

HCT NLP Week 4

问答摘要与推理
Seq2Seq（二）

Outline

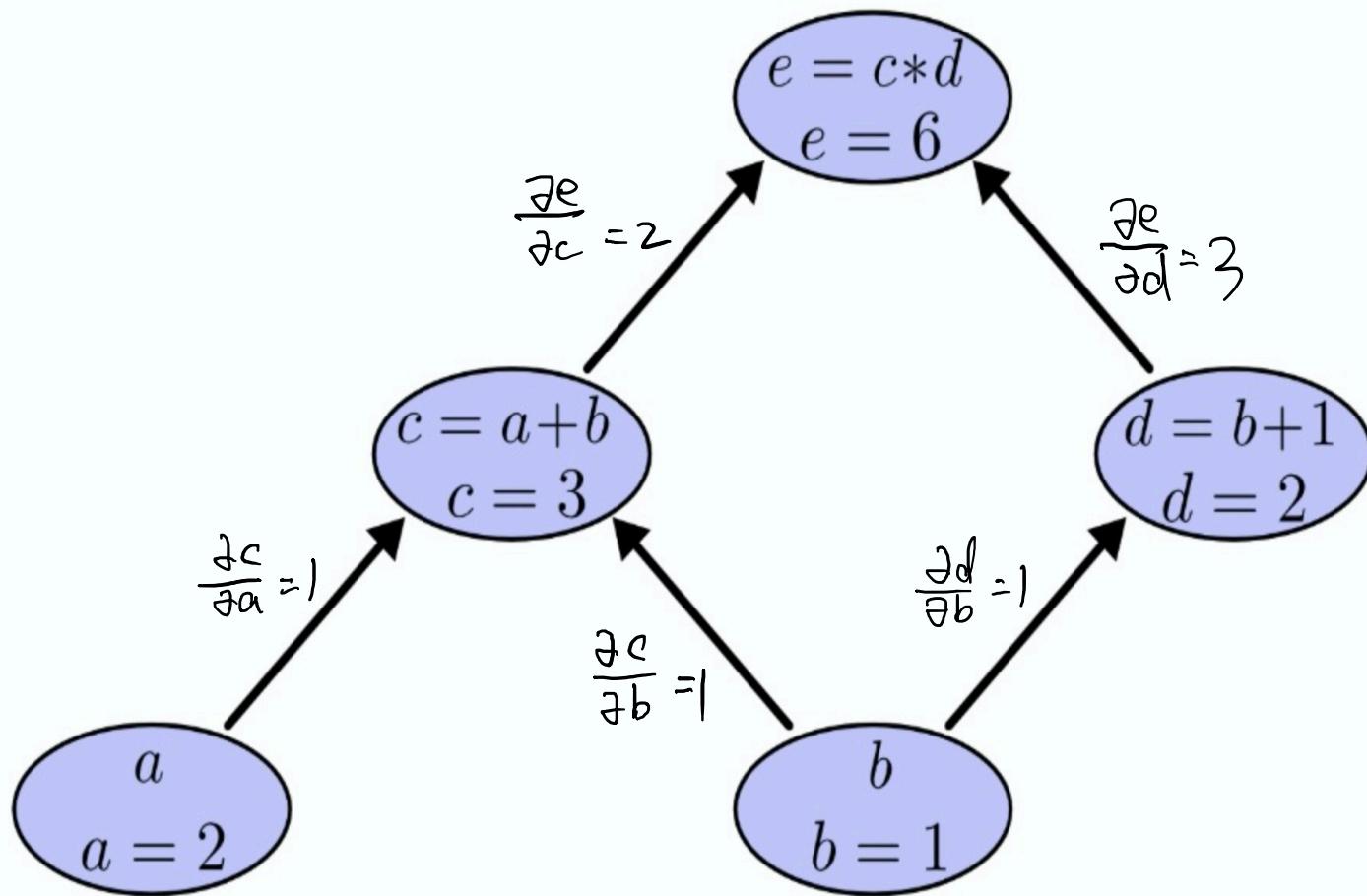
- 深度学习框架图计算理论
- Beam search
- 生成式文本摘要问题补充
- Baseline代码实践

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



$$e = (a+b) * (b+1)$$

$$c = a+b$$

$$d = b+1$$

$$e = c * d$$

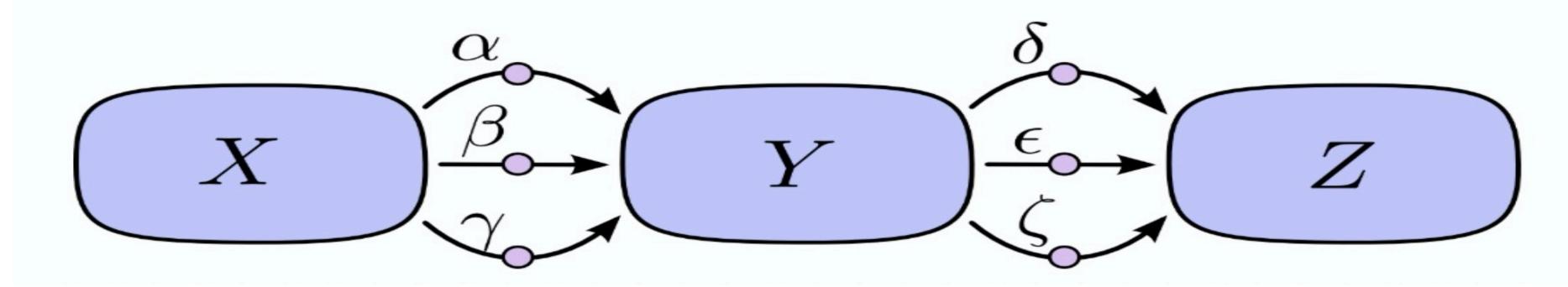
$$\frac{\partial e}{\partial b} = 1 \times 2 + 1 \times 3 = 5$$

Derivatives on it ?

图计算理论

Factoring Paths

$$\frac{\partial z}{\partial x} = \alpha\delta + \alpha\varepsilon + \alpha\varsigma + \beta\delta + \beta\varepsilon + \beta\varsigma + \gamma\delta + \gamma\varepsilon + \gamma\varsigma = (\alpha + \beta + \gamma)(\delta + \varepsilon + \varsigma)$$

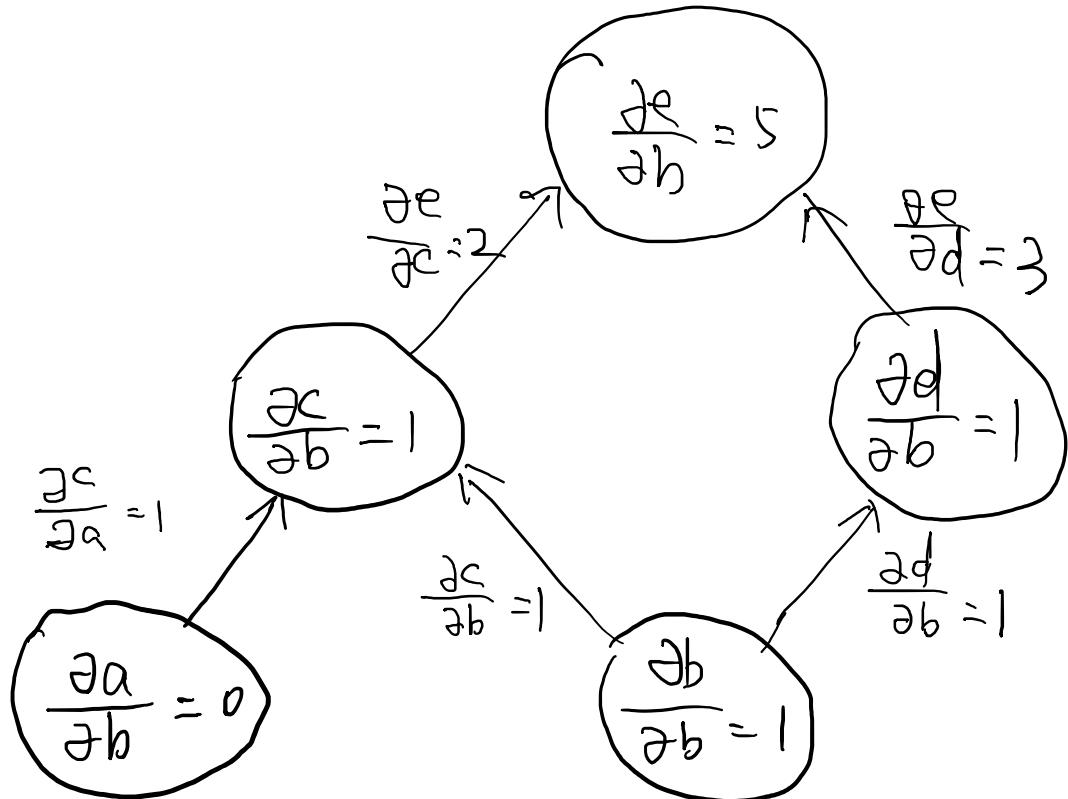


forward-mode differentiation

reverse-mode differentiation

图计算理论

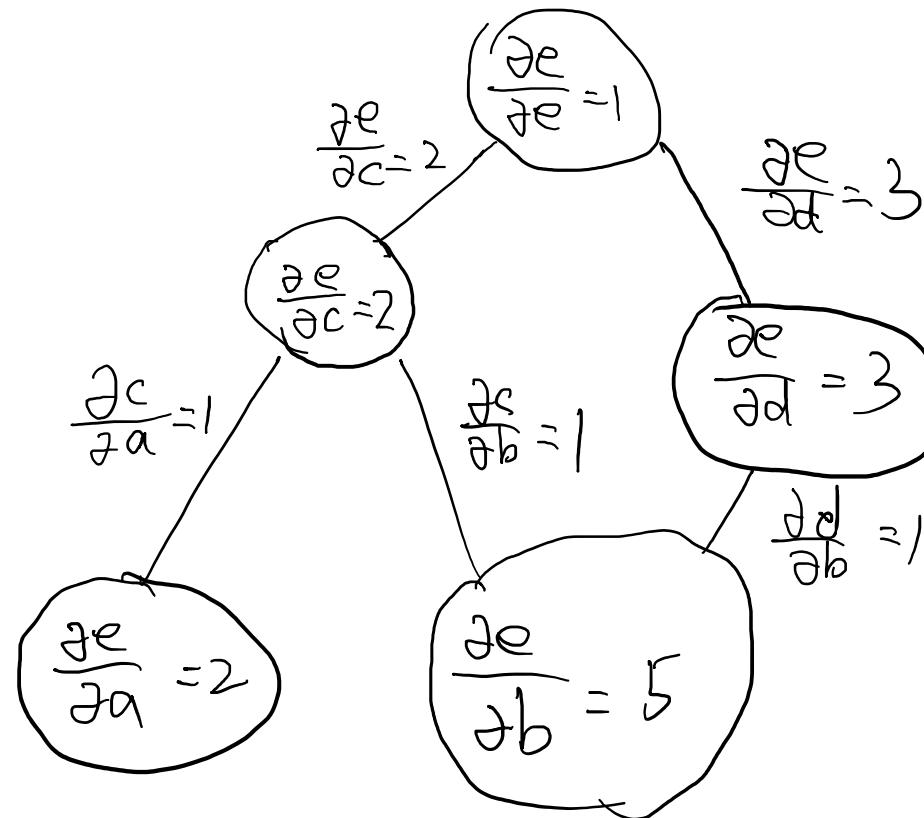
forward-mode differentiation



这张图从下往上推

图计算理论

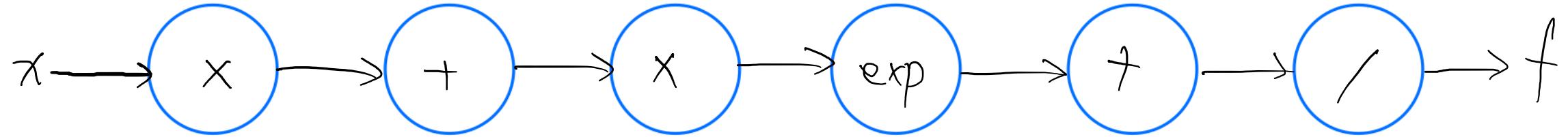
reverse-mode differentiation



这张图从上往下推

图计算理论

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



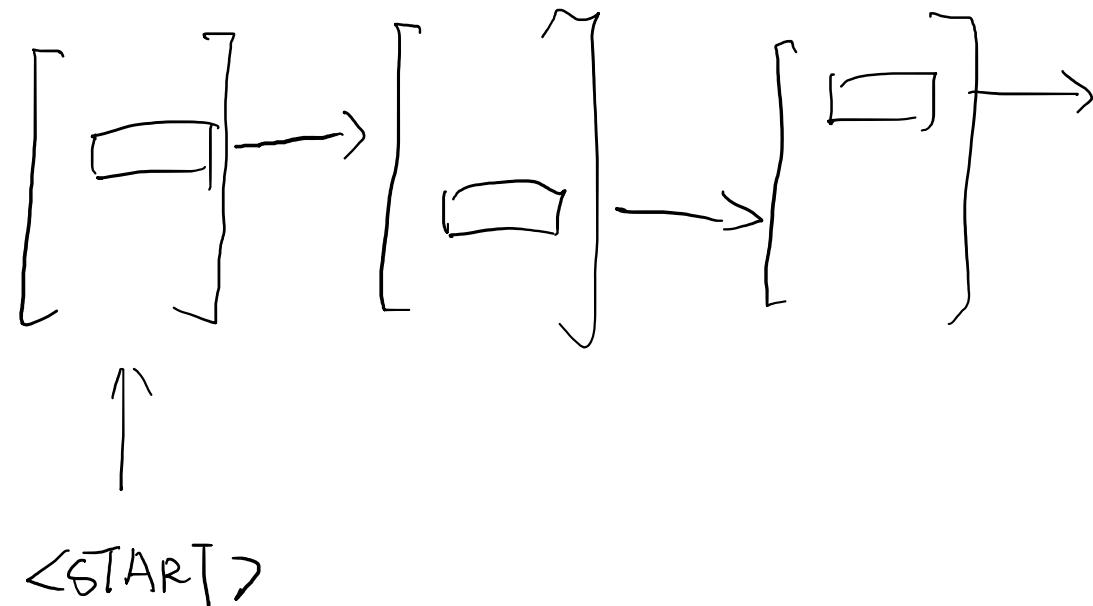
Outline

- 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

Beam search

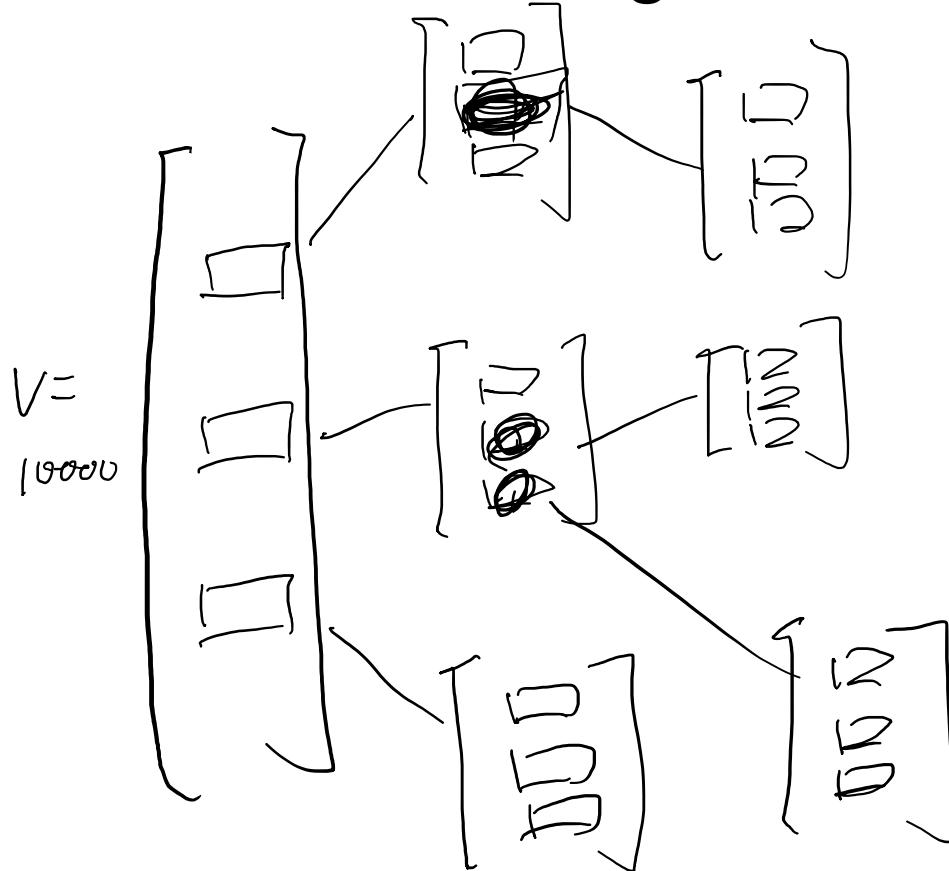
greedy search algorithm

行驶 没有 极端 的 感觉



Beam search

Beam search algorithm

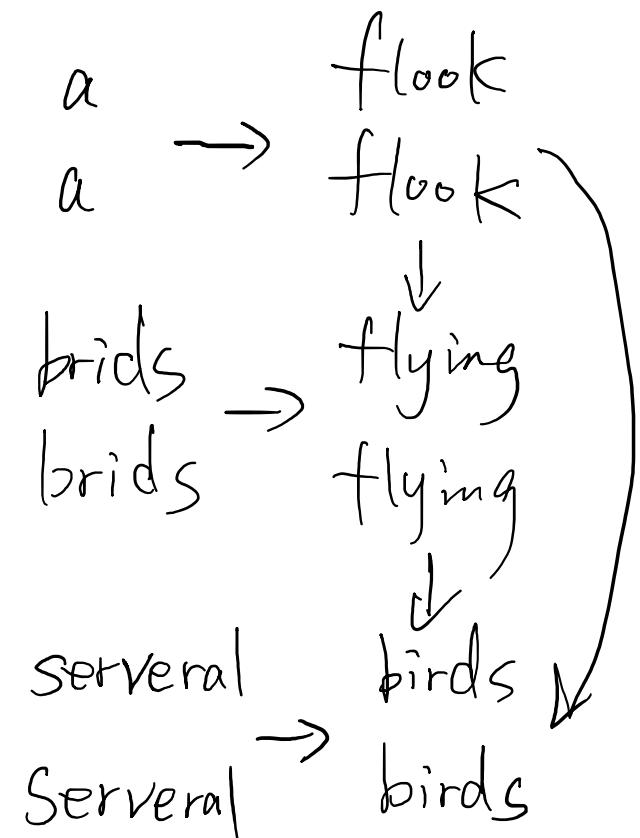


Beam search

Beam search algorithm

Beam search

Beam search algorithm



$B=6$, $G=3$, 每一组的 beam width 为 2
打分 scores

此时刻可能输出的词在上面的组里见过，
我们对这个词分数减 1，没出现，不惩罚。

diverse beam search

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

- 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

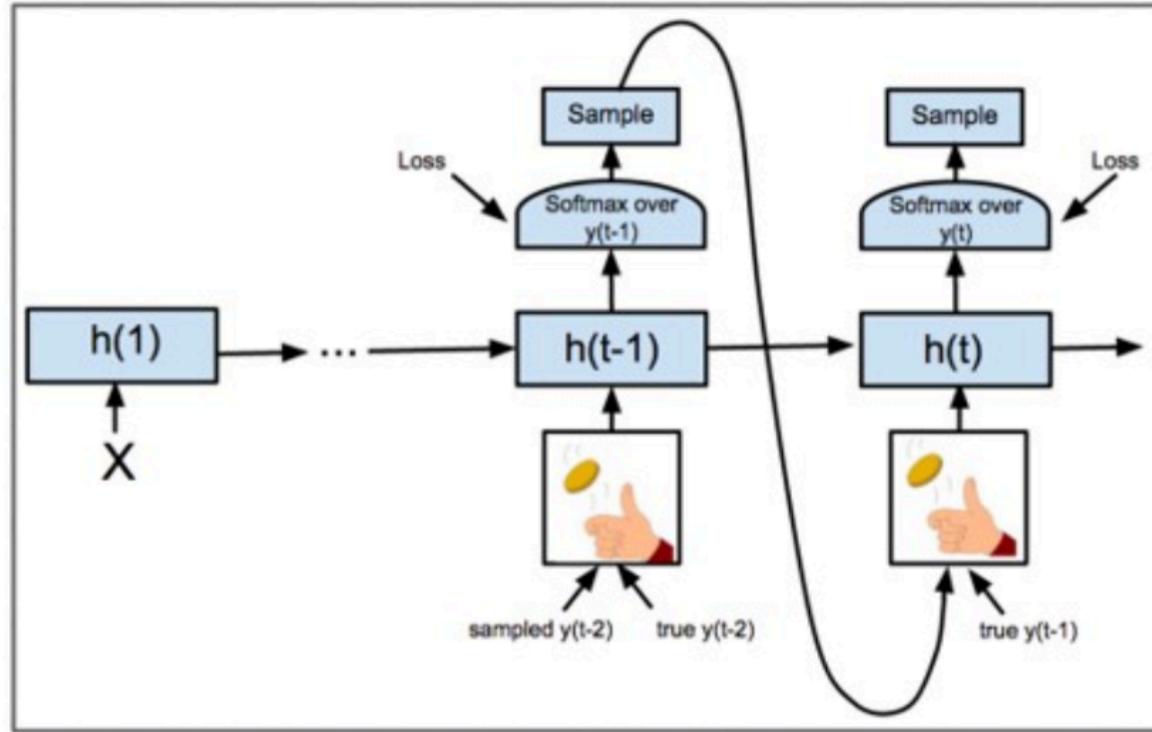
Scheduled Sampling

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](#)

Scheduled Sampling



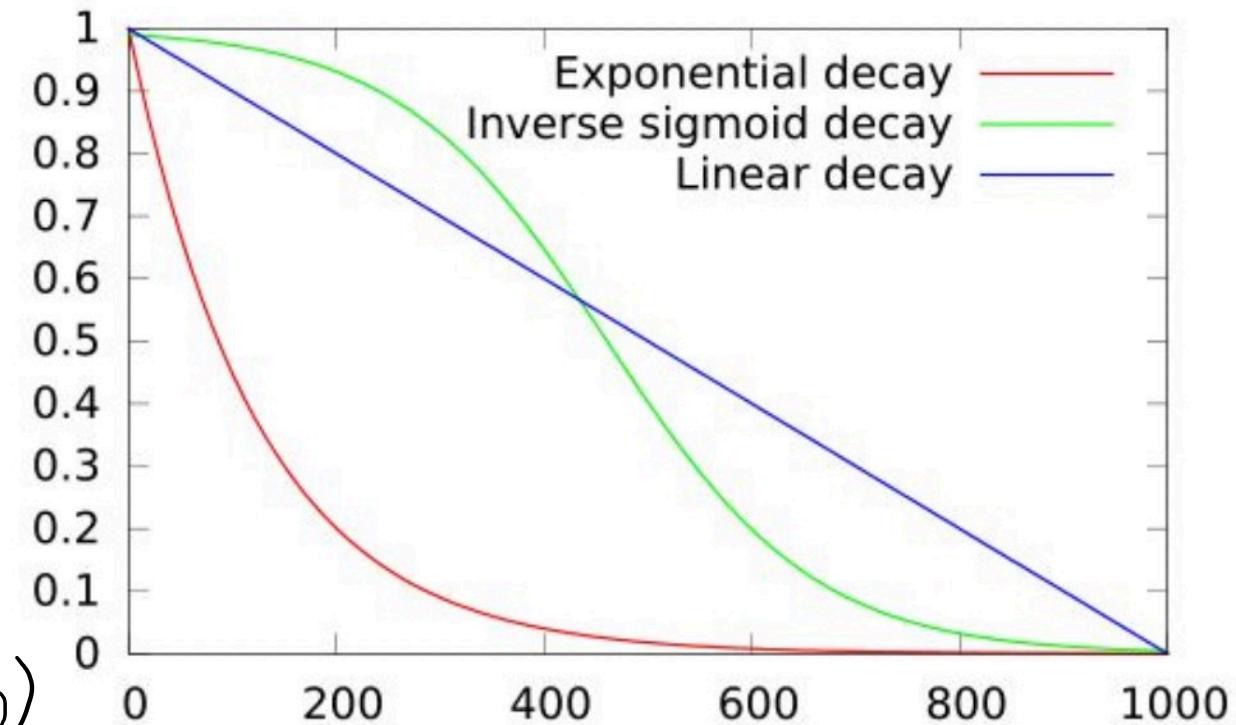
定义 $\epsilon \in [0, 1]$

$$\epsilon \cdot x_{t-1} + (1-\epsilon) \hat{x}_{t-1}$$

真值
预测

(1) 线性 decay

$$\epsilon_i = \max(\epsilon, K - c_i)$$



(2) 指数 decay $\epsilon_i = k^i$

(3) 反 sigmoid decay $\epsilon_i = k / (k + \exp(i/k))$

Datasets

CNN dailymail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.

单文档摘要：
Gigaword
LCSTS
Newsroom
Xsum

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) ▶

train_text.txt

短文本的内容，约100-200字

train_label.txt

短文本的摘要，约10-20字

ROUGE

recall-oriented understand for
gisting evalution

- ROUGE-N
- ROUGE-L
- ROUGE-S
- ROUGE-W
- ROUGE-SU

$$\text{Rouge}_N = \frac{\text{参考摘要和自动摘要共有 } n\text{-gram个数}}{n\text{-gram的个数} (\text{参考摘要})}$$

自动摘要 the cat was found under the bed (Y)
 参考摘要 the cat was under the bed (X₁)

$$\text{Rouge}_1(X_1, Y) = \frac{6}{6} = 1.0$$

$$\text{Rouge}_2(X_1, Y) = \frac{4}{7} = 0.8$$

ROUGE

recall-oriented understand for
gisting evalution

- ROUGE-N
- ROUGE-L
- ROUGE-S
- ROUGE-W
- ROUGE-SU

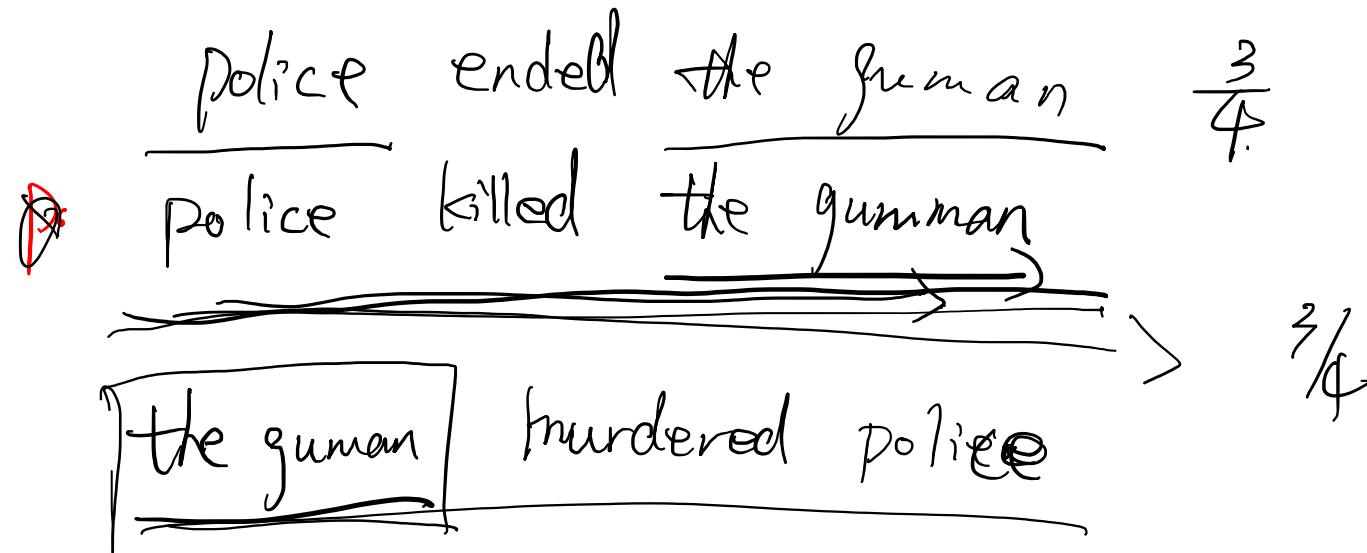
$$R_{lcs} = \frac{Lcs(x, y)}{m}$$

m为参考摘要长度

$$P_{lcs} = \frac{Lcs(x, y)}{n}$$

n为候选摘要长度

$$F_{lcs} = \frac{(1+\beta^2) R_{lcs} P_{lcs}}{R_{lcs} + \beta^2 P_{lcs}}$$



$$\rightarrow \beta=1$$

$$Rouge_L = \frac{3}{4}$$

S_2 优于 S_3

Initializing neural networks

$$\begin{bmatrix} w \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

(1) 固定方差 \leftarrow
 高斯分布
 均匀分布

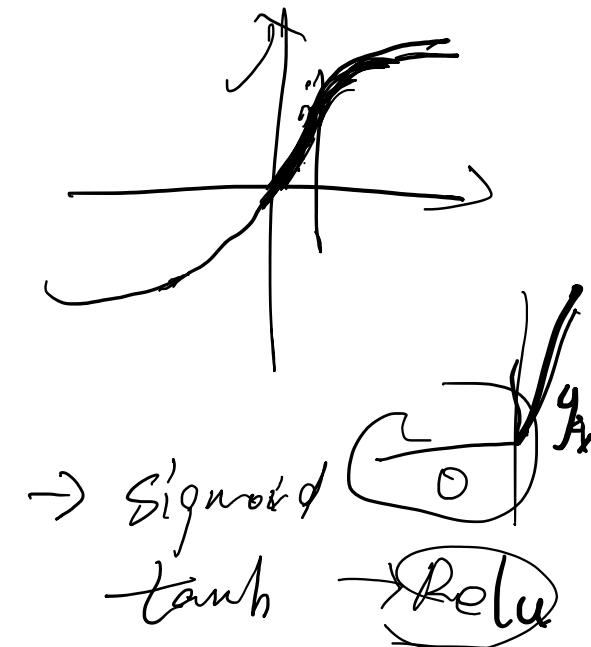
$$\mathcal{N}(0, \sigma^2)$$

(2) 梯度方差缩放

Xavier 初始化 (Glorot)

He 初始化

Relu



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

(Stochastic) Gradient Descent

Momentum

RMSprop

Adam

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

作业

Bye !