

BERT - Paper

Language Models . LM.

A language model is a probabilistic model which lets us assign probabilities to a sequence of words.

A LM allows us to compute the following:

$$P[\text{China}][\text{"Shanghai is a city in"}]$$

Next Token Prompt

Usually a NN is trained to predict these probabilities. A NN trained on a large corpora of text is known as LLM.

How to train a LM?

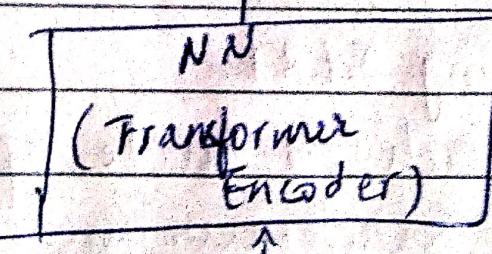
→ Target sequence → Before my bed lies a pool of
[10 tokens] moon bright [EOS]



Entropy loss



Output sequence (10 tokens) TK1 TK2 TK3 TK4 TK5 TK6 TK7 TK8
TK9 TK10



Input (10 tokens) → [sos] Before my bed lies a pool of
moon bright.

→ How to inference a language model?

Before my → Ask the LM to write
prompt first of the poem.

Output seq → Before my bed lies



NN
(Transformer Encoder)



Input seq → [sos] Before my bed

Append



Output seq → Before my bed

NN
(Transformer Encoder)



Input seq → [sos] Before my

Introducing BERT

BERT's arch is made up of layers of encoders
of the Transformer model.

• BERT's base

1) 12 encoder layers

2) ~~The~~ The size of the hidden size of the
feed forward layer is 3072

3) 12 attention heads

• BERT's large

1) 24 encoder layers

2) size of hidden size of feed forward layer
is 4096.

3) 16 Attention heads

Differences with Vanilla Transformer:

1) Embedding vector is 768 & 1024 for 2 models.

2) Positional embeddings are absolute and learnt
during training & limited to 512 positions

3) The linear layer head changes according to
the application.

4) The linear layer head changes according to
the application.

→ Used the wordpiece tokenizer, which also
allows sub-word tokens. The vocabulary
size is ~30,000 tokens.

BERT vs GPT/LLM

BERT stands for Bidirectional Encoder Representations from Transformers.

→ The importance of the left & the right context in BERT.

1. Unlike common LMs, BERT does not handle 'special tasks' with prompts, but rather, it can be specialized on a particular task by means of fine tuning.

⇒ BERT - pre-training

Masked Language Model (MLM)

Also known as Cloze task, It means that randomly selected words in a sentence are masked, and the model must predict the right word given the left & right context.

2. Unlike common LMs, BERT has been trained using the left context and the right context both.

3. " . . . ", BERT is not built specifically for text generation.

Rome is the capital of Italy, which is why it hosts many government buildings.

4. " . . . ", BERT has not been trained on the Next Token Prediction Task, but rather, on the Masked Language Model and Next Sentence Prediction task.

↓
randomly select one or more tokens and replace them with the special token [MASK]

Rome is the [MASK] of Italy, which is

Tasks in GPT/LLM vs BERT

Question Answering

in GPT/LLM

→ Prompt Engineering

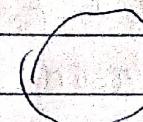
Question Answering in

BERT → Fine Tuning

Base-Trained BERT

↓

Fine Tune on QA



↓
capital

Left and right context in BERT

10% of the time it is not replaced.

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

\Rightarrow MLM \rightarrow training.

softmax $\begin{bmatrix} Q \\ (10, 768) \end{bmatrix} \times \begin{bmatrix} K^T \\ (768, 10) \end{bmatrix}$ None of the tokens is substituted by \rightarrow 0 to make it zero (0).
 $\sqrt{768}$
This is the reason it is Bidirectional Encoder

Target (1 token)

capital

loss \rightarrow Run backprop to update the weights

Each token "attends" tokens to its left and tokens to its right in a sentence.

\rightarrow Basically NO MASK.

TK1 TK2 TK3 TK4 TK5 TK6
TK7 TK8 TK9 TK10 TK11 TK12
TK13 TK14

BERT encoder

\rightarrow Masked Language Model (MLM): details

Rome is the capital of Italy, which is why it hosts many government buildings.

The pre-training procedure selects 15% of the tokens from the sentence to be masked.

When a token is selected to be masked

- 80% of the time is replaced with [MASK]

- 10% of the time it is replaced with a random token, e.g. 'zebra'

Input (14 tokens) Rome is the [mask] of Italy, which is why it hosts many government buildings

\Rightarrow Next Sentence Prediction (NSP)

Many downstream applications (e.g. choosing the right answer given 4 choices) require learning relationships between sentences rather than single tokens, that's why BERT has been pretrained on the Next Sentence Prediction Task.

Sentence A → Before my bed lies a pool of moon bright

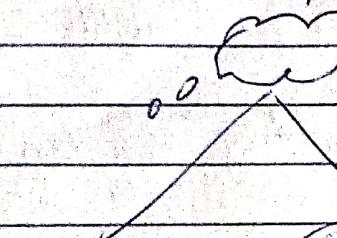
Sentence B → I could imagine that it's frost on the ground

• 50% of the time we select actual next sentence

• 50% of the time we select a random sentence from the text.

Sentence A = Before my bed lies a pool of moon bright

Sentence B = I look up and see the bright shining moon



Is Next

Not Next

⇒ Next Sentence Prediction (NSP) : segmentation embedding

Given the sentence A and the sentence B, how can BERT understand which tokens belong to the sentence A and which to the sentence B?

We introduce the segmentation embeddings!

We also introduce two special tokens [CLS] and [SEP]

Input [CLS] my dog is cute [SEP] he likes play #Hing [SEP]

Token Embeddings E_{CLS} E_m E_{dog} E_{is} E_{cute} E_(sep) E_{he} E_{likes} E_{#Hing} E_[SEP]

+ + + + + + + +

Segment Embeddings E_A E_A E_A E_A E_A E_A E_B E_C

E_A E_B E_C

+ + + + - + + +

Position Embeddings E₀ E₁ E₂ E₃ E₄ E₅ E₆ E₇

E₀ E₁ E₂

Next Sentence Prediction (NSP) ; Training Target Not Next

Linear layer (2 output features) + Softmax

Output (20 tokens) TK₁ TK₂ TK₃ TK₄ TK₅ TK₆ TK₇
TK₈ TK₉ TK₁₀ TK₁₁ TK₁₂ TK₁₃ TK₁₄

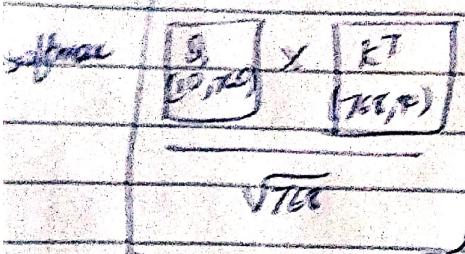
BERT encoder ← Sentence embedding

Input (20 tokens): [CLS] Before my bed lies a pool of moon bright [SEP] I look up and see the bright shining moon

[CLS] token in BERT

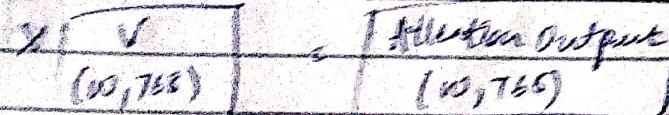
$$\text{Attention}(l, k, v) = \text{softmax}\left(\frac{QK^T}{\sqrt{dk}}\right)V$$

coming from the user such as hardware
software or billing



The [CLS] token always
interacts with all other
tokens, as we do not
use any mask.

So, we can consider the
[CLS] token as a token
that 'captures' the information
from all the other
tokens.



First Attention output \leftarrow [CLS] (1st row 1
1 column dot prod)

My wallet's also in this matter,
let's not say I don't know bill I have
nothing, I had me to bring charged Bob
bringing in just another instead of the
same wallet I tried returning until Bob says
not noticing but getting it back!
BUT WHERE SOETWERE BILLING

TRAINING

Target (1 token)

Wordpiece

Cross loss \leftrightarrow from backprop
softmax to update the
lower layers (3 output features)
+ softmax

* BERT Fine tuning
Text Classification

\rightarrow Assigning a label to a piece of text.

Eg. Imagine we are running an ISP &
receive complaints from customers.

We may want to classify requests

Output
tokens
(10 tokens)

TKE TKE TKE TKE TKE TKE TKE
TKG TKG TKG TKG TKG TKG TKG TKG

BERT Encoder

Input (10 words) [CLS] My wallet's also ...

Question Answering Task

Context

Answering ques is the task given a ~~text~~ context

Context: "Shanghai is a city in China, it is also a financial center, its fashion capital and industrial city".

Ques: "What is the fashion capital of China"

Answer: "Shanghai is a City in China, it is also a financial center, its fashion capital and industrial city".

2 problems

1. we need to find a way for BERT to understand which part of the input is the context, which one is the question.

2. we also need to find a way for BERT to tell us where the answer starts and where it ends in the context provided.

? Apply sentence A & sentence B logic
[CLS] ... [SEP] [SEP]

Sentence A
Question

Sentence B
Paragraph

