

Attention!

Yifeng Liu

CS 495

Preliminaries

- Definition of token
- Vector representation of token
- Softmax function

Today's Learning Objectives

Students will be able to:

- Describe the **basic mechanism** of attention in neural network
- Understand the **function** of components in scaled dot-product attention

Attention in Life



Attention in Life



Attention in context

- Try to predict the next word:
 - Paris is the capital of _____
- What are the most important words for your attention to predict?
 - Paris is the capital of (France)
- Each word(token) can be represented as a vector ($1 \times d$)
 - Hence we need to predict the vector for "France" mainly based on attentions to "Paris" and "capital"
 - But how to measure the extent of the "attention"?

Attention in context

- Inner product as attention
 - \mathbf{z} stands for the vector for “_____”
 - $\mathbf{X}=[\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5]$ stands for the vectors for “Paris is the capital of”
 - Compute the inner product of \mathbf{z} and $\mathbf{x}_1 \sim \mathbf{x}_5$ as $\mathbf{z}\mathbf{X}^T = [\mathbf{z}\mathbf{x}_1^T \dots \mathbf{z}\mathbf{x}_5^T]$
- Predict using the weighted average of given vectors (\mathbf{X})
 - $\mathbf{z}\mathbf{X}^T \cdot \mathbf{X}$
 - Normalize the weight using Softmax: $\text{Softmax}(\mathbf{z}\mathbf{X}^T) \cdot \mathbf{X}$
 - Suppose each element in \mathbf{z} and $\mathbf{X} \sim \mathcal{N}(0, 1)$
 - $\text{Var}(\mathbf{z}\mathbf{X}^T) = d$
 - Therefore, divide by \sqrt{d} to reduce the variance: $\text{Softmax}(\frac{\mathbf{z}\mathbf{X}^T}{\sqrt{d}}) \cdot \mathbf{X}$

$$\text{Softmax}(\mathbf{X})_i = \frac{e^{X_i}}{\sum_j e^{X_j}}$$

Attention in context

- How to learn with $\text{Softmax}(\frac{\mathbf{z}\mathbf{X}^T}{\sqrt{d}}) \cdot \mathbf{X}$?
 - Weighted average of given vectors by weights of inner products of given vectors
 - Weighted average of **representation** of given vectors by weights of inner products of **representation** of given vectors
 - $\mathbf{z} \rightarrow W_q \mathbf{z} = \mathbf{q}$, $\mathbf{X}^T \rightarrow W_k \mathbf{X}^T = \mathbf{K}^T$, $\mathbf{X} \rightarrow W_v \mathbf{X} = \mathbf{V}$
- $\text{Softmax}(\frac{\mathbf{q}\mathbf{K}^T}{\sqrt{d}}) \cdot \mathbf{V}$

Today's Learning Objectives

Students will be able to:

- ✓ Describe the **basic mechanism** of attention
 - $\text{Softmax}(\frac{qK^T}{\sqrt{d}}) \cdot V$
- ✓ Understand the **function** of components in scaled dot-product attention
 - To normalize the weight and control variance

Thank you very much!