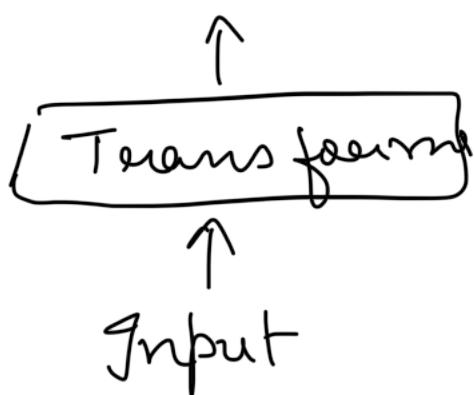


ELM - Week-3

L1 Intro to Language Modelling

- we have learnt about components of transformer architecture

Summary / Transferⁿ / Answer



- we can use this transformer to do many NLP tasks. we need to train a dataset for sep. model. And train the models from scratch.
- Parameters are randomly initialized for each task. Thus, taking long time for convergence. It is possible at many times, we do not have enough labelled data for many NLP tasks. Preparing labelled data is costly.
- But, we have many unlabelled

data across internet.

- If you ask some / any lay person based on some statement -
 - Task → statement followed by ques^m
 - but nobody was trained for this quesⁿ.
 - we develop a strong understanding of lang. w/o any explicit supervision
 - use the same idea.
- Can now my model develop basic understanding of lang. just by exposing to large amount of unlabelled raw text? (pre-training)
- How will it perform on downstream tasks with min. supervision?
(supervised fine-tuning)
- this 2-stage thing increases 92% accuracy.

12 Language Modelling

- pre-training = lang. modeling
- $V \Rightarrow$ vocab of lang. (unique words)

sentence
as a seq. x_1, \dots, x_n , $x_i \in V$
always there are some sentences
which are more probable than
others. We expect those sentences
to appear more frequently than
other sentences.

$$f: (x_1, \dots, x_n) \rightarrow [0, 1]$$

this funcⁿ is a lang. model

$$P(x_1, x_2, \dots, x_T) = P(x_1) P(x_1 | x_2)$$

$$\cdot \dots$$

$$= \prod_{i=1}^T P(x_i | x_1, \dots, x_{i-1})$$

if they are independent,

$$P(x_1, \dots, x_T) = \prod_{i=1}^T P(x_i)$$

- How do we make the model learn the lang.? Make it to learn to predict the next token in a seq.

If this idea works, then we have huge amount of such sequences to give model to work, ask it to predict $(k+1)$ th word.

- This task of predicting the next token in a sequence is called language modelling.
- But the words aren't indep. They do depend on previous words.

$$\prod_{i=1}^T P(x_i | x_1, \dots, x_{i-1})$$

it is the conditional distribution over the vocabulary.

L3 Transformers for Language Modelling

- How do we estimate these above conditional probabilities?
- Use autoregressive models where conditional probabs. are given by parameterized functions with a fixed no. of params like transformers.
- CLM (causal language model)

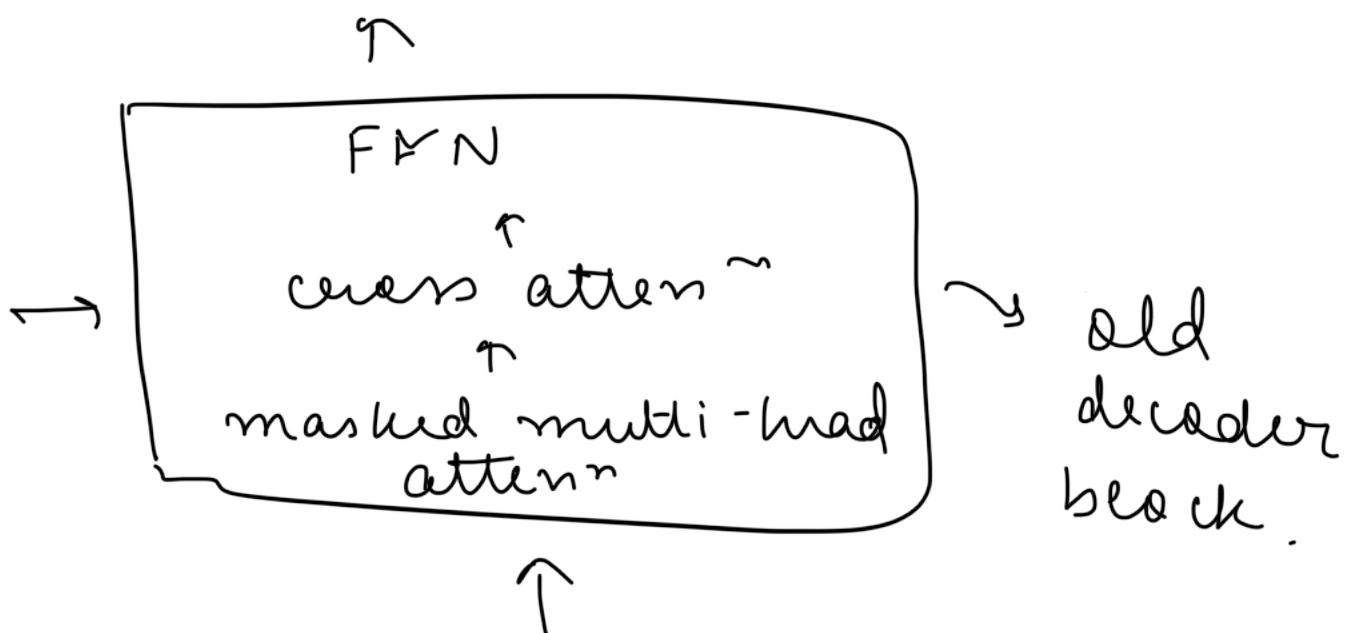
$$P(x_i | x_1, \dots, x_{i-1}) = f_\theta(x_i | x_1, \dots, x_{i-1})$$

we are looking for this f_θ .

- Can this be a transformer? Yes.
- Possibilities →
 1. encoder only → BERT
 2. decoder only → GPT
 3. encoder-decoder models. → TQ

L4 Causal Language Modelling

- we are looking into decoder only models.

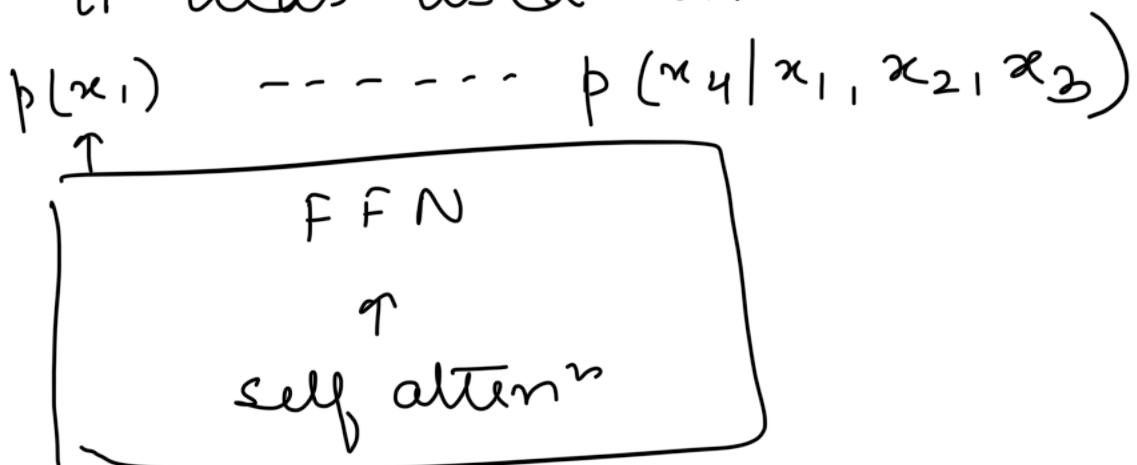


- input is a seq. of words. We want the model to see only the present & past inputs. We can achieve this by applying the task.

$$M = \begin{bmatrix} 0 & \infty & \infty & \infty \\ 0 & 0 & \infty & \infty \\ 0 & 0 & 0 & \infty \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

↳ this is neg.

But cross attenⁿ is not neg.
as it was used in enc. dec. model



$\langle q_0 \rangle \quad x_1 \quad x_2 \quad x_3$

output represent each term in
the chain rule. Probabs. are
calculated with the help of parameters
of model. The objective should be
to maximize the likelihood $L(\theta)$.

LS Gmeatice Pre-Trained Transformer

x is the input seq.

$$h_0 = X \in \mathbb{R}^{T \times d_{\text{model}}}$$

$$h_l = \text{Transformer}(h_{l-1}), \quad \forall l \in [1, n]$$

stack of such transforms (n)

$$P(x_i) = \text{softmax}(h_n[i] W_o)$$

$$L = \sum_{i=1}^T \log(P(x_i | x_1, \dots, x_n))$$

- input \rightarrow book corpus, long-range contiguous text (no shuffling of sentences)

- Byte-Pair encoding \rightarrow Tokenizer

- $V = 40478$ (vocab size)

- embedding dim = 768

- 12 decoder layers

- $T = \text{content size} = 512$

- attenⁿ heads = 12

- FFN layer size = $d_{\text{model}} \times 4$

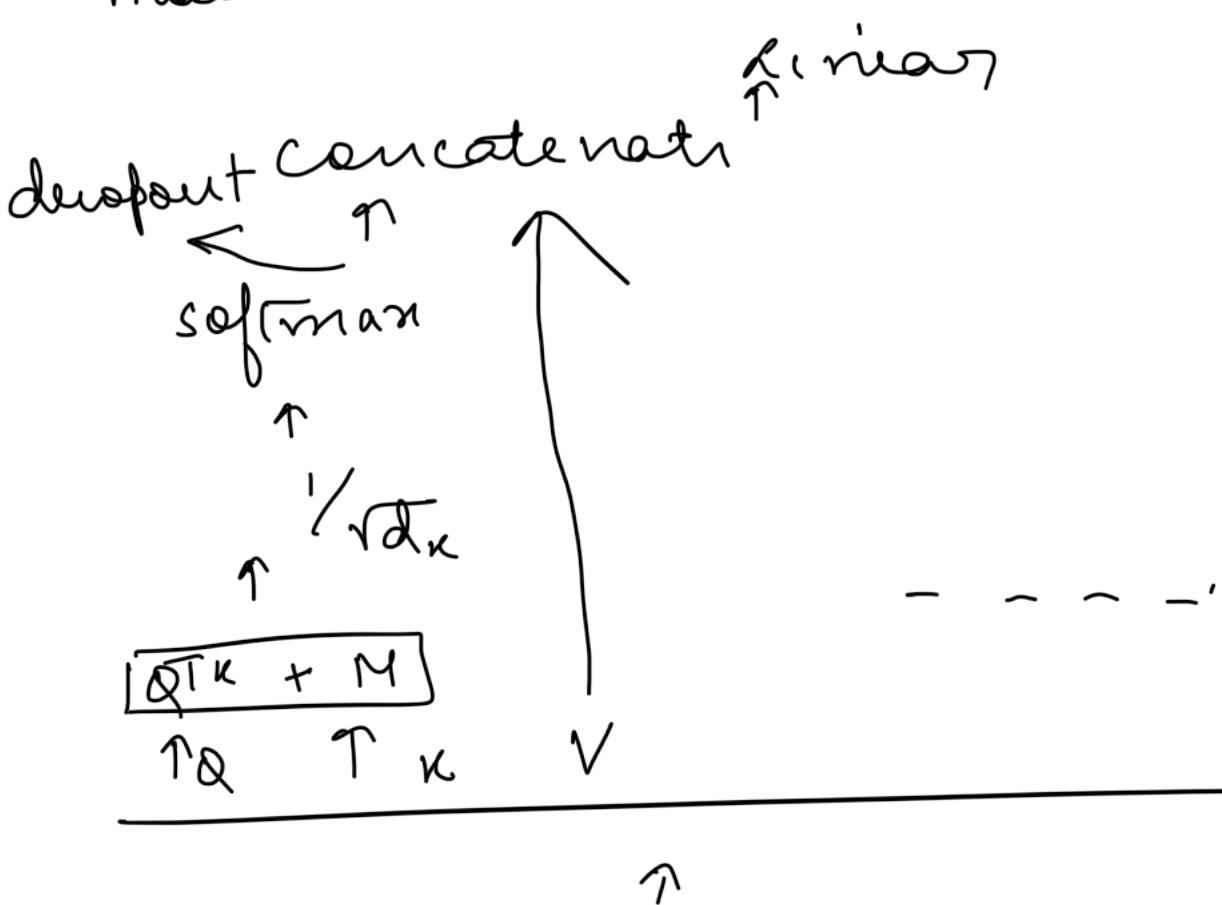
$$= 768 \times 4 = 3072$$

(64×12)

- GELU activation func → in FFN
- dropout, residual connection were used to increase the performance to reach to the convergence.

Input size $\Rightarrow T \times d_{\text{model}}$

- masked multi-head attention



- after linear layer

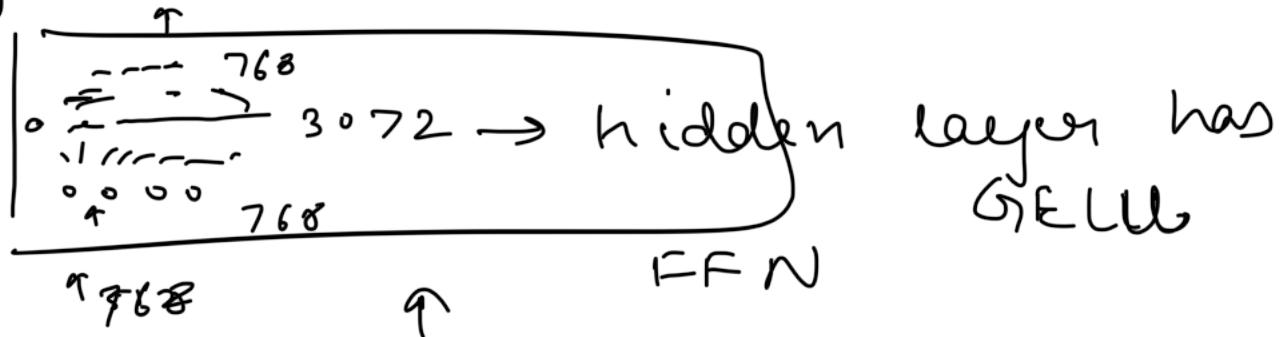
Residual \rightarrow layer Norm
connection
Linear

↳ allows the gradients to flow back.

$$H' = X + H$$

↳ skip / residual connection

- after MMHA,



$$\boxed{\text{MMHA}}$$

$$\begin{aligned} \text{- token emb.} &= V \times \text{emb-dim} \\ &= 40478 \times 768 = 31 \text{ m.} \end{aligned}$$

$$\begin{aligned} \text{- Posi"} \text{ emb.} &= \text{context len} \times \text{emb-dim} \\ &= 512 \times 768 = 0.3 \text{ m} \end{aligned}$$

31.3 m in input layer.

- the positional emb. are also learned, unlike original transformer which

uses fixed sinusoidal embeddings to encode the position.

- Attention parameters per block.

$$W_Q = W_K = W_V = (768 \times 64)$$

$$\text{Per attention} = 3 \times (768 \times 64)$$

$$12 \text{ heads} = 12 \times 3 \times 768 \times 64 \\ = 117 \text{ m}$$

$$\text{Attention} = 768 \times 768 = 0.6 \text{ m}$$

$$\text{for all 12 blocks} = 12 \times 2.3 \approx 27.6 \text{ m}$$

- FFN parameters per block.

$$2 \times (768 \times 3072) + 3072 + 768 \\ = 4.7 \text{ m} \times 12 \\ = 56.4 \text{ m}$$

$$\text{Total} = 117 \text{ m}$$

L6 Pre-Training & Fine-Tuning

- Pre-Training

$$\text{Batch size} = 64$$

$$\text{input} = (B, T, C) = (64, 512, 768)$$

↳ emb. dim.

Min.

$$L = - \sum_{x \in V} \sum_{\substack{j=1 \\ x \in \kappa_j}}^T y_j (\log(\hat{y}_j))$$

Optimizer → adam with cosine learning rate.

Strategy → Teacher forcing for quicker & stable convergence

- Fine-tuning → adapt the model for various downstream tasks (with a minimal Δ in the architecture).
- each sample of labelled data,
 $x_1, \dots, x_m \rightarrow \text{label } y$.
- initialize the params with the params learned by solving the pre-trained objective.
- at the output head, replace the pre-training LM head with a classification head (a linear layer W_y)
- Obj. is to predict the label of input seq.
 $\hat{y} = P(y | x_1, \dots, x_m) = \text{softmax}(W_y h_e^m)$

Min.

$$L = - \sum_{(x,y)} \log(\hat{g}_i)$$

- Eg. Sentiment Analysis

Text: ...

Sentiment: Positive $\hat{y} = \{+1, -1\}$

- Eg. Text Entailment / Contradiction

Text: ...

Hypothesis: ...

Entailment: True / False / Not Related

We use a delimiter token < \$ > to differentiate two sentences.

- the netwerk gets adapted as per the task.
- yes, it produces the same output seq. for a given prompt. \Rightarrow deterministic output.
- but, we don't want same text, or irrelevant text. Should be creative. Also, accelerate the gen' of tokens.