

Named Tensors: Transformer

FFN

$$\begin{aligned}
 X &\in \mathbb{R}^{\text{emb}:d_{\text{model}}} \\
 W^1 &\in \mathbb{R}^{\text{emb}:d_{\text{model}}, \text{hid}:d_{\text{ff}}} & b^1 &\in \mathbb{R}^{\text{hid}:d_{\text{ff}}} \\
 W^2 &\in \mathbb{R}^{\text{hid}:d_{\text{ff}}, \text{emb}:d_{\text{model}}} & b^2 &\in \mathbb{R}^{\text{hid}:d_{\text{ff}}} \\
 \text{FFN}(X; W, b) &= W^2 \underset{\text{hid}}{\cdot} \text{ReLU}(W^1 \underset{\text{emb}}{\cdot} X + b^1) + b^2
 \end{aligned}$$

Masked Attention

$$\begin{aligned}
 Q &\in \mathbb{R}^{\text{key}:d_k, \text{seq}':n}, K \in \mathbb{R}^{\text{key}:d_k, \text{seq}:n} \\
 V &\in \mathbb{R}^{\text{seq}:n, \text{val}:d_v}, M \in \mathbb{R}^{\text{seq}, \text{seq}'} \\
 \text{att}(Q, K, V, M) &= V \underset{\text{seq}}{\cdot} \underset{\text{seq}}{\text{softmax}} \left(\frac{Q \underset{\text{key}}{\cdot} K}{\sqrt{d_k}} + M \right)
 \end{aligned}$$

Multiheaded Self Attention

$$\begin{aligned}
 W^Q &\in \mathbb{R}^{\text{head}:h, \text{emb}:d_{\text{model}}, \text{key}:d_k} \\
 W^K &\in \mathbb{R}^{\text{head}:h, \text{emb}:d_{\text{model}}, \text{key}:d_k} \\
 W^V &\in \mathbb{R}^{\text{head}:h, \text{emb}:d_{\text{model}}, \text{val}:d_k} \\
 W^O &\in \mathbb{R}^{\text{head}:h, \text{val}:d_k, \text{emb}:d_{\text{model}}} \\
 X &\in \mathbb{R}^{\text{seq}:n, \text{emb}:d_{\text{model}}} \\
 \text{MHA}(X; W) &= \left[W^O \underset{\text{head}, \text{val}}{\cdot} \text{att}(Q, K, V, M) \right]_{\text{seq}' \rightarrow \text{seq}} \\
 Q &= W^Q \underset{\text{emb}}{\cdot} [X]_{\text{seq} \rightarrow \text{seq}'} \\
 K &= W^K \underset{\text{emb}}{\cdot} X \\
 V &= W^V \underset{\text{emb}}{\cdot} X \\
 M_{\text{seq}':i, \text{seq}:j} &= \begin{cases} 0 & j \leq i \\ -\infty & \text{otherwise} \end{cases}
 \end{aligned}$$

Layer Norm

$$\begin{aligned}
 X &\in \mathbb{R}^{\text{emb}:d_{\text{model}}} & \gamma, \beta &\in \mathbb{R}^{\text{emb}:d_{\text{model}}} \\
 \text{lnorm}(X; \gamma, \beta) &= \frac{X \underset{\text{emb}}{-} \text{mean}(X)}{\sqrt{\text{var}(X) + \epsilon}} \odot \gamma + \beta
 \end{aligned}$$

Position Encoding

$$\begin{aligned}
 X &\in \{0, 1\}^{\text{seq}:n, \text{vocab}:b} & \sum_{\text{vocab}} X &= 1 \\
 E &\in \mathbb{R}^{\text{vocab} \times v, \text{emb}:d_{\text{model}}} \\
 \text{embed}(X; E) &= (E \underset{\text{vocab}}{\cdot} X) \sqrt{d_{\text{model}}} + P \\
 P &\in \mathbb{R}^{\text{seq}:n, \text{hidden}:d_{\text{model}}} \\
 P_{\text{hidden}:i, \text{seq}:p} &= \begin{cases} \sin((p-1)/10000^{(i-1)/d_{\text{model}}}) & i \text{ odd} \\ \cos((p-1)/10000^{(i-2)/d_{\text{model}}}) & i \text{ even} \end{cases}
 \end{aligned}$$

Transformer

$$\begin{aligned}
 I &\in \{0, 1\}^{\text{seq}:n, \text{vocab}:b} & \sum_{\text{vocab}} X &= 1 \\
 X^0 &= \text{embed}(I) \\
 T^1 &= \text{lnorm}(\text{MHA}(X^0)) + X^0 \\
 X^1 &= \text{lnorm}(\text{FFN}(T^1)) + T^1 \\
 &\vdots \\
 T^L &= \text{lnorm}(\text{MHA}(X^{L-1})) + X^{L-1} \\
 X^L &= \text{lnorm}(\text{FFN}(T^L)) + T^L \\
 O &= \text{softmax}(W \underset{\text{vocab}}{\cdot} \underset{\text{emb}}{X^L})
 \end{aligned}$$