



任课教师: 孔雨秋

yqkong@dlut.edu.cn

目录(CONTENT)



- 01 基本循环神经网络
- 02 循环神经网络训练
- 03 小结和资源

参考课程



■ Andrew Ng, Recurrent Neural Network

https://www.bilibili.com/video/av88010803?p=1

■ 李宏毅,深度学习

https://www.bilibili.com/video/BV1JE411g7XF





基本循环神经网络 Basic RNN

循环神经网络(Recurrent Neural Network, RNN)

 \mathcal{X}



■处理序列数据

语音识别

音乐生成

情感分类

机器翻译

视频行为识别

命名实体识别

"Yesterday, Harry Potter met Hermione Granger."

"The quick brown fox jumped over the lazy dog"

 \rightarrow



"Voulez-vous chanter avec moi?"

"There is nothing to

like in this movie."

"Do you want to sing with me?"

---- Running

"Yesterday, Harry Potter met Hermione Granger."

循环神经网络



■ 触发字检测



Amazon Echo (Alexa)



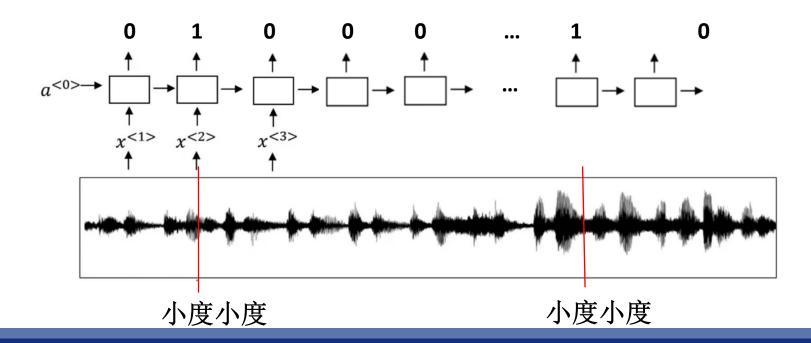
Baidu DuerOS (xiaodunihao)



Apple Siri (Hey Siri)



Google Home (Okay Google)





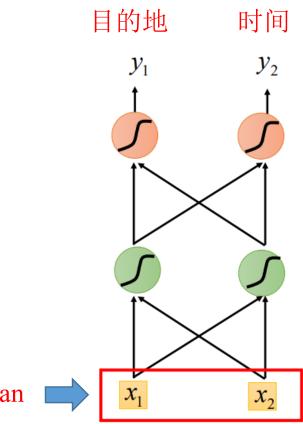
■前馈神经网络 V.S. 循环神经网络

Arrive Dalian on 1st December.

任务: 预测目的地是哪? 时间是?

输入:每次输入一个单词(特征表示)

输出:属于目的地或时间的概率





■前馈神经网络 V.S. 循环神经网络

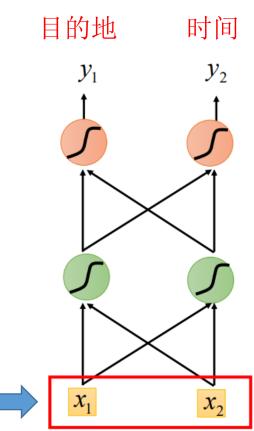
Arrive Dalian on 1st December.

任务: 预测目的地是哪? 时间是?

Arrive Dalian on 1st December. 目的地

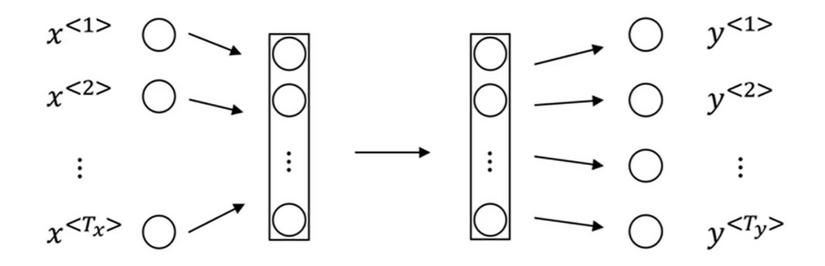
出发地 Leave Dalian on 1st December.

前馈神经网络不考虑上下文信息, 输入是固定的,输出就是固定的。





■前馈神经网络 V.S. 循环神经网络

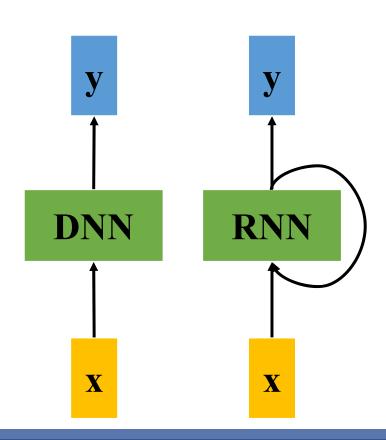


- □ 不同的训练样本,输入、输出序列的长度不同
- □ 在文本序列不同位置学习到的特征不共享



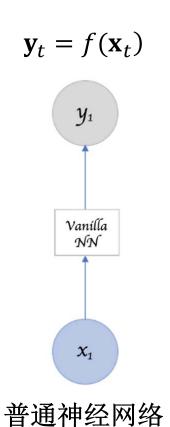
■ 前馈神经网络 V.S. 循环神经网络 Arrive Dalian on 1st December. 目的地 Leave Dalian on 1st December. 出发地

- □ 神经网络需要记忆功能,捕获历史时刻 的信息
- □ 将t时刻隐藏层的输出重新输入到隐藏 层中,作为t+1时刻的输入

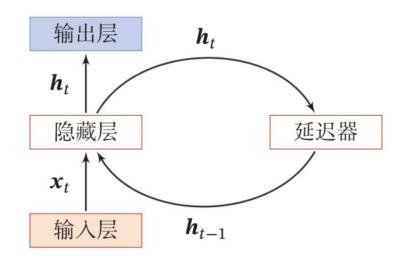




经验是智慧之父, 记忆是智慧之母。



$$\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t)$$



循环神经网络



■符号标记

X Harry Potter and Hermione Granger invented a new spell.

$$T_x = 9$$

$$\mathbf{x}^{<1>}$$

$$T_r = 9$$
 $\mathbf{x}^{<1>}$ $\mathbf{x}^{<2>}$ $\mathbf{x}^{<3>}$... $\mathbf{x}^{}$

$$\mathbf{x}^{< t>}$$

$$\mathbf{x}^{<9>}$$

 $y = 1 \quad 1 \quad 0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 0$

$$T_u = 9$$
 $\mathbf{y}^{<1>}$ $\mathbf{y}^{<2>}$ $\mathbf{y}^{<3>}$... $\mathbf{y}^{}$

$$\mathbf{v}^{<1>}$$

$$\mathbf{v}^{<2>}$$

$$\mathbf{X}^{(i)} = (\mathbf{x}^{<1>}, \mathbf{x}^{<2>}, \cdots, \mathbf{x}^{})$$

$$\mathbf{x}^{(i) < t >}$$

$$\mathbf{Y}^{(i)} = (\mathbf{y}^{<1>}, \mathbf{y}^{<2>}, \cdots, \mathbf{y}^{})$$

$$\mathbf{v}^{(i) < t >}$$

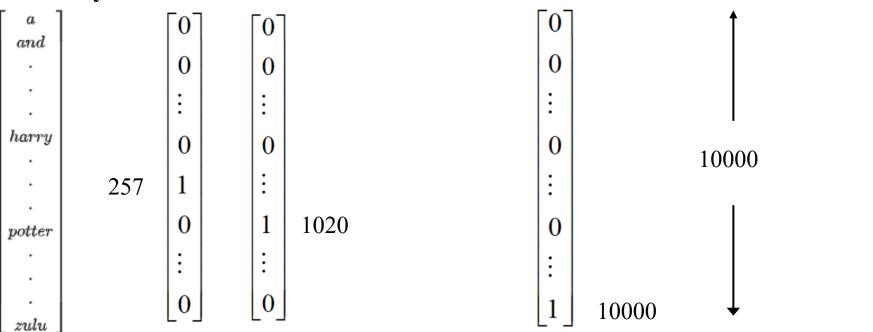


■符号标记

x Harry Potter and Hermione Granger invented a new spell.

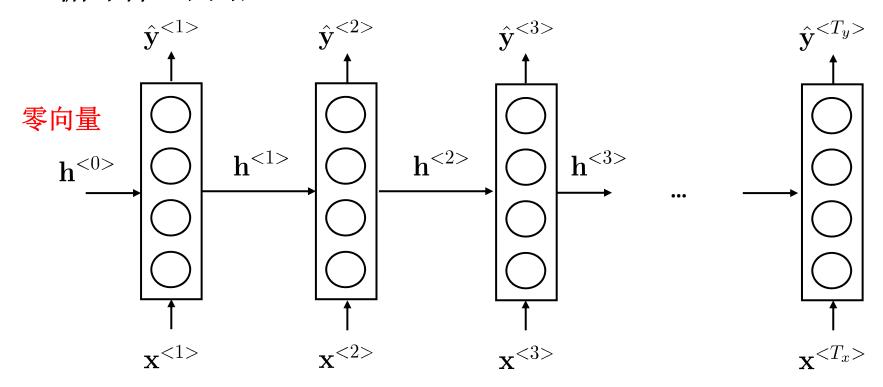
$$\mathbf{x}^{<1>}$$
 $\mathbf{x}^{<2>}$ $\mathbf{x}^{<3>}$... $\mathbf{x}^{}$... $\mathbf{x}^{<9>}$

Vocabulary



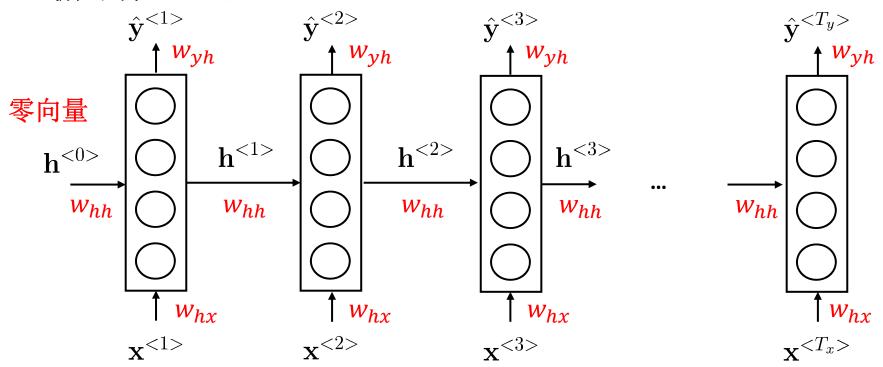
不在Vocabulary里的单词可用<UNK>标记





- □ 从左到右扫描序列数据
- □ 每个时间步共享参数

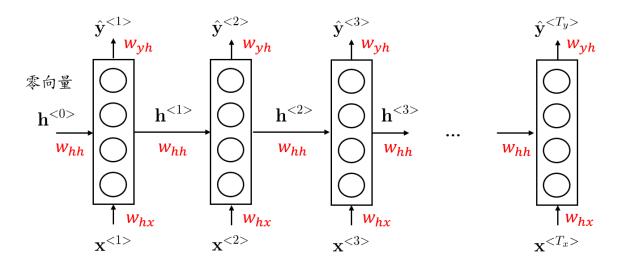


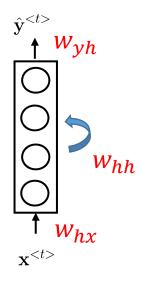


- □ 从左到右扫描序列数据
- □ 每个时间步共享参数



■循环神经网络





$$\mathbf{h}^{<0>} = \mathbf{0} \qquad \mathbf{h}^{<1>} = g(\mathbf{W}_{hh}\mathbf{h}^{<0>} + \mathbf{W}_{hx}\mathbf{x}^{<1>} + \mathbf{b}_h)$$
$$\hat{y}^{<1>} = g_2(\mathbf{W}_{yh}\mathbf{h}^{<1>} + \mathbf{b}_y)$$

$$\int \mathbf{h}^{\langle t \rangle} = g(\mathbf{W}_{hh} \mathbf{h}^{\langle t-1 \rangle} + \mathbf{W}_{hx} \mathbf{x}^{\langle t \rangle} + \mathbf{b}_h)$$

$$\hat{y}^{\langle t \rangle} = g_2(\mathbf{W}_{yh} \mathbf{h}^{\langle t \rangle} + \mathbf{b}_y)$$

激活函数: tanh

激活函数: sigmoid (二分类) softmax (多分类)



$$\int_{\mathbf{h}^{}} \mathbf{h}^{} = g(\mathbf{W}_{hh}\mathbf{h}^{} + \mathbf{W}_{hx}\mathbf{x}^{} + \mathbf{b}_{h}) \qquad \mathbf{h}^{} = g(\mathbf{W}_{h}[\mathbf{h}^{}, \mathbf{x}^{}] + \mathbf{b}_{h})$$

$$\hat{y}^{} = g_{2}(\mathbf{W}_{yh}\mathbf{h}^{} + \mathbf{b}_{y})$$

$$\begin{cases} \mathbf{h}^{} = g(\mathbf{W}_h[\mathbf{h}^{}, \mathbf{x}^{}] + \mathbf{b}_h) \\ \hat{y}^{} = g(\mathbf{W}_y\mathbf{h}^{} + \mathbf{b}_y) \end{cases}$$

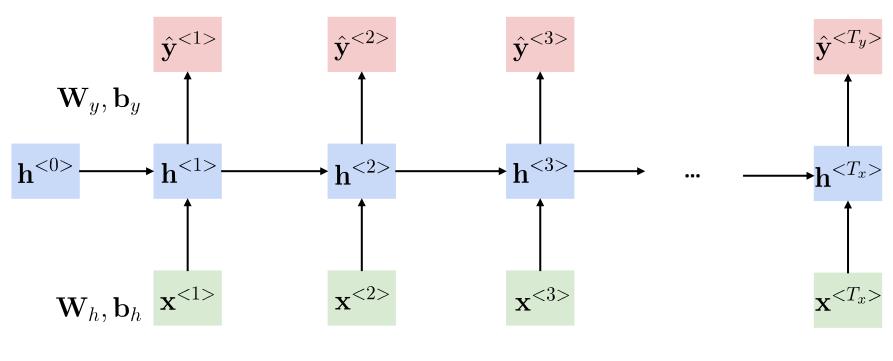
$$\mathbf{W}_{h} = [\mathbf{W}_{hh}; \mathbf{W}_{hx}]$$

$$(100,10100) \quad (100,1000)$$

$$[\mathbf{h}^{}, \mathbf{x}^{}] = \begin{bmatrix} \mathbf{h}^{} \\ \mathbf{x}^{} \end{bmatrix} = \begin{bmatrix} \mathbf{h}^{} \\ \mathbf{x}^{} \end{bmatrix} = \begin{bmatrix} \mathbf{h}^{} \\ \mathbf{x}^{} \end{bmatrix} = \mathbf{W}_{hh} \mathbf{h}^{} + \mathbf{W}_{hx} \mathbf{x}^{}$$





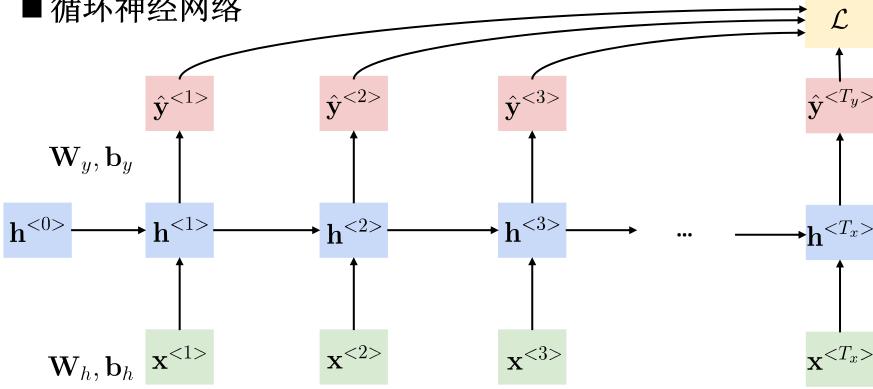


$$\mathcal{L}^{}(\hat{y}^{}, y^{}) = -y^{}\log \hat{y}^{} - (1 - y^{})\log(1 - \hat{y}^{})$$

$$\mathcal{L}(\hat{y}, y) = \sum_{t=1}^{T_y} \mathcal{L}^{}(\hat{y}^{}, y^{})$$







$$\mathcal{L}^{}(\hat{y}^{}, y^{}) = -y^{}\log \hat{y}^{} - (1 - y^{})\log(1 - \hat{y}^{})$$

$$\mathcal{L}(\hat{y}, y) = \sum_{t=1}^{T_y} \mathcal{L}^{}(\hat{y}^{}, y^{})$$

RNN



定理 6.1 - 循环神经网络的通用近似定理 [Haykin, 2009]:如果一个完全连接的循环神经网络有足够数量的 sigmoid 型隐藏神经元,它可以以任意的准确率去近似任何一个非线性动力系统

$$\mathbf{s}_t = \mathbf{g}(\mathbf{s}_{t-1}, \mathbf{x}_t), \tag{6.10}$$

$$\mathbf{y}_t = o(\mathbf{s}_t), \tag{6.11}$$

其中 \mathbf{s}_t 为每个时刻的隐状态, \mathbf{x}_t 是外部输入, $\mathbf{g}(\cdot)$ 是可测的状态转换函数, $\mathbf{o}(\cdot)$ 是连续输出函数,并且对状态空间的紧致性没有限制.



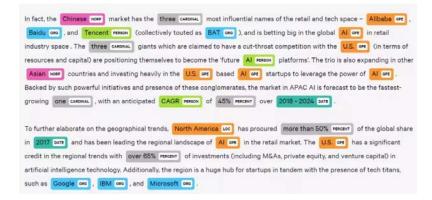
 \mathcal{X} "The quick brown fox jumped 语音识别 over the lazy dog" 音乐生成 "There is nothing to 情感分类 like in this movie." "Voulez-vous chanter "Do you want to sing 机器翻译 avec moi?" with me?" 视频行为识别 Running "Yesterday, Harry Potter "Yesterday, Harry Potter 命名实体识别 met Hermione Granger." met Hermione Granger."

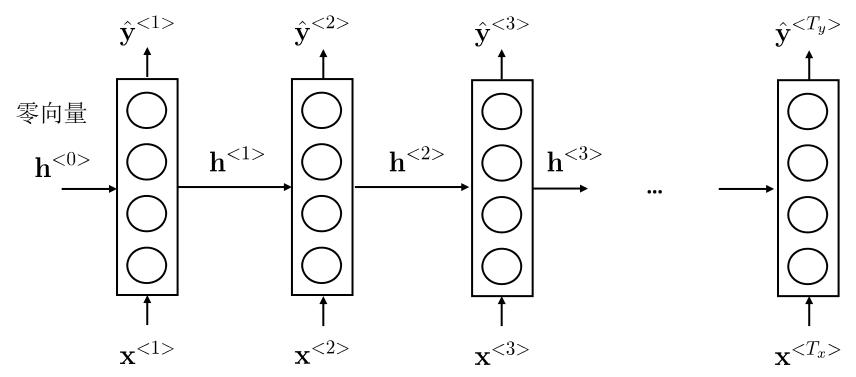


■ Many-to-many (命名实体识别)

$$T_x = T_y$$

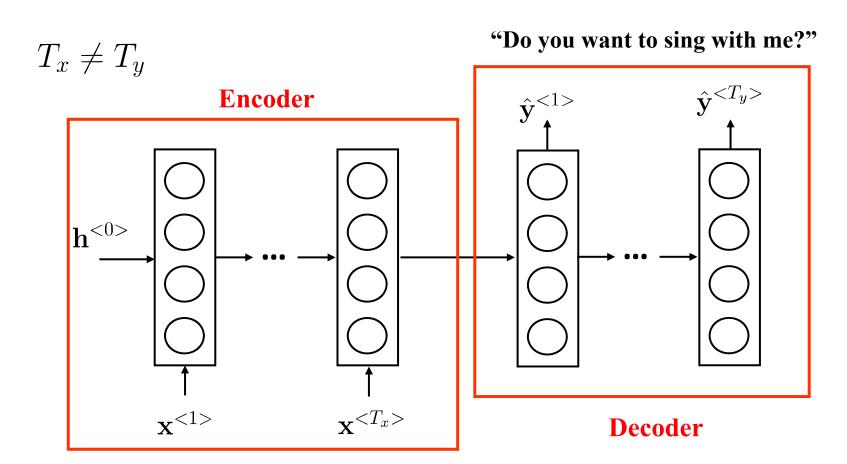
 $\hat{y}_t = g(\boldsymbol{h}_t), \quad \forall t \in [1, T]$







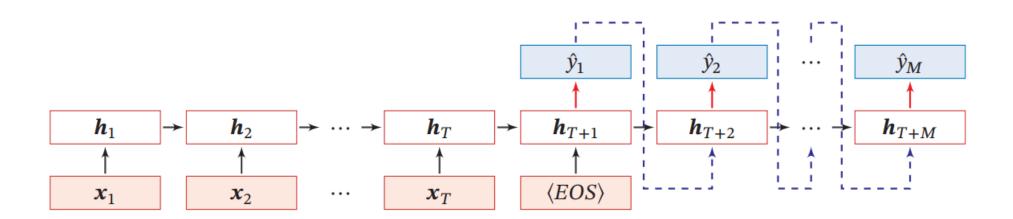
■ Many-to-many (机器翻译)



"Voulez-vous chanter avec moi?"

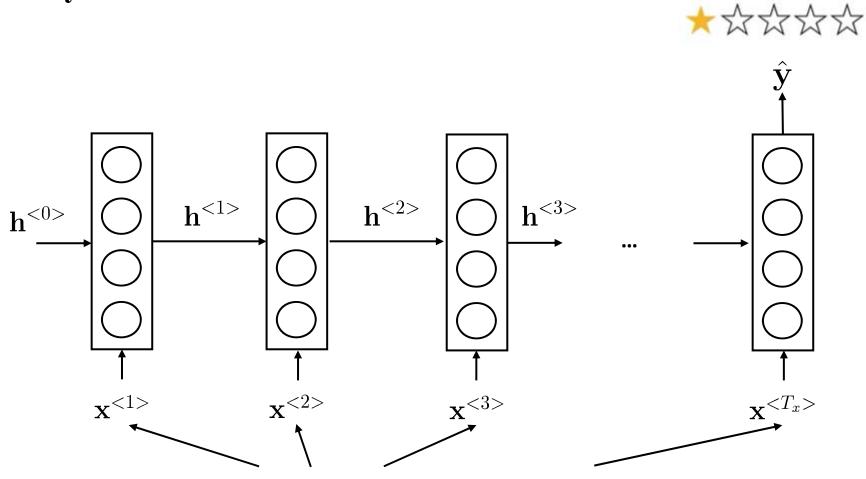


■ Many-to-many (机器翻译)





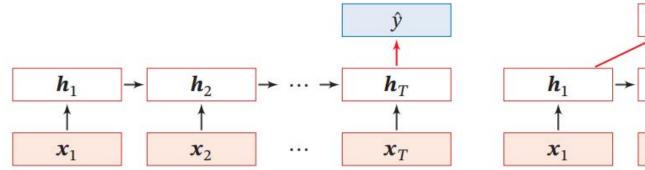
■Many-to-one (情感分类)



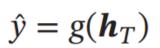
"There is nothing to like in this movie."

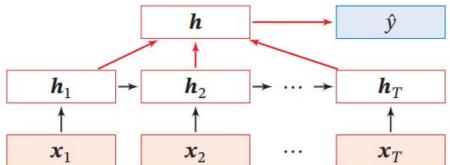


■Many-to-one (情感分类)



(a) 正常模式



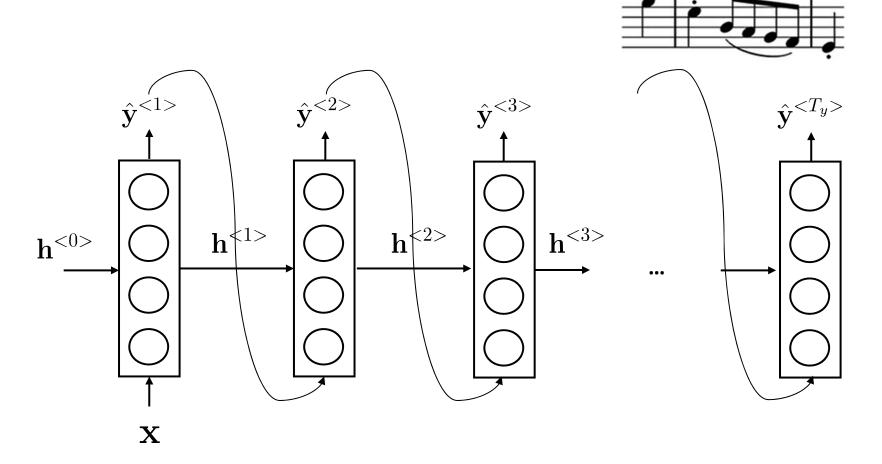


(b) 按时间进行平均采样模式

$$\hat{y} = g\Big(\frac{1}{T} \sum_{t=1}^{T} \boldsymbol{h}_t\Big)$$

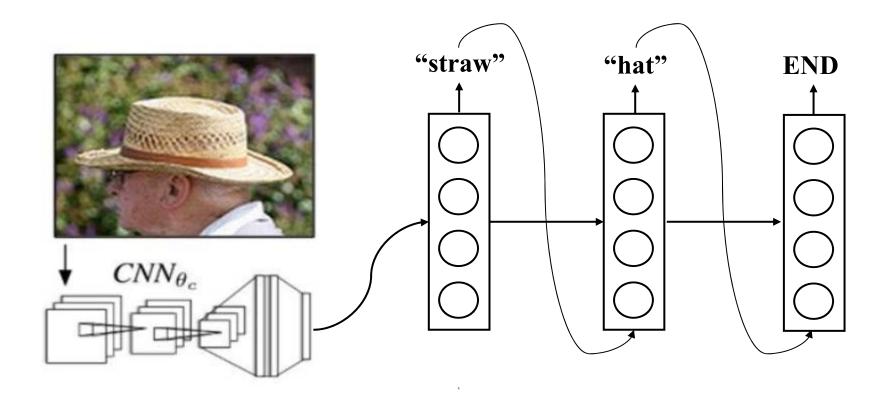


■One-to-many (音乐生成)



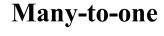


■One-to-many (图像加字幕)



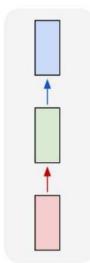


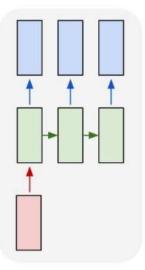
One-to-one One-to-many

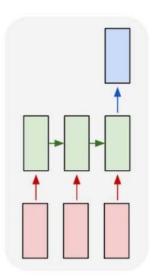


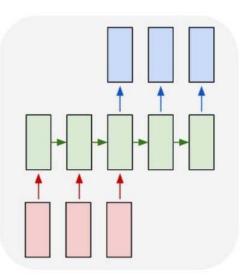
Many-to-many

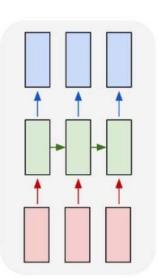
Many-to-many













- 语言模型 (Language model): 是自然语言处理的重要技术,常用于语音识别和机器翻译等。
- 语言模型给一个序列(句子)分配概率,找到最可能出现的词序列。 同时估计某个序列中各单词出现的概率。
- □ 语音识别系统: The apple and pair salad.

The apple and pear salad.

P(The apple and pair salad)=3.2e-13

P(The apple and pear salad)=5.7e-10

P(Sentence)=?
$$P(\hat{y}^{<1>}, \hat{y}^{<2>}, \cdots, \hat{y}^{})$$



■ 构建语言模型

□ 训练集: 大型文本语料库

□ 特征表示: one-hot向量

Cats average 15 hours of sleep a day.

$$\mathbf{y}^{<1>} \mathbf{y}^{<2>}$$
 ...

□ 用RNN构建词序列的概率模型

Vocabulary



<EOS

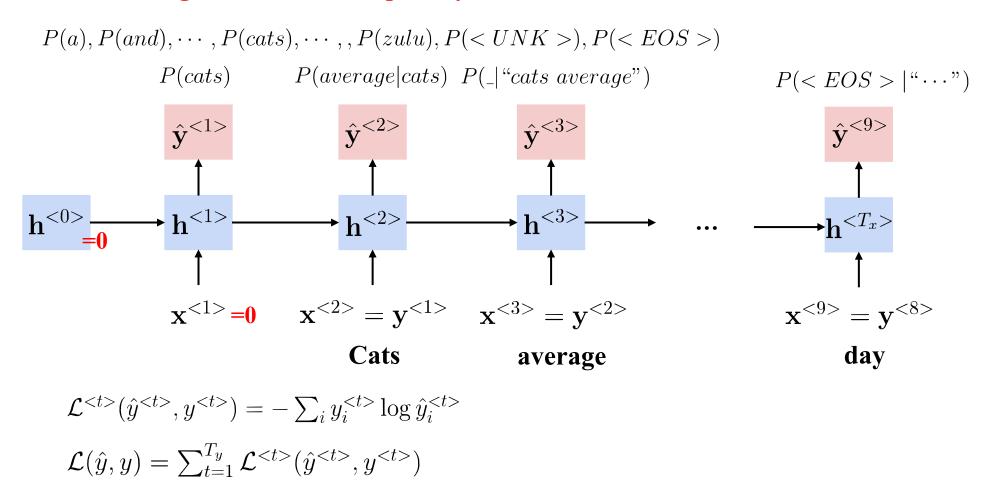
标点



Cats average 15 hours of sleep a day.

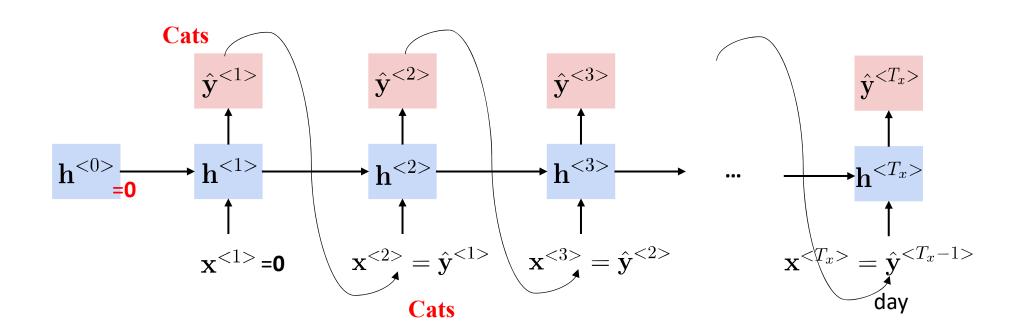


Cats average 15 hours of sleep a day.





- 用训练好的RNN生成句子
 - □ 在每个时间步下,根据输出的概率分布进行采样。
 - □将采样得到的单词作为下一个时间步的输入

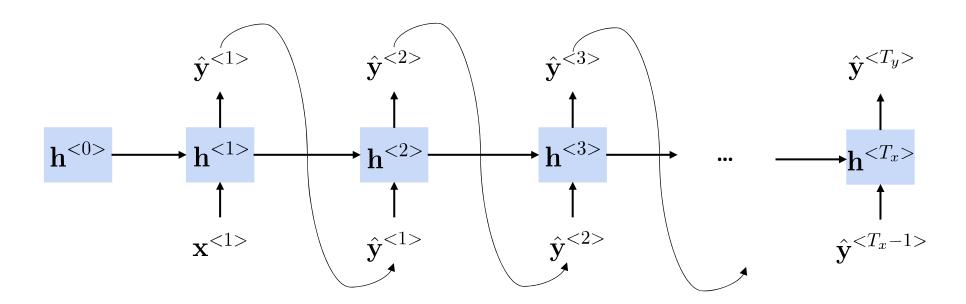




■ 字符级别语言模型

Vocabulary=[a, and, ..., zulu, <UNK>]

Vocabulary=[a, b, c, ..., z, ..., ?, .., 0, 1, 2, ...]



语言模型



■ 文本生成

News

President enrique peña nieto, announced sench's sulk former coming football langston paring.

"I was not at all surprised," said hich langston.

"Concussion epidemic", to be examined.

The gray football the told some and this has on the uefa icon, should money as.

Shakespeare

The mortal moon hath her eclipse in love.

And subject of this thou art another this fold.

When besser be my love to me see sabl's.

For whose are ruse of mine eyes heaves.

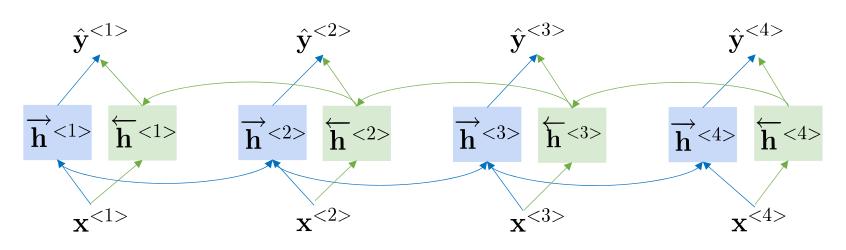




- 双向循环神经网络 (Bi-directional RNN)
 - □ 获取历史数据的信息 + 获取未来数据的信息

He said, "Teddy bears are on sale!"

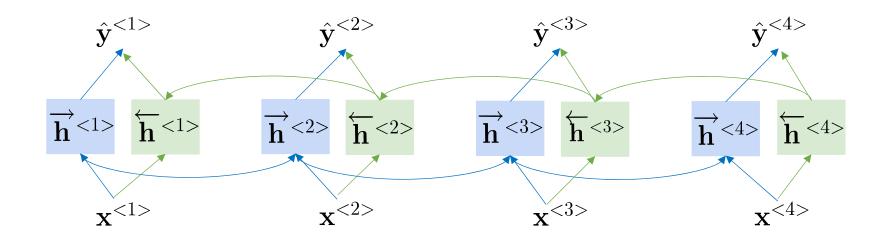
He said, "Teddy Rooseverlt was a great President!"



$$\hat{y}^{< t>} = g(\mathbf{W}_y[\overrightarrow{\mathbf{h}}^{< t>}, \overleftarrow{\mathbf{h}}^{< t>}] + \mathbf{b}_y)$$



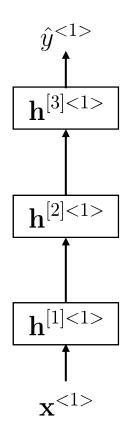
- 双向循环神经网络 (Bi-directional RNN)
 - □ 获取历史数据的信息 + 获取未来数据的信息
- ◆ 缺陷: 需要获取完整的序列(句子)才能进行预测,不能进行实时处理





■ 深层循环神经网络 (Deep RNN)

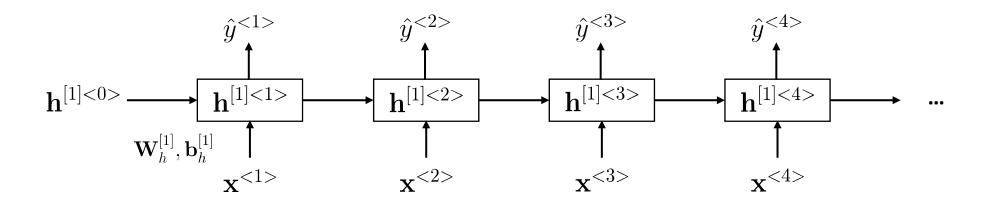
前馈神经网络





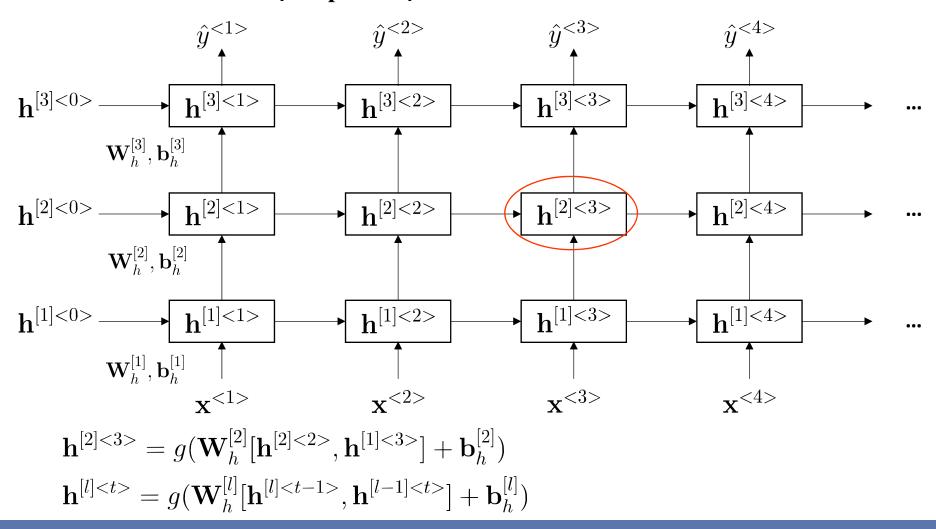
■ 深层循环神经网络 (Deep RNN)

循环神经网络

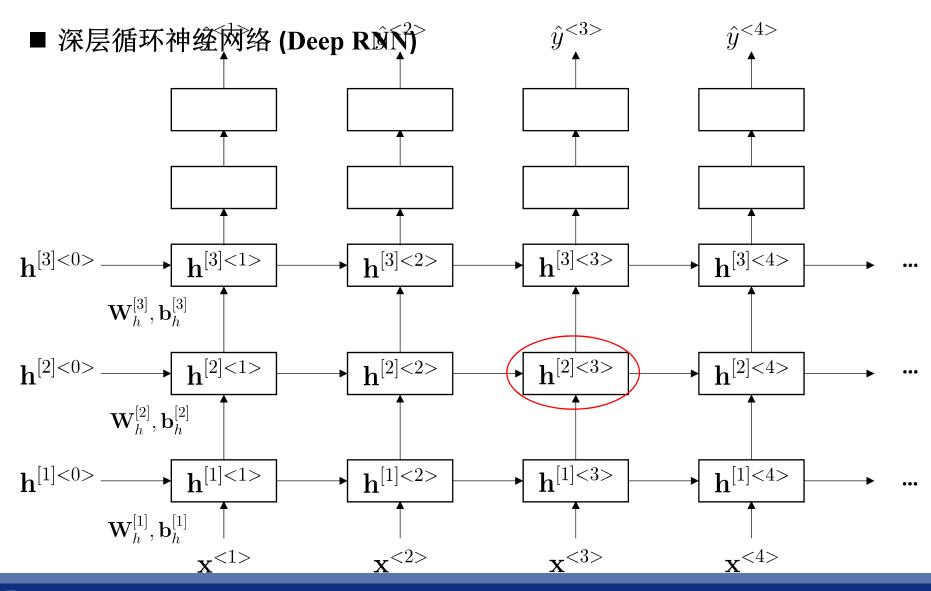




■ 深层循环神经网络 (Deep RNN)





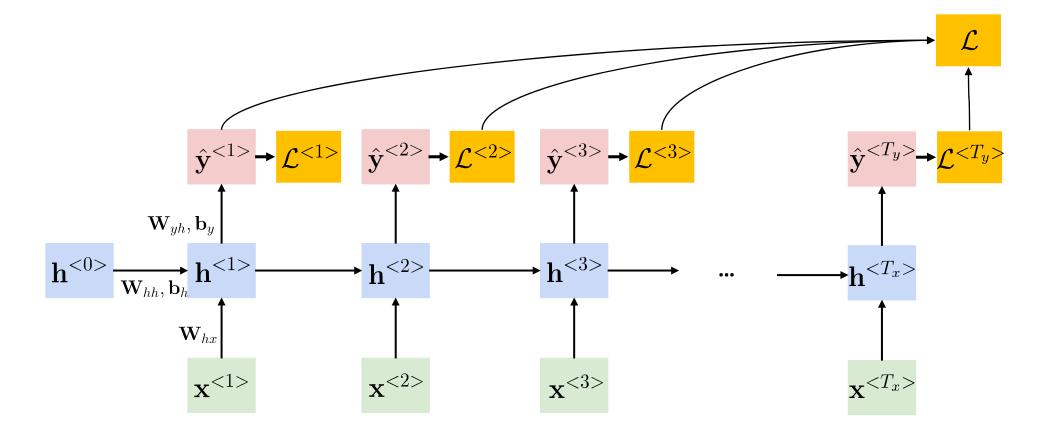






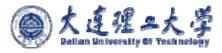
循环神经网络训练 RNN Training



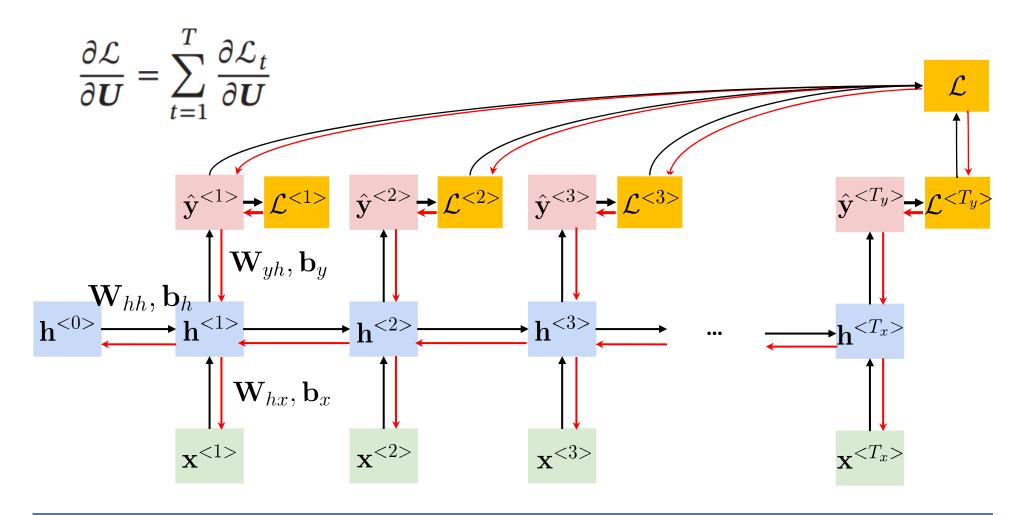


$$\mathcal{L}_t = \mathcal{L}(y_t, g(\boldsymbol{h}_t))$$

$$\mathcal{L} = \sum_{t=1}^{T} \mathcal{L}_t$$



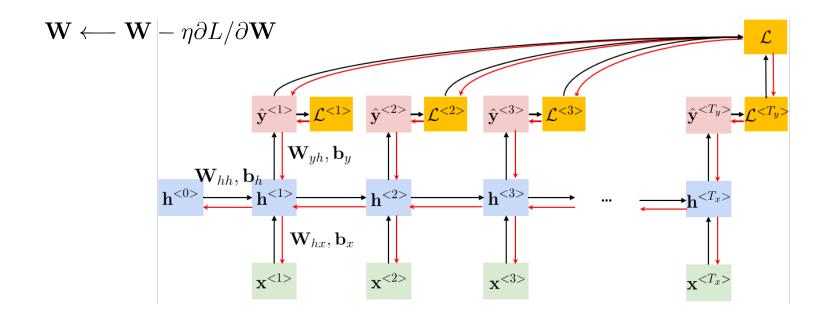
■ BPTT (Back-propagation through time): 随时间反向传播





■ BPTT (Back-propagation through time): 随时间反向传播

BPTT的中心思想和BP算法相同,本质是梯度下降法。需要寻优的参数有三个, \mathbf{W}_{hh} , \mathbf{W}_{hx} , \mathbf{W}_{yh} . 与BP算法不同的是, \mathbf{W}_{hh} , \mathbf{W}_{hx} 参数在寻优的过程中需要追溯之前历史数据,因此RNN有记忆功能。



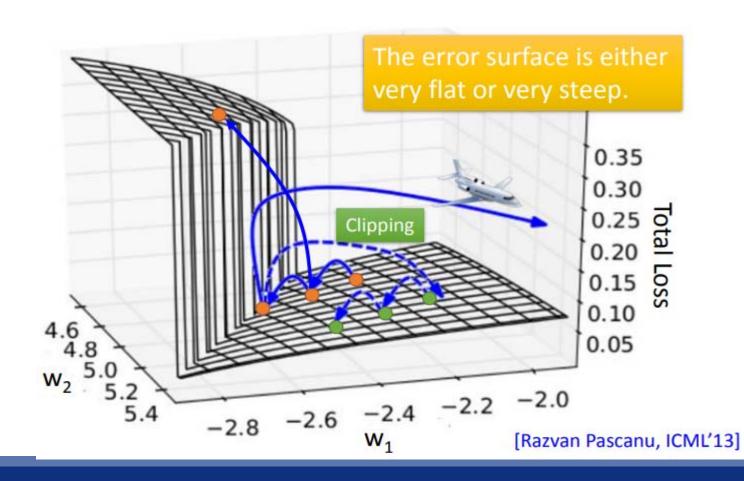


- BPTT (Back-propagation through time): 随时间反向传播
 - □ 训练RNN通常比较困难



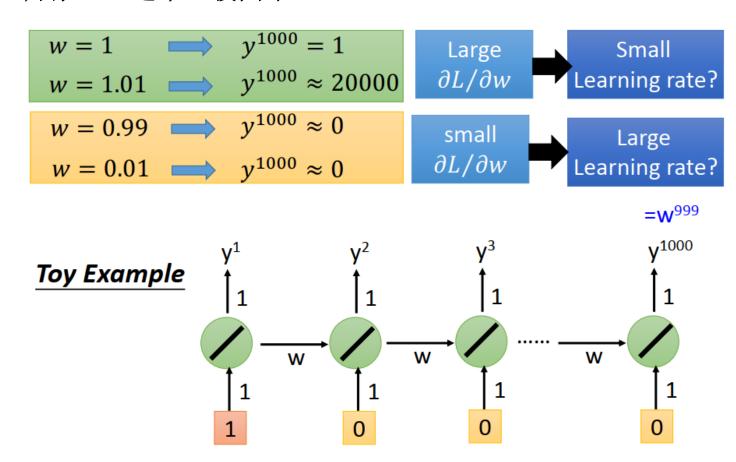


- BPTT (Back-propagation through time): 随时间反向传播
 - □ 训练RNN通常比较困难





- BPTT (Back-propagation through time): 随时间反向传播
 - □ 训练RNN通常比较困难



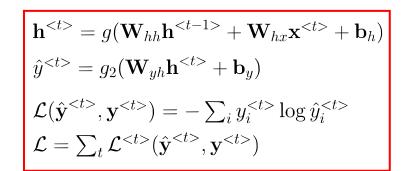


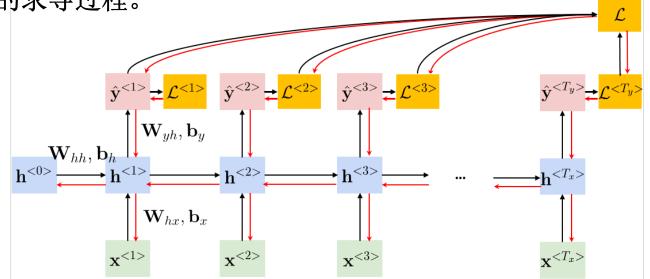
- BPTT (Back-propagation through time): 随时间反向传播
 - □ 对 W_{yh} 求偏导:

$$\frac{\partial \mathcal{L}^{\langle t \rangle}}{\partial \mathbf{W}_{yh}} = \frac{\partial \mathcal{L}^{\langle t \rangle}}{\partial \hat{y}^{\langle t \rangle}} \cdot \frac{\partial \hat{y}^{\langle t \rangle}}{\partial \mathbf{W}_{yh}}$$

其中嵌套着激活函数,是复

合函数的求导过程。

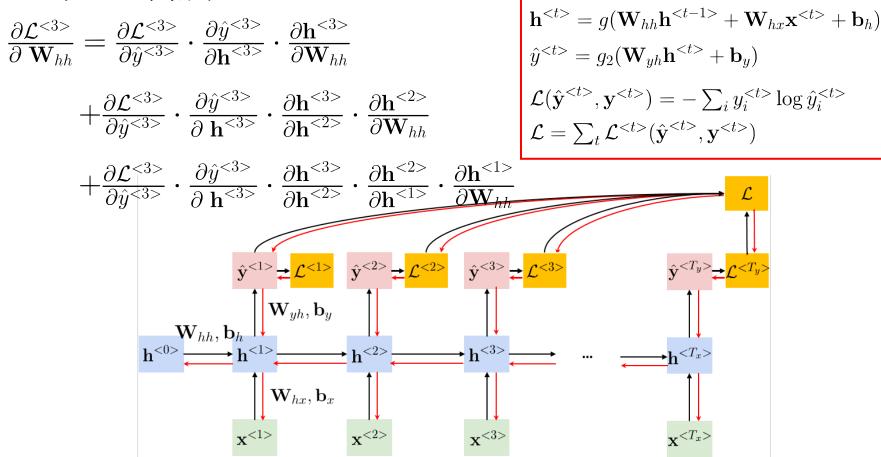






■ BPTT (Back-propagation through time): 随时间反向传播

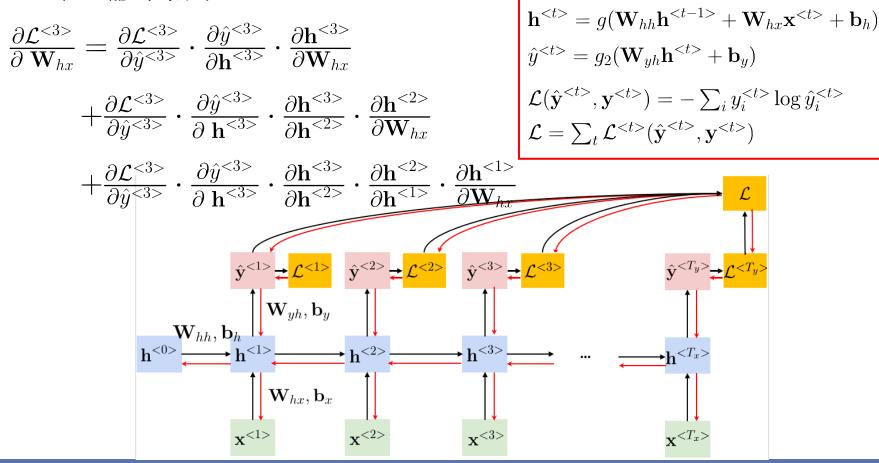
□ 对 W_{hh} 求偏导:





■ BPTT (Back-propagation through time): 随时间反向传播

□ 对 W_{hx} 求偏导:





BPTT (Back-propagation through time): 随时间反向传播

□ 对 W_{bb} 求偏导:

$$\frac{\partial \mathcal{L}^{< t>}}{\partial \mathbf{W}} = \sum_{k=0}^{t} \frac{\partial \mathcal{L}^{< t>}}{\partial \hat{\mathbf{L}}^{< t>}} \cdot \frac{\hat{y}^{< t>}}{\partial \mathbf{L}^{< t>}}$$

$$\prod_{j=k+1}^{t} (tanh)' \cdot \mathbf{W}_{hh}$$

$$\frac{\partial \mathcal{L}^{}}{\partial \mathbf{W}_{hh}} = \sum_{k=0}^{t} \frac{\partial \mathcal{L}^{}}{\partial \hat{y}^{}} \cdot \frac{\hat{y}^{}}{\partial \mathbf{h}^{}} \left(\prod_{j=k+1}^{t} \frac{\partial \mathbf{h}^{}}{\partial \mathbf{h}^{}} \right) \frac{\partial \mathbf{h}^{}}{\partial \mathbf{W}_{hh}}$$

 \square 对 \mathbf{W}_{hx} 求偏导:

$$\frac{\partial \mathcal{L}^{}}{\partial \mathbf{W}_{hx}} = \sum_{k=0}^{t} \frac{\partial \mathcal{L}^{}}{\partial \hat{y}^{}} \cdot \frac{\hat{y}^{}}{\partial \mathbf{h}^{}} \left(\prod_{j=k+1}^{t} \frac{\partial \mathbf{h}^{}}{\partial \mathbf{h}^{}} \right) \frac{\partial \mathbf{h}^{}}{\partial \mathbf{W}_{hx}}$$

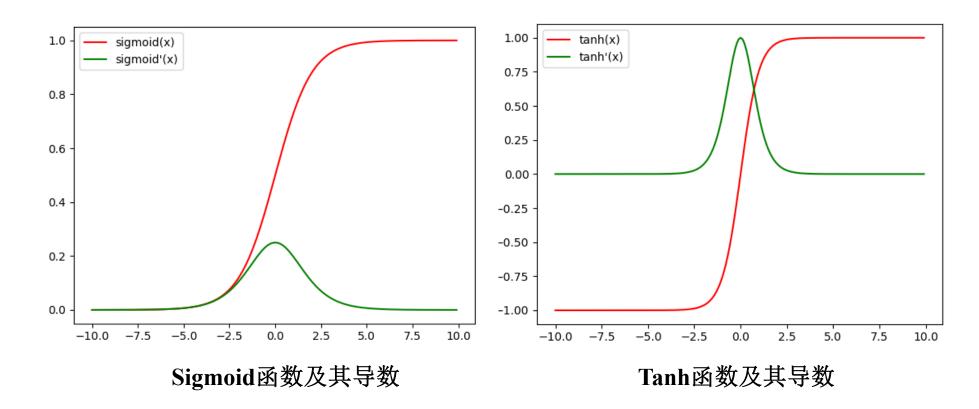
$$\prod_{j=k+1}^{t} (sigmoid)' \cdot \mathbf{W}_{hh}$$

$$\left(\prod_{j=k+1}^{t} \frac{\partial \mathbf{h}^{< j>}}{\partial \mathbf{h}^{< j-1>}}\right) \frac{\partial \mathbf{h}^{< k>}}{\partial \mathbf{W}_{hx}}$$

连乘会导致激活函数导数的连乘,进而会导 致"梯度消失"和"梯度爆炸"现象的发生

$$\mathbf{h}^{} = g(\mathbf{W}_{hh}\mathbf{h}^{} + \mathbf{W}_{hx}\mathbf{x}^{} + \mathbf{b}_h)$$
$$\hat{y}^{} = g_2(\mathbf{W}_{yh}\mathbf{h}^{} + \mathbf{b}_y)$$
$$\mathcal{L}(\hat{\mathbf{y}}^{}, \mathbf{y}^{}) = -\sum_i y_i^{} \log \hat{y}_i^{}$$
$$\mathcal{L} = \sum_t \mathcal{L}^{}(\hat{\mathbf{y}}^{}, \mathbf{y}^{})$$





Sigmoid和tanh的导数都在[0,1]范围内,随着时间序列的不断深入,小数的累乘就会导致梯度越来越小直到接近于0,这就是"梯度消失"现象。梯度消失就意味消失那一层的参数再也不更新,那么那一层隐层就变成了单纯的映射层,毫无意义了。



- 长程依赖问题
 - -- 梯度爆炸问题: 梯度削减
 - -- 梯度消失问题: 1. 选取更好的激活函数
 - 2. 改变传播结构

-- 记忆容量问题: 随着 h_t 不断累积存储新的输入信息,会发生饱和现象. 假设 $g(\cdot)$ 为 Logistic 函数,则随着时间 t 的增长, z_t 会变得越来越大,从而导致 h 变得饱和. 也就是说,隐状态 h_t 可以存储的信息是有限的,随着记忆单元存储的内容越来越多,其丢失的信息也越来越多.

$$\mathbf{h}^{} = g(\mathbf{W}_{hh}\mathbf{h}^{} + \mathbf{W}_{hx}\mathbf{x}^{} + \mathbf{b}_h)$$

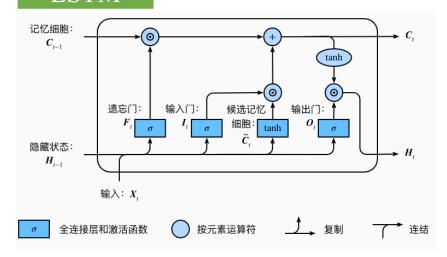


- 长程依赖问题
 - -- 梯度爆炸问题
 - -- 梯度消失问题
 - -- 记忆容量问题

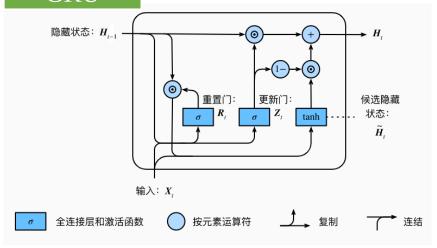
■ 如何解决?

引入门控机制来控制信息的累积 速度,包括有选择地加入新的信息, 并有选择地遗忘之前累积的信息.

LSTM



GRU







小结和资源

Take Home Message and Resource

小结



- ■RNN的结构
- ■语言模型
- Bi-RNN & Deep RNN
- ■RNN的训练
 - □随时间反向传播
 - □梯度爆炸
 - □梯度消失