

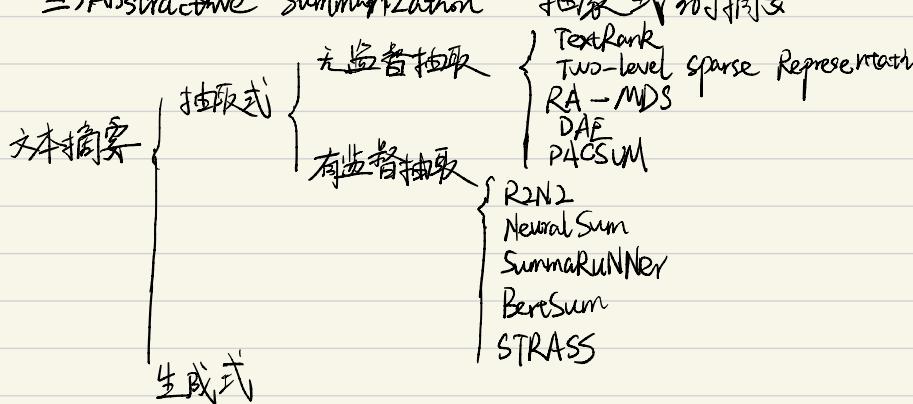
NLP Text summarization part 2

一般文本摘要有两种方式：

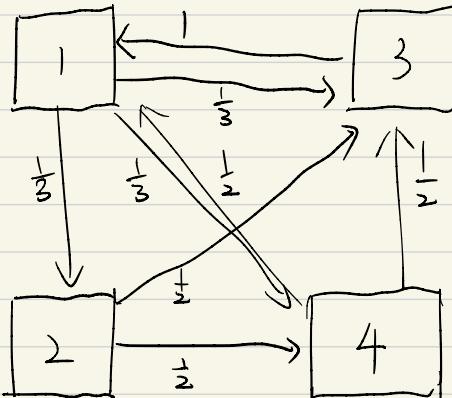
① 记号笔：所有的内容都来自原文 \Rightarrow extractive summarization

② 钢笔：写出来，就是看完之后做总结概括

\Rightarrow Abstractive Summarization



Google: PageRank 对网页排序



$$\begin{matrix} & 1 & 2 & 3 & 4 \\ 1 & 0 & 0 & 1 & \frac{1}{2} \\ 2 & \frac{1}{3} & 0 & 0 & 0 \\ 3 & \frac{1}{3} & \frac{1}{2} & 0 & \frac{1}{2} \\ 4 & \frac{1}{3} & \frac{1}{2} & 0 & 0 \end{matrix}$$

$$V = \begin{pmatrix} 0.25 \\ 0.25 \\ 0.25 \\ 0.25 \end{pmatrix}$$

$$AV = \begin{pmatrix} 0.37 \\ 0.08 \\ 0.33 \\ 0.20 \end{pmatrix}$$

$$A^2 V = A(AV) = A \begin{pmatrix} 0.37 \\ 0.08 \\ 0.33 \\ 0.20 \end{pmatrix} = \begin{pmatrix} 0.43 \\ 0.12 \\ 0.27 \\ 0.16 \end{pmatrix}$$

$$A^3 V = \begin{pmatrix} 0.35 \\ 0.14 \\ 0.29 \\ 0.22 \end{pmatrix}$$

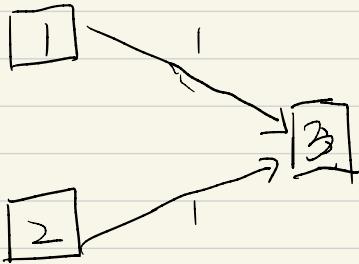
$$A^7 V = \begin{pmatrix} 0.38 \\ 0.12 \\ 0.29 \\ 0.19 \end{pmatrix}$$

$$A^8 V = \begin{pmatrix} 0.38 \\ 0.12 \\ 0.29 \\ 0.19 \end{pmatrix}$$

Converge

Two problems in PageRank

- Nodes with no outgoing edges (dangling nodes)



3 没有出边

$$V_1 = \begin{bmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{bmatrix} \quad \dots \quad V_2 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$V_1 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \frac{2}{3} \end{bmatrix}$$

damping factor \rightarrow solution

Text Rank 也用了 Page Rank 的思想

TextRank Process

1. 把所有篇文章中的句子连接起来
2. 把文本分割成一个个的句子
3. 对每个句子做向量化
4. 计算句子之间的相似度，并存储到矩阵中
5. 把相似度矩阵转化为图，其中句子是节点，相似度是节点间的权重
6. 取到 top-k 权重的句子，作为最终的摘要

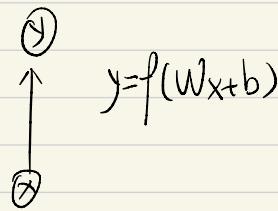
networkX \Rightarrow PageRank 算法

生成式算法：

Sequence Components:

Recurrent Neural Network:

Start from a single layer



classical RNN:

$$h_0 \rightarrow h_1$$

$$\begin{matrix} \uparrow \\ x_1 \quad x_2 \quad x_3 \quad x_4 \end{matrix}$$

$$h_1 = f(Ux_i + Wh_0 + b)$$

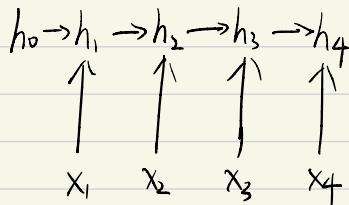
token: embedding 32D vector



$$h_0 \rightarrow h_1 \rightarrow h_2$$

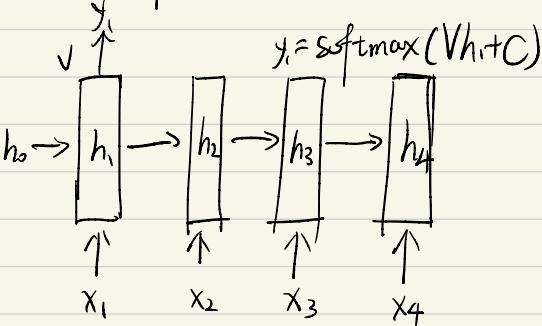
$$\begin{matrix} \uparrow & \uparrow \\ x_1 & x_2 \end{matrix} \quad x_3 \quad x_4$$

$$h_2 = f(Ux_2 + Wh_1 + b)$$

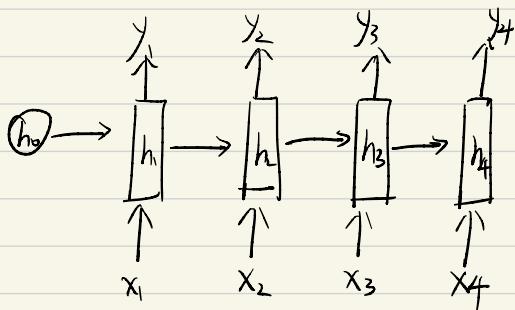


每一步使用的参数 U, W, b 都是一样的，也就是该步骤的参数都是共享的，这是 RNN 的主要特点
共享记忆

RNN Output



不同的 RNN 结构：N vs N



常用于命名实体识别
计算视频中每一帧的分类标签
因为要对每帧进行计算，因此输入和输出序列很长

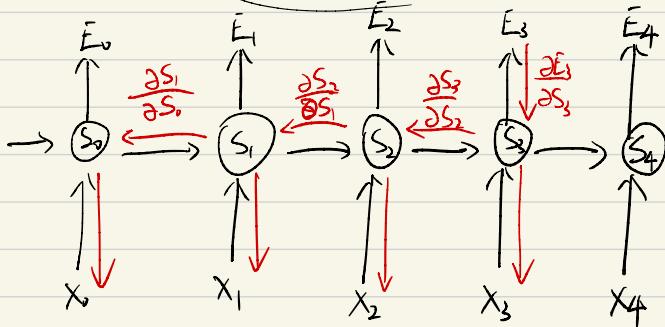
Backpropagate Through Time (BPTT) 反向传播

$$\frac{\partial E}{\partial W} = \sum_t \frac{\partial E_t}{\partial W} \quad \frac{\partial E_3}{\partial W} = \frac{\partial E_3}{\partial g_3} \frac{\partial g_3}{\partial s_3} \frac{\partial s_3}{\partial W} \quad \text{But } S_3 = \tanh(Ux_0 + WS_2)$$

But $S_3 = \tanh(WX_4 + WS_2)$ S_3 depends on S_2 , which depends on W and S_1 , and so on.

0.01

$$\frac{\partial E_1}{\partial W} = \sum_{k=0}^3 \left(\frac{\partial E_k}{\partial y_k} \right) \left(\frac{\partial y_k}{\partial S_k} \right) \left(\frac{\partial S_k}{\partial S_k} \right) \left(\frac{\partial S_k}{\partial W} \right)$$



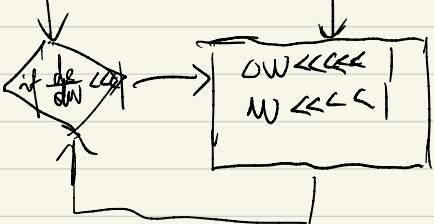
Vanishing Gradient 薄度消失

Back propagation

$$W = W + \Delta W$$

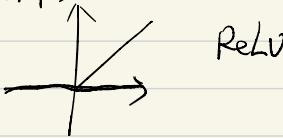
$$\Delta W = \eta \frac{de}{dW}$$

$$e = (\text{Actual output} - \text{Model output})^2$$



Solve Vanishing Gradient (梯度下降)

- ReLU activation function
- LSTM, GRU



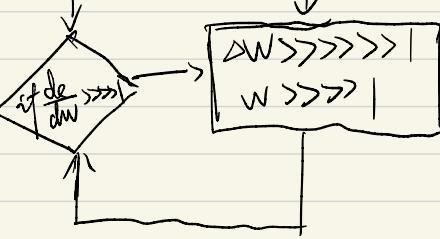
Exploding Gradient 梯度爆炸

反向传播
↓

$$W = W + \Delta W$$

$$\Delta W = n \frac{de}{dW}$$

$$e = (\text{Actual Output} - \text{Model Output})^2$$



Solve Exploding Gradient

①

梯度裁剪

Common solution: clipping gradient
if $\|g\| > V$ clip g: gradient

$\cdot V$: threshold

If norm of gradient exceeds some threshold, clip it!

②

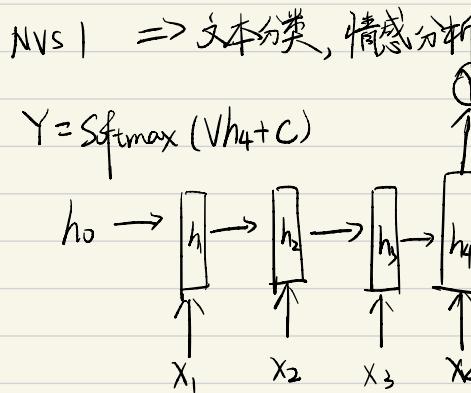
Truncated BPTT

k1: 前向传播的时间步。一般来说，这个参数影响模型训练的快慢，即权重更新的频率

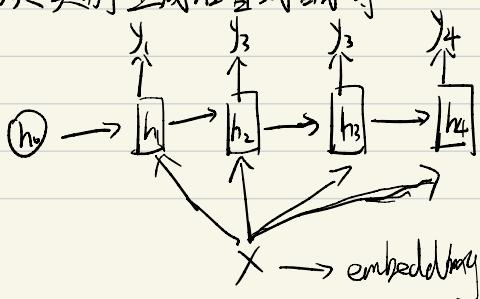
k2: 使用BPTT反向传播的时间步，一般来说，这个参数需要大一点，这样网络能更好的学习序列信息，但这个参数太大的话会导致梯度消失

③ RMSprop to adjust learning rate

Bidirectional RNN



VSN \Rightarrow 从图像生成文字 (image caption)，此时输入的X就是图像的特征，而输出的Y序列就是一段句子
从类别生成语音或音乐等



其实都是 RNN，只是输入和输出不一样了

图片通过 CNN 映射到向量空间中变成了 feature map

普通 RNN: `tf.keras.layers.SimpleRNN(4, input_shape=(3, 2))`

↑
h₀ hidden state 整个序列长度。

输入的部分加一个 embedding layer: onehot \Rightarrow 向量

两个 trick:
① 使用已经有的 embedding 作为参数
② embedding_lookup

`embedding_initializer = tf.keras.initializers.Constant()`

`embedding_lookup(..., feature_batch)`

多输出的 RNN:

`tf.keras.layers.SimpleRNN(4, input_shape=(), return_sequences=True)`

每个时间步增加层:

`tf.keras.layers.TimeDistributed(tf.keras.layers.Dense(...))`

多层次加: `model = tf.keras.Sequential()`

`model.add(tf.keras.layers.SimpleRNN(..., return_sequences=True))`

`model.add(...)`

`model.add(...)` (这一步一定要有，否则)

就没有向上传递了

双向RNN:

model.add(tf.keras.layers.Bidirectional())

RNN的优点与缺点:

Advantage: RNN can model sequence of data (i.e. time series) so that each sample can be assumed to be dependent on previous ones

· share parameters 参数共享

Disadvantage: ①梯度消失和梯度爆炸

②训练RNN很困难

③不能捕获很长的序列

LSTM

The Gates

Forget gate: Decides how much of past you should remember

$$f_t = \sigma(w_f \cdot [h_{t-1}, x_t] + b_f)$$

↓
bias

因为用了 sigmoid $f_f/f_i > f_c < 1$

Update Gate: decides how much of this unit is add to current state

$$t \rightarrow \sigma(w_i \cdot [h_{t-1}, x_t] + b_i) \rightarrow \text{range}(0, 1)$$

$$t' = \tanh(w_c \cdot [h_{t-1}, x_t] + b_c) \rightarrow \text{range}(-1, 1)$$

Memory pipe:

$$C_t = f_t C_{t-1} + i_t t$$

Output Gate: Decides which part of the current cell makes it to the output

$$O_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = O_t * \tanh(C_t)$$

GRU: LSTM的变种 减少LSTM的参数 减少了门，用一个门控制

LSTM & GRU 的对比：

GRU的参数更少，收敛速度更快，因此其实际花费时间要少很多，这大大加速了训练过程

从表现上讲，两者之间没有优劣定论，两者 performance 差距往往不大，建议调参带来的效果明显，不如在 LSTM or GRU 的激活函数 (ex. tanh 改为 tanh 等)

和权重初始化上下功夫

一般来说，先 GRU 作为基本单元，因为它收敛速度快，可加快试验进程，快速迭代

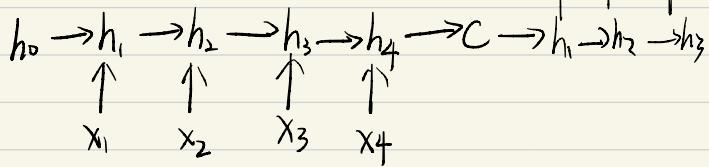
Seq2Seq 其实是 N vs M

N vs M

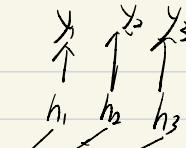
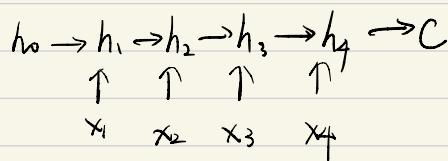
$$(1) C = h_4$$

$$(2) C = g(h_4)$$

$$(3) C = g(h_1, h_2, h_3, h_4)$$



Encoder-Decoder Architecture
translate, 语音, 文本摘要



变种

paper: Sequence to Sequence learning with Neural Networks

Input: use pre-trained embeddings or train our own embedding on our dataset

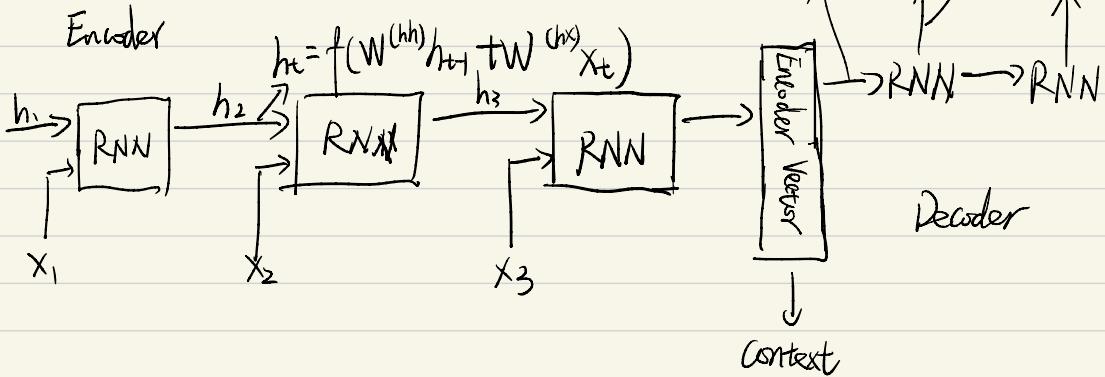
Context Vector

encoder \rightarrow context \rightarrow decoder

Usually size 256, 512 or 1024

$$y_t = \text{softmax}(W^S h_t)$$

Match in Seq2Seq



Attention mechanism

encoder 部分实际不能保存那么多信息，注意力机制会有选择地关注 input

Additive attention:

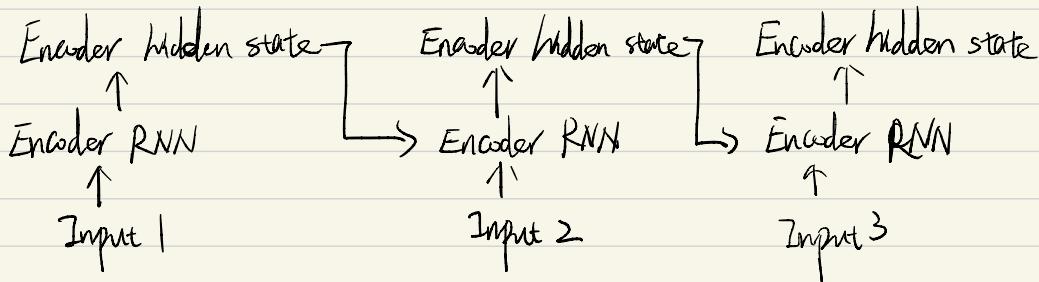
self-self with attention:

用当前的输入信息跟之前的每一个 hidden state 作一个打分，来看看多相关，打完分之后做一个 softmax，作出概率分布，加和算 context vector

缺点： $\Sigma \text{weight} \neq \text{encoder hidden state}$

加性 attention how to compute:

① Producing the encoder hidden states



②③④ Calculating Alignment Scores

additive

$$\text{Score}_{\text{alignment}} = \frac{W_{\text{combined}} \cdot \tanh(W_{\text{decoder}} \cdot H_{\text{decoder}} + W_{\text{encoder}} \cdot H_{\text{encoder}})}{\downarrow \text{字典相似} \quad \downarrow \text{学习相似} \quad \downarrow \text{学习相似}}$$

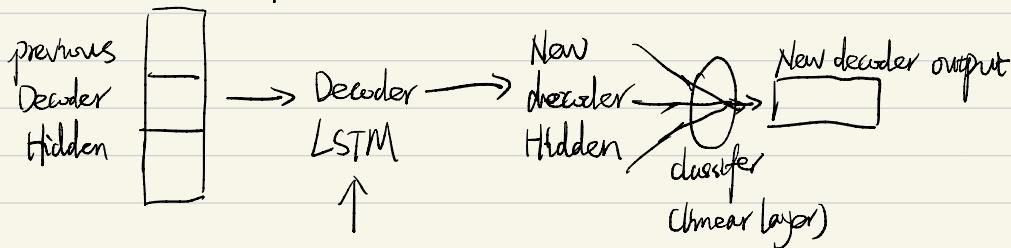
⑤ Softmax the alignment scores

$$\text{softmax}(\text{alignment score}) = \text{attention weight}$$

⑥ Calculating the Context Vector

$$\text{Encoder outputs} * \text{Attention weights} = \text{Context Vector}$$

⑦ Decoding the Output



Context Vector + Embedding of previous decoder output

multi attention for additive attention 有多于一个的：

$$\text{add: concat. score}_{\text{align}} = W \cdot \tanh(W_{\text{combined}} (\text{Encoder} + \text{Decoder}))$$

$$\text{Multi: Generativ: score}_{\text{align}} = W (\text{Encoder} \cdot \text{Decoder})$$

LSTM:

return_state : \rightarrow hidden & context State

首先将 Encoder 的 embedding

attention机制:

tensorflow-addons

seq2seq.LuongAttention

loss function: cross entropy

Teacher forcing: 第一个输出错了会导致下一个也出错, 后面会连续出错

Ground Truth: 强制改输入

cons:

训练时有 Ground Truth, 但是 predict 时是没有的 \Rightarrow exposure bias

解决方法: ① Beam Search

② Scheduled Sampling: 抽硬币的方法来决定是不是使用
Ground Truth

一开始高概率用 Ground Truth, 后面随着学习越来越好,
然后一点点下降 Ground Truth 概率