

Our tensor notation is discussed in Section ?? of Bayesuvius.

TA = Transformer Architecture

ℓ = maximum number of words in a sentence segment. $\alpha \in [\ell]$, $\ell \sim 100$

L = number of words in vocabulary, $\beta \in [L]$, $L \gg \ell$

$d = d_{\underline{q}} = d_{\underline{k}} = d_{\underline{v}} = 64$, hidden dimension per head, $\delta \in [d]$.

$n_{\underline{h}} = 8$, number of heads, $\nu \in [n_{\underline{h}}]$

$D = n_{\underline{h}}d = 8(64) = 512$, hidden dimension for all heads, $\Delta \in [D]$

$\Lambda = 6$, number of layers in plate (a.k.a., stack), $\lambda \in [\Lambda]$

reshaping

$$T^{\nu, \delta} \rightarrow T^{\Delta} \quad (T^{[n_{\underline{h}}], [d]} \rightarrow T^{[D]}) \quad (1)$$

$$T^{\Delta} \rightarrow T^{\nu, \delta} \quad (T^{[D]} \rightarrow T^{[n_{\underline{h}}], [d]}) \quad (2)$$

concatenation

$$T^{[n]} = (T^0, T^1, \dots, T^{n-1}) = (T^{\nu})_{\nu \in [n]} \quad (3)$$

Hadamard product (element-wise, entry-wise multiplication)

$$T^{[n]} * S^{[n]} = (T^{\nu} S^{\nu})_{\nu \in [n]} \quad (4)$$

Matrix multiplication ($T^{[n]} = T^{[n], [1]}$ is a column vector)

$$(T^{[n]})^T S^{[n]} = \text{scalar} \quad (5)$$

$$T^{[a], [b]} S^{[b], [c]} = \left[\sum_{\beta \in [b]} T^{\alpha, \beta} S^{\beta, \gamma} \right]_{\alpha \in [a], \gamma \in [c]} \quad (6)$$

Most treatments of TA, including the Transformers is All You Need paper, order the operations chronologically from left to right. So if A occurs before B , they write AB . This is contrary to what is done in Linear Algebra, where one orders the operations chronologically from right to left, and one writes BA . We will adhere to the Linear Algebra convention, since it is so prevalent and is the overwhelming precedent.

$$x^{\beta, \alpha} = \delta(\beta, \beta(\alpha)) \quad (x^{[L], [\ell]} \text{ has one hot columns.}) \quad (7)$$

$$e^{\delta, \alpha} = \sum_{\beta} E^{\delta, \beta} x^{\beta, \alpha} \quad (e^{[d], [\ell]} = E^{[d], [L]} x^{[L], [\ell]}) \quad (8)$$

$$Q^{\nu, \delta, \alpha} = \sum_{\delta'} W_{\underline{q}}^{\nu, \delta, \delta'} e^{\delta', \alpha} \quad (Q^{[D], [\ell]} = W_{\underline{q}}^{[D], [d]} E^{[d], [\ell]}) \quad (9)$$

$$K^{\nu,\delta,\alpha} = \sum_{\delta'} W_{\underline{k}}^{\nu,\delta,\delta'} e^{\delta',\alpha} \left(K^{[D],[\ell]} = W_{\underline{k}}^{[D],[d]} E^{[d],[\ell]} \right) \quad (10)$$

$$V^{\nu,\delta,\alpha} = \sum_{\delta'} W_{\underline{v}}^{\nu,\delta,\delta'} e^{\delta',\alpha} \left(V^{[D],[\ell]} = W_{\underline{v}}^{[D],[d]} E^{[d],[\ell]} \right) \quad (11)$$

$$B^{\nu,\alpha,\alpha'} = \frac{1}{\sqrt{d}} \sum_{\delta} K^{\nu,\delta,\alpha} Q^{\nu,\delta,\alpha'} \left(B^{[n_h],[\ell],[\ell]} = \left[\frac{1}{\sqrt{d}} (K^{\nu,[d],[\ell]})^T Q^{\nu,[d],[\ell]} \right]_{\nu \in [n_h]} \right) \quad (12)$$

$$A^{\nu,\delta,\alpha'} = \sum_{\alpha} V^{\nu,\delta,\alpha} \underbrace{\text{softmax}(B^{\nu,\alpha,\alpha'})}_{P(\alpha|\nu,\alpha')} \quad (13)$$

$$\sum_{\alpha \in [\ell]} P(\alpha|\nu, \alpha') = 1 \quad (14)$$

$$A^{\nu,\delta,\alpha'} \rightarrow A^{\Delta,\alpha'} \left(A^{[n_h],[d],[\ell]} \rightarrow A^{[D],[\ell]} \right) \quad (15)$$

Column vector notation:

$$B^{\nu,\alpha,\alpha'} = \frac{1}{\sqrt{d}} (K^{\nu,[d],\alpha})^T Q^{\nu,[d],\alpha'} \quad (16)$$

Important: Note that the softmax() makes the α component a probability, not the α' one!

For example, suppose $\nu = 1$ (one head), $\ell = 2$ (a 2 word segment), and $d = 3$ (hidden dimension is 3). The $Q^{[3],[2]}$, $K^{[3],[2]}$, $V^{[3],[2]}$ are 3×2 matrices (i.e, 2 3-dim column vectors). One uses the $Q^{[3],[2]}$ and $K^{[3],[2]}$ to arrive at a 2×2 matrix $P(\alpha'|\alpha)$ of probabilities. Then one uses that matrix of probabilities to replace

$$[V^{[3],0}, V^{[3],1}] \rightarrow [V^{[3],0}P(0|0) + V^{[3],1}P(1|0), V^{[3],0}P(0|1) + V^{[3],1}P(1|1)] \quad (17)$$

• Positional Embedding (a.k.a. encoding) Matrix

$E_{pos}^{[d],[\ell]}$

$$E_{pos}^{\delta,\beta} = \begin{cases} \sin\left(\frac{\beta}{10^{4\delta/d}}\right) = \sin\left(2\pi \frac{\beta}{\lambda(\delta)}\right) & \text{if } \delta \text{ is even} \\ \cos\left(\frac{\beta}{10^{4(\delta-1)/d}}\right) = \cos\left(2\pi \frac{\beta}{\lambda(\delta)}\right) & \text{if } \delta \text{ is odd} \end{cases} \quad (18)$$

$E_{pos}^{\delta,\beta}$ changes in phase by $\pi/2$ every time δ changes by 1. Its wavelength λ is independent of β , but increases rapidly with δ , from $\lambda(\delta = 0) = 2\pi * 1$ to $\lambda(\delta = d) = 2\pi * 10^4$.

Total Embedding equals initial embedding plus positional embedding: $E = E_0 + E_{pos}$

The purpose of positional embedding is to take $x^{\beta,\alpha}$ to $e^{\delta,\alpha} = \sum_{\beta} E_{pos}^{\delta,\beta} x^{\beta,\alpha}$ where $e^{\delta,\alpha}$ changes quickly as δ (i.e., position) changes.

- **ReLU**

For a tensor T of arbitrary shape

$$ReLU(T) = (T)_+ = \max(0, T) \quad (19)$$

max element-wise

- **Feed Forward neural net**

$$F(x^{[1],[\ell]}) = ReLU(x^{[1],[\ell]} W_1^{[\ell],[d]} + b_1^{[1],[d]}) W_2^{[d],[\ell]} + b_1^{[1],[\ell]} \quad (20)$$

$$F(x^{[\ell]}) = W_2^{[\ell],[d]} ReLU(W_1^{[d],[\ell]} x^{[\ell]} + b_1^{[d]}) + b_1^{[\ell]} \quad (21)$$

- **Softmax**

$\text{softmax}()$ takes a vector and returns a vector of probabilities of the same length

$$x^{[n]} \rightarrow P^{[n]} \quad (22)$$

where

$$P^\alpha = \frac{\exp(x^\alpha)}{\sum_{\alpha \in [n]} \exp(x^\alpha)} \quad \left(P^{[n]} = \frac{\exp(x^{[n]})}{\| \exp(x^{[n]}) \|_0} \right) \quad (23)$$

For example,

$$(1, 0, 0) \rightarrow (e, 1, 1)/norm \quad (24)$$

$$(10, 0, 0) \rightarrow (e^{10}, 1, 1)/norm \approx (1, 0, 0) \quad (25)$$

For any $a \in \mathbb{R}$,

$$(a, a, a) \rightarrow (1, 1, 1)/3 \quad (26)$$

- **Skip Connection (Add & Normalize)**

A **skip connection** is when you split the input to a **filter** into two streams, one stream goes through the filter, the other doesn't. The one that doesn't is then merged with the output of the filter via a **add & normalize** node. The reason for making skip connections is that the signal exiting a filter is usually full of

jumps and kinks. By merging that filter output with some of the filter input, one smooths out the filter output to some degree. This makes back-propagation differentiation better behaved.

The filter might be a Multi-Head Attention or a Feed Forward NN.

Add & Normalize just means $(A + B)/norm$ where A and B are the two input signals and “norm” is some norm of $A + B$ (for instance, $\|A + B\|_2$).

Normalization keeps the signal from growing too big and saturating the signal entering components upstream. Normalization can also involve subtracting the mean $\langle X \rangle$ of the signal X so as to get a signal $X - \langle X \rangle$ that has zero mean.

- **Redundancy**

For better results, the Decoder contains a pair of multi-head attentions in series, and Λ of those pairs in parallel.