道路工程深度学习技术 第七周 Transformer 原理及应用



- 1 Transformer 原理
- 2 Transformer 应用



1 Transformer 原理

处理 Sequence 数据的模型 Self-attention 结构层 Multi-head Self-attention Positional Encoding 其他结构层

② Transformer 应用



处理 Sequence 数据的模型

1 Transformer 原理

处理 Sequence 数据的模型

Multi-head Self-attention

Positional Encoding

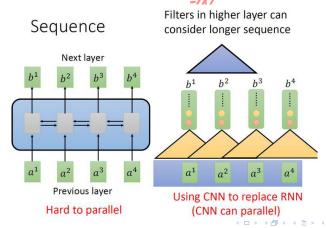
其他结构层

2 Transformer 应用



Sequence to Sequence model

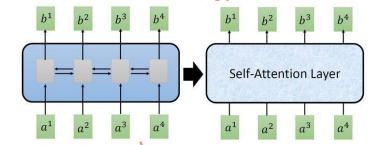
- Transformer 是一个 Sequence to Sequence model
- RNN 很不容易并行化 (hard to parallel)



处理 Sequence 数据的模型

Self-attention 结构层

- 输入和输出和 RNN 是一模一样的
- 输出 b₁-b₄ 可以并行化计算



东南大学 交诵学院

● Transformer 原理

处理 Sequence 数据的模型

Self-attention 结构层



Self-attention

https://arxiv.org/abs/1706.03762





q: query (to match others)

$$q^i = W^q a^i$$

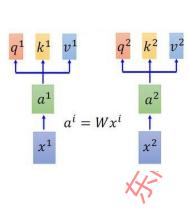
k: key (to be matched)

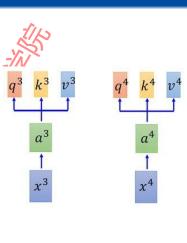
$$k^i = W^k a^i$$

v: information to be extracted

$$v^i = W^v a^i$$

Self-attention 结构层





注意力机制

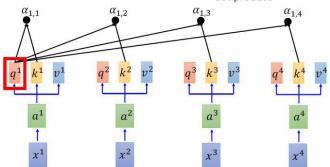


拿每個 query q 去對每個 key k 做 attention

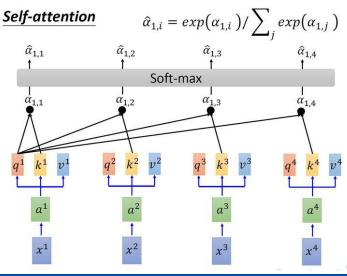
d is the dim of q and k

Scaled Dot-Product Attention: $\alpha_{1,i} = \underbrace{q^1 \cdot k^i}/\sqrt{d}$

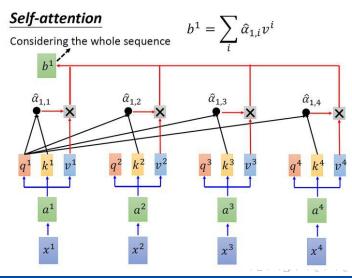
dot product



Softmax 标准化



Sequence 组合

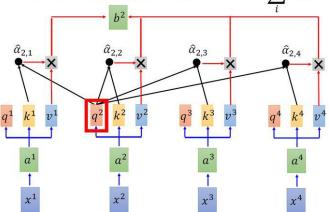


Sequence 组合

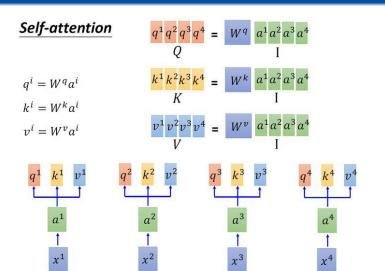
Self-attention

拿每個 query q 去對每個 key k 做 attention

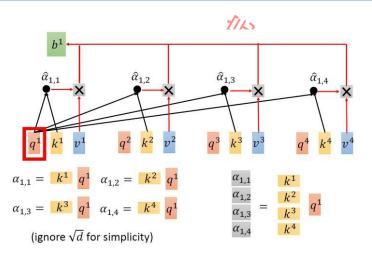
$$b^2 = \sum \hat{\alpha}_{2,i} v^i$$



Sequence 组合

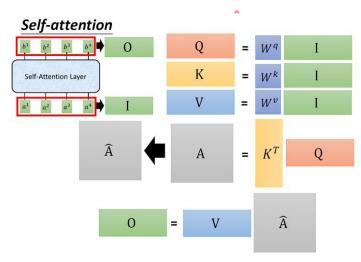


Sequence 组合流程



Sequence 组合流程

Self-attention 流程



1 Transformer 原理

处理 Sequence 数据的模型 Self-attention 结构层

Multi-head Self-attention

Positional Encoding 其他结构层

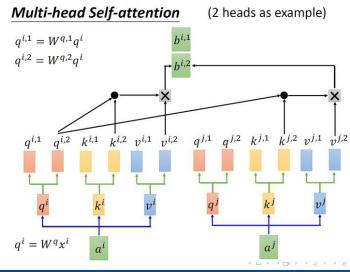
② Transformer 应用



Multi-head Self-attention

Multi-head Self-attention (2 heads as example) $q^{i,1} = W^{q,1}q^i$ $q^{i,2} = W^{q,2}q^i$ $q^{i,2}$ $k^{i,1}$ $k^{i,2}$ $v^{i,1}$ $v^{i,2}$ $q^{j,1}$ $q^{j,2}$ $k^{j,1}$ $k^{j,2}$ $v^{j,1}$ $v^{j,2}$ $q^i = W^q a^i$ a^{J}

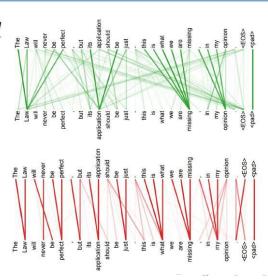
Multi-head Self-attention



Multi-head Self-attention

|Multi-head Self-attention 案例

Multi-head Attention



1 Transformer 原理

处理 Sequence 数据的模型 Self-attention 结构层

Positional Encoding

其他结构层

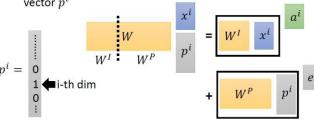
② Transformer 应用

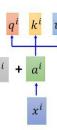


Multi-head positional encoding

Positional Encoding

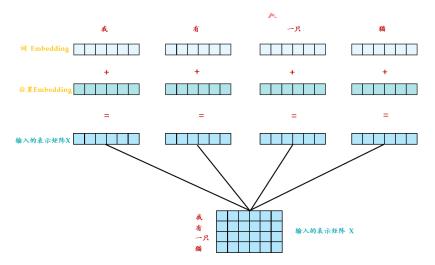
- · No position information in self-attention.
- · Original paper: each position has a unique positional vector e^i (not learned from data)
- In other words: each x^i appends a one-hot vector p^i





Positional Encoding

Multi-head positional encoding 案例



其他结构层

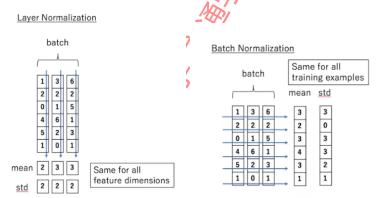
● Transformer 原理

处理 Sequence 数据的模型 其他结构层



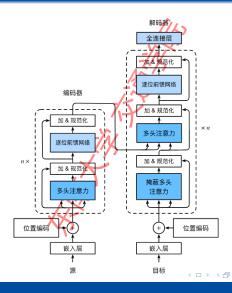
Layer normalization

- 对每一小批数据, 在批这个方向上做归一化
- LN 是每一个样本上计算均值与方差
- Batch normalization(BN)是批样本上计算均均值与方差



其他结构层

Decoder 层



^{其他结构层} 输出层

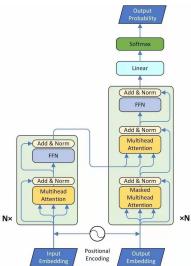
Which word in our vocabulary am is associated with this index? Get the index of the cell 5 with the highest value (argmax) log_probs 2 3 4 5 ... vocab size Softmax logits 0 1 2 3 4 5 ... vocab size Linear



Decoder stack output

其他结构层

整体模型

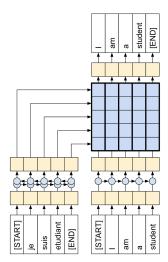




- ① Transformer 原理
- 2 Transformer 应用

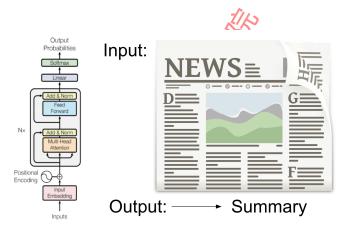


机器翻译

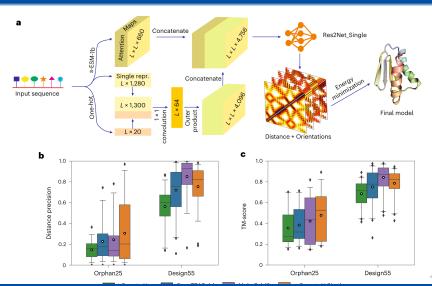




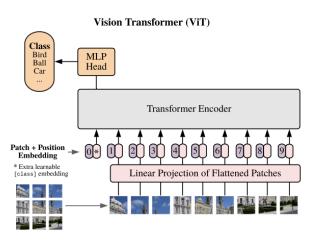
文本总结



序列数据分析

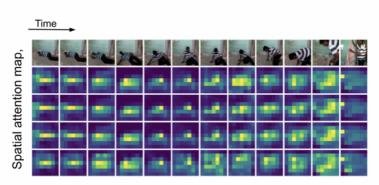


图像与视频处理



图像与视频处理





Visualization of the spatial attention maps in TokenLearner, over time. As the person is moving in the scene, TokenLearner pays attention to different spatial locations to tokenize.