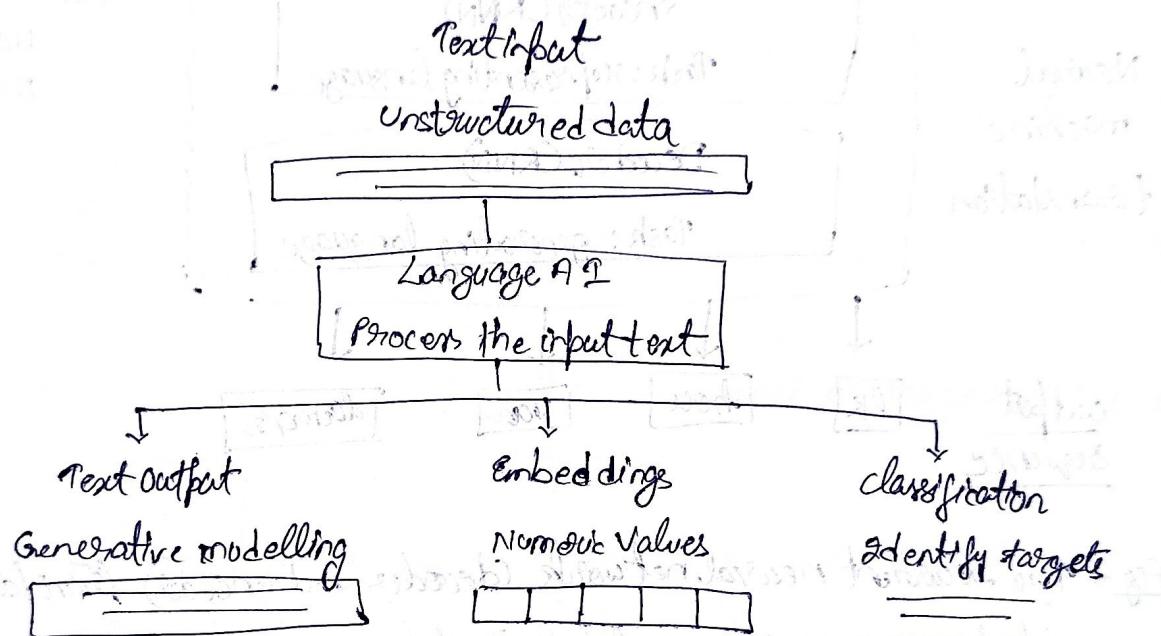


LLM (Large Language Model)

- Deep Learning is a subset of AI.
- Language AI is used to process and generate human language.
- * After 2018, the AI models uses transformer (convert one language to another).
- * Transformer can be Decoder, encoder or both encoder and decoder.

LLM - The deep learning based model used transformer to train on the natural language like text or image.

- The language AI is the subfield of AI used to understand process and generate the human language.



Bag of words - A method for representing unstructured text.

Steps
1) Tokenization - Splitting the sentences into individual words or subwords (tokens) is known as tokenization (tokens).

Input

This is a cute dog

↓ split into by whitespace

"that" "is" "a" "cute" "dog"

↓ create vocabulary

[that][is][a][cute][dog][my][cat]

Input

my cat is cute

↓

"my" "cat" "is" "cute"

numbers. This is why the tokenization comes into play.

Steps behind the tokenization -

- 1) Assign each token a unique number. These integers are then vectorized and feed into the language model.

2) The language model then processes the vectorized input and outputs the integers again, then the integers are subsequently converted back into text. That is the generated output.

There are 3 types of tokenization -

i) Byte Pair Encoding - Greedy merging of frequently adjacent character pairs to form subwords. Vocabulary is built by iteratively combining the most common pairs.

Eg: GPT, OpenNMT

ii) Word Piece - selects subwords to maximize the likelihood of the tokenized sequence under a language model. Often uses prefix indicators like for suffixes.

Eg: BERT, DistilBERT

iii) Unigram Language Model - Probabilistic model that assigns likelihoods to subword units and selects the most probable segmentation. Allows multiple segmentations per word.

Eg: XLNet, T5, Sentence Piece.

Eg: Byte pair encoding example: aa abc abc

Initial Corpus: ["aa", "abc", "abc"]

Step 1: Prefetching

"aa": ["aa", "aa", "aa", "-"]

"abc": ["ab", "bc", "bc", "-"]

"abc": ["a", "b", "c", "-"]

Step 2: Iterative merging

1. Iteration 1:

- Most frequent pair: ("aa", "aa")
 - Merge into: "aab"
 - Corpus update: "aab", "a", "-"
 - Vocabulary: ["a", "ab", "c", "-", "aab"]
2. Iteration 2
- Merge into: "abc"
 - Corpus Update:
 - ["a", "a", "-"]
 - ["abc", "-"]
 - ["abc", "a", "-"]
 - Vocabulary: ["a", "b", "c", "-", "aab", "abc"]

3. Iteration 3

(aa, ab)

Merge into: "aa"

• Corpus update:

- ["aa", "-"]
- ["abc", "-"]
- ["abc", "a", "-"]

• Vocabulary: ["a", "b", "c", "-", "aab", "abc", "aa"]

Final Tokenized Corpus

- ["aa", "-"]
- ["abc", "-"]
- ["abc", "a", "-"]

Final Vocabulary: ["a", "b", "c", "-", "aab", "abc", "aa"]

Final tokens: [aa, abc, a] = aabbccaa

- The decoder block in transformer consists of
 - i) Masked self-attention
 - ii) encoder attention
 - iii) encoder, new, network

- q19
- BERT is an encoder only model (introduced in 2018).
 - GPT 1 is an decoder only model. Trained with 4000 books and has 117 million parameters. (Proposed in 2018).

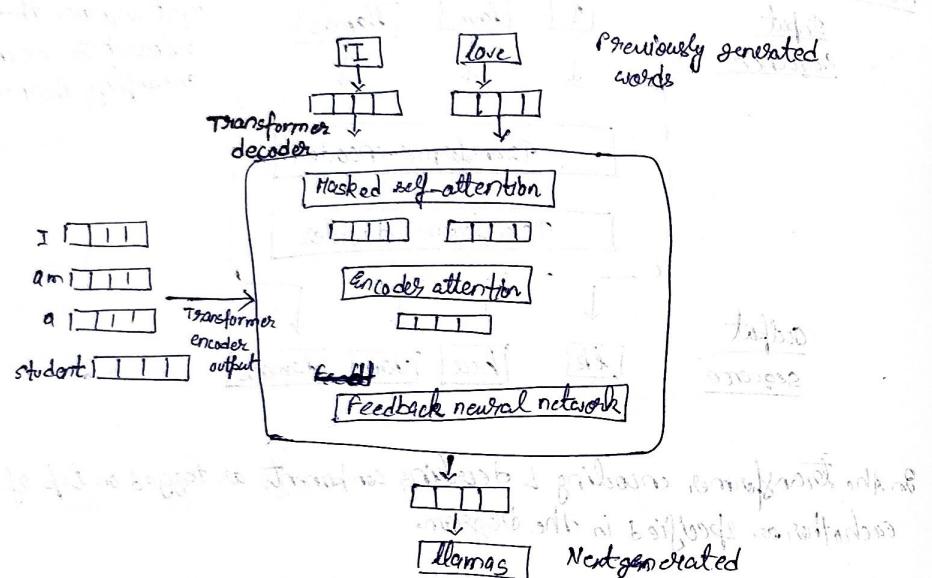


Fig. 9 → The decoder has an additional attention layer that attends to the output of the encoder.

BERT - (Bidirectional encoder Representations from transformer)

- It's an encoder only architecture, they focus on representing language.

Input sentence.

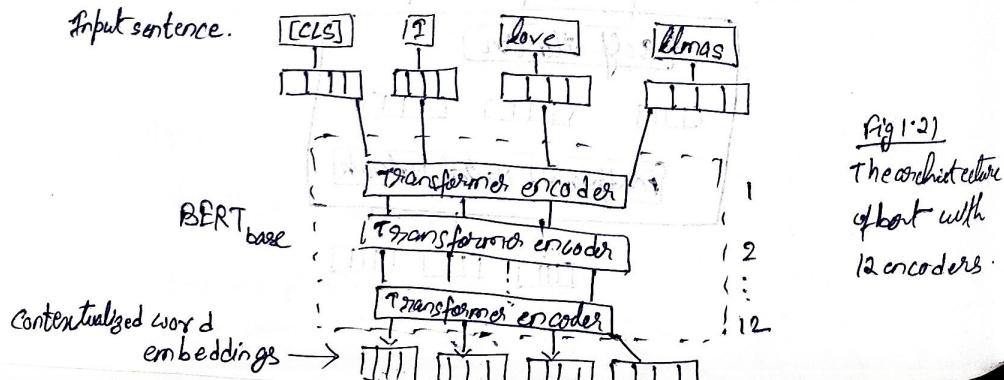


Fig 1.21
The architecture of BERT with 12 encoders.

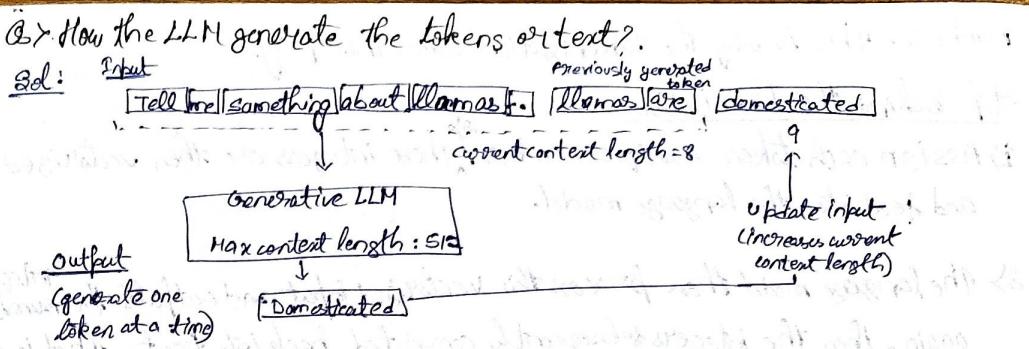


Fig 1.27 → The context length is the maximum context an LLM can handle

Chapter - 2 : Tokens and Embeddings

Tokentization and embeddings -

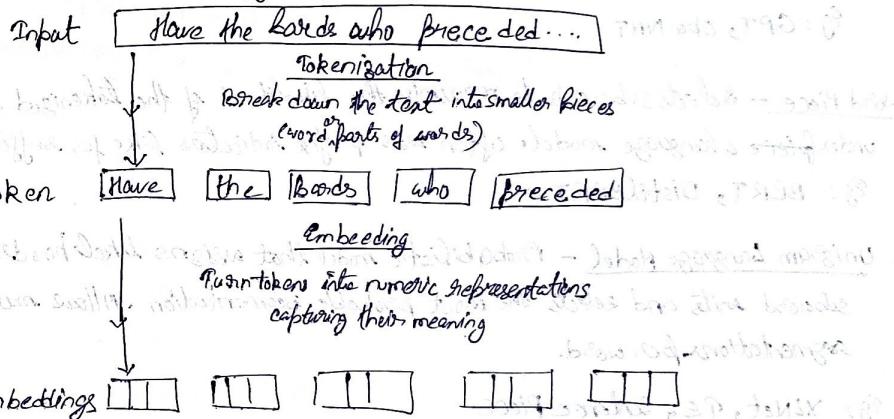


Fig 2.1 → Language models deals with text in small chunks called tokens. For the language model to compute language, it needs to turn tokens into numeric representations called embeddings.

Tokentization - is the process of splitting of sentence into small parts or sequence of tokens which can be either words, subwords or characters.

Q) Why tokentization is required?

Suppose, we want to pretrain our model and we have 10 GB of language data available.

However, the data is in text form but the language model processes the data in numbers. Therefore, we need to convert the text into the

Input

My cat is cute

↓
my [cat] [is] [cute]

↓
that is a cute dog my cat
0 1 0 1 0 1 1 ← vector representation

Bag of words

count individual words

Convert the text from english to dutch language

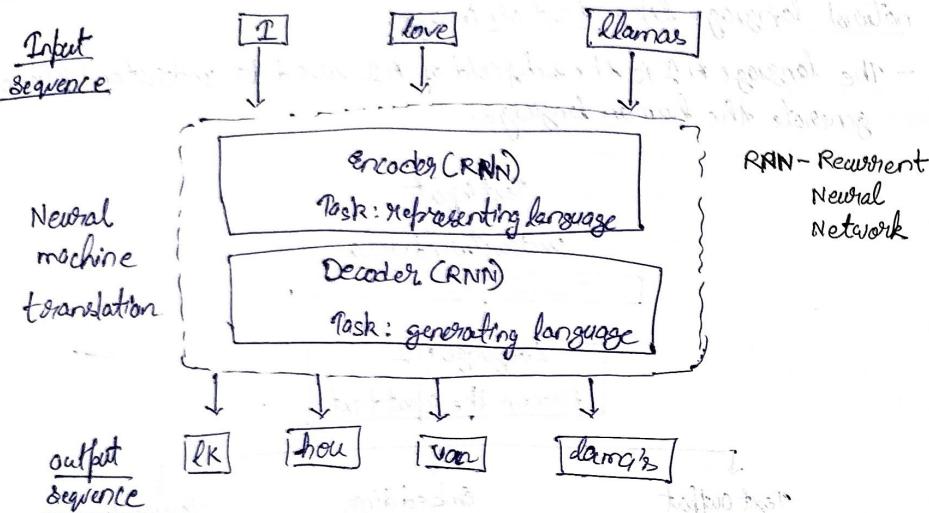


Fig - Two recurrent neural network (decoder and encoder) translating an input sequence from English to Dutch.

Embedding - process of converting tokens into vector forms which try to capture its meaning.

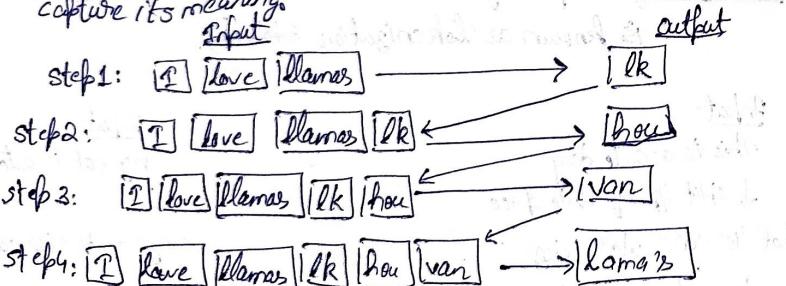


Fig 112 → Each previous o/p token is used as i/p to generate next token

embedding - allows us to measure the semantic similarities b/w 2 words.

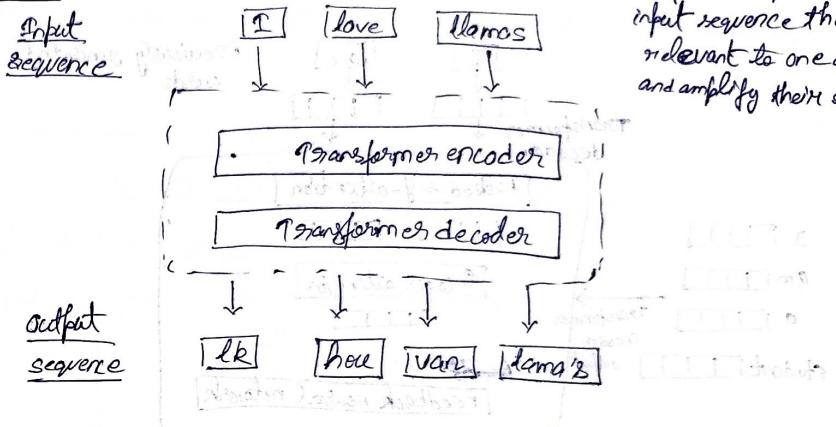
- Using various real metrics we can judge how close one word to another

Types of embedding

i) Document embedding ii) Sentence embedding

iii) Word embedding iv) Token embedding.

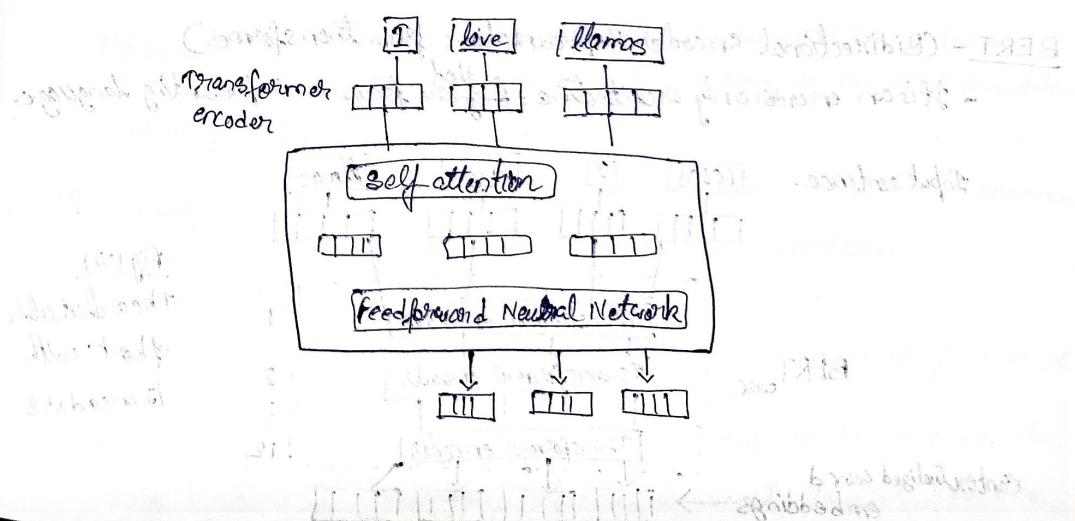
Attention - "Attention is all you need". It allows a model to focus on parts of the input sequence that are relevant to one another and amplify their signals.



In the transformer encoding & decoding components are tagged on top of each other as specified in the diagram.

- The encoder block in transformer consists of 2 parts -

- i) self attention and ii) feed forward network



Word Piece - The word piece tokenization was developed by google that the primary goal is to handle the unseen words more effectively than traditional word level tokenization.

Suppose the vocabulary is given as [UNK, a, aa, abc, ###b, ###c]

Cases [aa, abe, abc]

case 1: abc and ab are in vocabulary

- "aa":

- Longest prefix match : "aa" ∈ vocab

- Tokenized = ["aa"]

- "abe":

- Longest prefix match : "abc" ∈ vocab

- Tokenized = ["abc"]

- final tokenized output: [[aa], [abe], [abc]]

case 2: abc and ab are not in vocabulary

step by step tokenization of "abc":

1. start at position 0: Try longest match from vocab

- Try "aab" → not in vocab

- Try "ab" → not in vocab

- Try "a" → match

- emit: "a"

2. Position 1: Remaining string "bc" → try (aab) → not in vocab

- Try "###bc" → not in vocab

- Try "###b" → match

- emit: "b"

3. Position 2: Remaining string is "c".

- Try "###c" → match

- emit: "c"

Final Tokenization: ["a", "##b", "##c"]

Note: Word Piece doesn't look back or look ahead to match "###bc" unless it's explicitly in the vocabulary

Sentence Piece - is language independent subword tokenization method developed by google which is designed to train directly on raw text without requiring the pre-tokenization.

- unlike BPE and wordpiece, sentence piece treats the raw string of unique code characters for tokenization.

- The sentence piece use unigram language model that assigns a probability ' P_i ' to each subword ' t_i ' in the vocabulary to find the total cost.

$$\text{cost}(t_i) = -\log P_i$$

Eg: Suppose sentence is "abc". Use sentence piece algorithm to tokenize the "abc".

Token	Probability P	cost = - $\log P_i$
-------	---------------	---------------------

a	0.20	1.60
---	------	------

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
-----	------	------

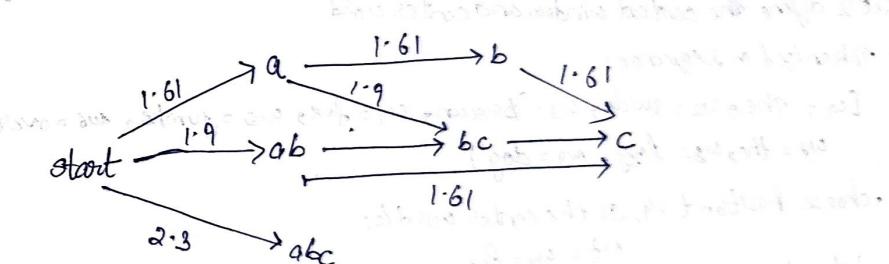


Fig 4: The cost lattice for "abc"

Segmentation: [aa, ab, ac]

Cost: $1.61 + 1.61 + 1.61 = 4.83$

Segmentation: [ab, ac]

Cost: $1.90 + 1.61 = 3.51$

Segmentation: [aa, abc]

Cost: $1.61 + 1.90 = 3.51$

Segmentation: [abc]

Cost: 2.30

Fig 3.10 - When generating text, it's important to cache the computational results of previous tokens instead of repeating the same calculation over and over again.

Transformer block - The transformer LLMs are composed of series of transformer blocks and each block processes its input then passes the results of its processing to the next block.

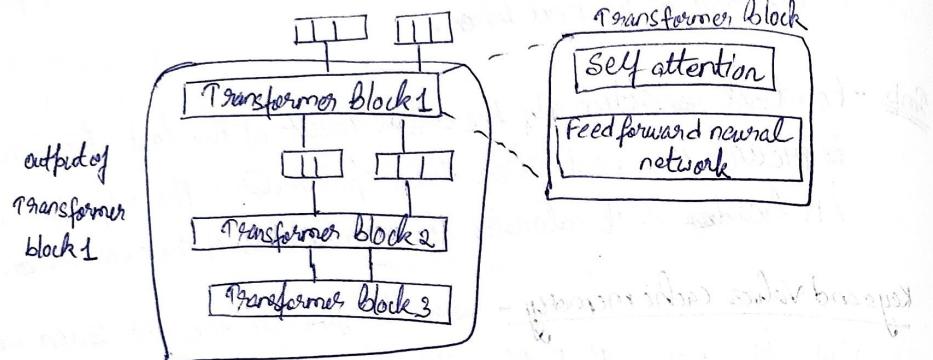


Fig 3.12. A transformer block is made up of a self-attention layer and a feedforward neural network.

The Shawshank - is a movie. When we search "The Shawshank" it will automatically complete the whole moviename "The Shawshank ~~Recall~~ Redemption".

19/9

Attention Layer - Attention is a technique that uses to process the input and generates the output according to the previous inputs.

- The self attention layer is mainly concerned with incorporating relevant information from other tokens/input tokens and positions.
- The feed forward layer ~~as~~ layer finds the majority of the models' processing capacity.
- The self attention mechanism is made of 2 major steps
 - i) Relevance scoring
 - ii) Feed forward or Combining information.

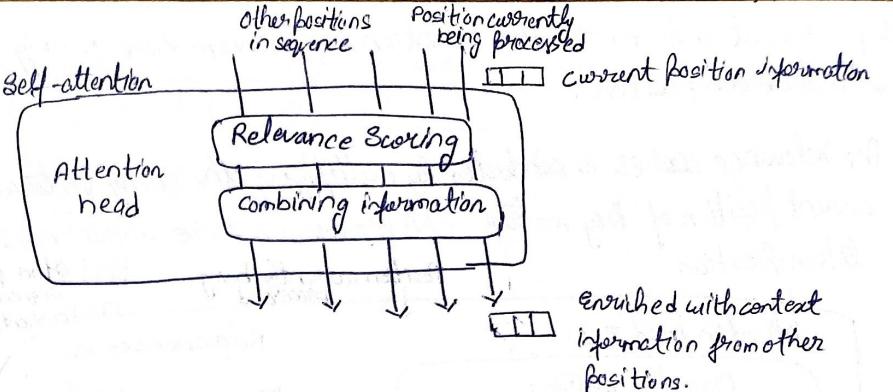


Fig 3.16. Attention is made up of 2 major steps: relevance scoring for each position, then a step where we combine the information based on those scores.

Attention Head - contains all the parallel applications of attention into the self attention layer to give the transformer more extensive attention capability.

Q) How attention is calculated?

- The attention layer in generative LLM used for processing the attention for a single position.
- The inputs to the layer are
 - i) Vector representation of the current position / information.
 - ii) Vector representation of previous tokens.
- The goal here is to produce a new representation of current position that incorporates relevant information from the previous tokens.
- In further the training process produces 3 major projection matrix to interact with the attention head for the calculations.
 - i) Query projection matrix
 - ii) A key projection matrix
 - iii) A value projection matrix.

Relevance Scoring

Relevance Scoring - In generated transformer every time it will generate only 1 token that means the processing will be always

Chapter 3 - Looking Inside of Large Language Models

Auto-regressive model - In machine learning, the models that consume their early predictions to make later predictions means the models first generate the ~~the~~ token & used to generate the second token are known as "auto-regressive models".

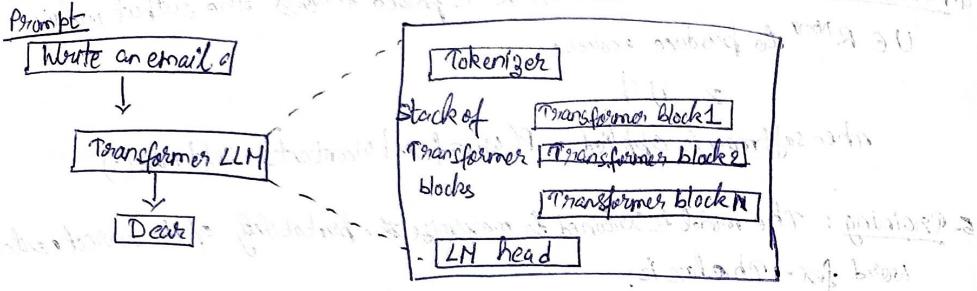


Fig 3.4 - A transformer LLM is made up of a tokenizer, a stack of transformer blocks, and a language modelling head.

Decoding strategy - ~~In the~~ The method of choosing a single token from the probability distribution means the probability score for each token in the vocabulary is called as "Decoding strategy".

- Choosing the highest scoring token every time is called as "Greedy Decoding".

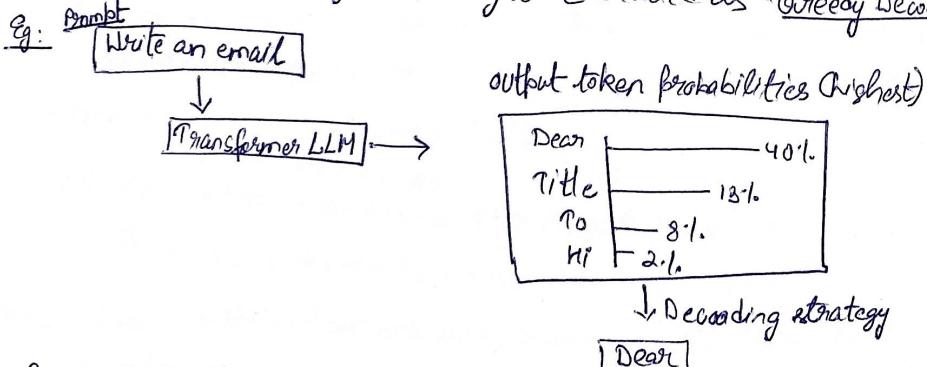


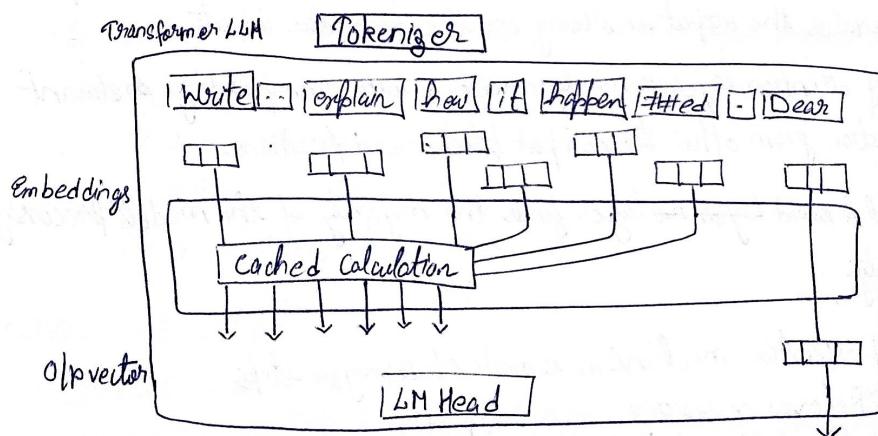
Fig 3.7 - The tokens with the highest probability after the model's forward pass. Our decoding strategy decides which of the tokens to output by sampling based on the probabilities.

Q) How the Transformer LLM works?

- Ans:
- First, the tokenizer in transformer LLM creates or converts the tokens which in further goes for token embedding means the vector representation.
 - Each generated token, the process flows once through each of the transformer blocks in the stack in order.
 - After that generated output goes to LM head which contains the probability distribution for the next token.

Prob - For text generation only the output result of the last stream is used to predict the next token that output vector is the only input into the LM head as it calculates the probabilities of the next token.

Keys and Values cache memory - When we generate the 2nd token we simply append the output token to the input and do the forward pass to the model. However, if we give the model the ability to ~~call~~ cache the result of the previous calculation, we no longer need to repeat the calculation of the previous strings. This time the only needed calculation is for the last string. This process is called as "Key Value cache" which speeds up the generation process significantly.
- This is known as the "attention mechanism".



optimal path via Viterbi Decoding

The Viterbi algorithm selects the path with minimum total cost.

$$\text{In this case } E[\text{abc}] = 2.30$$

Best segmentation $[\text{abc}]$, so total cost = 2.30

Embedding - It will find the semantic similarity b/w words using vectors matrix that can judge how close one word to another.

The embedding techniques has been used to convert the each token in input sequence into the numerical forms called as Input embedding.

- Each token is mapped to a vector using an embedding matrix

$$x = \text{embedding}(t)$$

t = token id

x = resulting vector

Word2Vec - It's the most embedding technique uses a sliding window to generate the semantic and syntactic relationship.

Eg: Consider the sentence: "The quick brown fox jumps over the lazy dog."

Let's define the context window and center word

- Tokenized in sequence:

$$[w_1 = \text{the}, w_2 = \text{quick}, w_3 = \text{brown}, w_4 = \text{fox}, w_5 = \text{jumps}, w_6 = \text{over}, w_7 = \text{the}, w_8 = \text{lazy}, w_9 = \text{dog}]$$

- choose position $t=4$, so the center word is:

$$w_t = w_4 = \text{fox}$$

Let, the context window size be $c=2$. Then the context words are

$$w_{t-2} = w_2 = \text{quick}, w_{t-1} = w_3 = \text{brown}, w_{t+1} = w_5 = \text{jumps}, w_{t+2} = w_6 = \text{over}$$

Cbow objective - Use the context words quick, brown, jumps and over to predict the center word fox.

embedding and prediction steps

1. One hot encoding - Each context word w_i is represented as a one-hot vector

$x_i \in \mathbb{R}^V$, where V is the vocabulary size.

2. embedding lookup : each one-hot vector is multiplied by the input embedding matrix $W \in \mathbb{R}^{N \times N}$ to obtain its dense embedding

$$e_i = W^T x_i$$

3. context aggregation : the embedding of four context words are averaged:

$$h = \frac{1}{4}(e_{\text{quick}} + e_{\text{brown}} + e_{\text{jumps}} + e_{\text{over}})$$

4. Prediction: the hidden vector $h \in \mathbb{R}^N$ is passed through the output matrix $U \in \mathbb{R}^{N \times V}$ to produce scores:

$$z = U^T h$$

then softmax is applied : $P(w_t = \text{fox} | \text{context}) = \text{softmax}(z)$

5. Training: The model is trained to maximize the probability of the correct center word fox. The loss is:

$$L = -\log P(\text{fox} | \text{quick, brown, jumps, over})$$

Result: after training, the rows of W become meaningful word embeddings.

Word that appear in similar context (like fox and dog) will have similar embeddings.

1 position at a time. So, the attention mechanism here is only concerned with those one position.

The relevance scores is conducted by multiplying the query vector with the current position of key matrix. This produces a score about each previous token position.

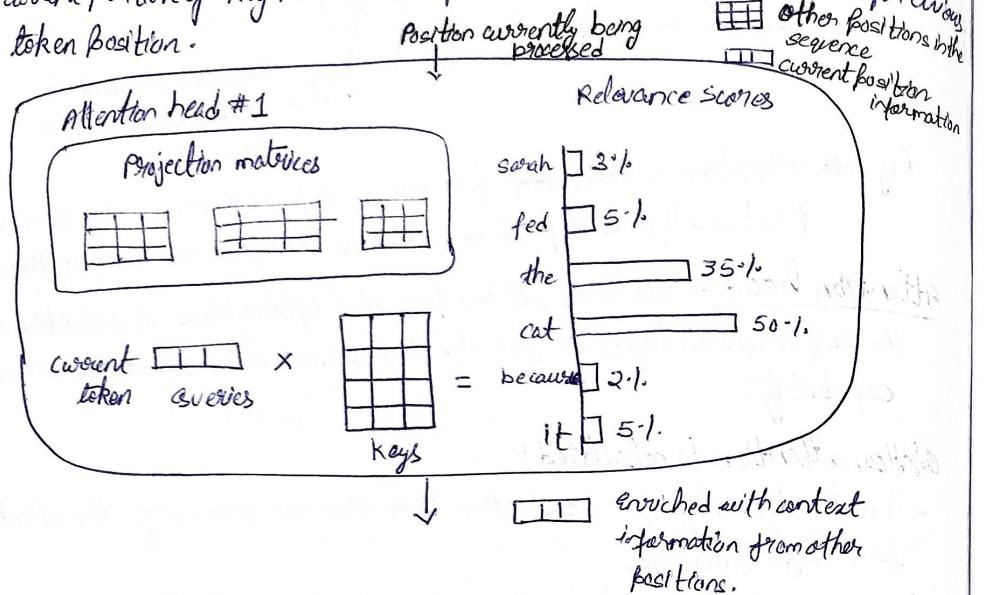


Fig 3.20 → Scoring the relevance of previous tokens is accomplished by multiplying the query associated with the current position with the keys matrix.

• Passing these scores by a softmax function normalizes these scores by sum of ~~to 1~~ to 1.

Combining information - After performing the relevance scores the attention mechanism multiplies the scores with the value vectors for each token that means ~~to 1~~.

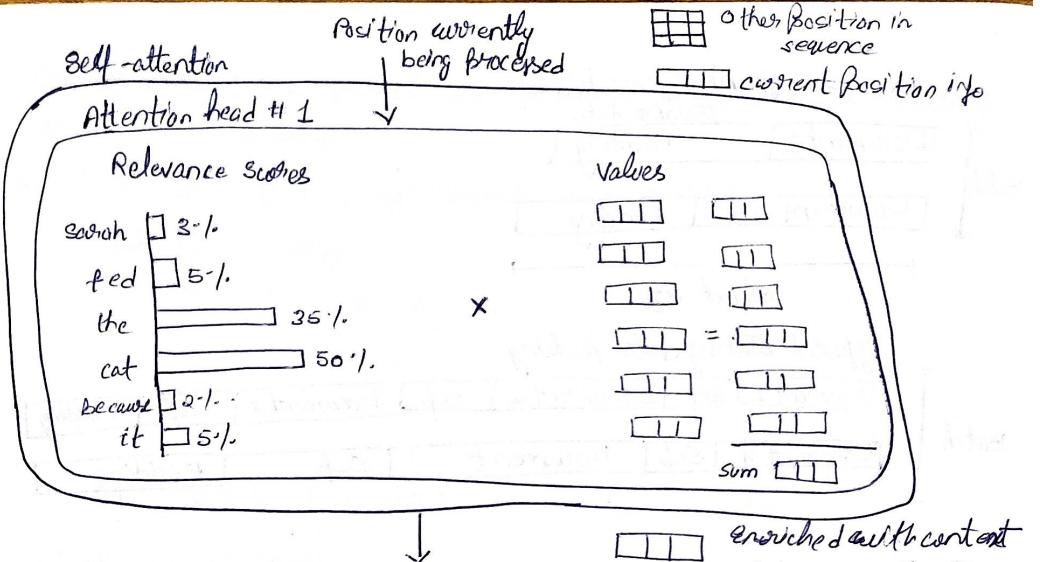


Fig 3.21 Attention combines the relevant info of previous positions by multiplying their relevance scores by their respective value vectors.

Local / sparse attention - To make more efficient attention, the local or sparse attention mechanism limits the context of previous tokens by only paying attention to a small no. of previous tokens / positions.

Positional Embedding - also called as RoPE (Rotary Position Embedding).

- The RoPE is a key component in Transformer LLM.

Q: Why positional embedding is required?

Sol: Let us take an example.

There are lot of documents are trained using language model, however the training set is much shorter than that content.

- That means it would be inefficient to allocate the entire content.

Ex: 4K ~~age~~ into a 10 word sequence.

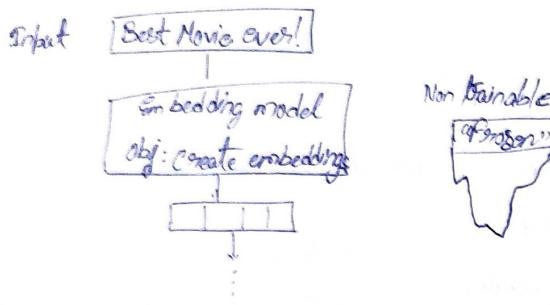
• So during model training the documents are packed together into each context in the training batch, that is known as Packing.

our model make the prediction.

Classification task that leverage embeddings -

Instead of directly using the representation model for the classification, the embedding model uses a step approach for the text classification.

i) In the first step convert the text input into the embeddings using embedding model.



ii) These embeddings in the 1st step serve as the input features to the classifier. The classifier is trainable and can perform classification.

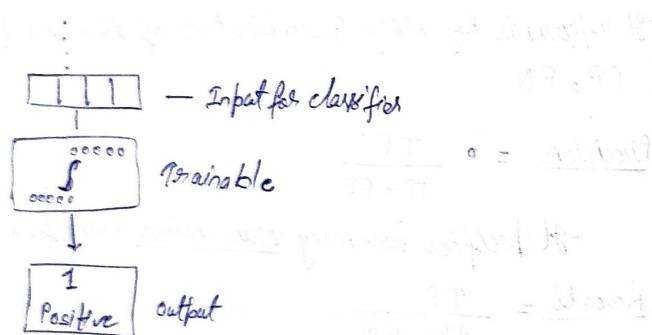


Fig 4:12

Benefit of this method is that -

We didn't need to fine tune our embedding model which can be costly.

1010125

Text classification using generative models -

The generative models takes as input some text and generate text in the form of sequence of tokens which is also known as sequence to sequence model. However, the task specific model takes the sequence of tokens as input but the model generates the numerical values instead of tokens.

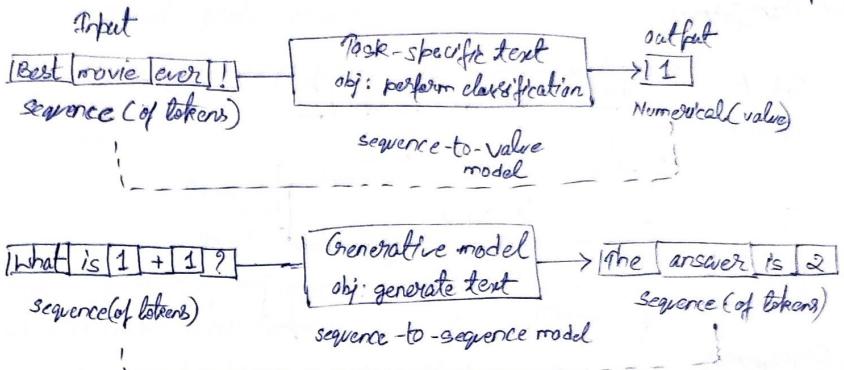


Fig 4:17 A task-specific model generates numerical values from sequences of tokens while a generative model generates sequences of tokens from sequence of tokens.

- These generative models are generally trained on a wide variety of task and usually don't perform out-of-box.

For eg: If we give a generative model a movie review without any context then it has no idea what to do.

- This guiding process is done mainly through the instruction or known as "prompt" that we provided to the model.

- Iteratively improving the prompt to achieve the desired result is called as prompt engineering.

- Prompt engineering allows prompts to be updated to improve the output generated by the model.

Text to Text

CH-4: Text Classification

Text classification - classification of text is used to extract the meaningful insights like entities and detecting the language.

- The text classification in the field of language models can be represented using:

- i) Text classification with representation model
- ii) Text classification using generative model

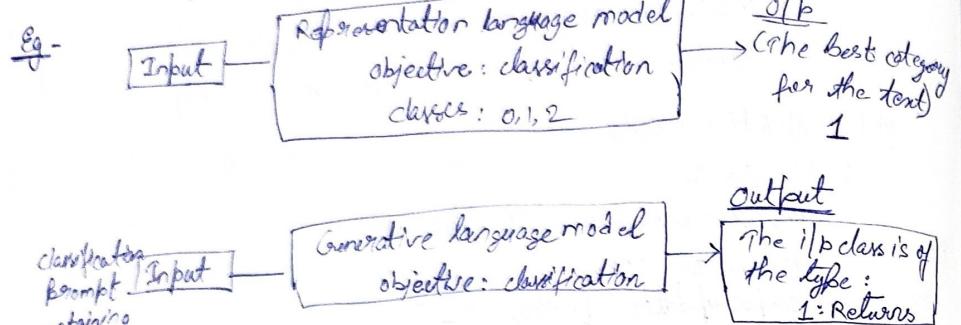


Fig 4:2. Although both representation and generative models can be used for classification, their approaches differ.

Text classification with representation models -

- i) Task specific model
- ii) Embedding model

- The task specific model and embedding model are used for the text classification where the task specific uses tokenizer and token embedding to generate the word embedding.

- On the other hand the embedding model uses the supervised ML techniques for the text classification.

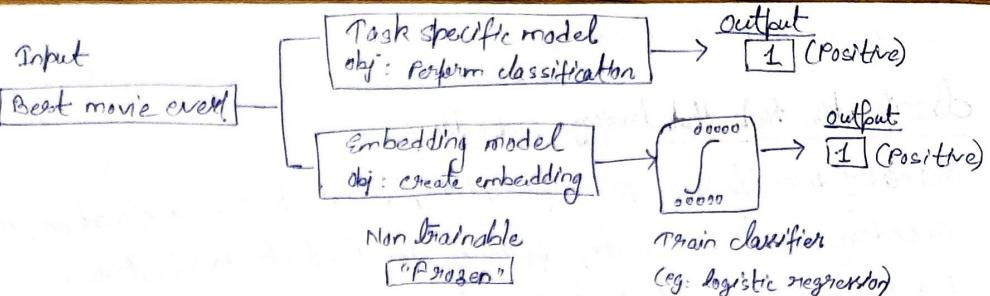


Fig 4:4. Perform classification directly with a task-specific model or indirectly with general purpose embeddings.

Confusion matrix

		Actual values	
		P	N
Predicted values	P	TP	FP
	N	FN	TN

↓
incorrectly classified
as negative

- The confusion matrix shows whether the model correctly or incorrectly predicted result correspond to the actual values

- It represents by using 4 combinations of true and false values i.e TP, TN, FP, FN

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

- It specifies how many true values identified actually.

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

- It defines out of all values which are actually predicted true.

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

Q - Accuracy find how many are true out of all predictions.

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

- It combines both recall and precision scores which measures how effective

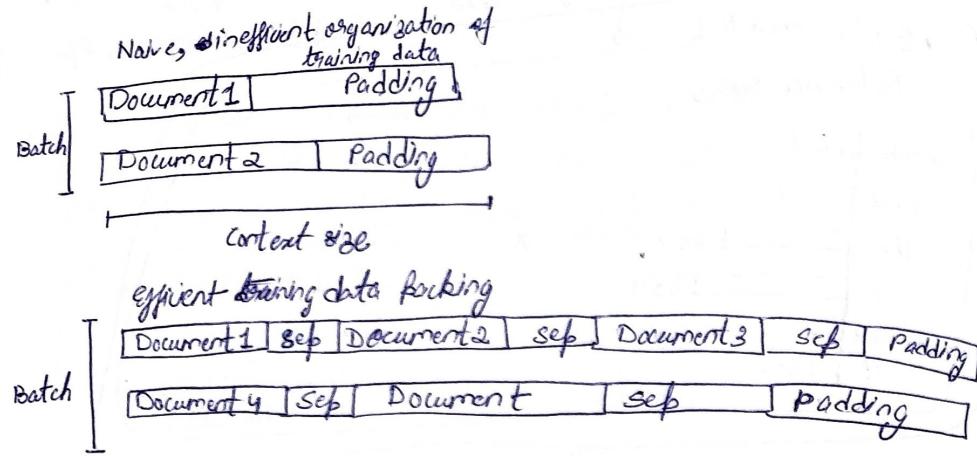


Fig 3.31. Packing is the process of effectively organizing short training documents into the context. It includes grouping multiple documents in a single context while minimizing the padding at the end of context.

Rotary embeddings -

Rotary embeddings are a method to encode positional information in a way that captures absolute and relative token position information. This is based on the idea of rotary vectors in their embedding.

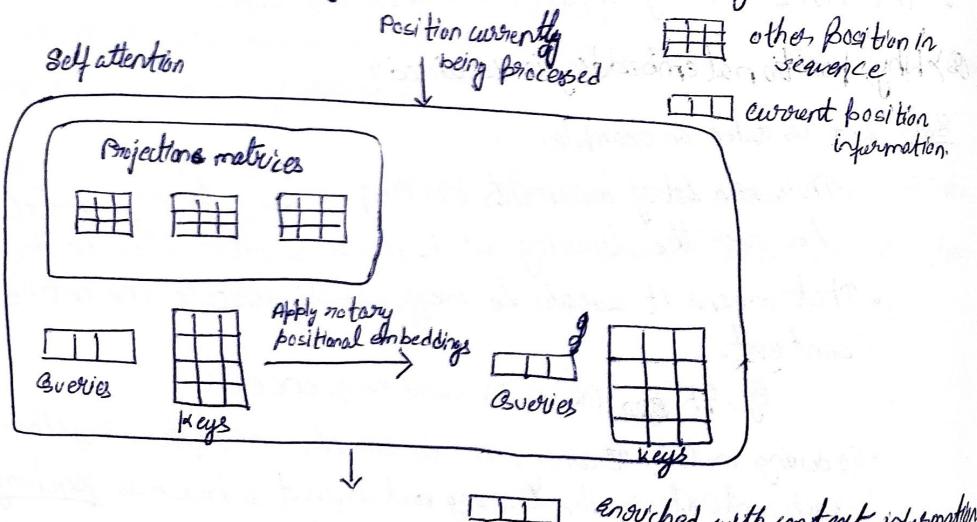


Fig 3.33. Rotary positional embeddings are added to the representation of tokens just before the relevance scoring step in all NLP models.

Ques 269
Q) "The dog chased the squirrel" use byte pair encoding to create vocabulary.
Ans: find the no. of tokens.

Step 1: Initial corpus : ["the", "dog", "chased", "the", "squirrel"]

Step 1: Preprocessing

"the": ["t", "h", "e", "-"]
 "dog": ["d", "o", "g", "-"]
 "chased": ["c", "h", "a", "s", "e", "d"]
 "the": ["t", "h", "e"]
 "squirrel": ["s", "q", "u", "i", "x", "r", "e", "l", "r"]

Step 2: Iterative merging

1. Iteration 1

Most frequent pairs: ("t", "h", "e")

Merge into: "the"

corpus updates: ["the", "-"]

["d", "o", "g", "-"]

["c", "h", "a", "s", "e", "d"]

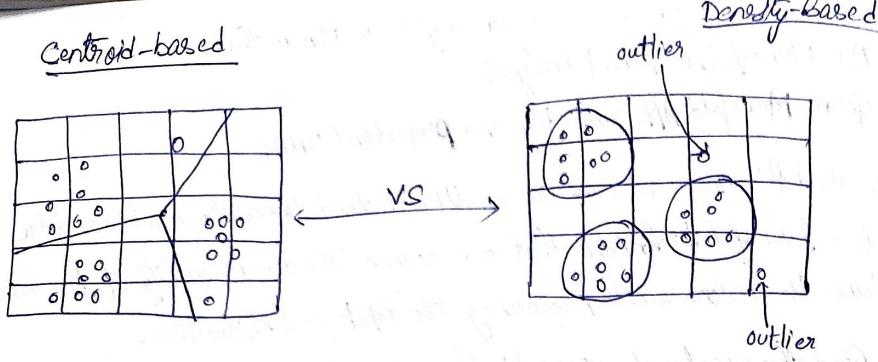
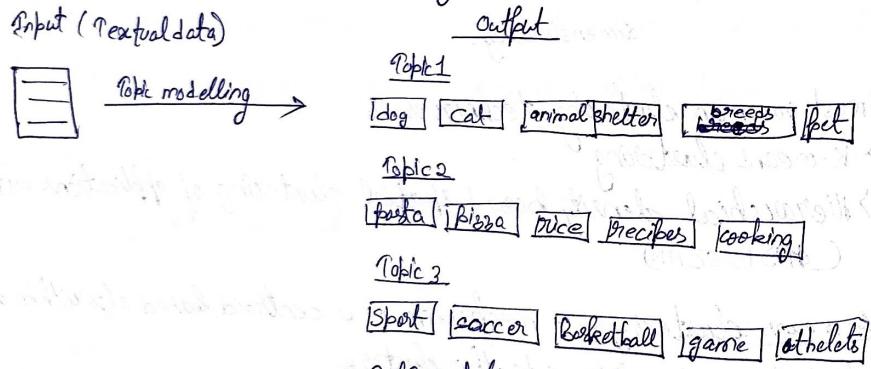


Fig: 5-7 : The clustering algorithm not only impacts how clusters are generated but also how they are viewed.

- The HDBSCAN algorithm allows dense clusters to be found without having to explicitly specify the no. of clusters, allows only dense clusters.
- It also ~~dictates~~ dictates the ~~outliers~~ outlier in the data which are the datapoints that doesn't belong to any ~~else~~ cluster.

21/10 Topic modelling -

The process of providing themes or topics to a cluster according to their idea of finding themes each known as topic modelling.



Representation can take many forms keywords
bi-grams, labels etc.

Fig: 5-9: Traditionally topics are represented by a number of keywords

- but can take another forms.
- Each topic is characterized by a probability distribution of words in a corpus vocabulary.

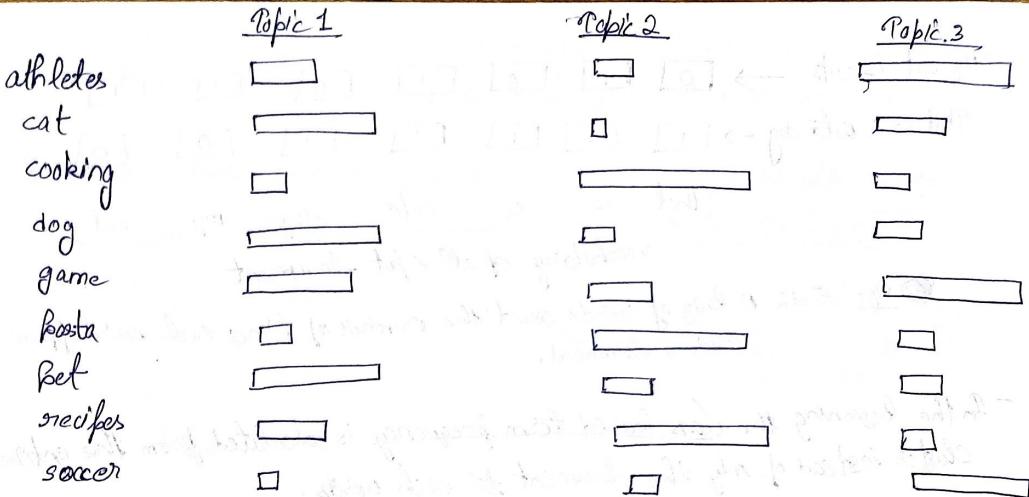


Fig: 5-10 : Keywords are extracted based on their distribution over a single topic .

BERT Topic : A Modular Topic Modelling Framework -

- The BERTopic modeling leverage the clusters of semantically similar text to extract various types of topics representations.
- It follows two steps approach -

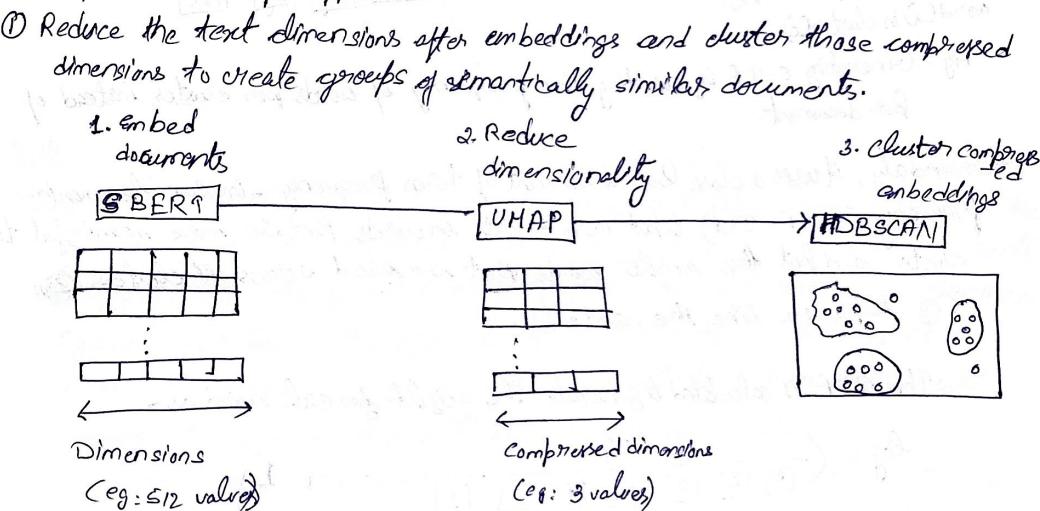


Fig 5-11: The first part of BERTopic's Pipeline is to create clusters of semantically similar documents.

- ② BERTopic model provides a distribution of words in a corpus vocabulary by leveraging a classic method known as Bag of Words.

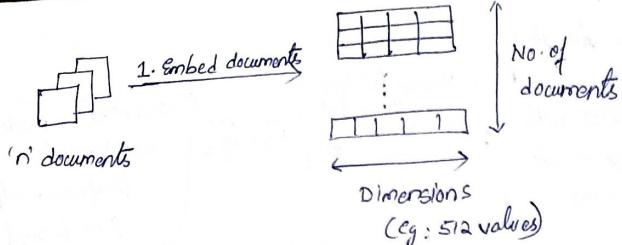


Fig 5.3 step 1 - We convert documents to embeddings using an embedding model.

- Here the textual data converting into embeddings means the numerical representation of text that attempt to capture its meaning.

Step 2 Reduce the dimensionality of embeddings using dimensionality reduction model.

- The dimensionality reduction process can be troublesome for many clustering techniques when there have larger volume of data as it gets more difficult to identify the meaningful clusters.

- As the dimensions in terms of no. increases there is an exponential growth in the no. of possible values within each dimension.

- So, the dimension reduction technique aim to preserve the global structure of high dimension data by finding low dimensional representations.

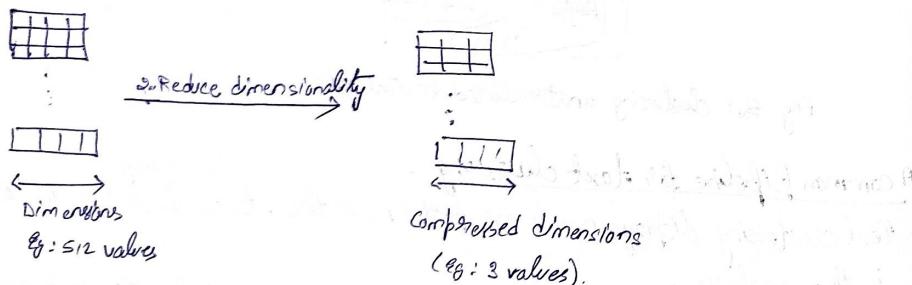


Fig 5. step 2 : The embeddings are reduced to a lower-dimensional space using dimensionality reduction.

- There are 2 most common dimensionality reduction method.
 - i) PCA (Principal Component Analysis)
 - ii) Uniform Manifold Approximation and Projection (UMAP)

* PCA is used for linear data whereas UMAP tends to handle non linear data.

* UMAP is a powerful technique that can reduce the dimensionality of the data into lower dimensions while preserving its topological structure.

Step 3 Find the groups of semantically similar documents with a cluster model.

- Clustering of the documents requires a set of clusters to be generated.

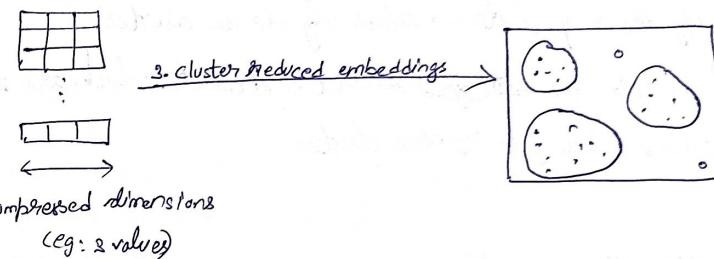


Fig 5.6 step 3 : We cluster the documents using the embeddings with reduced dimensionality.

2 most common clustering techniques are

- i) K-means clustering
- ii) Hierarchical density based spatial clustering of applications with noise (HDBSCAN)

* K-means clustering - mechanism is a centroid-based algorithm used to find the similar data into the clusters.

- It's widely used when clusters are spherical and of similar size.

* HDBSCAN - It is a hierarchical density based clustering method that calculate the no. of clusters to be found without forcing all datapoints to be part of a cluster.

Pert to text transfer transformer

- Instead of encoder-only model (BERT) and decoder-only model (chatGPT), the text-to-text transfer model (T5) contains both encoders and decoders for the classification of text.
- T5 model contains 12 encoder and 12 decoders that work for text-to-text transforms.

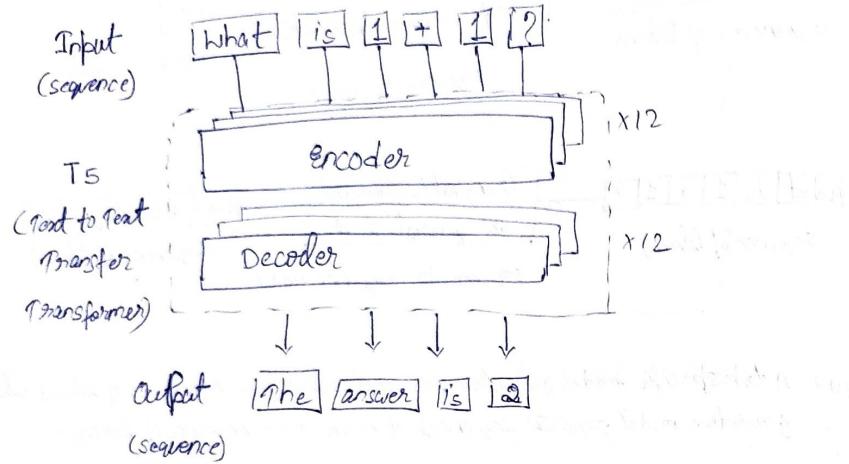


Fig: 4.19 The T5 architecture is similar to the original transformer model, a decoder-encoder architecture.

- In 1st step, the models first pre-trained with set of tokens known as masked language modelling.
- In 2nd step of training, the model is fine tuned with variety of task.

ChatGPT for classification -

- The ChatGPT uses open AI during the training process that generate the desired output to an input data in the prompt. every time it checks for the preference and final output will be accordingly preference - tuning. for e.g:

Generate output with instruction-tuned model

- (A) An abbreviation for the Master of Laws.
- (B) I am not familiar with ...
- (C) Large language models are artificial ...

Human labels rank the output
preference tuning

Generative model
obj: generate
preferred output
given prompt

Create preference Data

Fig: 4.23 Manually ranked preference data was used to generate the final model, chatGPT.

14/10

CHAPTER - 5

Text clustering and topic modelling

Text clustering - is an unsupervised technique that ~~aims~~ to group similar text based on their semantic context meaning and relationships.

- The resulting of clustering semantically similar documents not only facilitate efficient categorisation of large volume of data unstructured text but also allow for quick exploratory data analysis (EDA).

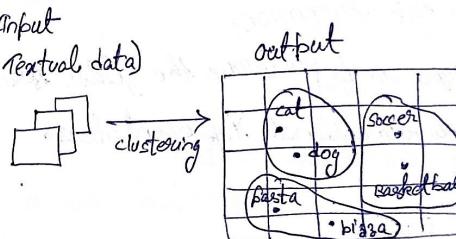


Fig 5.1 clustering unstructured textual data

A common pipeline for text clustering -

- Text clustering displays non linear patterns in the data which may contain the complexity in the analysis.
- To apply the text clustering from the graph based neural network to centroid based clustering technique we use 3 common steps and algorithms.

steps

- 1) Convert the input documents to embeddings with an embedding model.

Top-p - It is a nucleus sampling method which controls the subset of tokens (nucleus) that the LLM can handle, means it will consider tokens until it reaches the cumulative probability.

For eg: If we set top-p to 0.1, it will consider tokens until it reaches that value but if we say top-p to 1 then it will consider all tokens.

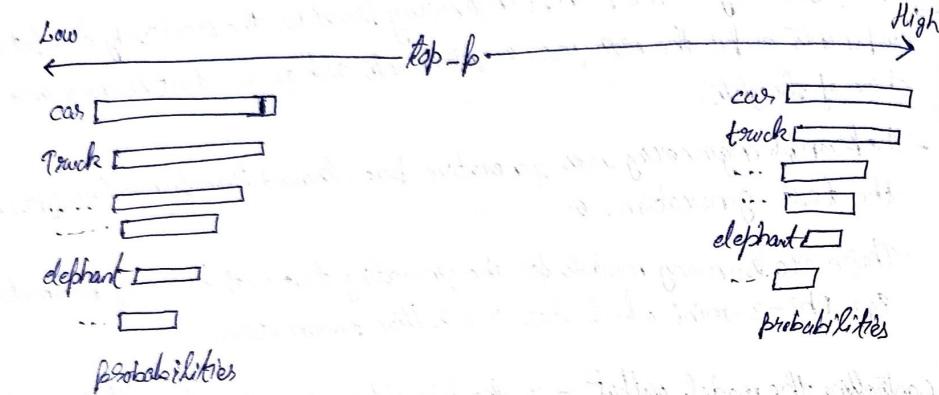


Fig 6.5 - A higher top-p increases the no. of tokens that can be selected to generate and vice versa.

In this figure, the lowering the value it will consider the fewer tokens and generally give less creative outputs while increasing the value the LLM to choose from more tokens.

Introduction to Prompt engineering

It is a tool that evaluate the output of a model as well as to design soft safety methods. This is known as Prompt optimization.

- The basic LLM model is a prediction machine that is based on a ~~set~~ certain input, where the prompt tries to predict the words that might give the appropriate or ~~set~~ relevant answers.

- In the prompt box we can ~~contai~~ continue adding or updating the elements until we can get the response what we are looking for.

- We can add additional examples.

Eg: User cases in more details, provide additional context etc to get the result.

- In BERT topic the re-ranker models are known as representation models.
- For eg: if we have millions of documents and hundred topics then the representation model only needs to be applied once for every topic instead for every document.

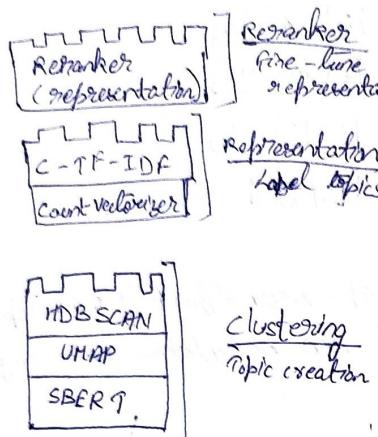
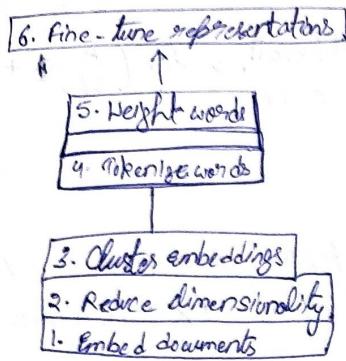
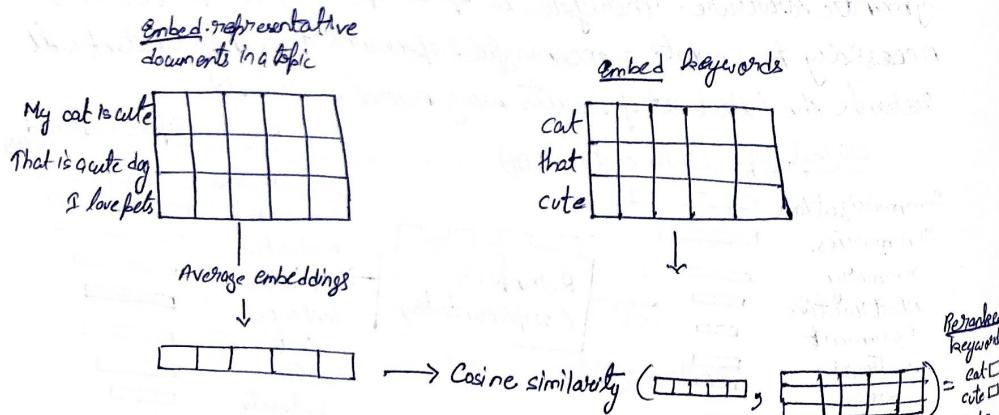


Fig 5.20 - The re-ranker (representation) block sits on the top of C-TF-IDF representation.

- A wide variety of representation model have developed for.
- Eg: KeyBERT inspired
- KeyBERT Inspired uses C-TF-IDF to extract the most represented documents per topic by calculating the similarity between a document's C-TF-IDF values and those topics related to that keywords.



5.22 KeyBERT inspired representation model

- In this figure the average document embeddings per topic is calculated and compare to the embeddings of candidate keywords to rerank the keywords.

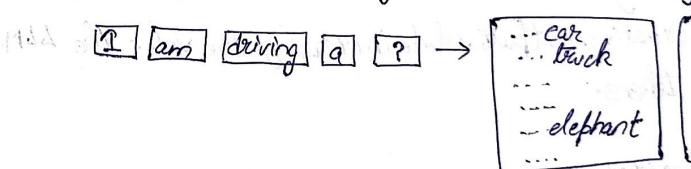
4.11 Chapter 6: Prompt engineering

Prompt engineering - The prompt engineering provides the creativity and potential components as per the user requirements such as in context learning and chain of thoughts.

- The prompt engineering uses generative pre-trained transformer (i.e gpt) for the text generation.
- There are so many models for the generating text with billion of parameters eg: Phi-3-mini which has 3.8 billion parameters.

Controlling the model output - In the prompt engineering the output can be generated according to the different response for the exact same prompt. each time an LLM needs to generate on a token.

Eg: There's a sentence given "I am driving a ?"



There are 2 important factors defined in the 5 function as

i) Temperature

ii) Top-k

- The temperature controls the randomness and creativity of the text that generated means it defines how likely it's to choose tokens that are less probable.

- It will generate the probability of any token with the same response every time bcoz it always chooses the most likely word.

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 input documents

Fig: 5-12 A bag of words count the number of times each word appears inside a document.

- In the beginning the cluster-based term frequency is calculated from the entire cluster instead of only the document for each words.

My cat is cute → class based Term Frequency (c-TF)
that is a cute dog → [1 2 1 2 1 1 1 1]
 ... → [9 6 3 18 10 5 10]
 ... → [9 1 3 2 0 5 0]
 $\text{Frequency (tf)} = \|t_{fx,c}\|$
 $\text{word}(x) \in \text{cluster}(c)$

that is a cute dog my cat

Fig: Generating c-TF by counting the frequency of words per cluster instead of per document.

Secondly, it was a cluster-based variant of term frequency - inverse document frequency (c-TF-IDF) to put more weights on words that are more meaningful to cluster, and less weights words that are used across all clusters (e.g.: stopwords like the, and, a).

The IDF is calculated to generate the weight for each word as -

$$\log \left(\frac{29}{\frac{1}{4} + 1} + 1 \right)$$

that is a cute dog my cat

Average frequency (A) divided by the total frequency of each word (C_{fx})

$$= \log \left(\frac{A}{C_{fx}} + 1 \right)$$

Fig 5-14: Creating a weighting scheme

- In the last the weighted value c-TF-IDF is calculated by multiplying the result of a weight i.e. IDF for each word with its frequency i.e. c-TF.

4. Create a cluster-based bag-of-words

Count Vectorizer

[1 2 1 2
1 6 3 18
1 1 3 2]

$\|t_{fx,c}\|$

5. Weights terms

c-TF-IDF

$$\|t_{fx,c}\| \times \frac{\text{IDF}}{\text{c-TF}}$$

Fig: The second part of BERT Topics pipeline is representing the topics, the calculation of the weights of term *x* is a cluster *c*.

24/10

Adding a special Lego Block - The pipeline in BERT topic has advantage like fast and modular but it follows some disadvantages as it still represents a topic through bag of words with they are without taking into common syntactic structure. Therefore, to speed up the process it is necessary to generate a meaningful representation. For that it will rerank the initial set of results using neural search like

original Topic (with c-TF-IDF)

summarization
summaries
summary
abstractive
document
extractive
rouge
documents
factual
evaluation

Reranker
(representation)

Ranked topic (in improved order)

summarization
summaries
abstractive
sentences
pert
metaphors
datasets
neural
model