

LLM from Scratch - Elias Hossain

① Self-Attention Mechanism

Self-attention mechanism that allows a model to weigh different parts of input sequence when computing representations for each token.

Given an input sequence X of length n with each token represented as an embedding vector x_i , Self attention computes a weighted sum of all token embeddings for each token.

- Query (Q): Represents the current token
- Key (K): Represents each token's importance in the context.
- Value (V): Represents the token's actual contribution to the output.

The attention scores are computed using the Scaled-dot product attention:

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

Where,

- Q, K, V are metrics derived from the input embeddings after linear transformation.
- d_k is the dimension of the key vectors.
- Softmax normalizes the scores to sum to 1.

P.T.O

No. _____
Date _____

Mathematical steps:

1/ Compute Queries, keys and values:
 $Q = XW_q$, $K = XW_k$, $V = XW_v$

Where, W_q , W_k , W_v are learnable weight matrices.

2/ Compute attention Scores:

$$A = \frac{QK^T}{\sqrt{d_k}}$$

3/ Apply Softmax:

$$\alpha = \text{Softmax}(A)$$

4/ Compute output:

$$Z = \alpha V$$

Causal (Masked) Attention

Causal Attention is a type of Self-attention where tokens cannot attend to future tokens. This is essential for autoregressive models like GPT, ensuring that predictions only depend on past information.

Theory— This attention is implemented using a masking matrix, which sets future token attention scores to $-\infty$ before applying the Softmax activation.

Mathematical Implementation

- ① Create a lower triangular mask M :

$$M_{ij} = \begin{cases} 0 & \text{if } j \leq i \\ \infty & \text{if } j > i \end{cases}$$

- ② Modify attention scores:

$$A_{\text{masked}} = \frac{QK^T}{\sqrt{d_k}} + M$$

- ③ Apply softmax and compute the final weighted sum:

$$\alpha = \text{Softmax}(A_{\text{masked}})$$

$$Z = \alpha V$$

This ensures that each token can only attend to itself and previous tokens.

P.T.O

Example of a masking matrix:

From a sequence length of $n=5$, the masking matrix M looks like this:

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

This matrix is added to the attention score matrix before applying Softmax, it ensures that future tokens have a weight of nearly zero.

BOOM

Multi-Head Attention:

Multi-head attention (MHA) extends Self-attention by applying multiple attention mechanisms in parallel, allowing the model to capture different types of relationships.

Theory

Instead of computing a single set of Q, K, V , MHA splits the input into multiple attention heads, each with its own learned weight matrices. The outputs are then concatenated and projected back into the original embedding space.

Mathematical Formulation:

① Compute separate projections:
 $Q_i = XW_Q^i, K_i = XW_K^i, V_i = XW_V^i$

② Compute attention output:
 $Z_i = \text{Softmax} \left(\frac{Q_i K_i^T}{\sqrt{d_K}} \right) V_i$

③ Merged all heads:

$$Z = \text{Concat}(Z_1, Z_2, \dots, Z_n)$$

④ Apply final linear transformation:

$$Z_{\text{final}} = ZW_O$$

Where W_O is a learnable weight matrix for projection back to the model's embedding space.