

第五章 问答摘要与推理-项目模型算法提升

HCT NLP Week 5

问答摘要与推理- 项目模型算法提升

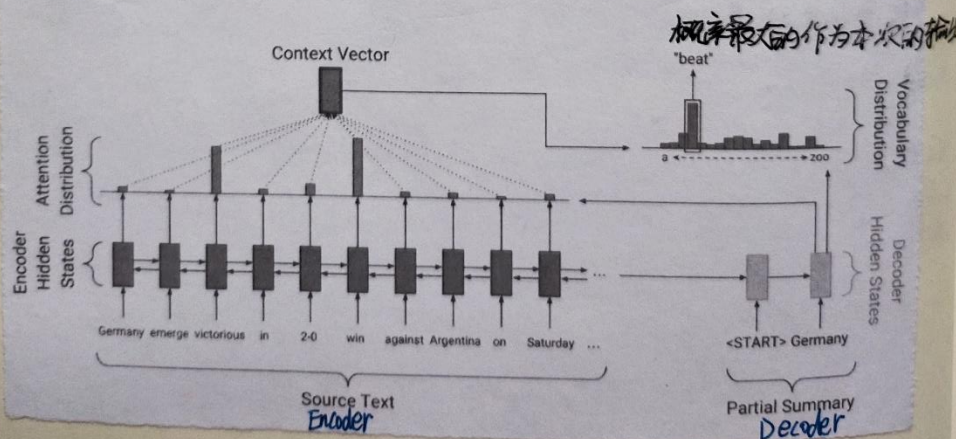
Outline

- OOV 和Word-repetition^{词语重复}解决
- Training Strategies
- 抽提式文本摘要基本方法
- 相关代码实践

Outline

- OOV 和Word-repetition解决
- Training Strategies
- 抽提式文本摘要基本方法
- 相关代码实践

典型的seq2seq

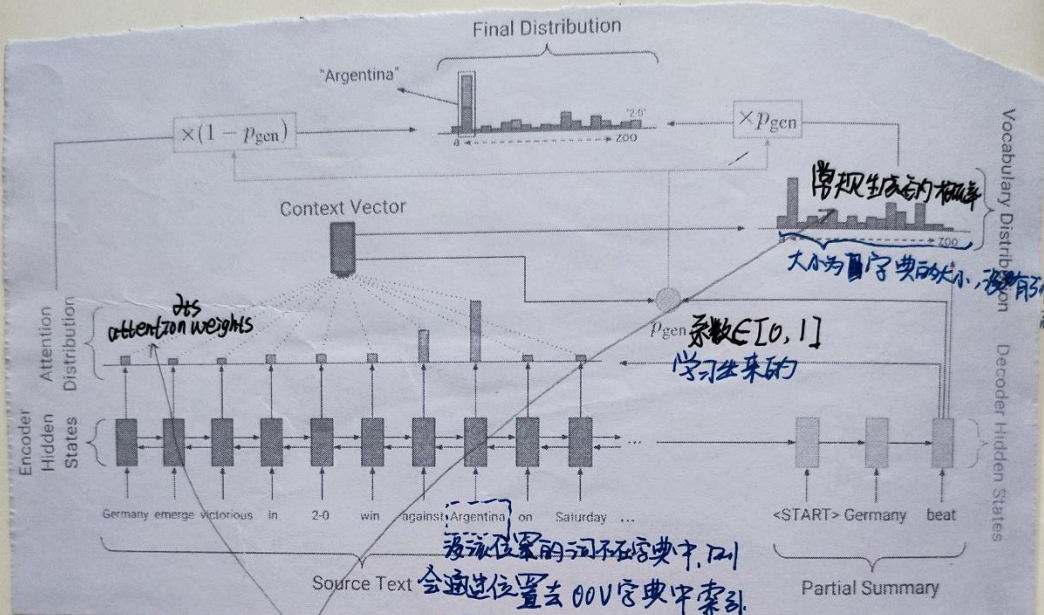


through time
从模型的路径上看



PGN (Pointer-Generator Networks)

用于解决 OOV 问题



$$p_{gen} = \sigma (w_h^T \times h_t + w_s^T \times s_t + w_x^T \times x_t + b)$$

\downarrow Sigmoid $\in (0,1)$ Context vector decoder hidden 输入的词向量 embedding 偏置

因为 w_1, w_2, w_3 也行只是写法问题
 w_h, w_s, w_x 为神经网络权重

$$p(w) = p_{gen} \cdot p_{vocab}(w) + (1 - p_{gen}) \cdot \sum_{i:w_i=w} \frac{w_i}{n}$$

\downarrow 常规生成的概率 \downarrow 即输入词对应的 attention weight

该位置的词不在字典中, 则
 会通过位置去 OOV 字典中索引
 然后代入进行计算, 避免
 该位置为 <unk>

常规生成的概率
 大小为字典的大小, 没有 OOV

p_{gen} 系数 $\in [0,1]$
 学习出来的

PGN 特点

1. pointer-generator network能够很容易的复制输入的文本内容，可以通过Pgen 来调节。
2. pointer-generator network能够从输入的文本内容中复制OOV词汇，这是最大的优点，这个也可以采用更小的词汇表vocabulary，较少计算量和存储空间。
3. pointer-generator network训练会更快，在seq2seq训练过程中用更少的迭代次数就能取得一样的效果。

Get To The Point: Summarization with Pointer-Generator Networks

Repetition Handling 解决词语重复问题

model generated summaries suffer from both word-level and sentence-level repetitions.

Temporal Attention

Intra-decoder Attention

这两种方法应用的很少

★ Coverage (机器翻译中用得最多)

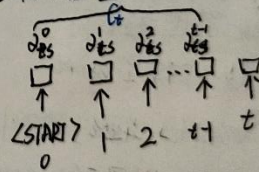
重

减少生成复词，对已生成的词做一个惩罚和

Coverage:

① $C_t^t = \sum_{s=0}^t d_{ts}$ 到 t 位置的 Coverage 等于从 0 位置到 t 位置的 attention weights 的和

其中 $C_0 = 0, C_1 = 1$; G 为一个非标准的分布



$$C_t^t = U^T \cdot \tanh(W_h h_t + W_c C_t^t + b)$$

(score) \downarrow tanh 激活函数 \downarrow encoder hidden state \downarrow decoder hidden state \downarrow coverage 偏置
 (Dense) \downarrow 权重 \downarrow 用梯度控制 \downarrow 时间 sigmoid

注：一层的全连接层为一个矩阵 [20]

trick: 在解码时加 trigram block

$$CovLoss_t = \sum \min(d_{ts}, C_t^t) \quad \text{Coverage loss}$$

$$Loss_t = -(\log p(w_t) + \lambda \sum \min(d_{ts}, C_t^t)) \quad \text{最终 loss}$$

\uparrow 固定系数值 0.5, 1

Outline

- OOV 和 Word-repetition 解决
- Training Strategies
- 抽提式文本摘要基本方法
- 相关代码实践

TRAINING STRATEGIES 训练策略

A. Word-Level Training 词语级别的训练

two different methods for avoiding the problem of exposure bias.

1) Cross-Entropy Training (XENT) 交叉熵 (loss 函数用交叉熵)

teacher forcing 方法 ~~会造成~~ exposure bias

2) Scheduled Sampling 避免 exposure bias

是一种解决训练和生成时输入数据分布不一致的方法。在训练早期该方法主要使用目标序列中的真实元素作为解码器输入，可以将模型从随机初始化的状态快速引导至一个合理的状态。随着训练的进行，该方法会逐渐更多地使用生成的元素作为解码器输入，以解决数据分布不一致的问题。该方法应用在模型的训练阶段，生成阶段不使用。

Scheduled Sampling for Sequence Prediction with Recurrent Neural Networks

TRAINING STRATEGIES

B. Sequence-Level Training 序列(句子)级别的训练

RL algorithms reinforcement learning 增强学习 (强化学习)

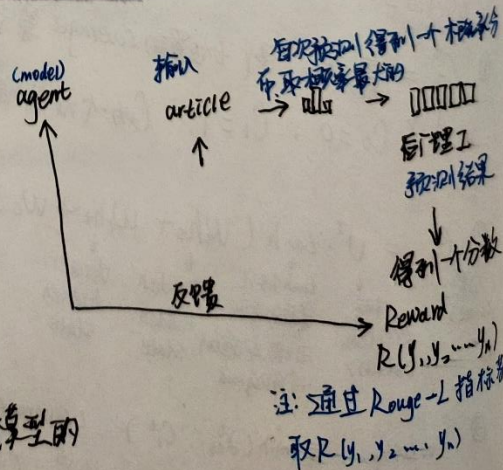
强化学习过程如下:

- ① 取一个模型 agent
- ② 初始化, 并进行预测
- ③ 用预测结果更新 agent

注: 强化学习资料:

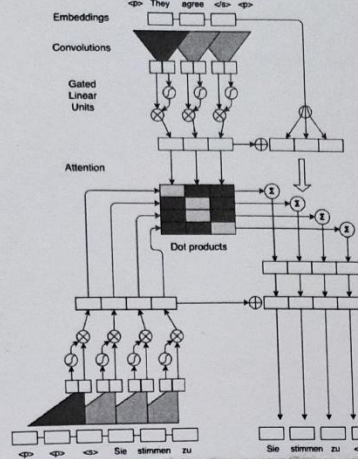
1. 网站: 周博磊

RL 不是通过反向梯度求导来更新模型的
(没有交叉熵, 没有 loss 概念)



Beyond RNN (现在用的不多) 之前在翻译网络中用的较多

CNN



Position Embedding

层叠CNN构成了hierarchical representation表示

融合了Residual connection、liner mapping的多层attention

采用GLU做为gate mechanism

进行了梯度裁剪和精细的权重初始化, 加速模型训练和收敛

Convolutional Sequence to Sequence Learning

14

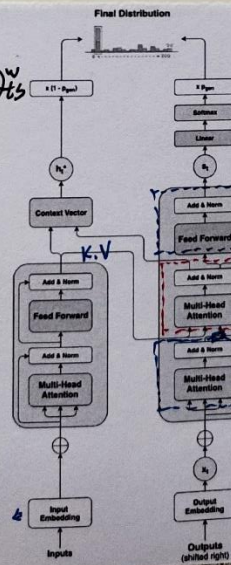
Beyond RNN为了解决GRU和LSTM不能并行计算的问题

CNN: Convolutional Neural Networks 卷积神经网络

Beyond RNN

用Transformer替换LSTM, gru等神经网络
短视频中有transformer的基础介绍

Encoder



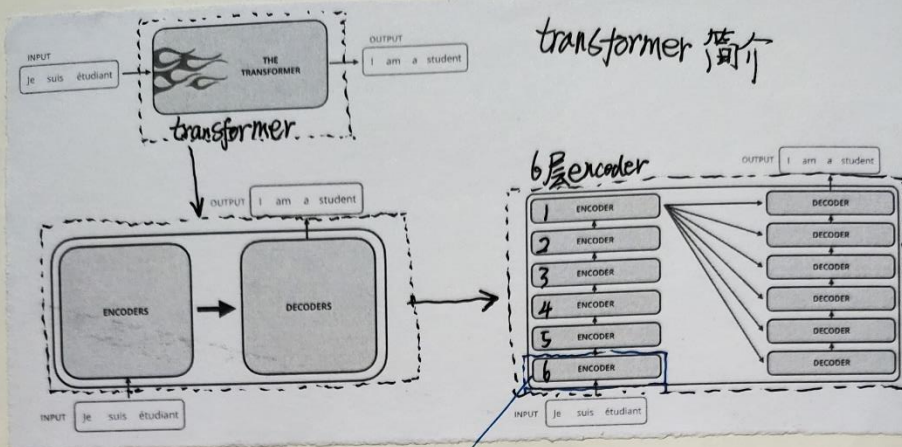
$P_{vocab} \times P_{gen}$

与Encoder中的结构相同
与Encoder中的结构有差异

Decoder

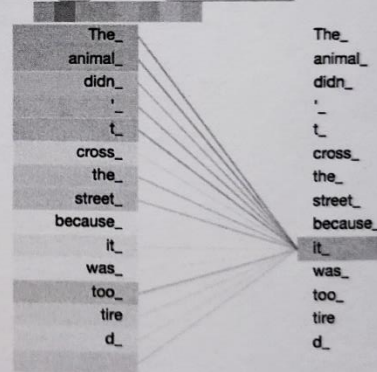
Transformers and Pointer-Generator Networks for Abstractive Summarization

15

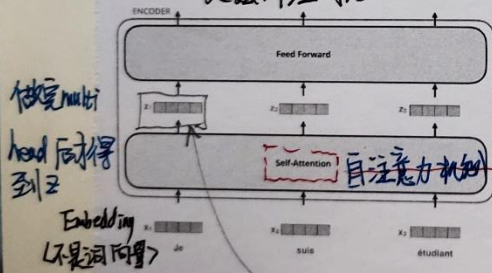


词与词之间的关系

Layer: 5 Attention: Input - Input



双层神经网络



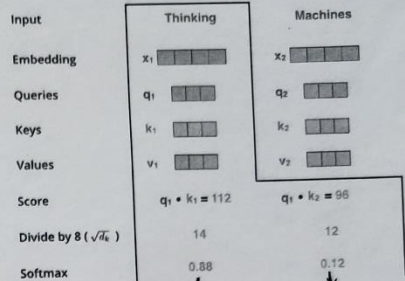
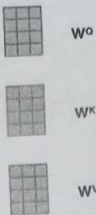
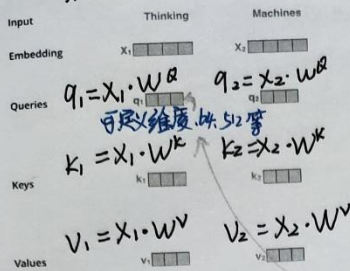
做完multi head 后得到 z1 z2 z3

Embedding 不是词向量

自注意力机制

输入: 2个词

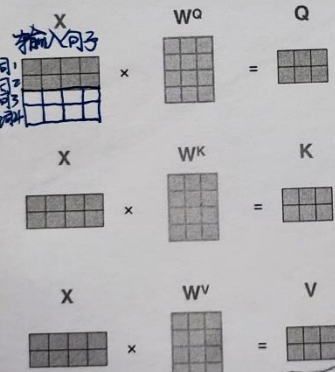
Self-Attention 计算



dim为8的
维度, 4x8

Thinking与Thinking
的关系
Thinking与
Machines的关系

由词
构成
句子



$$\text{softmax} \left(\frac{Q \cdot K^T}{\sqrt{d_k}} \right) \cdot V = Z \Rightarrow \text{Multi-headed} \Rightarrow \Sigma$$

Multi-headed

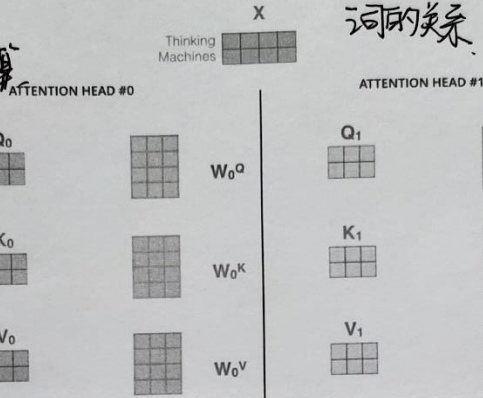
重复多次 self-attention 操作, 计算不同维度的词的关系

利用多个头进行计算

因为初始值

矩阵不同

计算结果也不同



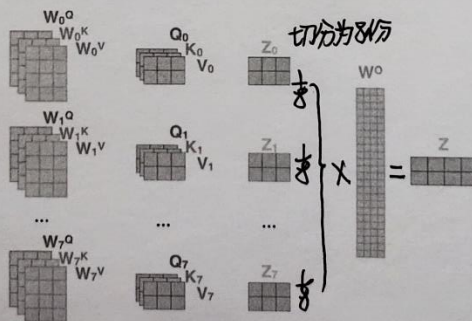
Multi-headed

- 1) This is our input sentence*
- 2) We embed each word*
- 3) Split into 8 heads. We multiply X or K with weight matrices
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer

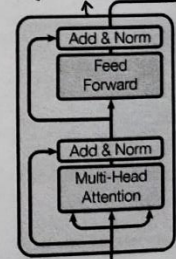
Thinking Machines

* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

R



Layer normalization



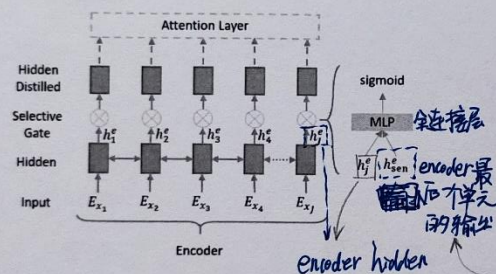
Other Studies

- 1) Network Structure and Attention
- 2) Extraction + Abstraction
- 3) Long Documents

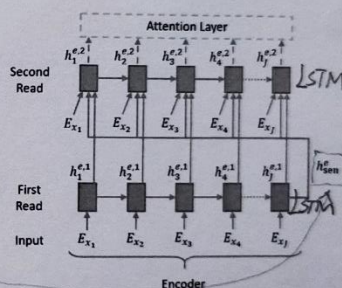
Improving Encoded Representations

对Encoder改进

Selective Encoding



Read-Again Encoding



Improving Decoder 对Decoder的改进思想

Embedding Weight Sharing 共享权重

Extraction + Abstraction

Extractor + Pointer-Generator Network

Key-Information Guide Network (KIGN)

Reinforce-Selected Sentence Rewriting

SUMMARY GENERATION

Diverse Beam Decoding

the top-B hypotheses may differ by just a couple tokens at the end of sequences, which not only affects the quality of generated sequences but also wastes computational resources

Outline

- OOV 和Word-repetition解决
- Training Strategies
- 抽提式文本摘要基本方法
- 相关代码实践

Text summarization

Extractive text summarization 抽提式文本摘要

Source Text: Peter and Elizabeth took a taxi to attend the night party in the city.

从原文中抽取

While in the party, Elizabeth collapsed and was rushed to the hospital.

Summary: Peter and Elizabeth attend party city. Elizabeth rushed hospital.

Abstractive text summarization

Source Text: Peter and Elizabeth took a taxi to attend the night party in the city.

While in the party, Elizabeth collapsed and was rushed to the hospital.

Summary: Elizabeth was hospitalized after attending a party with Peter.

Extractive Text Summarization

houchangtech.com

抽取式文本摘要步骤:

对词进行表征: 如词向量, one-hot, tfidf

① 输入文本表征的构建

topic

representation

方法 { Frequency-driven Approaches
Latent Semantic Analysis
Bayesian Topic Models

② 根据表征给句子打分

indicator

representation

方法 { Graph Methods
Machine Learning

③ 根据一定数量的句子组合

Text Summarization Techniques: A Brief Survey

23

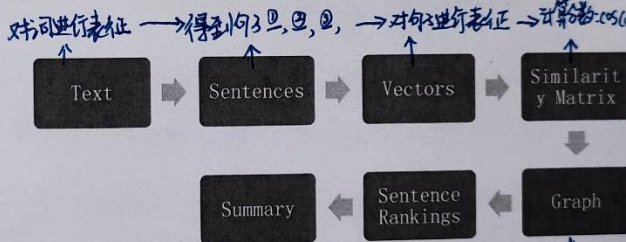
Graph Methods

TextRank算法

jupyter 1ecwe-5

有此图用代码

传统机器学习方法

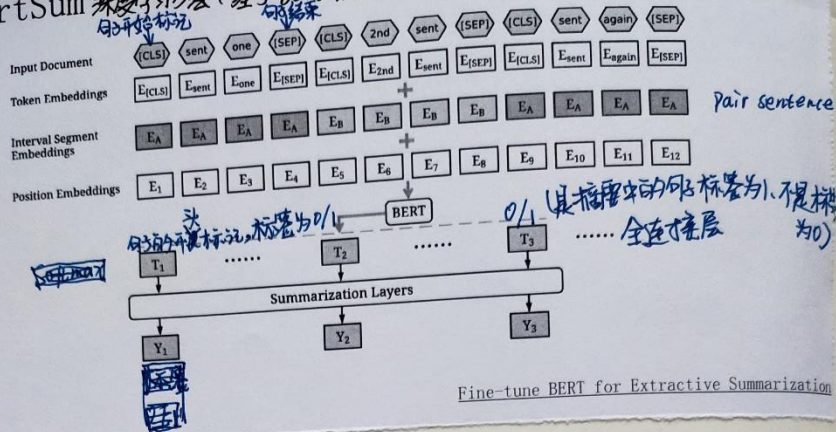


1. 第一步是把所有文章整合成文本数据
2. 接下来把文本分割成单个句子
3. 然后, 我们将为每个句子找到向量表示 (词向量)
4. 计算句子向量间的相似性并存放在矩阵中

5. 然后将相似矩阵转换为以句子为节点、相似性得分为边的图结构, 用于句子TextRank计算



BertSum 预训练方法 (基于 bert 改造的)



注: (1) bert是无监督学习

(2) bert是自编码结构

(3) bert是一个字一个字切分的 (切分的单位为字, 今天星期几-1)

(4) bert 模型缺点: 模型大

bert在训练时, 15%的词会被mask (掩盖) 掉, 其中 (这15%中) 的80%会完全被mask (类似空型), 10%的给错误词, 10%的给正确词

bert就是很多层 transformer 的编码器, 做相邻语句的断任务, 实现无监督训练

给定上一个词, 预测下一个词为自回归问题, bert为自编码问题

bert只有encoder, 没有decoder.

Outline

- OOV 和 Word-repetition 解决
- Training Strategies
- 抽提式文本摘要基本方法
- 相关代码实践

附：

- (1) transformer 模型 github: <https://github.com/huggingface/transformers>
- (2) 文章: BERT-Pre-training of Deep Bidirectional Transformers for Language Understanding:
<https://arxiv.org/pdf/1810.04805.pdf>
(翻译博文: https://blog.csdn.net/sinat_33741547/article/details/86311310)
- (3) bert 模型 github: <https://github.com/jihun-hong/Bert-Classifier/tree/master/src/models>