

## Large Language Model

copy list with random pointer

King - Man + Women

King Minus Man plus Women is Queen

Statements like this needs semantic meaning to be preserved

This can be done

i) Tokenize the given input into subwords such as King, woman, queen, plus, minus

ii) Through positional encoding order of the words are preserved

iii) The semantic meaning can be inferred from a cluster of similar sentences.

iv) This semantic meaning must be retrieved through methods like retrieval augmented generation (RAG)

## Chapter 1: Understanding the LLMs

What is lang AI?

What are AI system

AI is the process of developing intelligent computer program

so that it can understand (Human) AI and replicate it. The AI systems are pre-trained with a large amount of data

Language AI is the subfield of AI that focuses on developing technologies capable of understanding processing and generating human language

Language AI models are pre-trained on large corpus of data therefore they are called as Large language models

## Heaps and Priority Queue

The non-transformer models are simple LM which can process and structure the input using simple vector representation or neural networks

e.g. Bag-of-words, word2vec, attention etc

They do not perform well if input sentence is large. Transformer based models rely on transformer framework in order to process the input. The attention mechanism in transformers makes it suitable to process longer input sequence

ex - BERT, GPT, T5, switch

## Representing language as a bag of words

A sentence comprises of several words. Each word is again comprised of characters. Each word is separated from another by a whitespace.

vocabulary

ip → tokenization → Vector Embedding

Tokenization: A large ip is broken into smaller units called tokens. Every time we get an ip it must be broken into tokens and these are added into data dictionary called the vocabulary. This vocabulary is the reference for any language models

Only unique tokens are added to the vocabulary (BoW)

"Bag of words" → 3 tokens (word-level tokenization)

ip: That is a white dog

ip: My cat is cute

that [ ] is [ ] a [ ] white [ ] dog [ ] my [ ] cat [ ]  
1 2 1 2 10 1 1 1

7 vocabulary

NOTE:

All though repeated words are not added to vocabulary but the vocabulary keeps track of how many time a word is used in input

### Vector Representation of Bag-of-words

Each token must be converted to a numerical representation this numerical representation allow the mapping of the i/p to the vocabulary

The Embedding technique used in bag of words is called one-Hot Encoding where there is a 'zero' in the vocabulary if the word is not present in the input

e.g. i/p my cat is cute

that	is	a	cute	dog	my	cat
↓						
0	1	0	1	0	1	0

one-Hot encoding is a fixed length encoding technique

0	1	0
0	0	0
1	0	1

0's no meaningful content  
1's content

sparse matrix

One-Hot Encoding contains a large no of 0's which is computationally expensive. The biggest drawback is that it doesn't contain any semantic

BoW cannot prense the semantic meaning of word

### Advantages and Disadvantages of BoG

Advantages - 1) Easy to use  
2) Used as base for newer models

Disadvantages 1) If there are no whitespace the BoW can't generate

### Dense vectors

Instead of storing embedding as sparse matrix we try to store the embeddings as a dense vector the most popular dense vector technique are -

- 1) word2vec
- 2) glove

### Dense vector Representation

Dense vectors are embeddings of tokens which try to preserve the semantic meanings of words

Dense vectors are created using Neural Network ANN is adjust with some features and some weights

In Dense vector representation there are a large no of classes each word in the input has a similarity score for every class

	car	puppy	houses	apple	baby	
animal	0.1	-0.2	-0.5	-0.7	0.1	
newborn	-0.1	0.1	-0.2	-0.1	0.0	Number
human	-0.1	0.3	0.1	0.9	0.8	of
plural	0.4	0.4	0.4	0.5	-0.3	dimensions
fruit	-0.1	-0.5	-0.3	0.8	-0.1	

Embedding of words that are similar will be close to each other in dimensional space

o o	
cats dog	
puppy building	
apple house	
banana adult	
dense vector	baby

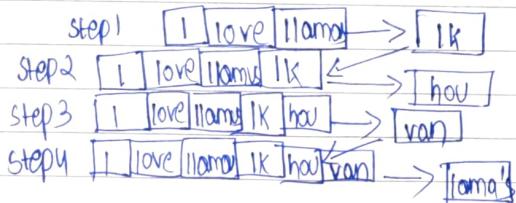
### Attention (RNN with Attention)

Although dense vectors can preserve the semantic meaning of the input sequence. It doesn't perform efficient when the context window is large.

To overcome this advantage RNN (Recurrent Neural Network) were proposed

RNN are Neural Networks which can model sequences as an input, this helps preserve the contextual meaning as well as the input sequence

### RNN - To capture semantic meaning

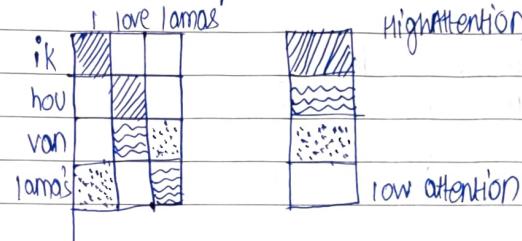


To remember a RNN in the  $i^{\text{th}}$  step takes the output of previous  $1 \text{ to } i-1$  steps along with the input.

The disadvantage of a recurrent neural network is that a large no of output words has to be passed everytime which is computationally very expensive. To reduce this computational complexity the attention mechanism was produced.

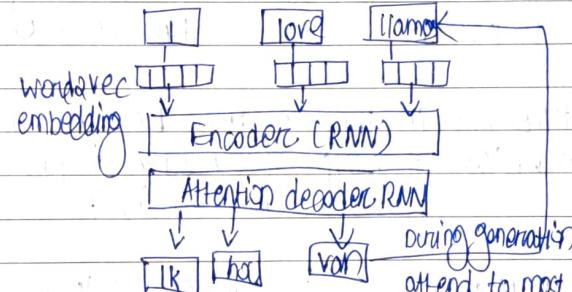
### Attention mechanism

Attention allows the model to focus on the specific parts of the input that are relevant to each other.



From the attention matrix the model can predict which is the most relevant part of the output. Therefore it can preserve both meaning and sequence.

### RNN + Attention Mechanism



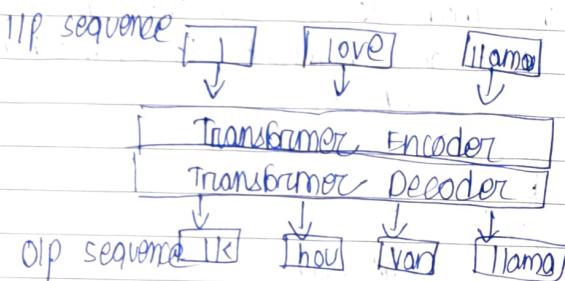
Using RNN + Attention we can maintain sequence for large context. This helps us in translation related work. Therefore these models are known as representation model.

## Transformers

A Transformer is a neural network that consists of 2 layers

- (1) Encoder layer
- (2) Decoder layer

A transformer model can be trained in parallel unlike RNN model.



## Encoders & Decoders

### Encoding

The encoder takes a sequence of tokens and transforms it into a series of contextualized embeddings.

These embeddings capture the semantic meanings as well as the sequences.

There are 2 layers -  
- self attention  
- feed forward  
In the decoder only model there is no encoder the i/p is directly given to the decoder.

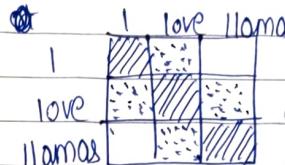
### Decoding

Decoders can be used in 2 scenarios -  
- Decoder only  
- Encoder & Decoder

In the encoder-decoder model the o/p of encoder is provided to decoder as input

## Self Attention

Self Attention is at the heart of transformed model the self attention is different from the previous ~~attention~~ attention.

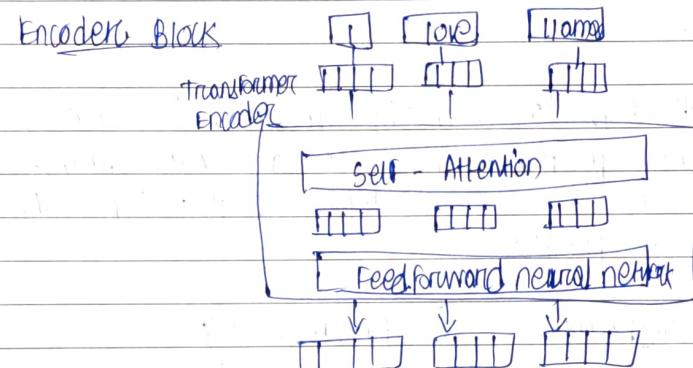


In SA the attention b/w every token in the i/p is calculated. No o/p is involved in case of self Attention

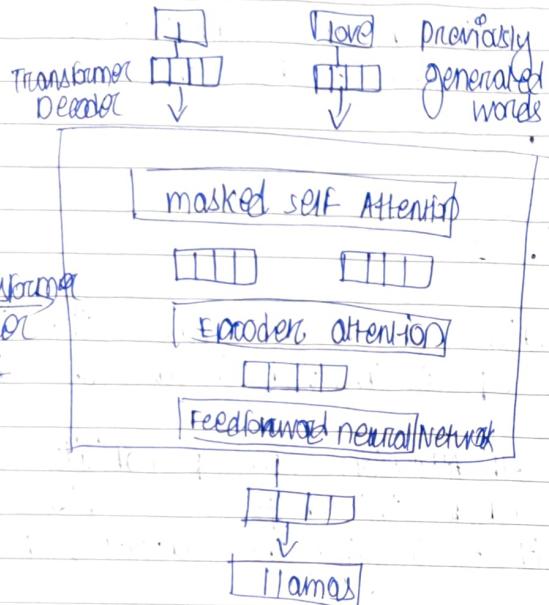
Because of the self Attention low probability token in inputs are discarded, this helps us to provide input to the model without changing the contextual meaning of the input.

This SA mechanism led to the creation of generative model.

### Encoder Block



## Decoder Layer



The decoder has an extra layer for encoder attention to pay self-attention to the encoder output masked.

The self-Attention layer in the decoder is masked to that future positions are not available during decoding otherwise the model will just copy the input instead of predicting

## Types of Transformers

- (I) Encoder only
- (II) Decoder only
- (III) Encoder - Decoder

DOMS

## Bidirectional Encoder Representations from Transformers

### Generative or completion LLM

User query  
(prompt)

Tell me something about llamas

Output  
(completion)

generative LLM  
Task: complete the ip

} Decoder only mode

Context window: represents the max no of tokens that a model can process.

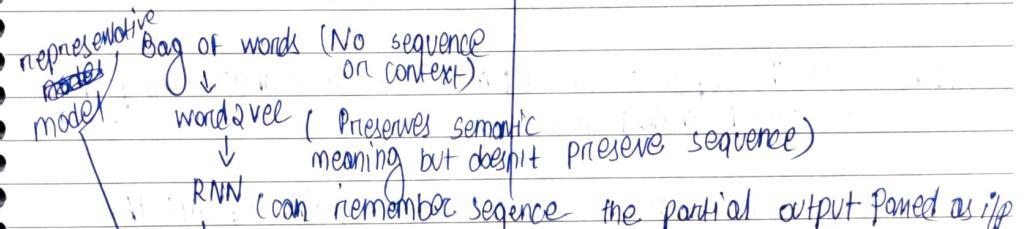
A large context window allows entire document to be passed through the LLM.

### Pretrained LLM

Before a LLM is released in the public it is trained on massive purpose of text data these are pretrained LLM.

### Fine-tuned LLM

A pretrained LLM that is available in the public can be finetuned on certain parameters to perform specific tasks



RNN+ Attention  $\rightarrow$  Transformers capable of representation as well as generation  
no need to pass old like RNN if can remember sequences through attention

DOMS

## Tokens and Embeddings

Input Have the bands who preceded

↓ Tokenisation

Breakdown the text into smaller pieces  
(words or parts of words)

Tokens Have the bands who preceded

↓ Embedding

        Turn tokens into numeric representations capturing their meaning

The LLM although they process and generate text can't understand characters in the text  
Therefore characters in the input to LLM must be converted to numeric value.

Converting the entire i/p sequence to numeric values doesn't make the conversion efficient

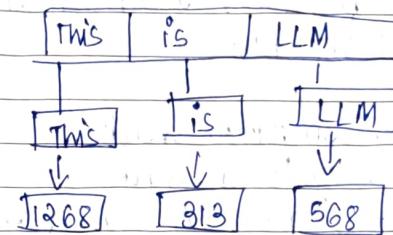
A better way to convert text into numerical values is to break the i/p into smaller units and create numerical representation of those units.

NOTE - All the open source LLMs are available in a python class called transformers

The input\_id variable stores the token ids of all the generated tokens.

## Input and Outputs of LLMs

when the i/p sequence given to the LLM it doesn't process the character tokens rather it processes numerical representations of those tokens.



In addition to being used to process the i/p text into a language model, tokenizers are used on the output of the language model to turn the resulting tokens into the o/p word for token associated with it.

Input Have the bands who preceded

Tokenizer: encode

Tokens Have the bands who preceded

Token ID 1075 1278 1278 3163 1050

↓ Language model

output TokenID: 1294

Tokenizer: decode ↓  
or7

when a model produces o/p it doesn't produce a sequence of text rather it produces a sequence of numerical values. These numerical values

must be converted to characters: This is done by mapping the token ID to its corresponding text.

### Properties of Tokens

i) The first token is ID 1 (`[S>]`) a special token indicating the beginning of the text

ii) Some tokens are complete words

iii) Some tokens are parts of words

iv) Punctuation characters are their own token

v) Space characters do not have its token

### What are Tokens?

The smallest units of text that a model processes is known as a token input sequence. It is broken into these small units.

Types of Tokenization - word

subword

character

Byte

### Tokenization methods

Word-level tokenization - The process of splitting input sequence of text into individual words

Advantages - simplifies the input by dividing them into words. Preserves the contextual meaning of each token. Works well with languages where DOMS

clear word boundaries is defined such as English

Disadvantages - similar words with slightly different meanings are stored in the vocabulary as diff tokens. This increases size of vocabulary.  
ex - apology, apologies, apologetic

ii) Word level tokenization has no mechanism to detect OOV

OOV (out of vocabulary) :- when the LLM model is pretrained the vocabulary is created. When the tokenizer finds a token which is not in the vocabulary it is called OOV token

Character-level tokenization - when input sequence is broken into individual characters including letter, numbers, digits, space

Advantages - → No OOV words issues

it works well with creative and noisy texts

it is more fine-grained (useful for spelling correction)

Disadvantages - if loses context when the input sequence is too long loss of semantic and contextual meaning because characters can't hold semantic meaning like words

Byte-level tokenization - breaks the input sequence into individual bytes that are used to represent unicode characters

This allows for universal representation including texts, numbers and emojis.

Advantages

- no unknown tokens DOMS

Size of vocabulary is very less - 1 byte = 256 tokens used, bcz of social media, coding, multilingual input

Disadvantages - the semantic and contextual meaning gets lost because of unicode representation

word-level	preserves semantic meaning	large vocab OOV
character level	more fine grained	it doesn't preserve semantic
Byte-level	universal representation of language	

Subword - level tokenization can produce words, character tokens, it is middle ground of character and words

it is useful for designing generative models like GPT

Advantages - it handles OOV words, reduces the vocabulary size

Improves generalization across word forms

Disadvantages - it may produce meaningless subwords

Requires pretraining to build this vocabulary

Complex to implement than word-level or character level

word tokens

I H a v e | t h e | n | b e r i d s | w h o | s p r e e d e d | - | -

subword tokens

I H a v e | t h e | n | b e r | d | s | w h o | s p r e | e | d | - | -

character tokens

H | a | v | e | | t | h | e | | n | b | e | r | d | s | | w | h | o | | s | p | r | e | e | d | - | -

byte tokens

H	a	v	e		space		t	h	e		space	-			
0	0	0	0	0	.	.	0	0	0	0	.	1	1	1	1
1	1	1	1	0	.	.	1	1	1	0	.	1	0	0	0
0	1	1	1	0	.	.	1	1	1	0	.	1	0	0	1
0	0	1	0	0	.	.	1	0	0	0	.	1	1	0	1
0	0	0	0	0	.	.	0	1	0	0	.	0	1	1	0
1	0	0	0	0	.	.	0	1	0	0	.	0	1	1	0
0	0	1	1	0	.	.	1	0	1	0	.	0	1	1	1
0	0	1	0	0	.	.	0	0	0	0	.	0	1	0	0
0	1	0	0	0	.	.	0	0	0	1	.	0	1	0	1

Problems solved by Tokenization

words vary in form (run, running)

Rare or new words (photosynthesis)

MisPELLING OR typos

Multilingual text

Model input size limits

How tokenization helps

Breaks into subwords  
captures shared meaning

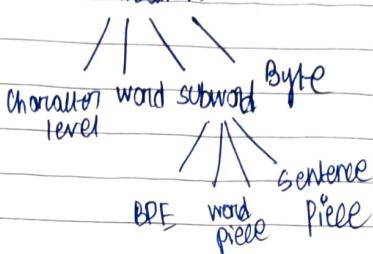
splits into known  
chunks  $\rightarrow$  avoids  
unknown tokens

subwords unit still makes  
sense (transformer)

uses shared vocabulary  
across languages

controls token count - fits  
within model window

Tokenization



Byte pair encoding is a subword level tokenization algorithm  
It begins with character level merges to form subwords

Required - corpus of words, No of merges

Output - Tokenized corpus, subword vocabulary

For each word in the corpus to split the word into  
characters  
preprocessing Append EOW symbol end for.  
Tokenization  
Initialize the vocabulary with unique characters & EOW  
for merge step = K to N

do  
count Frequency of all adjacent symbol pairs  
Identify the most frequent pair (a,b)  
Merge (a,b)  $\rightarrow$  ab and add to Frequency  
End for.

- I aa abc abc
- II 'a', 'a', 'b', 'b', 'c', 'c', 'a', 'b', 'c'
- III (a,b) - 2  
(b,c) - 2  
(a,b)  $\rightarrow$  ab

vocabulary  
'a', 'b', 'c',  
'\_'

Initial corpus: `["aa", "abc", "abc"]`

`"aa"` → `["a", "a", ""]`

`"abc"` → `["a", "bc", "c", ""]`

`"abc"` → `["a", "b", "c", ""]`

Step 1: Iterative Merging

Most frequent pair: `("a", "b")`

merge into: `"ab"`

corpus update: `["a", "b", ""]`

`["ab", "c", ""]`

`["ab", "b", "c", ""]`

`["ab", "bc", ""]`

Most frequent pair: `("ab", "c")`

merge into: `"abc"`

corpus update: `["a", "a", ""]`

`["abc", ""]`

`["abc", ""]`

Vocabulary: `["a", "b", "c", ""]`

Iteration 2:

Most frequent pair: `("a", "a")`

merge into: `"a"`

corpus update: `["aa", ""]`

`["abc", ""]`

`["abc", ""]`

Vocabulary: `["a", "b", "c", ""]`

Final Token: `["aa", ""]`

`["abc", ""]`

`["abc", ""]`

Final vocabulary: `["a", "b", "c", "-", "ab", "abc", "aa"]`

Byte pair encoding splits the input into words and splits those words into characters.

The unique characters are added to the vocabulary. From the input corpus it finds out the pair of most occurring characters (a-byte pair) and merges them.

This merged pair of characters is called as a ~~character~~ subwords and in future steps may be combined with other subwords.

For ex - `a, b` merged to `ab`

`ab, c` → `abc`

These subwords are added to vocabulary in each step and the final vocabulary as well as final tokenized corpus is returned.

Using Byte pair encoding find the final vocabulary and the final tokenized corpus

`hug pug pun bun hugis`  
`{hug, 10}, {pug, 15}, {pun, 12}, {bun, 14}, {hugis, 5}`

vocab

`hug, p, n, bys`

step 1: `hug -> ug`

`u, p -> 18, 16` → (`un` is now token)  
`p, u -> 17, 12` → (`ug` is now a token)

`hug pug p, un b, un hugis`  
`hug -> hug`

DOMS

Wordpiece - is a subword level tokenization method developed by google used by BERT

Like byte-pair encoding word piece also merges tokens to form the training vocabulary

Unlike byte-pair encoding which merges the most frequent pair of characters word piece follow a merging rule

$$\text{Score} = \frac{\text{frequency of pair}}{\text{frequency of first ele} \times \text{frequency of second ele}}$$

{un} {able}

{n} {u} {able}

By dividing the frequency of the pair by the product of the frequencies of each of its parts, the algorithm prioritizes merging the pairs where individual parts are less frequent in vocabulary.

it won't necessarily merge ("un", "able") even if the pairs occurs very frequently in the vocabulary because the two pairs un and able will likely each appear in a lot of other words and have a high freq

Wordpiece prioritizes merging the subwords that are less likely to appear as individual tokens throughout the vocabulary

Once the vocabulary is created using the merge rule wordpiece follows the given greedy strategy

1. Start at the beginning of the word
2. Find the longest substring that matches a token in the vocabulary
3. Add that token to the output
4. Move forward in the word and repeat until the entire word is tokenized

NOTE - This algorithm is for a single word and must be repeated until all words in the input sequence are tokenized

Wordpiece appends the symbol # before any subword to specify that the particular subword is part of another token.

Ex - Vocabulary = { "un", "#off", "#able", "#offable" }  
word = "unaffable"  
match = "un"  
Remaining = "affable"

M = "# off"

R = "able"

i = "# offable"

Final = { "un", "# off", "# able" }

unknown

NOTE → Wordpiece uses the UNK symbol to represent input subwords that do not match any token in the vocabulary.

Vocabulary = { UNK, a, aa, # off, # able }

aa → aa

abc → # off abc → # off c  
match match

vocabulary = {UNK, a, aa, ab, abc, ##b, ##c}

corpus =  $\{ "aa", "abc", "abc" \}$   
match match match

## Sentence Piece

is a subword level tokenization algorithm that is language independent it can work well with languages which do not have clear boundaries such as Chinese

unlike byte pair encoding / wordpiece that create the vocabulary by merging tokens sentence piece requires a large vocabulary to start.

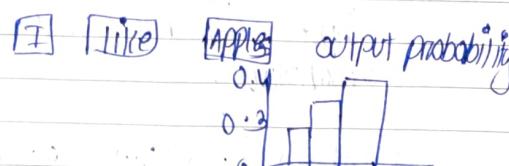
It introduces a special metacharacter to represent the spaces , the most used ex - is unigram model

## UNIGRAM MODEL

The unigram model is used in LLM models like T5, BigBird, XLNet.

INPUT TEXT  $\rightarrow$  BAG OF WORDS

I like apples  
INPUT TEXT



Since it produces language independent or universal tokens so it is called unigram

The unigram model uses viterbi algorithm to find the most probable segmentation of a token i.e the sequence of tokens that maximizes the product of their unigram probability

This probability is not the raw frequency rather it is a learned probability that the model learns based on how frequently the ~~the~~ unigram occurs in the vocabulary.

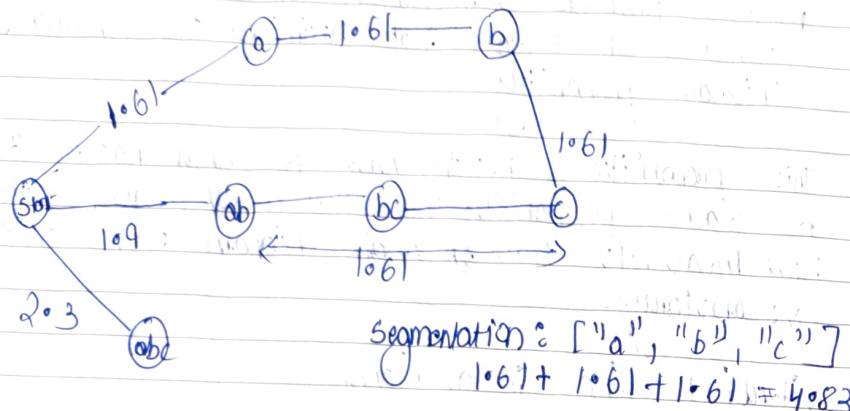
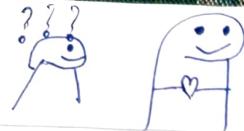
## viterbi decoding in a unigram

cost of a segment is the negative logarithm of the probability of its constituent tokens.

$$P_{abc} = P_a + P_b + P_c$$

The viterbi algorithm finds the most probable by finding segmentation which incurs the least cost.

token	probability	cost - logp
a	0.20	1.61
b	0.20	1.61
c	0.20	1.61
ab	0.15	1.90
bc	0.15	1.90
abc	0.10	2.30



Segmentation:  $\{["abc"]\}$   
cost:  $1.9 + 1.61 = 3.51$

Segmentation:  $\{["a"], ["bc"]\}$   
cost:  $1.61 + 1.9 = 3.51$

Segmentation:  $\{["abc"]\}$   
cost: 2.30

Given the vocabulary and n-gram, the most probable Segmentation is abc.

The probability of subwords present in the vocabulary are not fixed, they change over time depending on the input corpus therefore these probabilities are called as learned probabilities in the n-gram model.

The probability of a segment (abc)

$P_{abc}$  is the sum of probabilities of its constituent segments

$$P_{abc} = P_a + P_b + P_c \text{ or } P_{ab} + P_c \text{ or } P_a + P_{bc}$$

The cost of a segment  $C_{abc}$  is the negative log of probability  
 $C_{abc} = -\log P_{abc}$

These probabilities change over time based on the likelihood of a subword appearing on the training corpus

Ex corpus - aa abc abc

Token	probability
a	0.1
aa	0.3
abc	0.6

1.  $[\text{aa}] [\text{abc}] [\text{abc}]$  These are only valid
2.  $[\text{a}, \text{a}] [\text{abc}] [\text{abc}]$  segmentation using the given vocabulary.

$$\text{Segmentation 1: } 0.3 \times 0.6 \times 0.6 = 0.108$$

$$\text{Segmentation 2: } 0.1 \times 0.1 \times 0.6 \times 0.6 = 0.0036$$

Normalised probability -  $\frac{0.108}{0.1116}$

$$\underline{0.0036}$$

$$\underline{0.1116}$$

once the normalised probability is calculated unigram uses the EM (Estimation-Maximization) algorithm to update the probability. Expectation

Finding the Expected Count for each subword (E-word)

Count in Seg-1

<u>Subword</u>	<u>Count in segmentation-1</u>	<u>Expected count</u>
a	0	$0 \times 0.968 + 2 \times 0.032$
aa	1	$1 \times 0.968 + 0 \times 0.032 = 0.968$
abc	2	$2 \times 0.968 + 2 \times 0.032 = 2.000$

Maximization step

The maximization step the probability the subwords are updated on if they appear on the input corpus

<u>Subword</u>	<u>expected count</u>	<u>updated probability</u>	<u>(calculated)</u>
a	0.064	0.0211	$0.064 / 3.032$
aa	0.968	0.3193	$0.968 / 3.032$
abc	2.000	0.6596	$2.000 / 3.032$

As seen the probability of tokens 'aa' and 'abc' increased and 'a' decreased because the

a did not occur

aa occurred once increased slightly

abc occurred twice so it increased

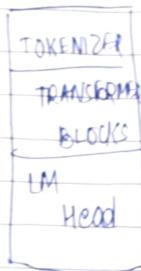
## Final Tokenization using Unigram

Given the initial learned probabilities and the subwords unigram model finds the most probable segmentation using the ~~Viterbi~~ Viterbi algorithm then using EM algo

so for the given vocabulary the best tokenization is "aa abe"

[Prev MyCopy]

Input - Write an email to Sarah



Vocabulary	
Dear	0.1.
A	0.1.
not	0.1.
This	0.1.

The i/p comes into the tokenizer there is broken individual tokens. These tokens then pass through transfer block which changes the probability of tokens in the vocabulary once the probability distribution of the vocabulary changes LM had selects the next token based on the decoding strategy

LM had selected one token from the vocabulary this is called as decoding. The easiest decoding strategy is to select that token which has highest probability this is known as greedy decoding.

Softmax Function transforms array of values to its corresponding probability distribution such that all the values remain between 0 to 1

$$\text{softmax} = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

$$Q = [0.0, 1.0, 0.0]$$

$$e^{2.0} = 7.39, e^{1.0} = 2.72, e^{0.0} = 1.0$$

$$\sum e^{x_j} = 7.39 + 2.72 + 1.0 = 11.11$$

DOMS



$$\text{softmax}(2.0) = \frac{7.39}{11.11} = 0.659 \quad \text{softmax}(1.0) = \frac{2.72}{11.11} = 0.242$$

$$\text{softmax}(0.0) = \frac{1.0}{11.11} = 0.091$$

Let us calculate attention weights for the word crossed in the i/p sentence *the chicken crossed the road*

Let us assume dimension  $d_A = 4$  therefore the scaling factor  $\sqrt{d_A} = 2$

$$Q_{\text{crossed}} = [1, 2, 1, 0]$$

key vectors for the tokens -

$$K_{\text{The}} = [1, 1, 0, 0]$$

$$K_{\text{chicken}} = [2, 1, 1, 0]$$

$$K_{\text{crossed}} = [1, 2, 1, 0]$$

$$K_{\text{The}} = [1, 1, 1, 0, 0]$$

$$K_{\text{Road}} = [1, 1, 1, 1, 0]$$

dot product  $Q_{\text{crossed}}$  by each  $K_i$

$$Q \cdot K_{\text{The}} = 1 \times 1 + 2 \times 1 + 1 \times 0 + 0 \times 0 = 3$$

$$Q \cdot K_{\text{chicken}} = 1 \times 2 + 2 \times 1 + 1 \times 1 + 0 \times 0 = 5$$

$$Q \cdot K_{\text{crossed}} = 1 \times 1 + 2 \times 2 + 1 \times 1 + 0 \times 0 = 6$$

$$Q \cdot K_{\text{Road}} = 1 \times 1 + 2 \times 1 + 1 \times 1 + 0 \times 0 = 4$$

$$\text{Scaled The} = \frac{3}{2} = 1.5$$

$$\text{Scaled chicken} = \frac{5}{2} = 2.5$$

$$\text{Scaled crossed} = \frac{6}{2} = 3.0$$

$$\text{Scaled Road} = \frac{4}{2} = 2.0$$

$$\text{softmax}(n_i) = \frac{e^{n_i}}{\sum e^{n_j}}$$

$$z = [1.5, 0.5, 3, 2]$$

$$e^z = [4.48, 1.18, 20.08, 7.38]$$

$$\sum e^z = 44.14$$

$$\text{softmax}(z) = [0.10, 0.276, 0.455, 0.167]$$

$$\text{crossed } \& \text{ me} = \frac{4.48}{44.14} = 0.10 \quad \times$$

$$\text{crossed } \& \text{ chicken} = \frac{1.18}{44.14} = 0.028$$

$$\text{crossed } \& \text{ crossed} = \frac{20.08}{44.14} = 0.46 \quad \checkmark$$

$$\text{crossed } \& \text{ road} = \frac{7.38}{44.14} = 0.17$$

Crossed gives highest attention itself followed by chicken & road gives least attention the token Gho.

We are assuming the value of v=1

### Positional Embedding

Using the FFNN and the attention mechanism the model can understand contextual relationship b/w the i/p tokens but it can't remember sequence of i/p tokens therefore it uses positional embedding to differentiate and remember the sequence of each i/p token.

most used embedding is ~~not~~ ROPE (Rotary Positional Embedding)

ROPE doesn't add positional embedding vectors.

ROPE rotates the query and key vectors for each token

" " " " " " " " using trigonometric functions. each tokens vector is rotated by an angle based on its position.

The dot product between rotated vectors encode relative position

### Decoding in ROPE

query's token  
Each token's vector is rotated based on its position

During decoding, the model computes these rotated vectors to compute attention scores.

The angle b/w vectors how far apart tokens are.

### Relative Position Decoding

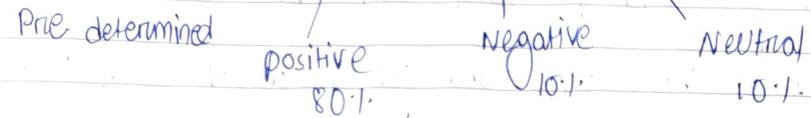


n-m = 2 tokens apart

## Text Classification Using LLMs (Ch-4)

### Text classification

Input - This is a good movie



Given an input statement on text, the process of classifying it into pre-determined terms is called as Text Classification.

It is useful in scenarios where there is a large amount of unorganized text that must be classified & organised. One of the examples is called sentiment analysis.

Text classification is useful in SA, review classification, document summarification.

Traditionally Text Classification is done using ML and DL models but since the introduction of transformers NN, JC is been done using LLM.

Text Classification - ML → word2vec, Glove  
DL → LTM, GRU

Input → LLM - +ve.

- - ve  
— neutral

Metrics - A classification model performance can be evaluated by using certain evaluation metrics.

There are 4 basic metrics/ primary metrics from which evaluation metrics are defined - TP, TN, FP, FN

		Actual Values	
		+ve	-ve
Predicted Values	+ve	TP	FP
	-ve	FN	TN
+ve review incorrectly classified	+ve		-ve review incorrectly classified as +ve
	-ve		-ve review correctly classified as -ve

Confusion Matrix - From confusion matrix we can derive the evaluation metrics such as accuracy, precision, recall, F1 score.

$$\text{Accuracy} \rightarrow \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Recall} \rightarrow \frac{TP}{TP + FN}$$

$$\text{Precision} \rightarrow \frac{TP}{TP + FP}$$

$$\text{F1 score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Support - No. of samples

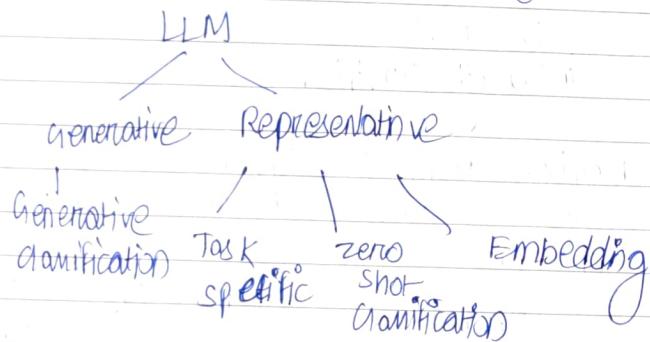
(Masks, weight) Avg

Precision - measures how many of the items found are relevant, which indicates the accuracy of the relevant results

Recall - refers to how many relevant classes are found which indicates its ability to find all relevant results

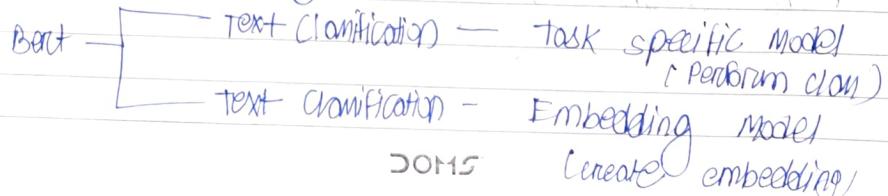
Accuracy - refers to how many correct predictions the model makes out of all predictions, which indicates overall correctness of the model

F1 score - balances both precision & recall to create a model's overall performance (categorised as best model)



### TC using Representative LLM

Classification with predefined representation models generally comes in 2 flavours either using a task specific model or an embedding model (when training dataset are labelled, used)



During Training if labelled datasets are not provided then representative models can do a 3rd type of classification

### Task specific Model

It is a fine-tuned variant of the base representative model such as "bert-base-uncased-finetuned-SST-2" which is bert fine-tuned on the labelled dataset Stanford Sentiment Treebank

Steps in TSM →

- 1) The input text is tokenized into subword units (Model i/p)
- 2) Model forward pass - The forward pass processes the tokens & finds semantic meaning
- 3) Classified Head - Based on semantic meaning of tokens the model has classification layer that outputs probabilities based on semantic meaning
- 4) Prediction - The class with highest probability is chosen

Embedding Model - In case we don't have specific model ie freely available or if we find user defined control on classifier in that case we can use a embedding model for text classification

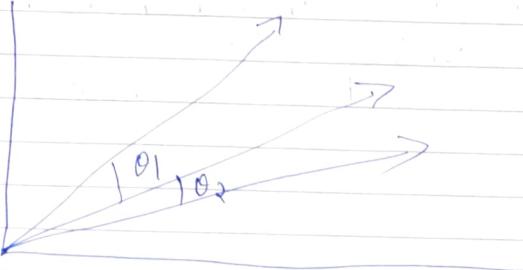
Steps →

- 1) Load an Embedding model - for example sentence transformers / all - npemb - base - VQ
- 2) Convert text to ~~string~~ Embeddings - The model generates a vector for each piece of text
- 3) Train a classifier - Feed these vectors into a standard classifier
- 4) Predict - for new input, I create the embedding and then run it through the classifier

### Zero shot classification

When there are no labelled datasets, then text classification can be done using representative LLM's through zero shot learning

- Step 1 - Create a text discipline of each label
- Step 2 - Embed the label description and your documents using the same embedding model
- 3 - Compute the cosine similarity b/w each document embedding and each label's embedding



DOMS

In cosine similarity, the similarity between 2 embeddings is found using trigonometric function cosine the angle.  $\cos(\theta)$  determines which 2 embeddings are closer

Ex - A very +ve movie review  
A very -ve movie review

Text to classify - This movie was an absolute masterpiece  
compute similarity with each label, whichever is higher gets assigned

### Text classification using generative LLMs

Plan - T5 for classification -

T5 → Text through Text Transfer Transformer

Plan - T5 is a fine-tuned variant of T5 esp for text classification

- 1) Load the model in your python program
- 2) Create a prompt - Is following sentence true or false
- 3) Generate output - the model's textual output might be the word true / false
- 4) Convert output to a label - Map - true → 1  
false → 0

DOMS

## Text Clustering & Topic Modeling

Text clustering - Given a large no of text inputs, the process of text clustering is to group similar documents into a cluster based on their content and semantic meaning. These groups are called as clusters and contain similar documents.

Clustering helps us in organising a large amount of unorganised input.

### Applications

Information retrieval - When the text is organised into clusters it is easier to retrieve data items and other related data items.

Recommender system - When the user selects a particular data item all other data items in the cluster can be provided as recommendation.

Spam filtering - When a message is labelled as spam all other (similar) messages are also added to the spam cluster.

Text summarization - Group similar sentences in a long text can help create a summary of the text.

Text classification using a ~~closed~~ <sup>Closed</sup> source model (Ch 10 P1)

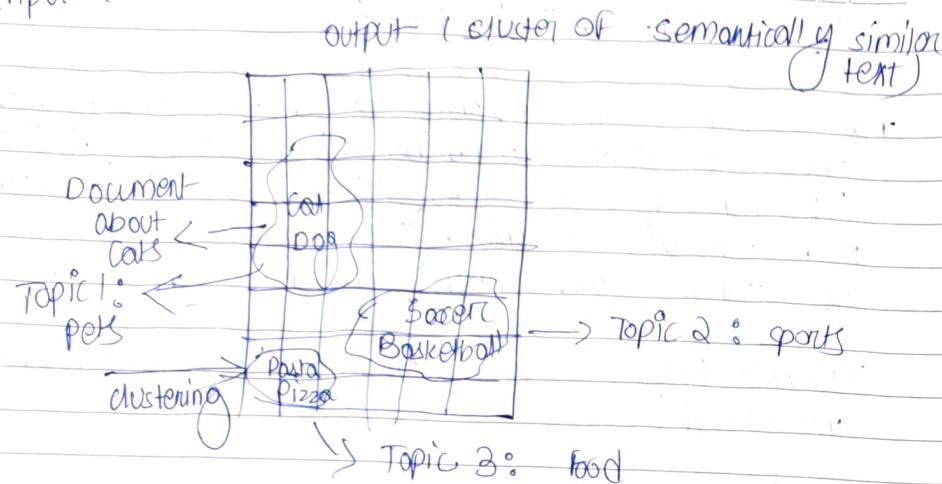
A closed source model provides API for task specific jobs like JC

Ex -

Document - The acting was superb and I laughed the whole time

Expected op : "1" (box its +ve)

Input (Textual data)



## Topic Modelling

Every cluster is created must be given an abstract definition which describes all the data items in that cluster. This abstract definition of a cluster is known as topic and process of doing so is called Topic Modelling.

A Topic helps in identifying a cluster as each topic is associated with several keywords.

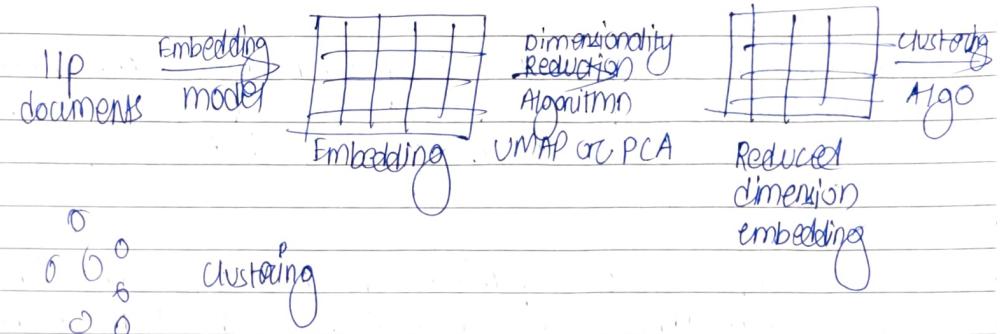
In this chapter we study the basic pipeline for text clustering and also discuss a fine-tuned model specifically made for topic modelling called as BERT Topic.

A common pipeline for text clustering

Traditionally Text clustering is performed using ML & DL methods but since introduction of transformer architecture text clustering is being done using LLMs

There are 3 steps which are executed for text clustering using LLM.

1. Convert the input documents into embedding with an embedding model
2. Reduce the dimensionality of embeddings with a dimensionality reduction model
3. Find groups of semantically similar documents with a cluster model

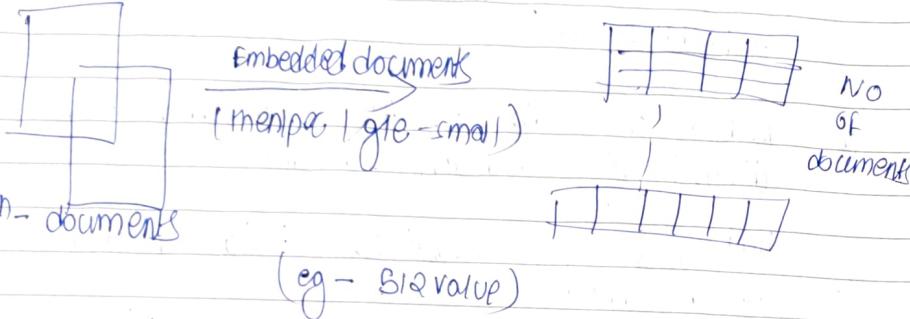


### Step 1) Embedding the input documents

Textual info / data is converted to its corresponding numerical representations known as embeddings. These embeddings can remember the context of meaning of words.

Embeddings for TC can be created using a general purpose embedding model called as "tinyper / gpt - small".

Embeddings that are produced by  $\mathbf{q}$  have a very high dimension



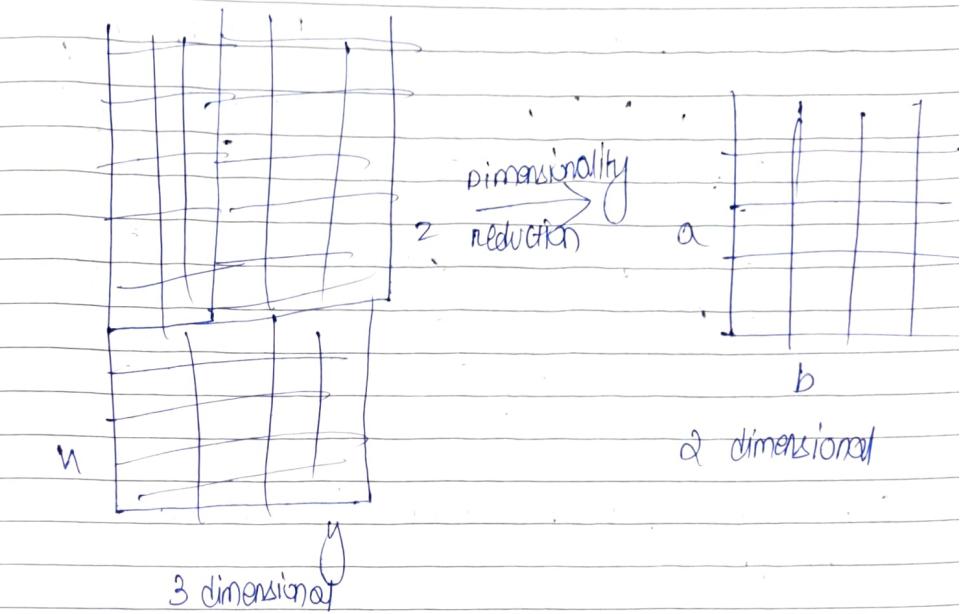
Step 2) Reducing dimensionality of embedding - An embedding of dimension  $D$  takes  $O(d^2)$  time for inferring the embedding vector when dimension  $d$  increases the surface increases

It is computationally expensive, therefore the dimensions of embeddings created in step 1 must be reduced

Dimensionality reduction is the process of reducing the dimensions and still getting OLP as if no dimension is reduced

2 of the dimensionality reduction is called ->

3 dimensional space  
2 dimensional space



Step 3) Cluster the Reduced Embedding - A clustering algo takes the reduced or compressed dimension embedding & produces clusters out of them

Centroid Based

In this approach one data item is chosen as centroid all other data item that surrounding the centroid within a threshold distance is included in cluster

These approaches can't detect outliers

density Based

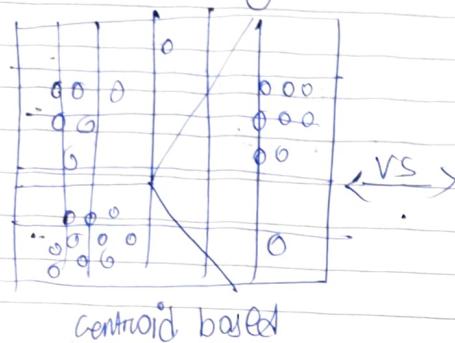
In this approach words have similar embeddings are put into a cluster

These approach can detect outliers

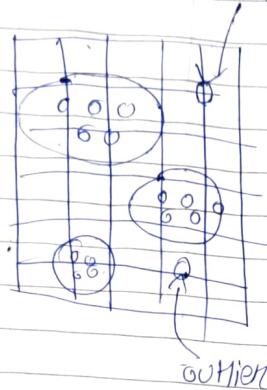
NOTE - Outliers are those data items which don't belong to any cluster

DOMS

Ex - K mean clustering



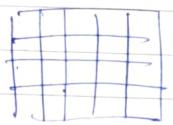
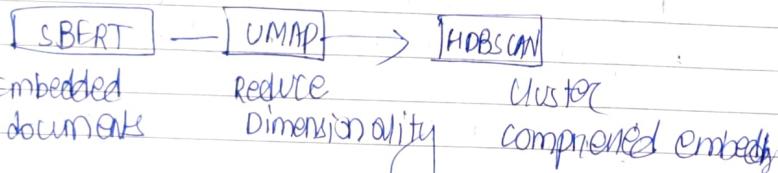
Ex - HDBSCAN output



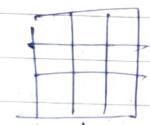
DBA are useful when we don't know no of clusters to be produced

Topic modelling using BERTopic

1st stage -



dimensional  
512 values



compressed dimensions  
3 values

BERTopic is a fine tuned version of BERT is specifically trained for topic modeling

In the 1<sup>st</sup> stage BERTopic uses SBERT in the embedding algo then it uses UMAP for reducing the dimension of the embedding

Finally it uses a density based approach like HDBSCAN to create the clusters

In stage 2 BERTopic uses the one word embeddings to count the frequency of each frequency which is called as term frequency.

My cat is cute → 0 1 0 1 0 1 1  
That is a cute dog → 1 1 1 1 1 0 0

that is a cute dog my cat,  
vocabulary of all tip documents

Term frequency calculates the frequency of each term per document but to give a topic to a cluster it must calculate the frequency of each term per document per cluster

Therefore instead of TF, BERTopic calculates CTF (term frequency) which is the frequency of each term in a document and in all documents across all clusters

my cat is cute → 1 2 1 2 1 1 1  
That is a cute dog → 1 6 3 1 8 0 5 0  
→ 1 1 3 2 0 3 0  
DOMS that is a cute dog my

frequency of word  $w$  in cluster  $C = \text{tf}_{Cw}, \text{cf}_w$

The class term frequency gives the frequency of all the terms used somewhere in any document of the cluster. This also calculates the frequency of words like "the", "a", "is" which can never be topic of cluster.

Therefore BERTopic must eliminate words that are not useful. To do this BERTopic calculates the inverse document frequency for each term to put more weight on words that are meaningful and puts less weight on words that are common across all clusters.  
ex - 'the', 'is', 'a'

$$TDF = \log \left( \frac{A + 1}{cf_n} \right) \quad A = \text{Avg freq} \\ cf_n = \text{total freq of each word}$$

$$= \log \left( \frac{24}{9.972219} + 1 \right) \text{ each } 1$$

that is a WTC dog? my cat.

After finding weight to each term in a cluster, the C-TF for a term is multiplied with its IDF value to find its weight.

$$\text{Weight of term } n \text{ in a clause} = \text{tf}_n, c \times \log \left( \frac{A}{\text{cf}} + 1 \right)$$

After assigning weight to each term in a cluster BERTopic finds out topic which has highest weight

This term is chosen as the topic of the cluster

## Final pipeline

## clustering



## Topic representation

Putting the two steps together clustering and representing topics results in the full pipeline of BERTopic. With this pipeline we can cluster semantically similar documents from the clusters generate topics represented by several keywords. The higher word's weight in a topic the more representative it is of that topic.

One of the major advantages of BERTopic, pipeline is it is Modular as we can replace or change any unit without effect other units.

## Using generate models for Text clustering

Generative models like GPT can also be used for text clustering. We have to provide the input documents & create a prompt for the generative model to cluster text. Along with the documents, important keywords must also be provided.

LLM

Keywords

[attention, text, nlp, transformer] + documents

I have a topic that contains following document  
[DOCUMENTS]

The topic is described by following keyword  
Based on the above info, give a short label  
of topic

I pretrained transformer model like nlp