# HCT NLP Week 5

问答摘要与推理-
项目模型算法提升

| Outline | • OOV 和 Word-repetition 解决 (词语重复) <br> • Training Strategies <br> • 抽提式文本摘要基本方法 <br> • 相关代码实践 |
|---|---|
| Outline | • OOV 和 Word-repetition 解决 <br> • Training Strategies <br> • 抽提式文本摘要基本方法 <br> • 相关代码实践 |

典型的seq2seq

# 问题

The encoder is not well trained via back propagation through time.

从模型的路径上看，encoder到实际输出有一定距离，从此限制了反向传播。

## OOV（Out-of-vocabulary未登录词 ）

摘要总结的结果有的时候并不准确，比如摘要的结果可能输出德国队以2-1比分击败阿根廷，但是实际比分是2-0，出现这个的原因是out-of-vocabulary words（OOV）的出现

## Word-repetition问题 词语级别的重复（不是句子重复）

摘要结果会出现repeat重复的信息，比如重复出现德国队击败阿根廷队



CopyNet

对S4单元的输入和输出进行改进

生成概率

生成的    copy的

copy部分的概率

生成部分

(b) Generate-Mode & Copy-Mode

Prob("Jebara")=Prob("Jebara", g)+Prob("Jebara", c)

Softmax

Vocabulary    copy部分    Source

对输出进行改进

对输入进行改进

Decoder

hi    Tony Jebara

$s_1$ $s_2$ $s_3$ $s_4$

<eos>  hi    Tony

Attentive Read

$S_4$    神经网络    DNN    M

Embedding for "Tony"

Selective Read for "Tony"

Encoder

$h_1$ $h_2$ $h_3$ $h_4$ $h_5$ $h_6$ $h_7$ $h_8$

"Tony"

为输入单元h8的hidden state

M

hello  ,  my  name  is  Tony Jebara  .

(a) Attention-based Encoder-Decoder (RNNSearch)

相对于Attentive Read

乘数

输入词的one-hot形式

(c) State Update

Incorporating Copying Mechanism in Sequence-to-Sequence Learning

输入词的one-hot形式乘以系数1，再与h8的hidden state 相加

# PGN ( Pointer - Generator ~~Attent~~ Networks)

用于解决 OOV 问题



**Final Distribution**

"Argentina"

$\times (1 - p_{gen})$   $\times p_{gen}$

常规生成的概率

Context Vector

大小为整个字典的大小, 没有OOV

Vocabulary Distribution

Attention Distribution

αᵗₛ attention weights

$p_{gen}$ 系数 ∈ [0,1]

学习出来的

Encoder Hidden States

Decoder Hidden States

Germany emerge victorious in 2-0 win against Argentina on Saturday ...

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>

<START> Germany beat

Source Text

Partial Summary

$$P_{gen} = \sigma ( W_h^T \times h_t + W_s^T \times S_t + W_x^T \times x_t + b )$$

- Sigmoid ∈ (0,1)
- Context vector
- decoder hidden
- 偏置
- 输入词的词向量 embedding

因为 $W_1, W_2, W_3$ 也行, 只是写法问题
$W_h^T, W_s^T, W_x^T$ 为神经网络的权重.

$$P_{(w)} = P_{gen} \cdot P_{vocab}(w) + (1 - P_{gen}) \sum_{i:w_i=w} a_{ts}^w$$

- 正常规模的字典的概率
- 即输入词对应的 attention weight

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 解决词语级复重复问题

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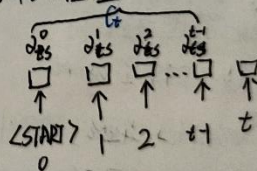
*Temporal Attention*

*Intra-decoder Attention*

} 这两种方法之间的很少

☆*Coverage*(机器翻译中用得多) 减少生成复词, 对已生成的词做一个惩罚 和重罚

**Coverage :**

① $C_i^t = \sum\limits_{t=0}^{t-1} a_{ts}$  到 $t$ 位置的 coverage 等于从 0 位置到 $t-1$ 位置的 attention weights 的 总和

其中 $C_0 = 0$ , $C_1 = 1$;  $C_i$ 为一个非标准的分布

② $C_i^t = U^T \cdot \tanh ( U_s h_s + U_t h_t + W_c C_i^t + b)$

分数 (score)
权重 连接层 Dense()
tanh 激活 表现力在 用概率控制到 时用sigmod
encoder hidden state
decoder hidden state
coverage 偏置

注: 一层的全连接层为一个矩阵 [图]

trick: 在预测时加 tri-gram block

③ $CovLoss_t = \sum min(a_{ts}^t, C_i^t)$  ( coverage loss )

$Loss_t = -\log P(w_t) + \lambda \sum min ( a_{ts}, C_i^t )$  ( 最终 loss )

↑ 自由调节的系数值,如 0.5, 1

## 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函数采用交叉熵)

teacher forcing 方法 ~~直接~~ 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

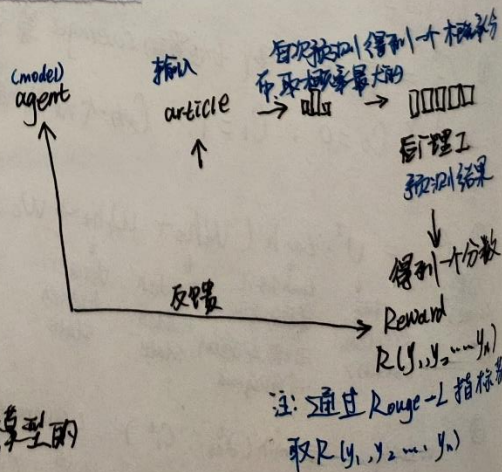**RL algorithms** reinforcement learning 增强学习 (强化学习)

强化学习过程如下：
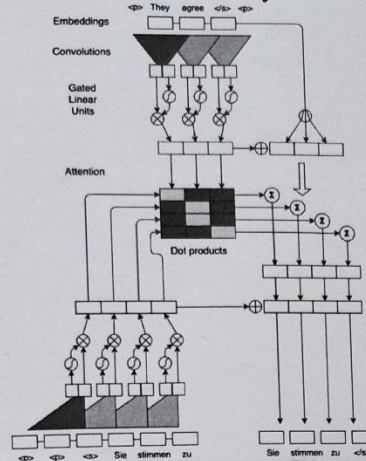
① 取一个模型 agent

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：

b站：同博客

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）

(model)
agent

抽取
article

自己预测得到一个概率
布再取概率最大的

后处理工
预测结果

得到一个分数
Reward
$R(y_1, y_2, \cdots y_n)$

反馈

注：通过 Rouge-L 指标数

取 $R(y_1, y_2, \cdots y_n)$

# Beyond RNN （现在用的少了）（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
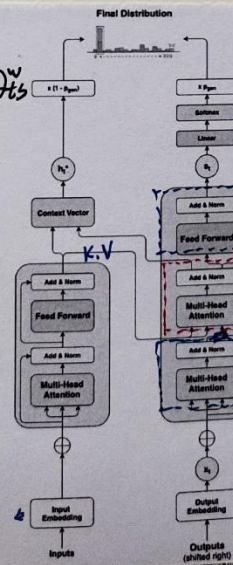
Convolutional Sequence to Sequence Learning

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

$(1-P_{gen})\sum_i a_i^t$  $P_{vocab} \times P_{gen}$

用Transformer替换
LSTM、gru等神经网络

注视频中有transformer
的基础介绍



与Encoder中的结构相同

与Encoder中的结构有差异

Decoder

Q

K、V

Encoder

此处的embedding不
是词向量

Transformers and Pointer-Generator Networks for Abstractive Summarization

transformer 简介

6层encoder

渐5词之间的关系

双层神经网络

做完multi
head 后得
到Z

自注意力机制

Embedding
（不是词向量）

Layer: 5 Attention: Input - Input

The_
animal_
didn_
'_
t_
cross_
the_
street_
because_
it_
was_
too_
tire
d_

The_
animal_
didn_
'_
t_
cross_
the_
street_
because_
it_
was_
too_
tire
d_

由
构
句

输入：2个词          Self-Attention 计算

Input

Embedding          Thinking $x_1$          Machines $x_2$

Queries          $q_1 = x_1 \cdot W^Q$          $q_2 = x_2 \cdot W^Q$          $W^Q$

可定义维度、64、512等

Keys          $k_1 = x_1 \cdot W^K$          $k_2 = x_2 \cdot W^K$          $W^K$

Values          $v_1 = x_1 \cdot W^V$          $v_2 = x_2 \cdot W^V$          $W^V$

$d_k$为定义的维度、64维

| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ($\sqrt{d_k}$) | 14 | 12 |
| Softmax | 0.88 | 0.12 |

Thinking与Thinking 的关系          Thinking与 Machine的关系

---

输入句子
由词构成句子

$X$          $W^Q$          $Q$

$X$          $W^K$          $K$          $\Rightarrow$          $softmax\left( \dfrac{Q \times K^T}{\sqrt{d_k}} \right)$          $V$

$X$          $W^V$          $V$

$= Z \Rightarrow Multi\text{-}headed \Rightarrow \boxed{Z}$

# Multi-headed

重复多次 self-attention 操作，计算一个词与个词的关系。

利用多个部分进行计算
因为初始化
矩阵不同，
计算结果也不同



ATTENTION HEAD #0

$Q_0$  $W_0^Q$

$K_0$  $W_0^K$

$V_0$  $W_0^V$

ATTENTION HEAD #1

$Q_1$  $W_1^Q$

$K_1$  $W_1^K$

$V_1$  $W_1^V$

Thinking Machines — X

---

# Multi-headed

1) This is our input sentence*

2) We embed each word*

3) Split into 8 heads. We multiply X or R 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

layer normalization

Thinking Machines — X

$W_0^Q$ $W_0^K$ $W_0^V$

$Q_0$ $K_0$ $V_0$

$Z_0$

切分为8份

$W^O$

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

$W_1^Q$ $W_1^K$ $W_1^V$

$Q_1$ $K_1$ $V_1$

$Z_1$  × = Z
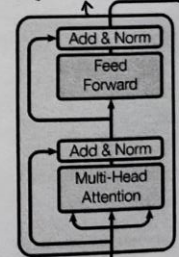
R

...

$W_7^Q$ $W_7^K$ $W_7^V$

...

$Q_7$ $K_7$ $V_7$

...

$Z_7$
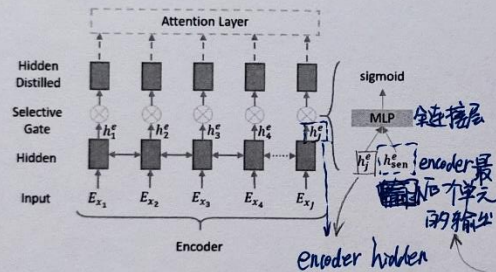
Add & Norm
Feed Forward
Add & Norm
Multi-Head Attention

# Other Studies

1) Network Structure and Attention

2) Extraction + Abstraction

3) Long Documents

# Improving Encoded Representations

对于Encoder的改进

*Selective Encoding*                                    *Read-Again Encoding*

链接层
encoder最后一个单元的输出
encoder hidden

# 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 <u>hospital.</u>

### 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.

# xtractive Text Summarization

抽提式文本摘要步骤:

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

对句子进行表征

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

② 根据表征给句子打分

indicator representation

抽取摘要

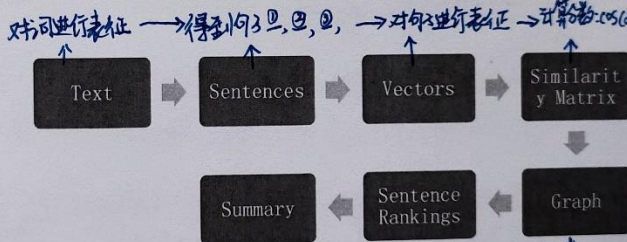方法 { Graph Methods → 句的构
Machine Learning

③ 根据一定数量的句子组合

Text Summarization Techniques: A Brief Survey

23

---

# Graph Methods

(此方法适合长文章)

TextRank算法

jupyter lecture-5

有此图圃代码

传统机器学习方法

对词进行表征 → 得到句子①,②,③ → 对句子进行表征 → 计算句:cos(①,②)分

①,②,③表示句子1,②③.

1. 2. 3
1 0.50.3 x
2 x x x
3 x x x x

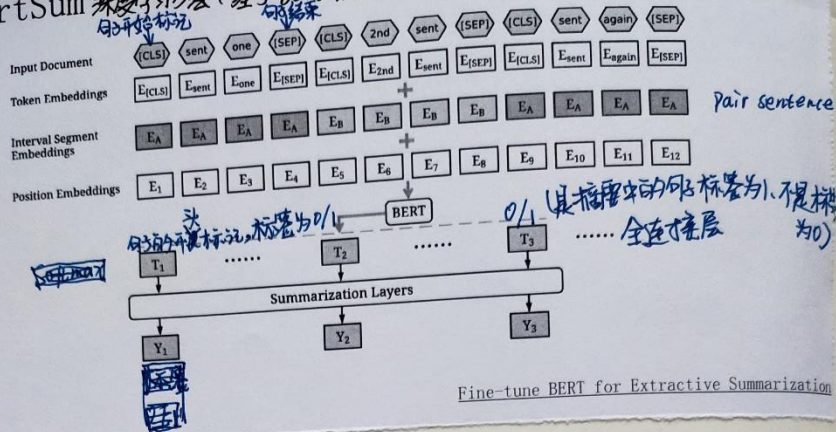Text ⇒ Sentences ⇒ Vectors ⇒ Similarity Matrix

Summary ⇐ Sentence Rankings ⇐ Graph

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

依据 PageRank�阳 迭查拣的

关辖话

最终,谁最重要分数最高

**BertSum** 深度学习方法（基于 bert 改造的）

Input Document: [CLS] sent one [SEP] [CLS] 2nd sent [SEP] [CLS] sent again [SEP]

Token Embeddings: $E_{[CLS]}$ $E_{sent}$ $E_{one}$ $E_{[SEP]}$ $E_{[CLS]}$ $E_{2nd}$ $E_{sent}$ $E_{[SEP]}$ $E_{[CLS]}$ $E_{sent}$ $E_{again}$ $E_{[SEP]}$

Interval Segment Embeddings: $E_A$ $E_A$ $E_A$ $E_A$ $E_B$ $E_B$ $E_B$ $E_B$ $E_A$ $E_A$ $E_A$ $E_A$    Pair sentence

Position Embeddings: $E_1$ $E_2$ $E_3$ $E_4$ $E_5$ $E_6$ $E_7$ $E_8$ $E_9$ $E_{10}$ $E_{11}$ $E_{12}$

BERT

$T_1$ …… $T_2$ …… $T_3$ …… 全连接层

Summarization Layers

$Y_1$ $Y_2$ $Y_3$

Fine-tune BERT for Extractive Summarization

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？令头星其 | -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（类似完形填空）, 10% 的给错误词, 10% 的给正确词

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

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, 而没有 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