

### 第三章 问答摘要与推理-Seq2Seq (二)

#### HCT NLP Week 4

#### 问答摘要与推理 Seq2Seq (二)

##### Outline

- 深度学习框架图计算理论
- Beam search
- 生成式文本摘要问题补充
- Baseline代码实践

##### Outline

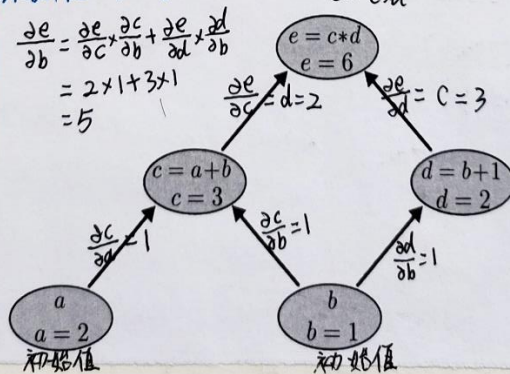
- 深度学习框架图计算理论
- Beam search
- 生成式文本摘要问题补充
- Baseline代码实践

#### 图计算理论

Computational Graphs

计算方向是从下往上

$$\begin{aligned}\frac{\partial e}{\partial b} &= \frac{\partial e}{\partial c} \times \frac{\partial c}{\partial b} + \frac{\partial e}{\partial d} \times \frac{\partial d}{\partial b} \\ &= 2 \times 1 + 3 \times 1 \\ &= 5\end{aligned}$$



$$e = (a+b) \times (b+1)$$

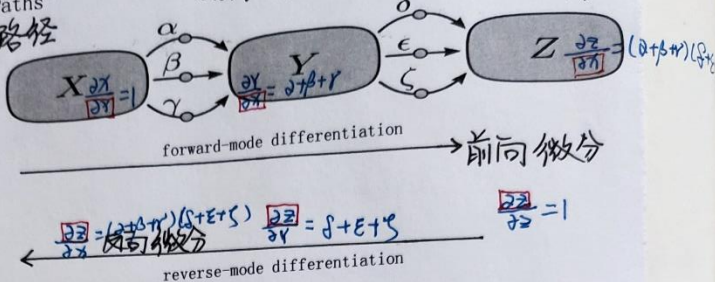
拆解:  $c = a+b$ ,  $d = b+1$   
 $e = c \times d$

# 图计算理论

Factoring Paths  
因式分解路径

$$\frac{\partial z}{\partial x} = \alpha \delta + \alpha \epsilon + \alpha \zeta + \beta \delta + \beta \epsilon + \beta \zeta + \gamma \delta + \gamma \epsilon + \gamma \zeta$$

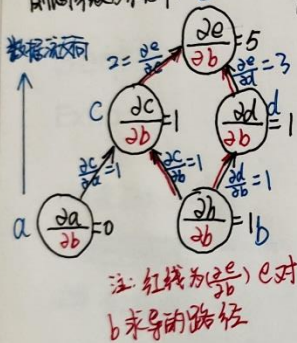
$$= (\alpha + \beta + \gamma) \times (\delta + \epsilon + \zeta) \text{ 有9条路径}$$



微分求导的一种方式

前向微分和反向微分是两种方式构建图，

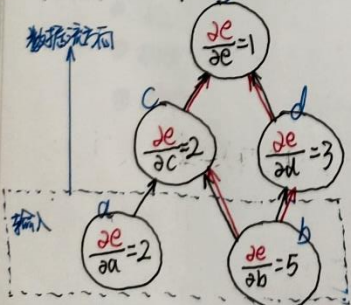
前向微分示例:



$$e = c \times d = (a+b) \times (b+1)$$

1. 前向微分时, 分母不变 (如都为  $\frac{\partial}{\partial b}$ ), 反向微分时  $\delta$  不变 (如都为  $\frac{\partial e}{\partial e}$ )
2. 前向微分, 一次只能对一个参数求导 (如  $b$ ), 反向微分一次构建图能对所有参数求导
3. 前向微分, 一个输入影响每个节点, 反向微分, 每个节点影响最终输出。

反向微分示例:



在构建图前, 已经计算出每个节点的值 (记该函数为最基函数)

注: tensorflow, pytorch 都是采用反向微分

若有100万个参数, 正向(前向)微分要计算100万次才能算出对所有参数的微分, 反向微分只需要计算1次就能算出对所有参数的微分, (则反向微分时的时间为正向的  $\frac{1}{100万}$ )



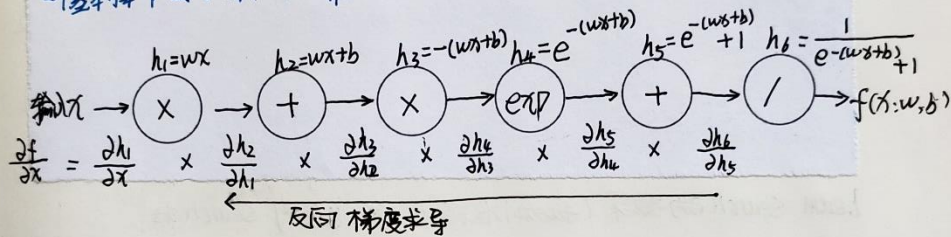
设有:  $z = \sigma(wx + b)$ .  $y = wx + b$ . 神经网络基本流程: 输入  $(x, y) \rightarrow \hat{y} \rightarrow \text{loss}$  或  $L(y, \hat{y})$

反向梯度求导:  $\frac{\partial L}{\partial w} = \frac{\partial z}{\partial w} \times \frac{\partial L}{\partial z}$   $\frac{\partial L}{\partial b} = \frac{\partial z}{\partial b} \times \frac{\partial L}{\partial z}$

## 图计算理论

$$f(x; w, b) = \frac{1}{\exp(-(wx + b)) + 1}$$

一个基本操作 (如加、减) 为一个算子



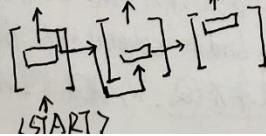
## Outline

- 深度学习框架图计算理论
- Beam search
- 生成式文本摘要问题补充
- Baseline 代码实践

greedy search algorithm 贪婪搜索算: 即在每步运算中只输出

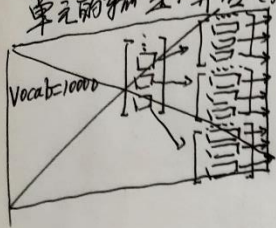
概率最大的那个值, 并代入下一步进行运算

缺点: 容易生视局部最优 (整体不是最优)

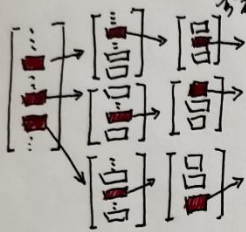


每步预测都只输出概率最大的一个

Beam search algorithm: 即将概率最大的几个 (如 3 个) 做为本单元的输出, 并传入下一个单元, 做为下一个单元的输出.



- ① beam search 搜索中, 若 beam-size = 1, 则 beam search 变为 greedy search
- ② beam search 是在广度优先下进行搜索空间的优化
- ③ beam search 比 greedy search 慢 (因为代入下一步的步数多)



设 beam-size = 3, 则每层选取 3 个概率最大的做为输出. 假设同表大小为 10000, 则每层要做 30000 次搜索.

设  $B = \text{beam-size}$ , 词典大小为  $V$ , 序列长度为  $L$ , 则 beam search 的算法复杂度为:  $O(B \cdot V \cdot L)$

beam search 的效果 (输出准确度) 比 greedy search 好

beam search 中  $B$  越大, 则输出选择越多 (输出多样性越多), 但计算的时间也越长, 同时消耗的内存也越大

beam search 的缺点: 多样性不明显 (输出的差别较小)

为了解决 beam search 的缺点, 可使用 **Diverse Beam Search**

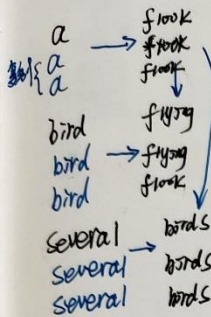
Diverse Beam Search 原理如下: 设 beam-size = 3



- ① 假设上一步输出为  $a$ , bird, several.
- ② 下一步计算时, 假设  $a$  作为输入时 flock 的概率最大, 则 bird 和 several 作为输入时, 在计算 flock 时进行一定的惩罚 (如减 1 分), 使  $a$  作为输入时 flock 作为输出的概率与 bird 和 several 作为输入 flock 作为输出的概率不同, 以此来增加输出的多样性

借鉴 diverse beam search 增加多样性的思想, 通过复制多份输入(上一步输出)的方式同样可以增加输出多样性

在代码实现中



复制后再传入下一步进行计算

注: diverse beam search 中惩罚的实现位置.

① 实现 score 函数, 初始化 score 如 [10, 10, 10, ...] 初始都为 10

② 在此步计算中, 若可能生成的词在前面的 ~~组合~~ 组合出现过, 则对这个词的分值减 1. 若没出现过, 则不减 (不惩罚)

## 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

下  
司  
至  
下  
为  
并  
据



# Outline

- 深度学习框架图计算理论
- Beam search
- 生成式文本摘要问题补充
- Baseline代码实践

## Scheduled Sampling

解决模型泛化能力不足的问题

A method for avoiding the problem of 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。

① 定义参数:  $\epsilon \in [0, 1]$  在第  $t$  步训练中真实值为  $x_t$ , 预测为  $\hat{x}_t$

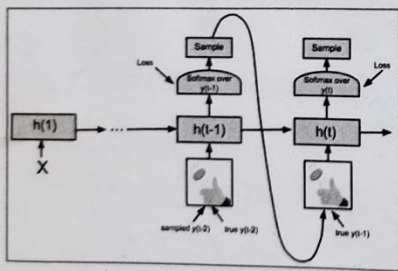
将  $\epsilon \cdot x_t + (1-\epsilon) \hat{x}_t$  做为下一步 ( $t+1$ ) 的输入,

② 开始的时候  $\epsilon$  取值较大, 随着训练步数的增多,  $\epsilon$  的值

取得较小.

Scheduled Sampling for Sequence Prediction with Recurrent Neural Networks

## Scheduled Sampling



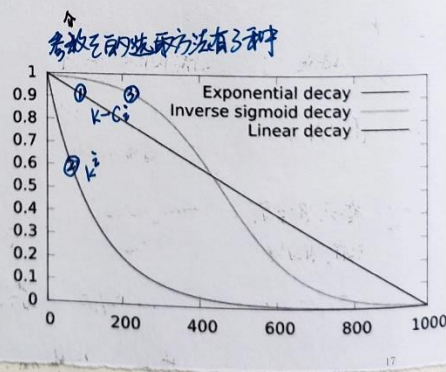
③ 指数衰减:  $\epsilon_i = k^i$

④ 逆 Sigmoid  $\epsilon_i = \frac{k}{k + e^{\frac{i}{k}}} \quad (k \geq 1)$

代码中 input 与 target 说明

input < START > 后厂理工学院 < END > target

①  $\epsilon_i = \max(\epsilon, k - C_i)$  衰减速率 线性衰减



## Datasets

数据集的选择很关键

CNN dailymail数据集: 文本摘要最具代表性的数据集

First highlight: Argentina coach Sabella believes Messi's habit of being sick during games is down to nerves.

First 2 sentences: Argentina coach Alejandro Sabella believes Lionel Messi's habit of throwing up during games is because of nerves. The Barcelona star has vomited on the pitch during several games over the last few seasons and appeared to once again during Argentina's last warm-up match against Slovenia on Saturday.

数据组成形式

单文本摘要数据集:

Gigaword

LCSTS

Newsroom

Xsum

train\_text.txt

train\_label.txt

短文本的内容, 约100-200字

短文本的摘要, 约10-20字

新浪微博摘要数据集 (679898 条数据)

["干杯 大哥!" 外卖小哥点头那一刻 泪目] 近日, 一男生和外卖小哥之间的互动, 在网上刷屏。视频中, 外卖小哥坐在电瓶车上啃干粮, 男生请他帮忙拧瓶盖。当外卖小哥拧开第二瓶水时, 男生说"干杯, 大哥, 天气很热, 加油!"。外卖小哥这才反应过来, 点头致谢 (央视) □

ROUGE 评价模型的如以, 基于召回模型

的评价指标  
recall-oriented understand for  
gisting evaluation

语言模型中的 n-gram

- ROUGE-N **RN**
- ROUGE-L **RL**: 1 (longest common subsequence)  
最长公共序列。
- ROUGE-S
- ROUGE-W
- ROUGE-SU

注: 分类问题可用准确率衡量  
生成式问题 (摘要生成) 不适合用  
准确度来衡量

BLEU 为基于准确度的评价指标。

参考摘要即参考答案 (标准答案)

注: Pyrouge 包 (pip 安装)

生成摘要问题通常用 rouge 评价

在参考摘要  
和自动摘要中  
同时发现

自动摘要即生成的摘要

ROUGE-N =  $\frac{\text{参考摘要和自动摘要中共有的 n-gram 个数}}{\text{参考摘要中 n-gram 的个数}}$  通常取 rouge-1, rouge-2

初: 参考摘要: the cat was under the bed (x)

生成摘要: the cat was found the under the bed (y)

则  $\text{Rouge-1}(x, y) = \frac{6}{6} = 1.0$  (蓝色标注)

$\text{Rouge-2}(x, y) = \frac{4}{5}$  (红色标注)



ROUGE-L:  $LCS(x, y)$ :  $x, y$  两个数据集中最长子序列相同的个数

$R_{LCS} = \frac{LCS(x, y)}{m}$ : 计算召回 (最长子序列的召回)  
 $m$  为参考答案的长度

$P_{LCS} = \frac{LCS(x, y)}{n}$ : 准确率  
 $n$  为生成答案的长度

$$F_{LCS} = \frac{(1 + \beta^2)RP}{R + \beta^2P}$$

$S_1$ : police killed the gunman 参考

$S_2$ : police ended the gunman 战

$S_3$ : the gunman murdered police 生成  
 句子的方向  $\rightarrow$

当  $\beta=1$  时, Rouge-L 为召回  $R_{LCS}$ , 则此时:

$$Rouge-L_{S2} = \frac{3}{4}$$

$$Rouge-L_{S3} = \frac{2}{4}$$

则生成的序列  $S_2$  优于  $S_3$

## Initializing neural

networks 网络初始化, 初始化一些权重、参数

- 基于固定方差初始化方法
  - 符合高斯分布的初始化 (期望为0)
  - 均匀分布初始化

- 基于差值缩放的参数初始化方法

- Xavier 初始化方法 (Glorot)
  - 根据每层神经元数量自动计算适当的参数初始化方法

$\downarrow$  适用于 sigmoid, tanh 作为激活函数。

- He 初始化, 但不适用于 ReLU, 适合 ReLU 作为激活函数。



期望  $N(0, \sigma^2)$   
 ① 若差值太小, 会使神经元输出太小, 经过多层后信号消失  
 ② 若差值太大, 采用 sigmoid 作为激活函数时, 丧失非线性能力, 容易出现梯度消失





## Batch size

1. Batch size是用于在每次迭代中训练模型的数据数量。一般的设置是32, 64, 128, 256, 512。
2. 选择正确的Batch size对于确保cost function和参数值的收敛, 以及模型的泛化能力。
3. Batch size决定更新的频率。Batch size越小, 更新就越快。
4. Batch size越大, 梯度越精确。也就是说, 在迭代计算的时候更容易跳过局部区域。
5. 比较大Batch size, 往往GPU memory是不够用的, 就需要通过并行计算的方式解决。

## Choice of optimizer

优化器

(Stochastic) Gradient Descent

Momentum

RMSprop

Adam 最常用, 会自动调节学习率, 但

Adam容易造成不收敛, 此时需要  
更换其他优化器

## Outline

- 深度学习框架图计算理论
- Beam search
- 生成式文本摘要问题补充
- Baseline代码实践见jupyter

附：

(1) 后厂理工 ai github: <https://github.com/HouchangX-AI>

(2) 微分代码 github:

[https://github.com/ZhaoYi1031/automatic\\_differentiation/blob/master/autodiff\\_test.py](https://github.com/ZhaoYi1031/automatic_differentiation/blob/master/autodiff_test.py)

(3) 后厂理工代码参考:

<https://github.com/HouchangX-AI/Question-and-answer-summary-and-reasoning>

(4) 优化器: <https://deeplearning.ai/ai-notes/optimization/>

(5) 网络参数初始化: <https://www.deeplearning.ai/ai-notes/initialization/>

(6) 华为云: <https://auth.huaweicloud.com/authui/login.html#/login>

(7) python argparse 讲解: <https://zhuanlan.zhihu.com/p/28871131>