, for example:

- "cat" = [0, 0, 1, 0, ..., 0]
- You do this for all 4 context words

Now you have 4 one-hot vectors.

3. Projection Layer (Hidden Layer)

Each one-hot vector is multiplied by a shared weight matrix to get a dense word embedding.

If embeddings are of size 100, each context word becomes a 100-d vector.

Now:

- You have 4 embeddings of size 100
- These are averaged to get one vector (size 100)

This is where averaging happens — not on output, but on input embeddings.

4. Prediction Layer (Output Layer)

This average vector is passed through another weight matrix to get a score for every word in the vocabulary (say, 10,000 words).

Then a softmax is applied to turn scores into probabilities.

The word with the highest probability should be:

"sits" (your target word)

5. Training (Backward Propagation)

If the predicted word is wrong:

- The model calculates error
- Uses backpropagation to update the weights (embedding matrix + output layer)

Embedding Matrix
= Vocab size x embed size

✓ Yes — you choose the embedding size manually.
It's a hyperparameter that you set before training the model.

Simple Summary:

Step	Action
1 Window Size	Choose context words around the target
2	Convert context words to one-hot vectors
3	Map to embeddings using a shared weight matrix
4	Average embeddings into one vector
5	Feed to output layer → predict target word
6	Update weights via backpropagation

1. window size (context words)
(calculate vocab)
2. One (vectorization of words)
w.r.t vocab
3. Calculate Embed Matrix
(dense vectorization) (100d v)
→ Vocab X
4. Average - Softmax (probability)
5. Training - Back propagate

Multi word
prediction

Skip Gram Intuition

softmax

🧠 What Skip-gram Does

Given the center (target) word, predict the surrounding context words.
The reverse of CBOW!

⚠ Step-by-Step (with Intuition)

Let's say you have this sentence:

→ "The cat sits on the mat"

You choose a window size = 2, and pick the target word = "sits"

✓ 1. Context Generation

With `window size = 2`, your model tries to predict the 2 words before and 2 after the target word.

So:

```
python
Target Word: "sits"
Context Words: ["The", "cat", "on", "the"]
```

Unlike CBOW, where we used multiple inputs to predict one output,
in Skip-gram we use one input (target word) to predict multiple outputs (each context word).

✓ 2. Convert Target Word → One-Hot Vector

You convert "sits" to a one-hot vector based on the vocabulary.

E.g., if vocab size = 10,000, "sits" = [0, 0, 1, 0, ..., 0]

✓ 3. Projection Layer (Hidden Layer)

The one-hot vector is multiplied by a weight matrix (embedding matrix)

→ Produces an embedding vector (say, size 100)

This embedding now represents the meaning of "sits"

⌚ Simple Summary:

Step	Action
1	Choose a target word
2	Identify context words within window
3	Convert target word to one-hot vector
4	Map to embedding using weight matrix
5	Use embedding to predict each context word (one at a time)
6	Backpropagate loss and update weights

✓ 4. Output Layer — Predict Context Words

Now, we use this embedding vector to predict each context word (one by one):

- Predict: "The"
- Predict: "cat"
- Predict: "on"
- Predict: "the"

Each prediction is a softmax over the whole vocabulary (10,000 words).

So you get 4 softmax outputs from one input.

→ $\text{vocab} \times \frac{\text{embedding size}}{\text{hyperparameter}}$

STEPS

1. window size (ie., 4 words to be predicted)
2. Target word to ONE
3. Embedding matrix

✓ 5. Training (Backward Propagation)

For each context word, the model:

- Calculates prediction error
- Uses backpropagation to adjust:
 - Embedding of the target word
 - Output layer weights for each context word

Q: when to apply CBOW vs. SkipGram

CBOW = Small Data Sets

SkipGram = Large Data Sets

Google word2vec
• 3 billion words
• Feature representation
of 300D vector

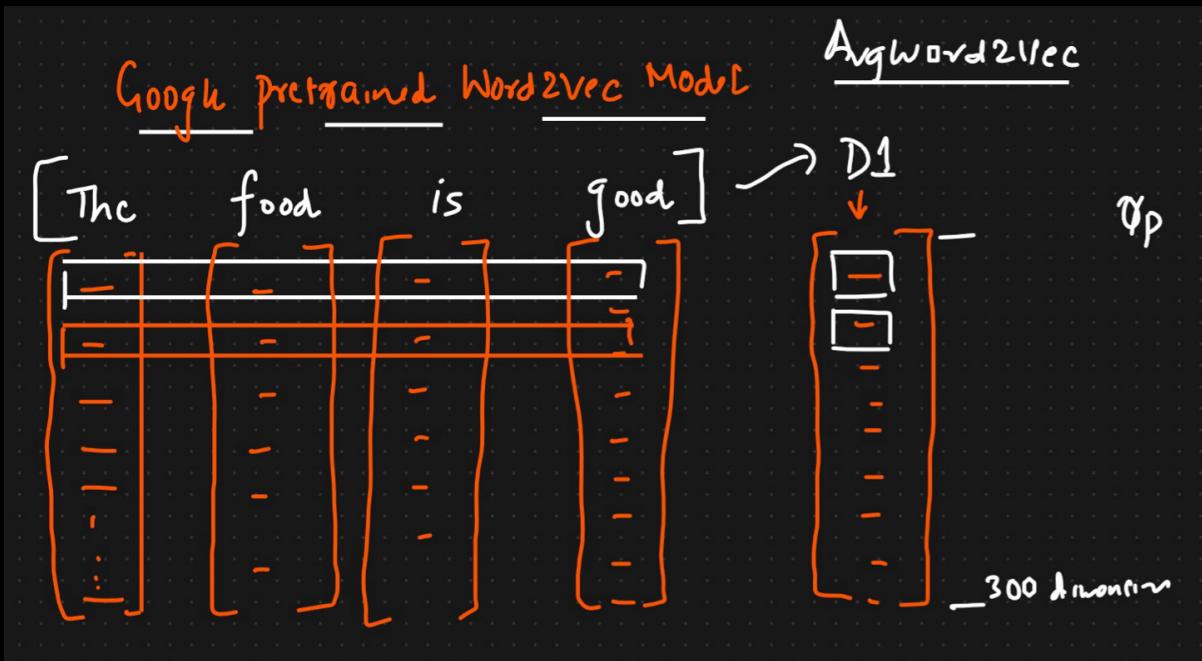
BENEFITS OF WORD2VEC :

- NO sparsity → Dense Matrix
- NO OOV (out of vocab) All words covered 3B
- fast, scalable & pretrained.
- Semantic Aware (most import)

Avg WORD2VEC

Given a sentence, it calculates the average of Word2Vec vectors of all the words in that sentence.

This gives you one single vector representing the whole sentence or document.



3. Use It in ML Models

This single vector can now be used as input to:

- A classifier (Logistic Regression, SVM, etc.)
- A clustering algorithm
- Similarity search, etc.

Q, what is ^{Toolbox} GENSIM & ^{Stanford} GLOVE

Gensim

A popular Python library for working with word embeddings and topic modeling.

Used to:

- Load Word2Vec, GloVe, and FastText models
- Train your own Word2Vec or Doc2Vec
- Find similar words, sentence vectors, etc.

* Think of Gensim as the toolbox to use or train embedding models.

GloVe (Global Vectors for Word Representation)

A word embedding technique developed by Stanford.

Key idea:

- Learns word meanings by counting how often words co-occur in large text corpora
- Combines global statistics (like matrix factorization) with local context (like Word2Vec)

* Unlike Word2Vec (predictive), GloVe is more count-based.

BEST PRACTICES

To avoid Data leakage

STEP I; Data preprocessing

- Remove signs (using RE) `re.sub('[a-zA-Z]', '')`
- Lowercase
- Tokenisation `word_tokenize`
- Lemmatization / Stemming (`PorterStemmer`)
`WordNet`
- Append Corpus

STEP II; Train Test split - To avoid Data leakage

`y = get_dummies (labels)`

`from sklearn.feature_extraction import`

`CountVectorizer, TfidfVectorizer`

`hyperparameters, max_features, ngram_range`

`x_train = CountVectorizer().fit_transform(x_train)`

`x_test = cv_transform(x_test).toarray`

STEP III ; Predictions from ML

`model = sklearn.naive_bayes import
MultinomialNB`

We can also use Random Forest classifier

Naive Bayes = Sparse Vectors

Random Forest = Dense Vectors.