

Transformer with GLU Variant Pseudocode

Zheling Zhang

March 6, 2024

Algorithm 1 Transformer model with Gated Linear Unit (GLU) variant

Require: x , a sequence of token IDs.

Ensure: y , the output sequence after passing through the Transformer model with a GLU variant.

- 1: Hyperparameters: N , the number of layers; d_{model} , the dimensionality of token embeddings; d_{ff} , the dimensionality of the feedforward layer; h , the number of attention heads.
 - 2: Parameters θ include all following parameters:
 - 3: $W_e \in R^{d_{\text{vocab}} \times d_{\text{model}}}$, the token embedding matrix.
 - 4: $W_p \in R^{\text{max_position} \times d_{\text{model}}}$, the positional embedding matrix.
 - 5: For each layer $l \in [1 \dots N]$:
 - 6: $W_Q^l, W_K^l, W_V^l \in R^{d_{\text{model}} \times (d_{\text{model}}/h)}$, attention parameter matrices for each head.
 - 7: $W_O^l \in R^{(d_{\text{model}}/h) \times d_{\text{model}}}$, output projection matrix for multi-head attention.
 - 8: $\gamma^l, \beta^l \in R^{d_{\text{model}}}$, parameters for layer normalization before and after the multi-head attention, respectively.
 - 9: $W_1^l \in R^{d_{\text{model}} \times d_{\text{ff}}}$, $W_2^l \in R^{d_{\text{ff}} \times d_{\text{model}}}$, weights for the feedforward network.
 - 10: $W_g^l \in R^{d_{\text{ff}} \times d_{\text{model}}}$, weights for the gating mechanism in the GLU variant.
 - 11: $b_1^l, b_2^l, b_g^l \in R^{d_{\text{model}}}$, biases for the feedforward network and the gating mechanism.
 - 12: $W_u \in R^{d_{\text{model}} \times d_{\text{vocab}}}$, the unembedding matrix.
 - 13: Initialize the output sequence y to an empty list.
 - 14: Compute the embedded input sequence $W_e^x = W_e[x] + W_p[\text{pos}]$, where pos is the position sequence.
 - 15: **for** each layer $l \in [1 \dots N]$ **do**
 - 16: Apply layer normalization: $X_{\text{norm}} = \text{LayerNorm}(E_x, \gamma^l, \beta^l)$.
 - 17: Calculate self-attention: $Z = \text{MultiHeadAttention}(X_{\text{norm}}, W_Q^l, W_K^l, W_V^l, W_O^l)$.
 - 18: Apply residual connection and layer normalization: $X_{\text{norm}} = \text{LayerNorm}(X_{\text{norm}} + Z, \gamma^l, \beta^l)$.
 - 19: Apply the first feedforward projection: $F1 = W_1^l X_{\text{norm}} + b_1^l$.
 - 20: Apply the GLU variant: $G = \text{GLU_Variant}(F1, W_g^l, b_g^l)$.
 - 21: Apply the second feedforward projection: $F2 = W_2^l G + b_2^l$.
 - 22: Apply residual connection: $E_x = X_{\text{norm}} + F2$.
 - 23: Append the result to the output sequence y .
 - 24: Compute the unembedded output: $P = \text{softmax}(W_u y)$.
 - 25: **return** P
 - 26: **function** GLU_VARIANT(X, W_g, b_g)
 - 27: Split the input matrix into two equal parts: $A, B = \text{split}(X, 2, \text{axis} = -1)$.
 - 28: Apply a non-linearity to the first part: $A = \tanh(A)$.
 - 29: Apply the gating mechanism to the second part: $B = \text{sigmoid}(W_g B + b_g)$.
 - 30: Element-wise multiplication of the results: $G = A \times B$.
 - 31: **return** G
-