



北京交通大学《深度学习》课件

第八讲 循环神经网络 I

北京交通大学 《深度学习》课程组





8.1 网络记忆能力

8.2 循环神经网络 (Recurrent Neural Network)

8.3 随时间反向传播 (BPTT)

8.4 双向循环神经网络



8.1 网络记忆能力

8.2 循环神经网络 (Recurrent Neural Network)

8.3 随时间反向传播 (BPTT)

8.4 双向循环神经网络



8.1 网络记忆能力

- 全连接神经网络和卷积神经网络
 - 只能单独的处理一个个的输入，前一个输入和后一个输入是完全没有关系的。
- 某些任务需要能够更好的处理**序列**的信息，
 - 即前面的输入和后面的输入是有关系的。
 - 比如，当理解一句话意思时，孤立的理解这句话的每个词是不够的，需要处理这些词连接起来的整个**序列**；
 - 当处理视频的时候，不能只单独的去分析每一帧，而要分析这些帧连接起来的整个**序列**。



8.1 网络记忆能力

实例1：语言模型

给定一句话前面的部分，预测接下来最有可能的一个词是什么

我昨天上学迟到了，老师批评了_____。

分词结果：

我 昨天 上学 迟到了， 老师 批评了 _____。



8.1 网络记忆能力

我昨天上学迟到了，老师批评了_____。

- 2-Gram模型

在语料库中，搜索『了』后面最可能的一个词

- 3-Gram模型

搜索『批评了』后面最可能的词

- 4-Gram, 5-Gram, ..., N-Gram

模型的大小和N的关系是指数级的，会占用海量的存储空间

而循环神经网络RNN理论上可以往前(后)看任意多个词



8.1 网络记忆能力

实例2: 槽填充(Slot Filling)问题

例如在一个自动订票系统中, 用户输入“我想八月二号到达台北”, 需要填充槽(slot),

目的地: 台北

到达时间: 2020.08.02

如何使用前馈神经网络(Feedforward network)解决此问题?



8.1 网络记忆能力

首先使用1-of-N encoding方法将每一个单词表示成一个向量。

1-of-N Encoding lexicon = {apple, bag, cat, dog, elephant}

The vector is lexicon size.

Each dimension corresponds
to a word in the lexicon

The dimension for the word
is 1, and others are 0

apple = [1 0 0 0 0]

bag = [0 1 0 0 0]

cat = [0 0 1 0 0]

dog = [0 0 0 1 0]

elephant = [0 0 0 0 1]



8.1 网络记忆能力

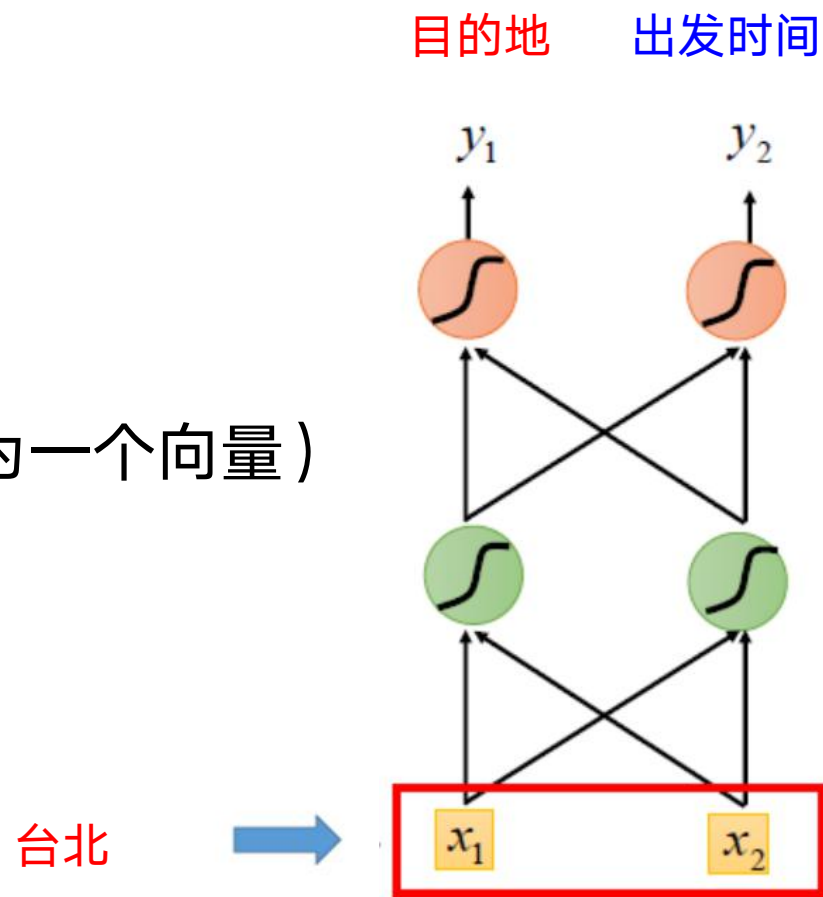
训练前馈神经网络，使得输入“台北”时，目的地的概率最大。

示例应用

通过前馈网络解决槽填充问题？

输入：一个单词 （每个单词都表示为一个向量）

输出：表示输入词属于槽的概率分布





8.1 网络记忆能力

若另一个用户输入“我想在八月二号**离开**台北”，需要填充槽(slot)，

出发地：台北

出发时间：2020.08.02

对于前馈神经网络，使得输入“台北”时，出发地的概率最大。

Problem: 对于同一个输入，输出的概率分布应该也是一样的，不可能出现既是目的地的概率最高又是出发地的概率最高。

因此，我们希望神经网络拥有“记忆”（memory）的能力，能够根据之前的信息（如**到达**或**离开**）得到不同的输出。



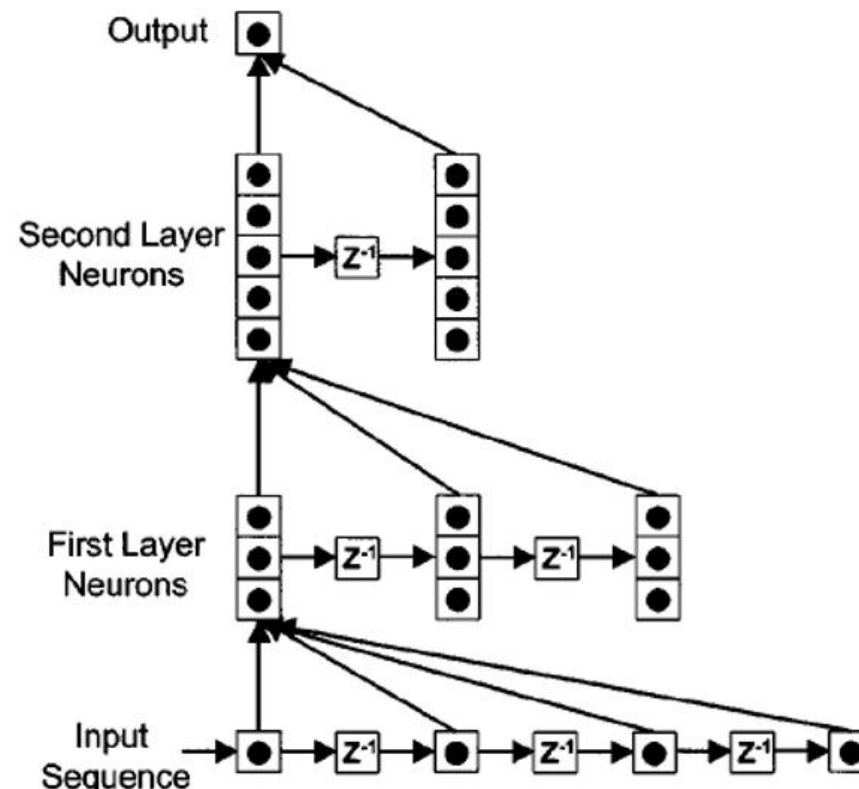
8.1 网络记忆能力

如何给网络增加记忆能力(除RNN外)?

- 延时神经网络 (Time Delay Neural Network, TDNN)
建立一个额外的延时单元, 用来存储网络的历史信息 (可以包括输入、输出、隐状态等)

$$\mathbf{h}_t^{(l)} = f\left(\mathbf{h}_t^{(l-1)}, \mathbf{h}_{t-1}^{(l-1)}, \dots, \mathbf{h}_{t-K}^{(l-1)}\right)$$

这样, 前馈网络就具有了短期记忆的能力





8.1 网络记忆能力

- 自回归模型（Autoregressive Model, AR）:一类时间序列模型，用变量 y_t 的历史信息来预测自己

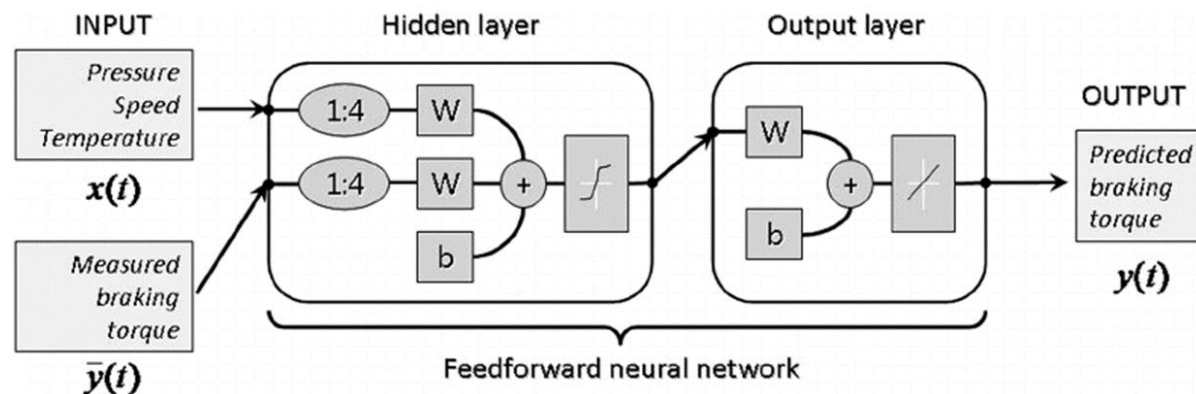
$$y_t = w_0 + \sum_{k=1}^K w_k y_{t-k} + \epsilon_t$$

$\epsilon_t \sim N(0, \sigma^2)$ 为第 t 个时刻的噪声

- 有外部输入的非线性自回归模型（Nonlinear Autoregressive with Exogenous Inputs Model, NARX）

$$y_t = f(x_t, x_{t-1}, \dots, x_{t-K_x}, y_{t-1}, y_{t-2}, \dots, y_{t-K_y})$$

其中 $f(\cdot)$ 表示非线性函数，可以是一个前馈网络， K_x 和 K_y 为超参数。





8.1 网络记忆能力

8.2 循环神经网络 (RNN)

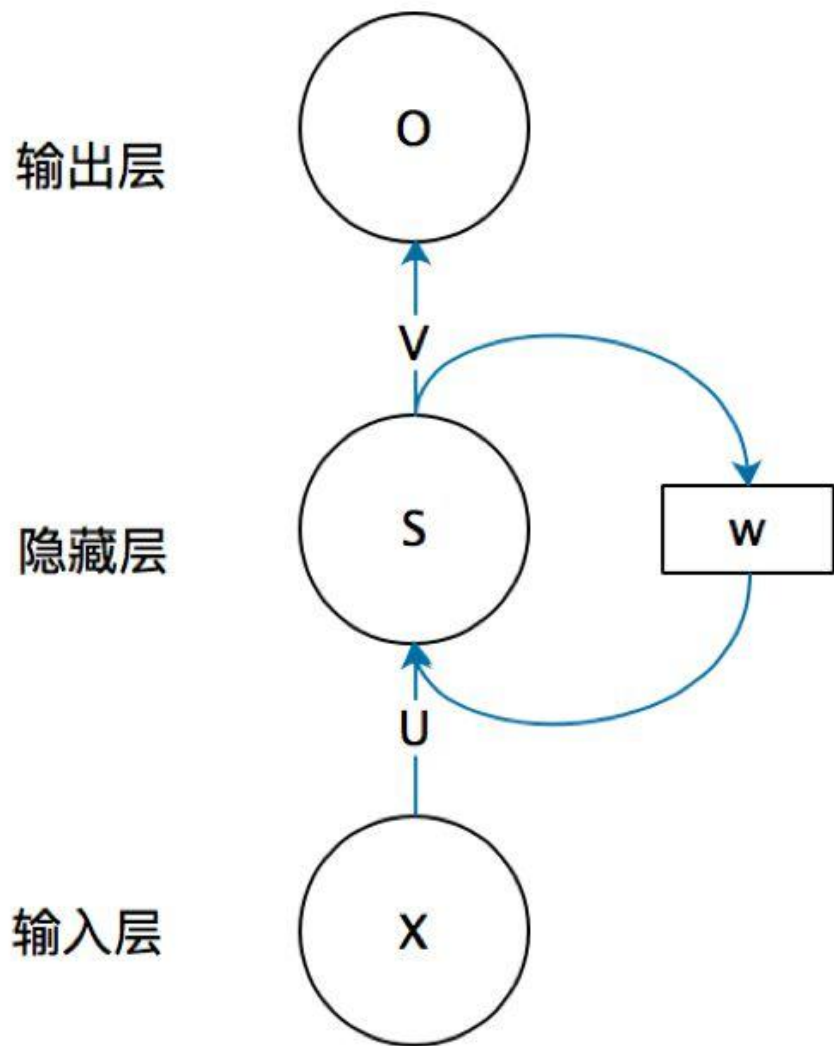
8.3 随时间反向传播 (BPTT)

8.4 双向循环神经网络



8.2 循环神经网络 (RNN)

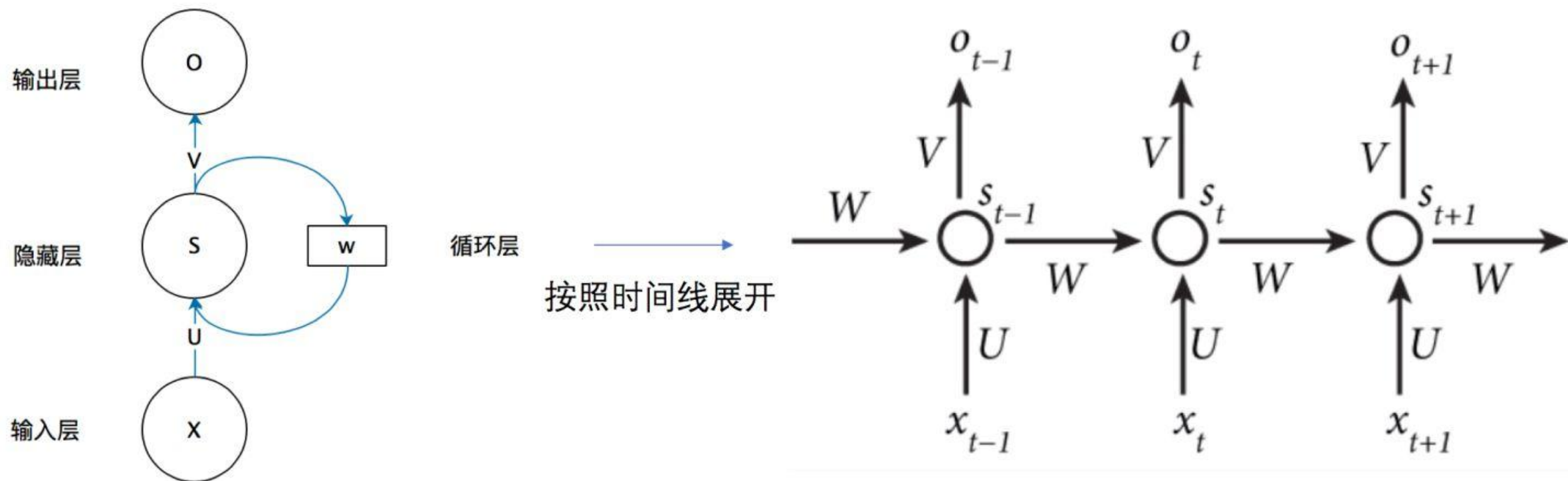
一个简单的循环神经网络由输入层、一个隐藏层和一个输出层组成。



- 去掉有 W 的带箭头的圈，它就变成普通的全连接神经网络
- x 是一个向量，它表示输入层的值
- s 是一个向量，它表示隐藏层的值， s 不仅仅取决于当前这次的输入 x ，还取决于上一次隐藏层的值
- o 也是一个向量，它表示输出层的值
- U 是输入层到隐藏层的权重矩阵
- V 是隐藏层到输出层的权重矩阵
- W 是隐藏层上一次的值作为这一次的输入的权重矩阵

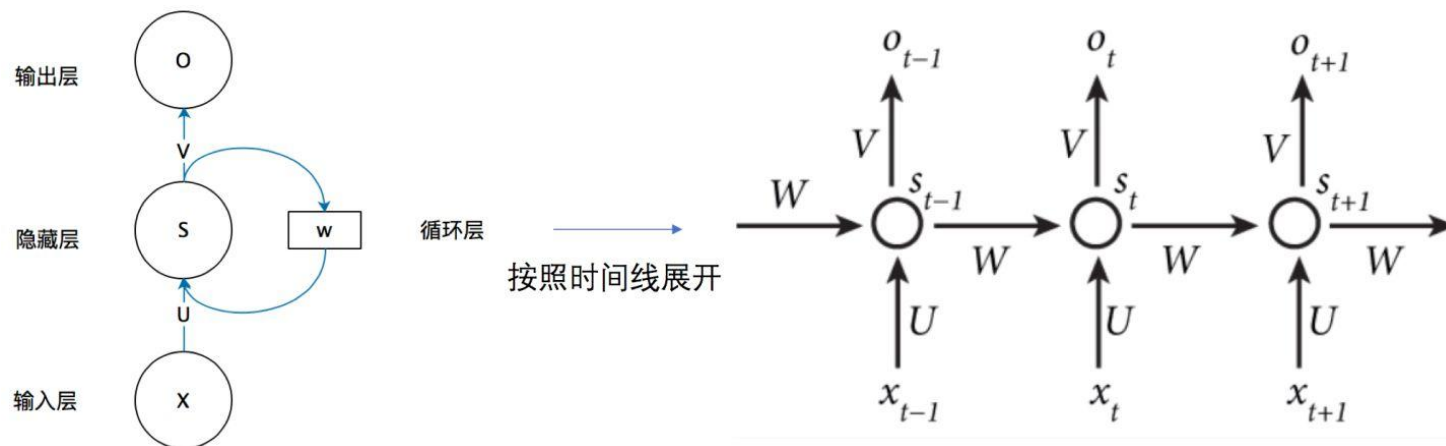


8.2 循环神经网络 (RNN)





8.2 循环神经网络 (RNN)



网络在 t 时刻接收到输入 x_t 之后，隐藏层的值是 s_t ，输出值是 o_t 。 s_t 的值不仅仅取决于 x_t ，还取决于 s_{t-1} 。

循环神经网络的计算：

$$o_t = g(Vs_t) \quad (\text{式1})$$

$$s_t = f(Ux_t + Ws_{t-1}) \quad (\text{式2})$$



8.2 循环神经网络 (RNN)

循环神经网络可以往前看任意多个输入值的原因

$$o_t = g(Vs_t) \quad (\text{式1})$$

$$s_t = f(Ux_t + Ws_{t-1}) \quad (\text{式2})$$

如果反复将（式2）带入（式1），我们将得到

$$\begin{aligned} o_t &= g(Vs_t) \\ &= g(Vf(Ux_t + Ws_{t-1})) \\ &= g(Vf(Ux_t + Wf(Ux_{t-1} + Ws_{t-2}))) \\ &= g(Vf(Ux_t + Wf(Ux_{t-1} + Wf(Ux_{t-2} + Ws_{t-3})))) \\ &= g(Vf(Ux_t + Wf(Ux_{t-1} + Wf(Ux_{t-2} + Wf(Ux_{t-3} + \dots))))) \end{aligned}$$

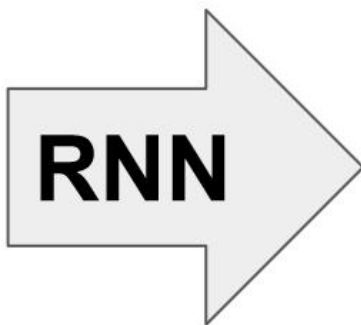
即输出值 o_t ，受前面历次输入值 x_t ， x_{t-1} 、 x_{t-2} 、 x_{t-3} ...的影响。



8.2 循环神经网络 (RNN)

RNN是一个序列模型

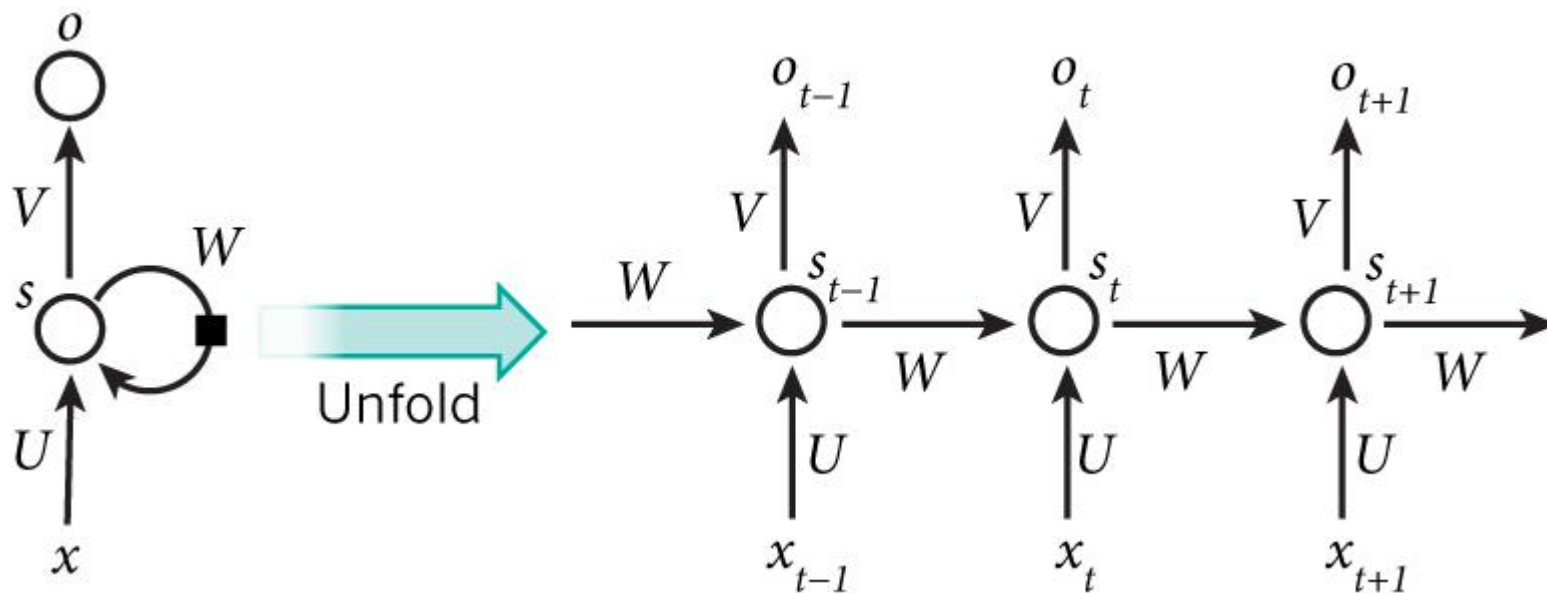
- 如何将可变长度序列作为输入？
- 如何预测可变长度序列作为输出？





8.2 循环神经网络 (RNN)

神经网络模型通过训练“学”到的东西蕴含在“权值”中。基础的神经网络只在层与层之间建立权连接，而RNN最大的不同之处在于层之间的神经元也建立了权连接。



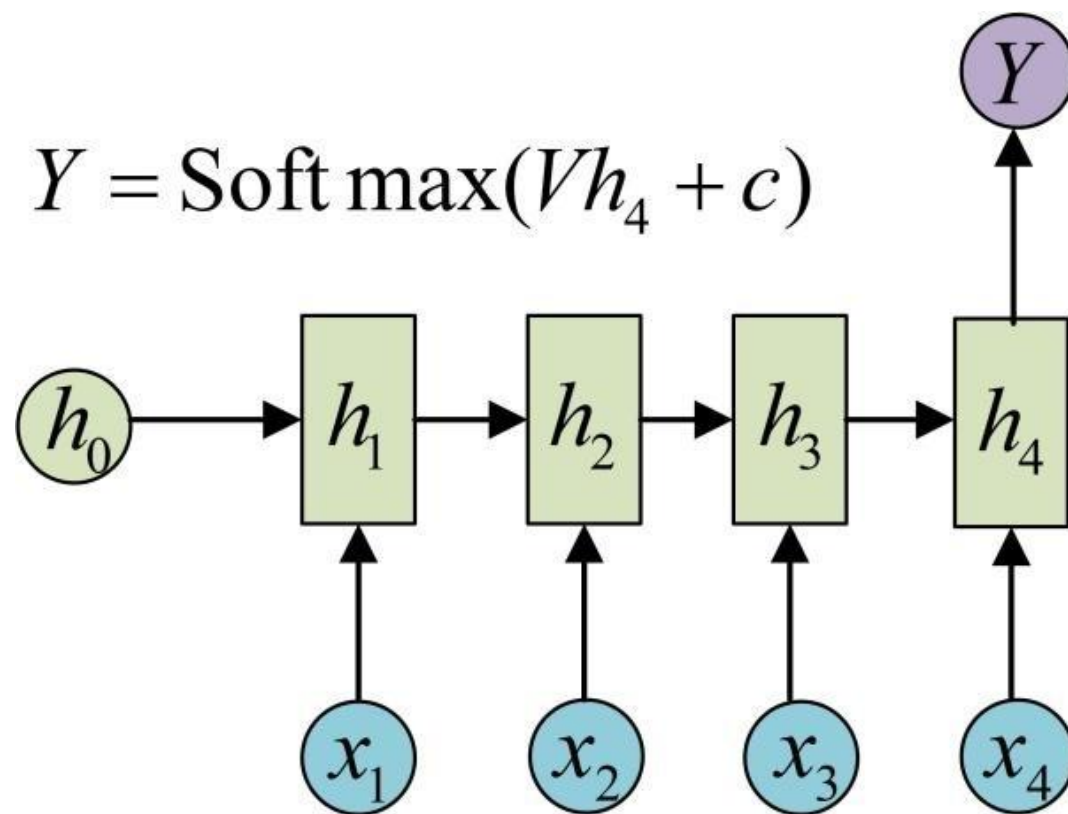
RNN的标准结构



8.2 循环神经网络 (RNN)

Many-to-One模型

RNN的标准结构在实际中并不能解决所有问题，例如，输入为一串文字，输出为分类类别。

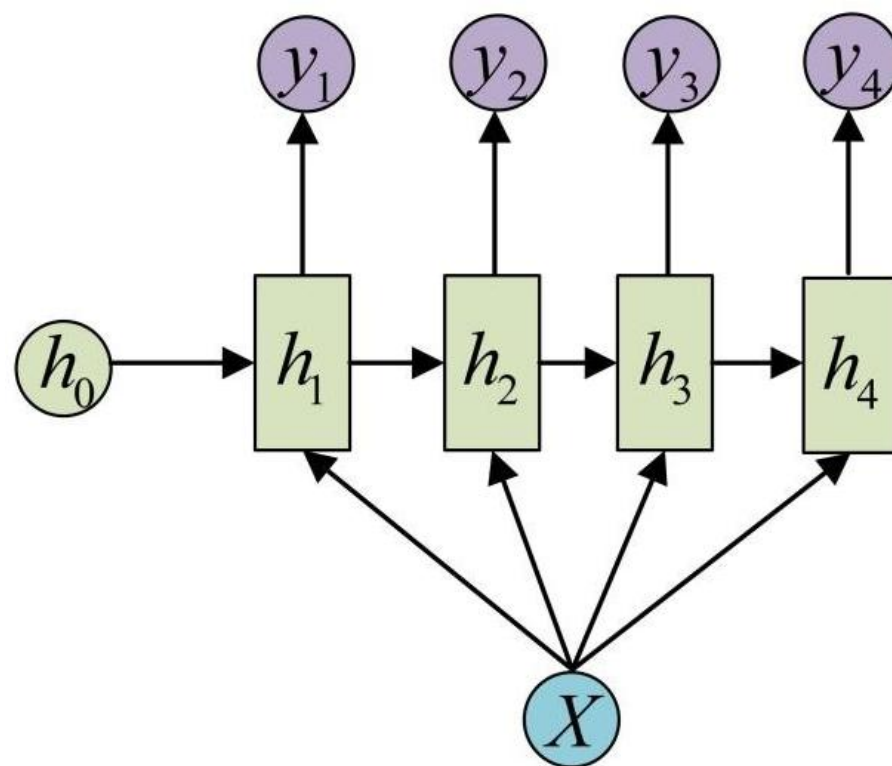
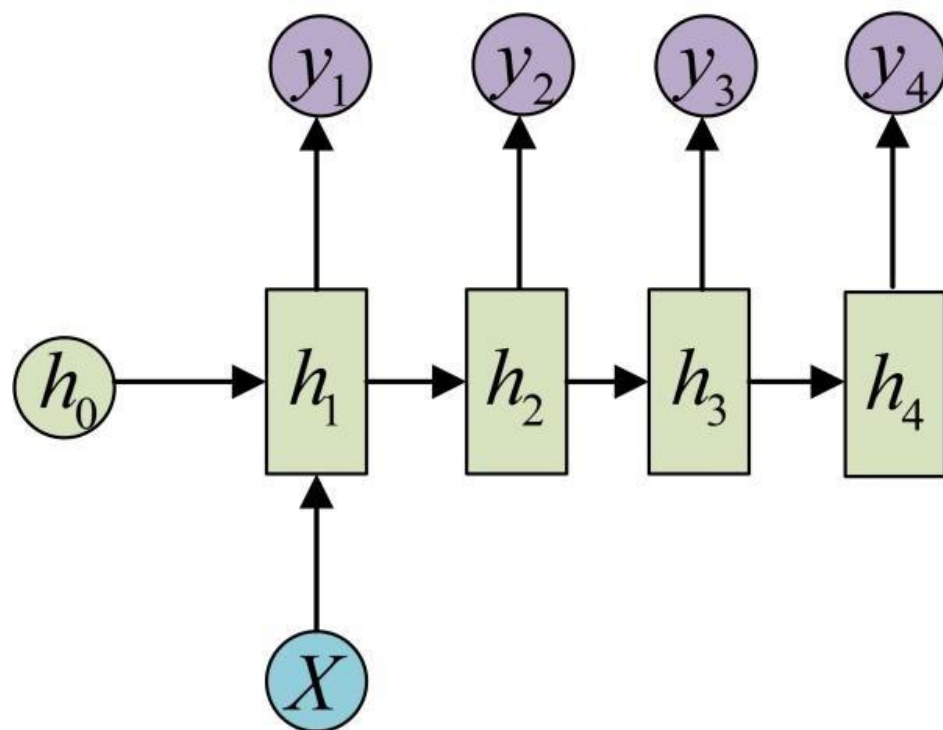




8.2 循环神经网络 (RNN)

One-to-Many模型

如输入为一张图片，输出为图片的文字描述

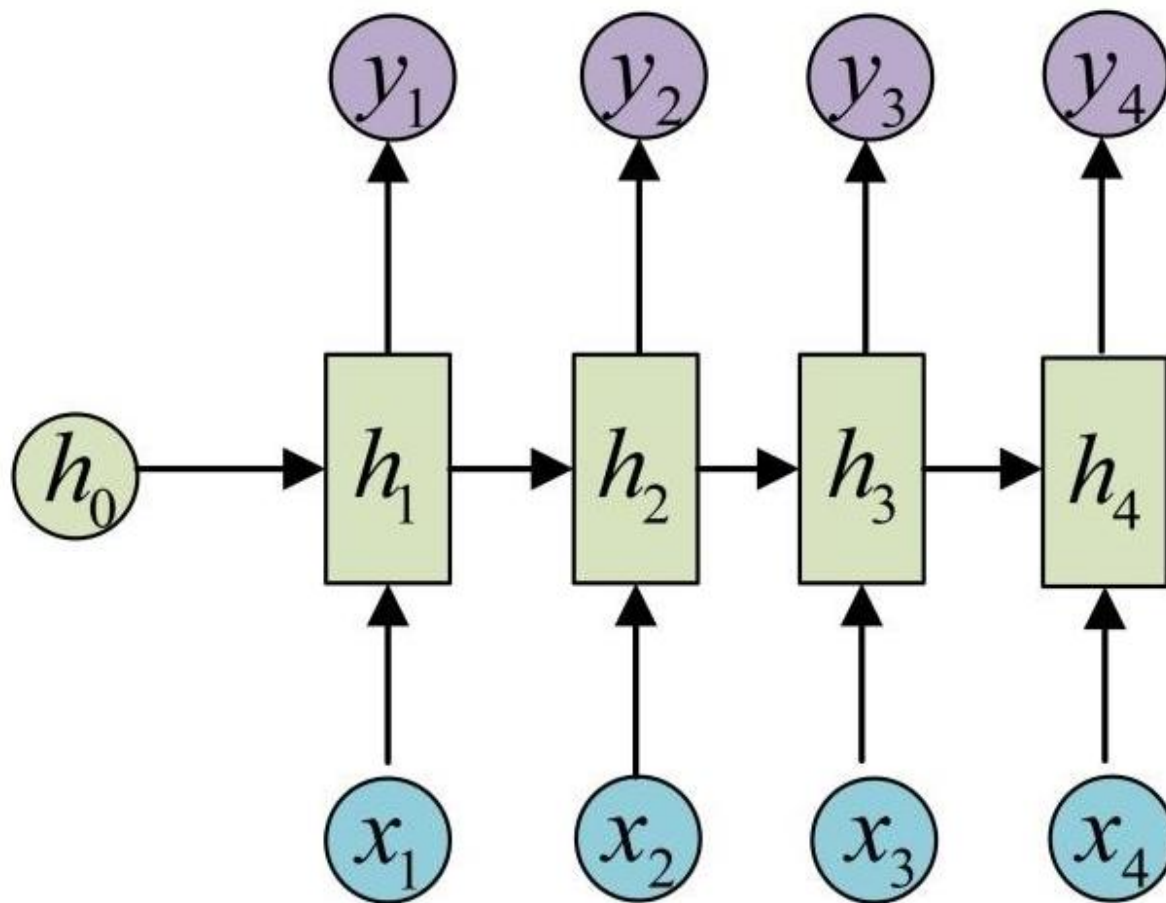




8.2 循环神经网络 (RNN)

Many-to-Many模型

输入和输出的序列个数相同，如输入为视频序列，输出为每一帧对应的标签

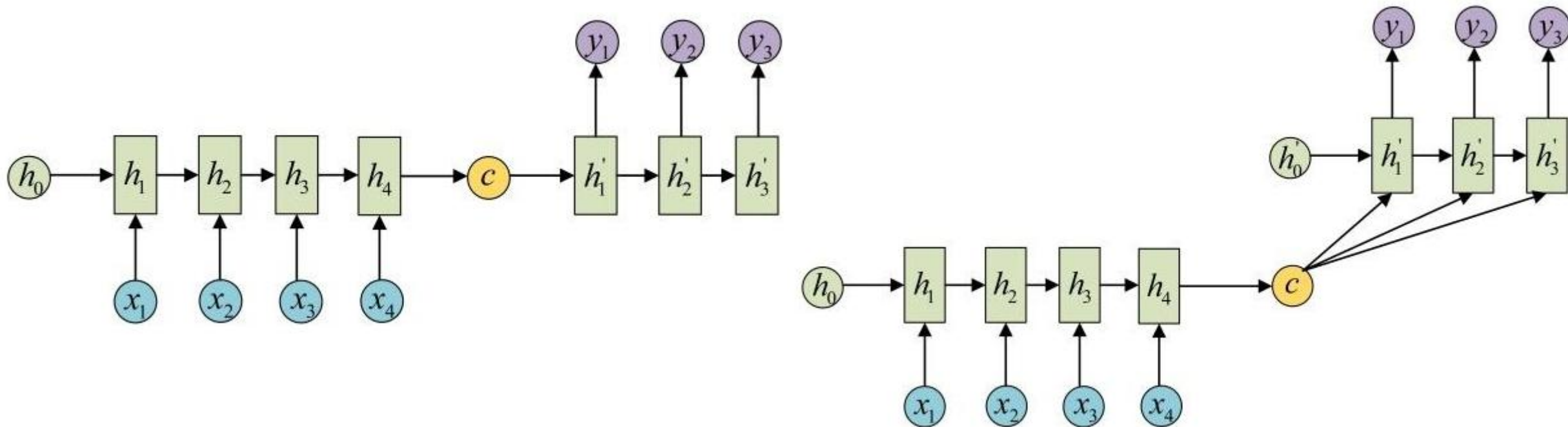




8.2 循环神经网络 (RNN)

Many-to-Many模型

输入和输出的序列个数不同，如机器翻译中，源语言和目标语言的句子往往并没有相同的长度。

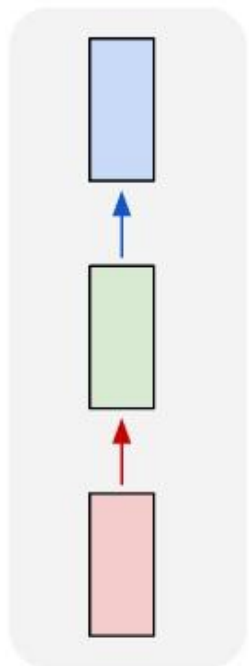


这种结构又叫Encoder-Decoder模型，也可以称之为Seq2Seq模型

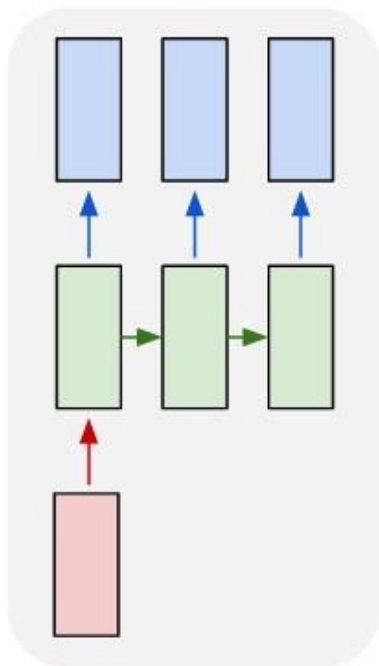


8.2 循环神经网络 (RNN)

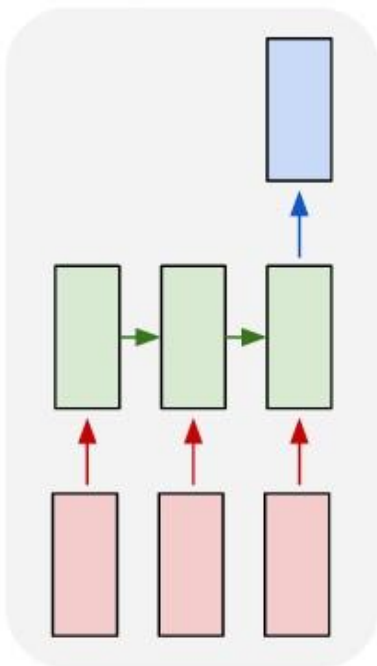
one to one



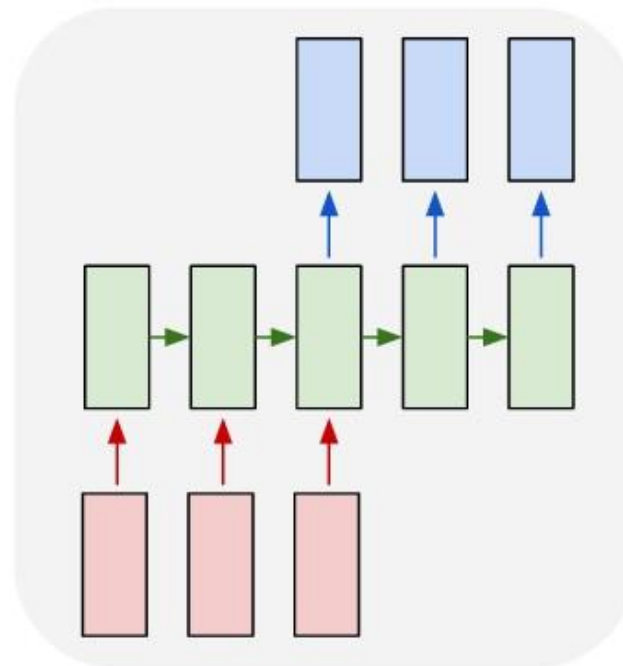
one to many



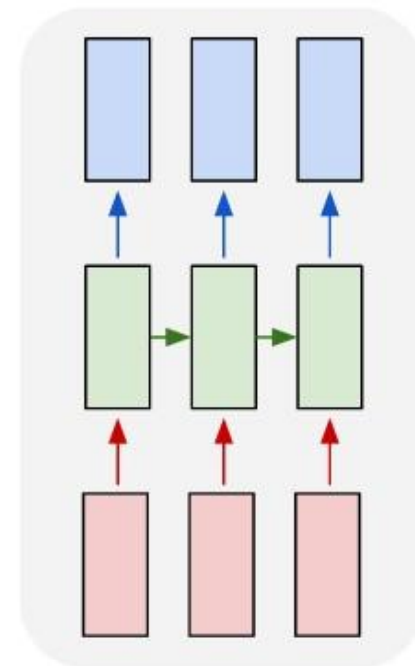
many to one



many to many



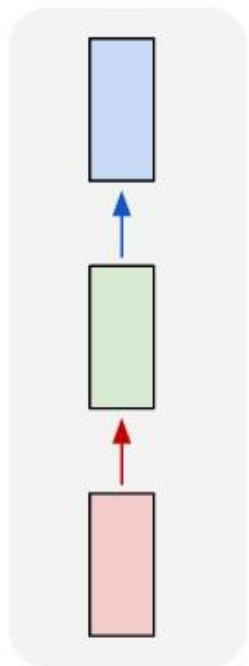
many to many



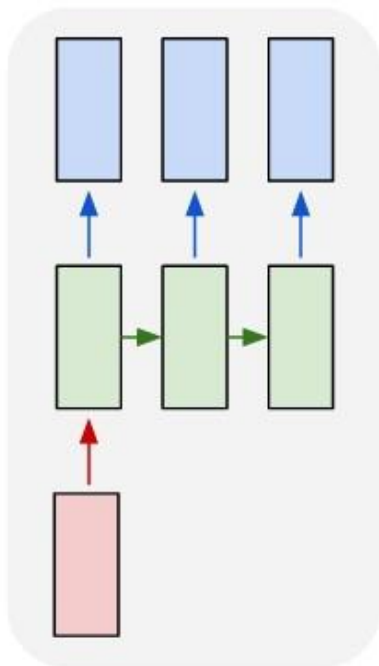


8.2 循环神经网络 (RNN)

one to one

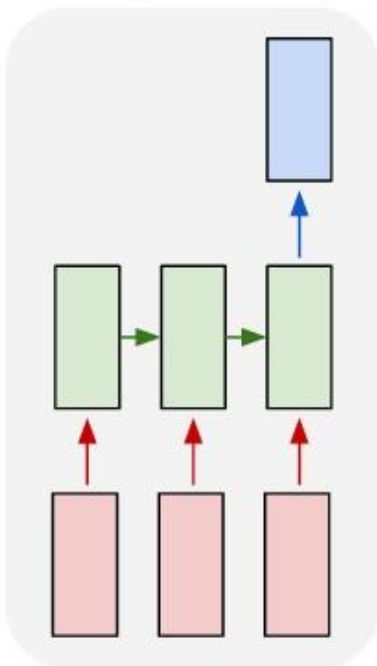


one to many



看图说话

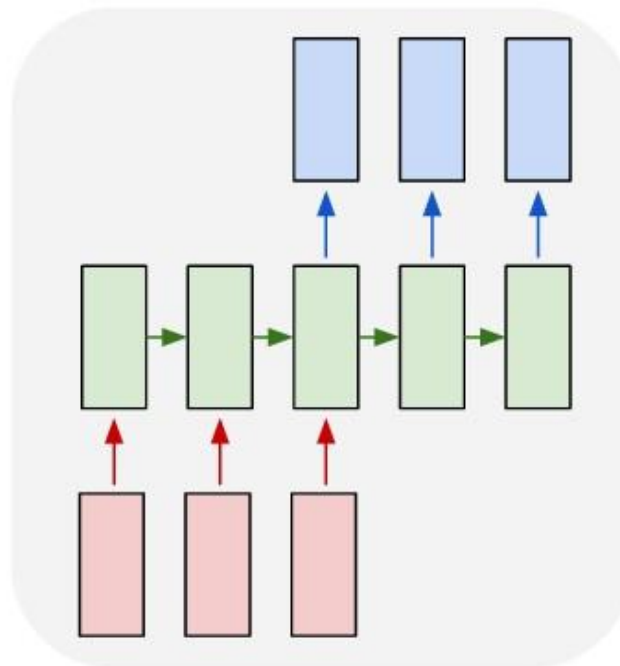
many to one



情感分析

文本分类

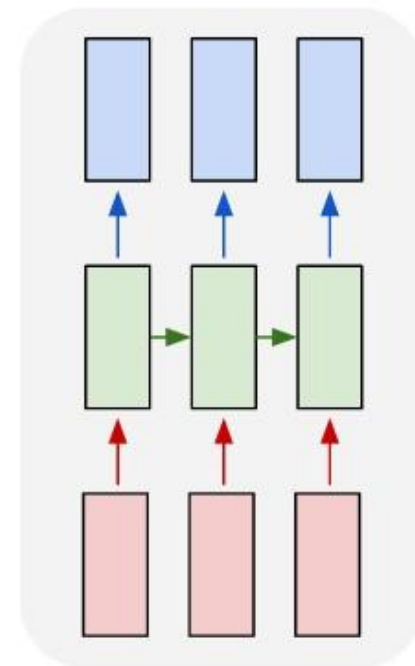
many to many



自动问答

机器翻译

many to many

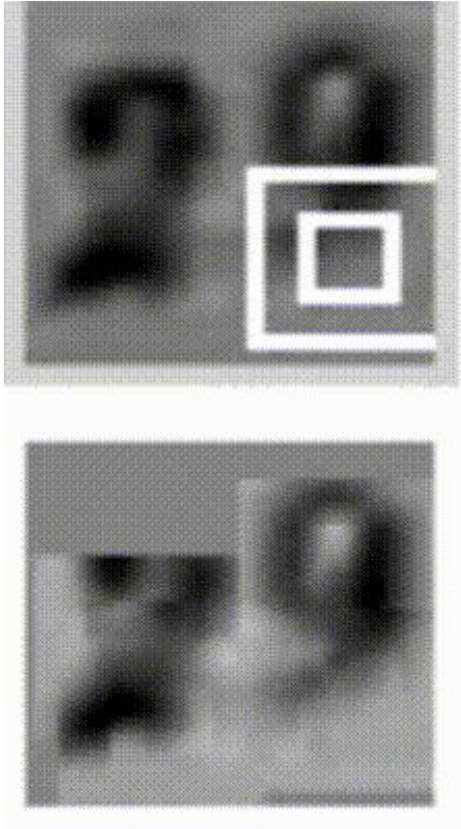


对视频的每一帧
打标签

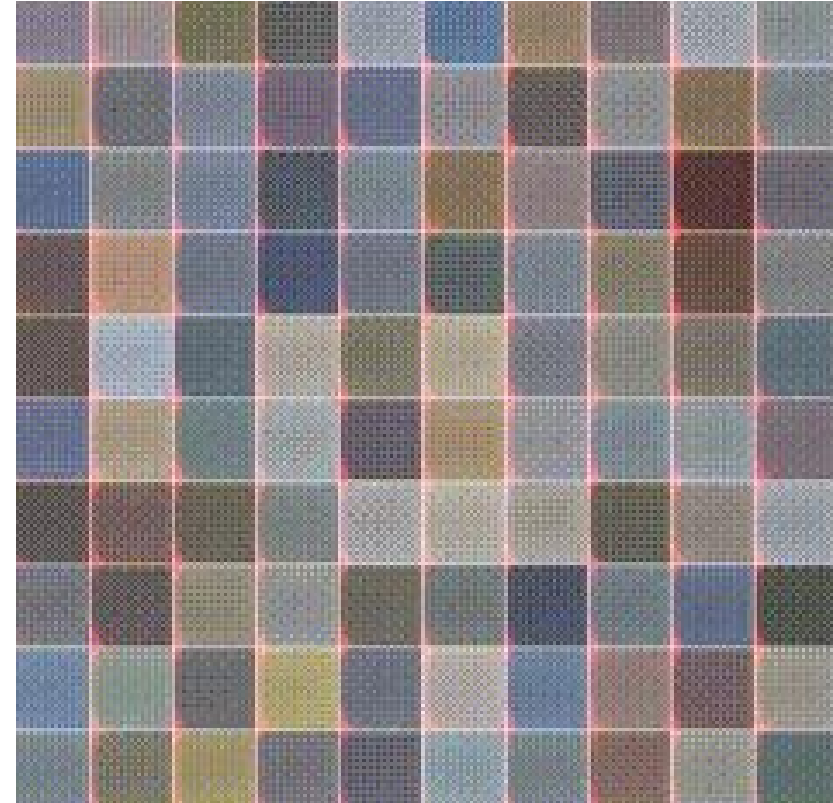


8.2 循环神经网络 (RNN)

RNN模型应用举例



RNN读门牌号



RNN绘制门牌号



8.2 循环神经网络 (RNN)

RNN模型应用举例



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"girl in pink dress is jumping in air."



"black and white dog jumps over bar."



"young girl in pink shirt is swinging on swing."

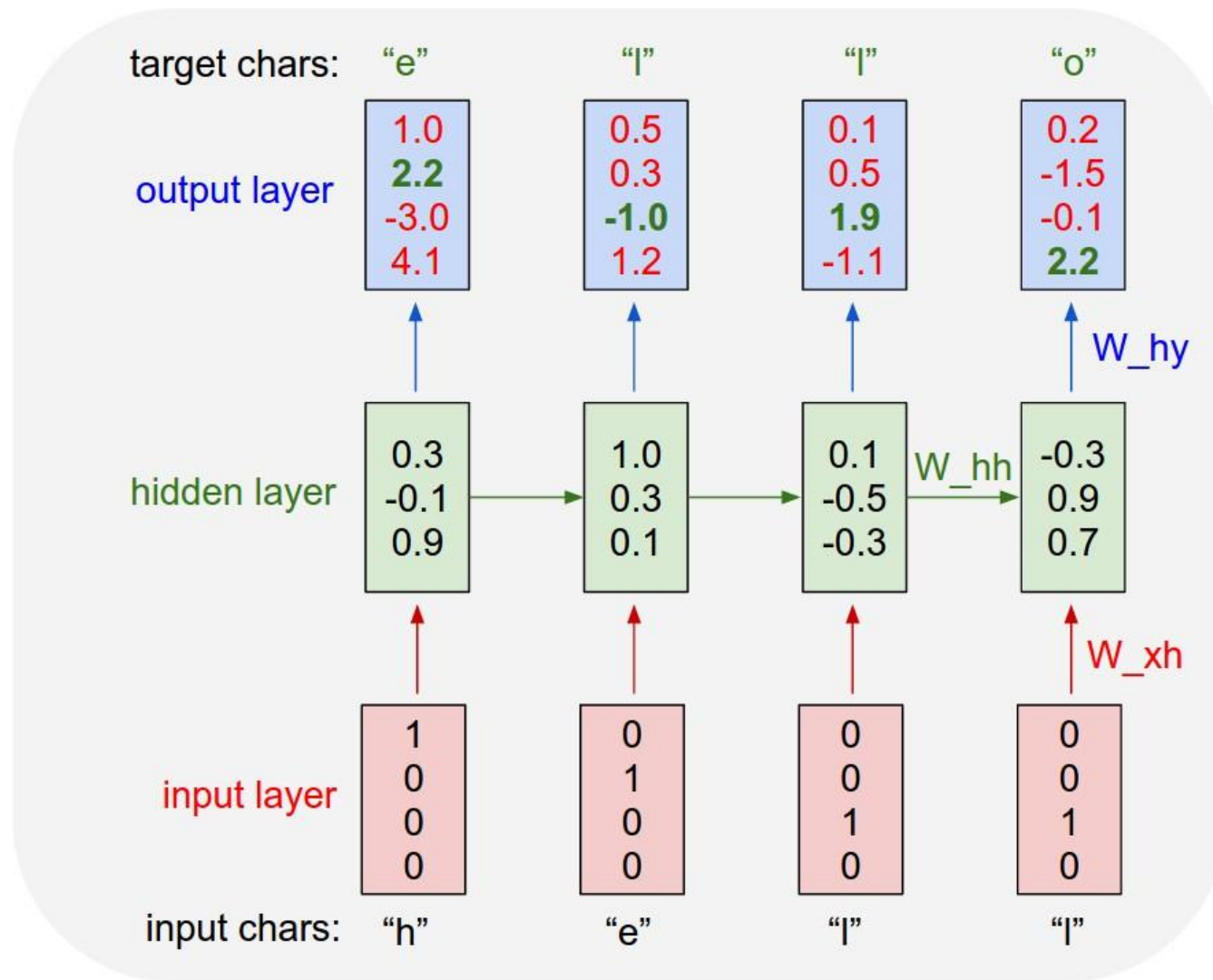


"man in blue wetsuit is surfing on wave."



8.2 循环神经网络 (RNN)

RNN模型应用举例



字符级别RNN



8.2 循环神经网络 (RNN)

RNN模型应用举例

PANDARUS:

Alas, I think he shall be come approached and the day
When little strain would be attain'd into being never fed,
And who is but a chain and subjects of his death,
I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul,
Breaking and strongly should be buried, when I perish
The earth and thoughts of many states.

DUKE VINCENTIO:

Well, your wit is in the care of side and that.

Second Lord:

They would be ruled after this chamber, and
my fair nudes begun out of the fact, to be conveyed,
Whose noble souls I'll have the heart of the wars.

Clown:

Come, sir, I will make did behold your worship.

VIOLA:

I'll drink it.

Shakespeare



8.2 循环神经网络 (RNN)

RNN模型应用举例

[玉辇临遐远，林台向晚开。寒霜万里雨，古木白云台。杞梓长晴属，藩旌旋北来。朝山犹似盖，天仆日萧纷。]

[银阁为君贵，昭原气极功。不知能似面，合合不甘穷。]

[从容峰岸啸初开，曾道神功月会频。只为钱陵庄叟旧，不曾知我有何人。]

[梧桐叶暗如丝拍，一片寒波鬓似丝。与子近来弹铗嫁，羞他画领学蛾眉。]

[渴戍平千重帝城，郊园初见杜陵僧。壮仙若望煮舟马，石上萧遮两袖僧。]

RNN写诗



8.1 网络记忆能力

8.2 循环神经网络 (RNN)

8.3 随时间反向传播 (BPTT)

8.4 双向循环神经网络

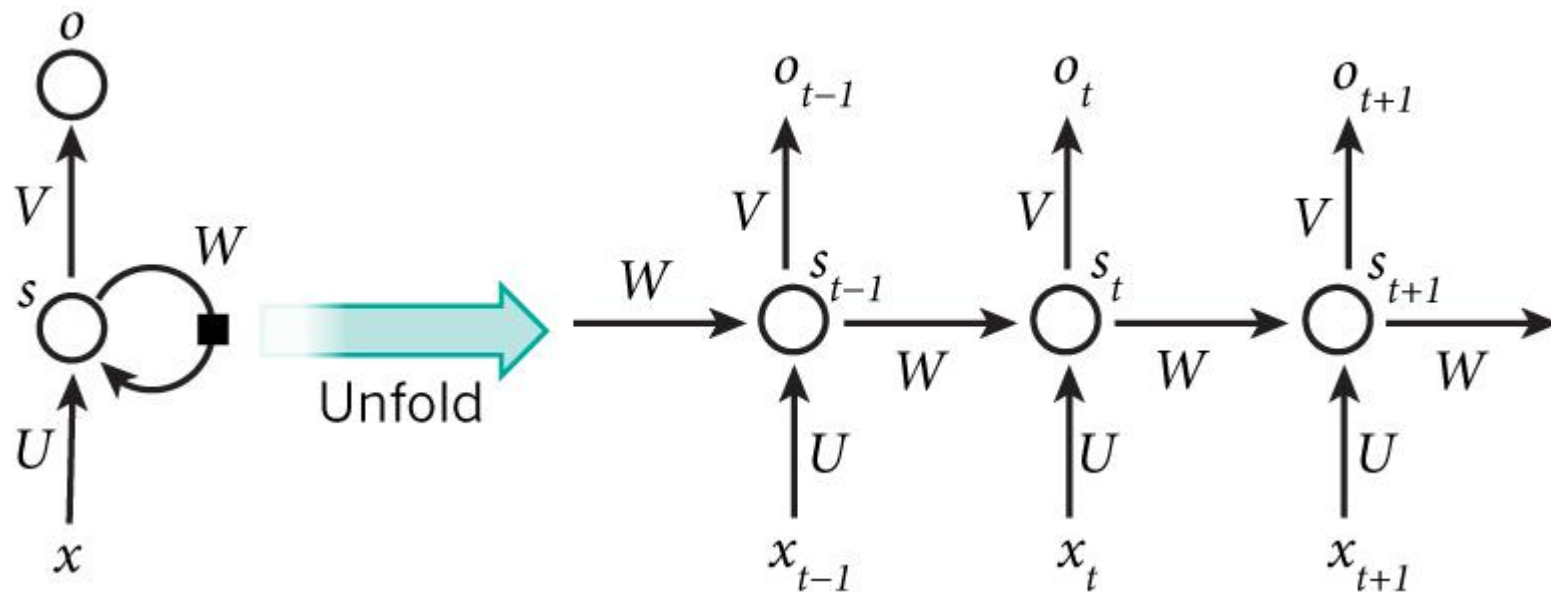


8.3 随时间反向传播 (BPTT)

- BPTT (Back-Propagation Through Time) 算法是常用的训练RNN的方法，因为RNN处理时间序列数据，所以要基于时间反向传播，故称作随时间反向传播。
- BPTT的中心思想是沿着需要优化的参数的负梯度方向不断寻找更优的点直至收敛（和BP算法相同）。
- BPTT算法本质上还是BP算法，BP算法本质是梯度下降法，那么求各个参数的梯度便成了此算法的核心。



8.3 随时间反向传播 (BPTT)



需要寻优的参数有三个，分别是 U 、 V 、 W 。与BP算法不同的是，权重矩阵 W 和 U 的寻优过程需要追溯之前的历史数据（BPTT算法的重点）。

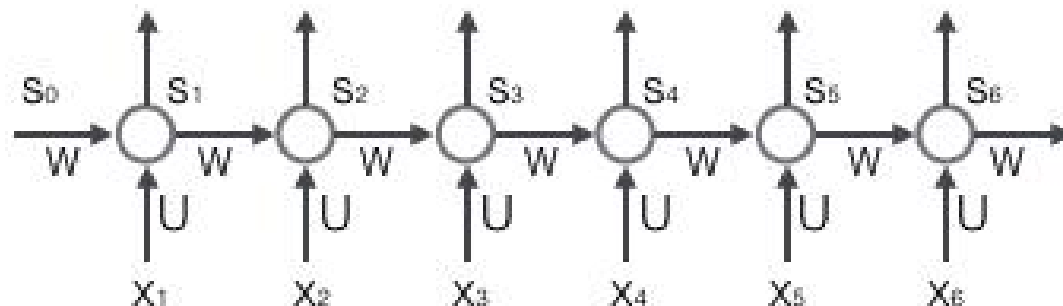


8.3 随时间反向传播 (BPTT)

BPTT算法是针对**循环层**的训练算法，它的基本原理和BP算法是一样的，也包含同样的三个步骤：

1. 前向计算每个神经元的输出值；
2. 反向计算每个神经元的**误差项** δ_j 值，它是误差函数E对神经元 j 的**加权输入** net_j 的偏导数；
3. 计算每个权重的梯度。

最后再用**随机梯度下降**算法更新权重。



RNN的循环层



8.3 随时间反向传播 (BPTT)

A. 前向计算

使用前面的**式2**对循环层进行前向计算：

$$\mathbf{s}_t = f(U\mathbf{x}_t + W\mathbf{s}_{t-1})$$

注意，上面的 \mathbf{s}_t 、 \mathbf{x}_t 、 \mathbf{s}_{t-1} 都是向量，用**黑体字母**表示；而 U 、 V 、 W 是矩阵，用大写字母表示。**向量**的下标表示时刻，例如， \mathbf{s}_t 表示在 t 时刻向量 \mathbf{s} 的值。



8.3 随时间反向传播 (BPTT)

A. 前向计算

我们假设输入向量 \mathbf{x} 的维度是 m ，输出向量 \mathbf{s} 的维度是 n ，则矩阵 U 的维度是 $n \times m$ ，矩阵 W 的维度是 $n \times n$ ，将上式展开成矩阵如下：

$$\begin{bmatrix} s_1^t \\ s_2^t \\ \vdots \\ s_n^t \end{bmatrix} = f \left(\begin{bmatrix} u_{11} & u_{12} & \dots & u_{1m} \\ u_{21} & u_{22} & \dots & u_{2m} \\ & & \ddots & \\ & & & u_{n1} & u_{n2} & \dots & u_{nm} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} + \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1n} \\ w_{21} & w_{22} & \dots & w_{2n} \\ & & \ddots & \\ & & & w_{n1} & w_{n2} & \dots & w_{nn} \end{bmatrix} \begin{bmatrix} s_1^{t-1} \\ s_2^{t-1} \\ \vdots \\ s_n^{t-1} \end{bmatrix} \right)$$



8.3 随时间反向传播 (BPTT)

B. 误差项的计算

BPTT算法将第 l 层 t 时刻的误差项 δ_t^l 值沿两个方向传播，一个方向是其传递到上一层网络，得到 δ_t^{l-1} ，这部分只和权重矩阵 U 有关；另一个方向是将其沿时间线传递到初始 t_1 时刻，得到 δ_1^l ，这部分只和权重矩阵 W 有关。

我们用向量 net_t 表示神经元在 t 时刻的加权输入，因为：

$$\begin{aligned}\text{net}_t &= U\mathbf{x}_t + W\mathbf{s}_{t-1} \\ \mathbf{s}_{t-1} &= f(\text{net}_{t-1})\end{aligned}$$

因此：

$$\frac{\partial \text{net}_t}{\partial \text{net}_{t-1}} = \frac{\partial \text{net}_t}{\partial \mathbf{s}_{t-1}} \frac{\partial \mathbf{s}_{t-1}}{\partial \text{net}_{t-1}}$$



8.3 随时间反向传播 (BPTT)

我们用 a 表示列向量，用 a^T 表示行向量，上式的第一项是向量函数对向量求导，其结果为Jacobian矩阵：

$$\begin{aligned}\frac{\partial \text{net}_t}{\partial s_{t-1}} &= \begin{bmatrix} \frac{\partial \text{net}_1^t}{\partial s_1^{t-1}} & \frac{\partial \text{net}_1^t}{\partial s_2^{t-1}} & \cdots & \frac{\partial \text{net}_1^t}{\partial s_n^{t-1}} \\ \frac{\partial \text{net}_2^t}{\partial s_1^{t-1}} & \frac{\partial \text{net}_2^t}{\partial s_2^{t-1}} & \cdots & \frac{\partial \text{net}_2^t}{\partial s_n^{t-1}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial \text{net}_n^t}{\partial s_1^{t-1}} & \frac{\partial \text{net}_n^t}{\partial s_2^{t-1}} & \cdots & \frac{\partial \text{net}_n^t}{\partial s_n^{t-1}} \end{bmatrix} \\ &= \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1n} \\ w_{21} & w_{22} & \cdots & w_{2n} \\ \vdots & \vdots & & \vdots \\ w_{n1} & w_{n2} & \cdots & w_{nn} \end{bmatrix} \\ &= W\end{aligned}$$



8.3 随时间反向传播 (BPTT)

同理，上式第二项也是一个Jacobian矩阵：

$$\frac{\partial \mathbf{s}_{t-1}}{\partial \mathbf{net}_{t-1}} = \begin{bmatrix} \frac{\partial s_1^{t-1}}{\partial net_1^{t-1}} & \frac{\partial s_1^{t-1}}{\partial net_2^{t-1}} & \dots & \frac{\partial s_1^{t-1}}{\partial net_n^{t-1}} \\ \frac{\partial s_2^{t-1}}{\partial net_1^{t-1}} & \frac{\partial s_2^{t-1}}{\partial net_2^{t-1}} & \dots & \frac{\partial s_2^{t-1}}{\partial net_n^{t-1}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial s_n^{t-1}}{\partial net_1^{t-1}} & \frac{\partial s_n^{t-1}}{\partial net_2^{t-1}} & \dots & \frac{\partial s_n^{t-1}}{\partial net_n^{t-1}} \end{bmatrix}$$

其中， $\text{diag}[\mathbf{a}]$ 表示根据向量 \mathbf{a} 创建一个对角矩阵，即

$$\text{diag}(\mathbf{a}) = \begin{bmatrix} a_1 & 0 & \dots & 0 \\ 0 & a_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & a_n \end{bmatrix}$$

$$= \begin{bmatrix} f'(net_1^{t-1}) & 0 & \dots & 0 \\ 0 & f'(net_2^{t-1}) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & f'(net_n^{t-1}) \end{bmatrix}$$

$$= \text{diag}[f'(\mathbf{net}_{t-1})]$$



8.3 随时间反向传播 (BPTT)

最后，将两项合在一起，可得：

$$\begin{aligned}\frac{\partial \text{net}_t}{\partial \text{net}_{t-1}} &= \frac{\partial \text{net}_t}{\partial s_{t-1}} \frac{\partial s_{t-1}}{\partial \text{net}_{t-1}} \\ &= W \text{diag}[f'(\text{net}_{t-1})] \\ &= \begin{bmatrix} w_{11}f'(net_1^{t-1}) & w_{12}f'(net_2^{t-1}) & \dots & w_{1n}f'(net_n^{t-1}) \\ w_{21}f'(net_1^{t-1}) & w_{22}f'(net_2^{t-1}) & \dots & w_{2n}f'(net_n^{t-1}) \\ & \cdot & & \\ & \cdot & & \\ w_{n1}f'(net_1^{t-1}) & w_{n2}f'(net_2^{t-1}) & \dots & w_{nn}f'(net_n^{t-1}) \end{bmatrix}\end{aligned}$$



8.3 随时间反向传播 (BPTT)

上式描述了将 δ 沿时间往前传递一个时刻的规律，有了这个规律，我们就可以求得任意时刻 k 的误差项 δ_k

$$\begin{aligned}\delta_k^T &= \frac{\partial E}{\partial \text{net}_k} \\ &= \frac{\partial E}{\partial \text{net}_t} \frac{\partial \text{net}_t}{\partial \text{net}_k} \\ &= \frac{\partial E}{\partial \text{net}_t} \frac{\partial \text{net}_t}{\partial \text{net}_{t-1}} \frac{\partial \text{net}_{t-1}}{\partial \text{net}_{t-2}} \cdots \frac{\partial \text{net}_{k+1}}{\partial \text{net}_k} \\ &= W \text{diag}[f'(\text{net}_{t-1})] W \text{diag}[f'(\text{net}_{t-2})] \cdots W \text{diag}[f'(\text{net}_k)] \delta_t^T \\ &= \delta_t^T \prod_{i=k}^{t-1} W \text{diag}[f'(\text{net}_i)]\end{aligned}\tag{式3}$$

(式3) 就是将误差项沿时间反向传播的算法。



8.3 随时间反向传播 (BPTT)

循环层将误差项反向传递到上一层网络，与普通的全连接层完全一样

循环层的加权输入 net^l 与上一层的加权输入 net^{l-1} 关系如下：

$$\begin{aligned} net_t^l &= U a_t^{l-1} + W s_{t-1} \\ a_t^{l-1} &= f^{l-1}(net_t^{l-1}) \end{aligned}$$

上式中的 net_t^l 是第 l 层神经元的加权输入（假设第 l 层是循环层）； net_t^{l-1} 是第 $l-1$ 层神经元的加权输入； a_t^{l-1} 是第 $l-1$ 层神经元的输出； f^{l-1} 是第 $l-1$ 层的激活函数。



8.3 随时间反向传播 (BPTT)

$$\begin{aligned}\frac{\partial \text{net}_t^l}{\partial \text{net}_t^{l-1}} &= \frac{\partial \text{net}_t^l}{\partial a_t^{l-1}} \frac{\partial a_t^{l-1}}{\partial \text{net}_t^{l-1}} \\ &= U \text{diag}[f'^{l-1}(\text{net}_t^{l-1})]\end{aligned}$$

所以

$$\begin{aligned}(\delta_t^{l-1})^T &= \frac{\partial E}{\partial \text{net}_t^{l-1}} \\ &= \frac{\partial E}{\partial \text{net}_t^l} \frac{\partial \text{net}_t^l}{\partial \text{net}_t^{l-1}} \\ &= (\delta_t^l)^T U \text{diag}[f'^{l-1}(\text{net}_t^{l-1})]\end{aligned}\tag{式4}$$

(式4) 就是将误差项传递到上一层算法。

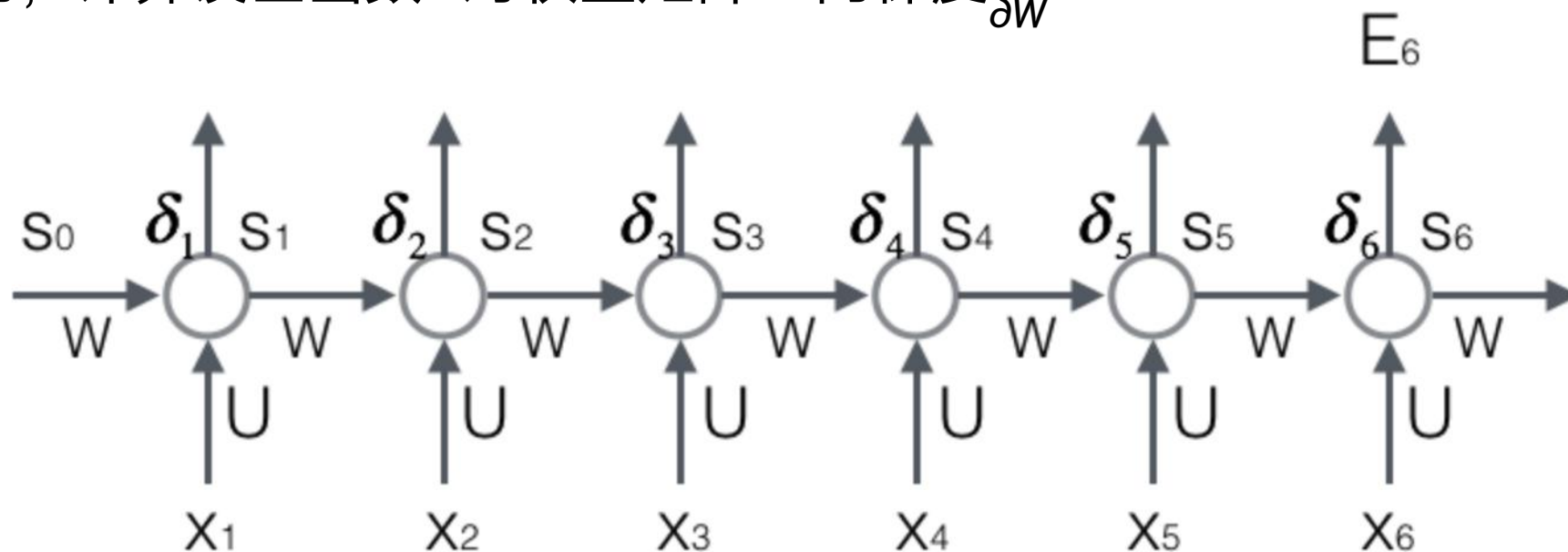


8.3 随时间反向传播 (BPTT)

C. 权重梯度的计算

BPTT算法的最后一步：计算每个权重的梯度。

首先，计算误差函数E对权重矩阵W的梯度 $\frac{\partial E}{\partial W}$



前两步中已经计算得到的量包括每个时刻 t 循环层的输出值 s_t ，以及误差项 δ_t 。



8.3 随时间反向传播 (BPTT)

类比之前介绍到的全连接网络的权重梯度计算算法，我们可以得到循环神经网络中梯度的计算：只要知道了任意一个时刻的误差项 δ_t ，以及上一个时刻循环层的输出值 s_{t-1} ，就可以按照下面的公式求出权重矩阵在t时刻的梯度 $\nabla_{W_t} E$ ：

$$\nabla_{W_t} E = \begin{bmatrix} \delta_1^t s_1^{t-1} & \delta_1^t s_2^{t-1} & \dots & \delta_1^t s_n^{t-1} \\ \delta_2^t s_1^{t-1} & \delta_2^t s_2^{t-1} & \dots & \delta_2^t s_n^{t-1} \\ \vdots & \vdots & \ddots & \vdots \\ \delta_n^t s_1^{t-1} & \delta_n^t s_2^{t-1} & \dots & \delta_n^t s_n^{t-1} \end{bmatrix} \quad (\text{式5})$$

使用反向传播算法的前馈神经网络随机梯度下降训练过程

输入：训练集 $D = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$ ，验证集 V ，学习率 α ，正则化系数 λ ，网络层数 L ，神经元数量 $\{M\}_{l=1}^L$ 。

输出： \mathbf{W} \mathbf{b}

```
1  随机初始化  $\mathbf{W}$   $\mathbf{b}$ ;  
2  repeat  
3      对训练集  $D$  中的样本随机重排序;  
4      For  $n = 1 \dots N$  do  
5          从训练集  $D$  中选取样本  $(\mathbf{x}^{(n)}, y^{(n)})$ ;  
6          前馈计算每一层的净输入  $\mathbf{z}^{(l)}$  和激活值  $\mathbf{a}^{(l)}$ ，直到最后一层;  
7          反向传播计算每一层的误差  $\boldsymbol{\delta}^{(l)} = f'_l(\mathbf{z}^{(l)}) \odot ((\mathbf{W}^{(l+1)})^\top \boldsymbol{\delta}^{(l+1)})$ ; //最后一层的误差为  $\boldsymbol{\delta}^{(L)} = f'_L(\mathbf{z}^{(L)}) \odot \frac{\partial L(\mathbf{y}, \mathbf{y})}{\partial \mathbf{y}}$   
          //计算每一层的梯度  
8           $\frac{\partial L(\mathbf{y}, \mathbf{y})}{\partial \mathbf{W}^{(l)}} = \boldsymbol{\delta}^{(l)} (\mathbf{a}^{(l-1)})^\top$ ;  
9           $\frac{\partial L(\mathbf{y}, \mathbf{y})}{\partial \mathbf{b}^{(l)}} = \boldsymbol{\delta}^{(l)}$ ;  
          //更新参数  
10          $\mathbf{W}^{(l)} \leftarrow \mathbf{W}^{(l)} - \alpha (\boldsymbol{\delta}^{(l)} (\mathbf{a}^{(l-1)})^\top + \lambda \mathbf{W}^{(l)})$ ;  
11          $\mathbf{b}^{(l)} \leftarrow \mathbf{b}^{(l)} - \alpha \boldsymbol{\delta}^{(l)}$ ;  
12     end;  
13 until 模型在验证集  $V$  上的错误率不再下降;
```



8.3 随时间反向传播 (BPTT)

下面可以简单推导一下式5。

我们知道：

$$\text{net}_t = Ux_t + Ws_{t-1}$$

$$\begin{aligned} \begin{bmatrix} \text{net}_1^t \\ \text{net}_2^t \\ \vdots \\ \text{net}_n^t \end{bmatrix} &= Ux_t + \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1n} \\ w_{21} & w_{22} & \dots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n1} & w_{n2} & \dots & w_{nn} \end{bmatrix} \begin{bmatrix} s_1^{t-1} \\ s_2^{t-1} \\ \vdots \\ s_n^{t-1} \end{bmatrix} \\ &= Ux_t + \begin{bmatrix} w_{11}s_1^{t-1} + w_{12}s_2^{t-1} \dots w_{1n}s_n^{t-1} \\ w_{21}s_1^{t-1} + w_{22}s_2^{t-1} \dots w_{2n}s_n^{t-1} \\ \vdots \\ w_{n1}s_1^{t-1} + w_{n2}s_2^{t-1} \dots w_{nn}s_n^{t-1} \end{bmatrix} \end{aligned}$$



8.3 随时间反向传播 (BPTT)

因为对 W 求导与 u_{x_t} 无关，我们不再考虑。现在考虑对权重项 w_{ji} 求导。通过观察上式我们可以看到 w_{ji} 只与 net_j^t 有关，所以：

$$\begin{aligned}\frac{\partial E}{\partial w_{ji}} &= \frac{\partial E}{\partial net_j^t} \frac{\partial net_j^t}{\partial w_{ji}} \\ &= \delta_j^t s_i^{t-1}\end{aligned}$$

按照上面的规律就可以生成式5里面的矩阵。



8.3 随时间反向传播 (BPTT)

我们已经求得了权重矩阵 W 在 t 时刻的梯度 $\nabla_{W_t} E$ ，最终的梯度 $\nabla_W E$ 是各个时刻的梯度之和：

$$\begin{aligned}\nabla_W E &= \sum_{i=1}^t \nabla_{W_i} E \\ &= \begin{bmatrix} \delta_1^t s_1^{t-1} & \delta_1^t s_2^{t-1} & \dots & \delta_1^t s_n^{t-1} \\ \delta_2^t s_1^{t-1} & \delta_2^t s_2^{t-1} & \dots & \delta_2^t s_n^{t-1} \\ . & . & . & . \\ \delta_n^t s_1^{t-1} & \delta_n^t s_2^{t-1} & \dots & \delta_n^t s_n^{t-1} \end{bmatrix} + \dots + \begin{bmatrix} \delta_1^1 s_1^0 & \delta_1^1 s_2^0 & \dots & \delta_1^1 s_n^0 \\ \delta_2^1 s_1^0 & \delta_2^1 s_2^0 & \dots & \delta_2^1 s_n^0 \\ . & . & . & . \\ \delta_n^1 s_1^0 & \delta_n^1 s_2^0 & \dots & \delta_n^1 s_n^0 \end{bmatrix} \quad (\text{式6})\end{aligned}$$

式6就是计算循环层权重矩阵 W 的梯度的公式。



8.3 随时间反向传播 (BPTT)

为什么最终的梯度是各个时刻的梯度之和呢？

还是从这个式子开始：

$$\text{net}_t = Ux_t + Wf(\text{net}_{t-1})$$

因为 Ux_t 与 W 完全无关，可以把它看做常量。现在，考虑式子加号右边的部分，

因为 W 和 $f(\text{net}_{t-1})$ 都是 W 的函数，因此要用到导数乘法运算：

$$(uv)' = u'v + uv'$$

因此：

$$\frac{\partial \text{net}_t}{\partial W} = \frac{\partial W}{\partial W} f(\text{net}_{t-1}) + W \frac{\partial f(\text{net}_{t-1})}{\partial W}$$



8.3 随时间反向传播 (BPTT)

最终需要计算的是 $\nabla_W E$:

$$\begin{aligned}\nabla_W E &= \frac{\partial E}{\partial W} \\ &= \frac{\partial E}{\partial \text{net}_t} \frac{\partial \text{net}_t}{\partial W} \\ &= \delta_t^T \frac{\partial W}{\partial W} f(\text{net}_{t-1}) + \delta_t^T W \frac{\partial f(\text{net}_{t-1})}{\partial W}\end{aligned}\quad (\text{式7})$$

先计算式7加号左边的部分。 $\frac{\partial W}{\partial W}$ 是矩阵对矩阵求导，其结果是一个四维张量 (tensor)，如下所示：



8.3 随时间反向传播 (BPTT)

$$\frac{\partial W}{\partial W} = \begin{bmatrix} \frac{\partial w_{11}}{\partial W} & \frac{\partial w_{12}}{\partial W} & \dots & \frac{\partial w_{1n}}{\partial W} \\ \frac{\partial w_{21}}{\partial W} & \frac{\partial w_{22}}{\partial W} & \dots & \frac{\partial w_{2n}}{\partial W} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial w_{n1}}{\partial W} & \frac{\partial w_{n2}}{\partial W} & \dots & \frac{\partial w_{nn}}{\partial W} \end{bmatrix}$$

$$= \begin{bmatrix} \begin{bmatrix} \frac{\partial w_{11}}{\partial w_{11}} & \frac{\partial w_{11}}{\partial w_{12}} & \dots & \frac{\partial w_{11}}{\partial_{1n}} \\ \frac{\partial w_{11}}{\partial w_{21}} & \frac{\partial w_{11}}{\partial w_{22}} & \dots & \frac{\partial w_{11}}{\partial_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial w_{11}}{\partial w_{n1}} & \frac{\partial w_{11}}{\partial w_{n2}} & \dots & \frac{\partial w_{11}}{\partial_{nn}} \end{bmatrix} & \begin{bmatrix} \frac{\partial w_{12}}{\partial w_{11}} & \frac{\partial w_{12}}{\partial w_{12}} & \dots & \frac{\partial w_{12}}{\partial_{1n}} \\ \frac{\partial w_{12}}{\partial w_{21}} & \frac{\partial w_{12}}{\partial w_{22}} & \dots & \frac{\partial w_{12}}{\partial_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial w_{12}}{\partial w_{n1}} & \frac{\partial w_{12}}{\partial w_{n2}} & \dots & \frac{\partial w_{12}}{\partial_{nn}} \end{bmatrix} & \dots \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix} = \begin{bmatrix} \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} & \begin{bmatrix} 0 & 1 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} & \dots \end{bmatrix}$$



8.3 随时间反向传播 (BPTT)

接下来，我们知道 $s_{t-1} = f(\text{net}_{t-1})$ ，它是一个**列向量**。让上面的四维张量与这个向量相乘，得到了一个三维张量，再左乘行向量 δ_t^T ，最终得到一个矩阵：

$$\begin{aligned}\delta_t^T \frac{\partial W}{\partial W} f(\text{net}_{t-1}) &= \delta_t^T \frac{\partial W}{\partial W} s_{t-1} \\ &= \delta_t^T \begin{bmatrix} \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & & & \\ 0 & 0 & \dots & 0 \end{bmatrix} & \begin{bmatrix} 0 & 1 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & & & \\ 0 & 0 & \dots & 0 \end{bmatrix} & \dots \end{bmatrix} \begin{bmatrix} s_1^{t-1} \\ s_2^{t-1} \\ \cdot \\ \cdot \\ s_n^{t-1} \end{bmatrix} \\ &= \delta_t^T \begin{bmatrix} \begin{bmatrix} s_1^{t-1} \\ 0 \\ \cdot \\ \cdot \\ 0 \end{bmatrix} & \begin{bmatrix} s_2^{t-1} \\ 0 \\ \cdot \\ \cdot \\ 0 \end{bmatrix} & \dots \end{bmatrix}\end{aligned}$$



8.3 随时间反向传播 (BPTT)

$$\begin{aligned} &= [\delta_1^t \quad \delta_2^t \quad \dots \quad \delta_n^t] \begin{bmatrix} \begin{bmatrix} s_1^{t-1} \\ 0 \\ \cdot \\ \cdot \\ 0 \end{bmatrix} & \begin{bmatrix} s_2^{t-1} \\ 0 \\ \cdot \\ \cdot \\ 0 \end{bmatrix} & \dots \end{bmatrix} \\ &= \begin{bmatrix} \delta_1^t s_1^{t-1} & \delta_1^t s_2^{t-1} & \dots & \delta_1^t s_n^{t-1} \\ \delta_2^t s_1^{t-1} & \delta_2^t s_2^{t-1} & \dots & \delta_2^t s_n^{t-1} \\ \cdot & & & \\ \cdot & & & \\ \delta_n^t s_1^{t-1} & \delta_n^t s_2^{t-1} & \dots & \delta_n^t s_n^{t-1} \end{bmatrix} \\ &= \nabla_{W_t} E \end{aligned}$$



8.3 随时间反向传播 (BPTT)

接下来计算式7加号右边的部分：

$$\begin{aligned}\delta_t^T W \frac{\partial f(\text{net}_{t-1})}{\partial W} &= \delta_t^T W \frac{\partial f(\text{net}_{t-1})}{\partial \text{net}_{t-1}} \frac{\partial \text{net}_{t-1}}{\partial W} \\&= \delta_t^T W f'(\text{net}_{t-1}) \frac{\partial \text{net}_{t-1}}{\partial W} \\&= \delta_t^T \frac{\partial \text{net}_t}{\partial \text{net}_{t-1}} \frac{\partial \text{net}_{t-1}}{\partial W} \\&= \delta_{t-1}^T \frac{\partial \text{net}_{t-1}}{\partial W}\end{aligned}$$



8.3 随时间反向传播 (BPTT)

于是，我们得到了如下递推公式：

$$\begin{aligned}\nabla_W E &= \frac{\partial E}{\partial W} \\ &= \frac{\partial E}{\partial \text{net}_t} \frac{\partial \text{net}_t}{\partial W} \\ &= \nabla_{W_t} E + \delta_{t-1}^T \frac{\partial \text{net}_{t-1}}{\partial W} \\ &= \nabla_{W_t} E + \nabla_{W_{t-1}} E + \delta_{t-2}^T \frac{\partial \text{net}_{t-2}}{\partial W} \\ &= \nabla_{W_t} E + \nabla_{W_{t-1}} E + \dots + \nabla_{W_1} E \\ &= \sum_{k=1}^t \nabla_{W_k} E\end{aligned}$$

这样，我们就证明了：最终的梯度 $\nabla_W E$ 是各个时刻的梯度之和。



8.3 随时间反向传播 (BPTT)

同权重矩阵W类似，我们可以得到权重矩阵U的计算方法：

$$\nabla_{U_t} E = \begin{bmatrix} \delta_1^t x_1^t & \delta_1^t x_2^t & \dots & \delta_1^t x_m^t \\ \delta_2^t x_1^t & \delta_2^t x_2^t & \dots & \delta_2^t x_m^t \\ \vdots & \vdots & \ddots & \vdots \\ \delta_n^t x_1^t & \delta_n^t x_2^t & \dots & \delta_n^t x_m^t \end{bmatrix} \quad (\text{式8})$$

式8是误差函数在t时刻对权重矩阵U的梯度。和权重矩阵W一样，最终的梯度也是各个时刻的梯度之和：

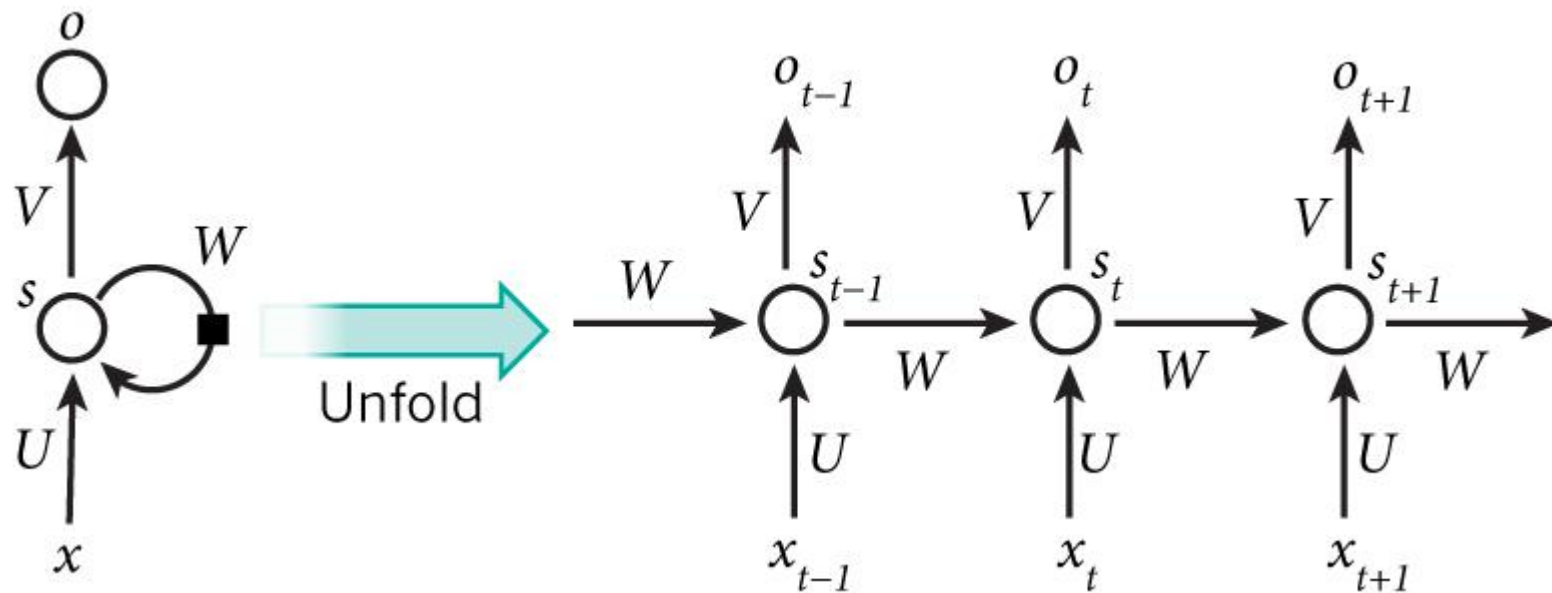
$$\nabla_U E = \sum_{i=1}^t \nabla_{U_i} E$$



8.3 随时间反向传播 (BPTT)

权重矩阵V的计算

参数V相对简单，只需关注目前的状态。RNN的损失也是会随着时间累加的，所以不能只求t时刻的偏导。



$$\frac{\partial E^{(t)}}{\partial V} = \frac{\partial E^{(t)}}{\partial o^{(t)}} \cdot \frac{\partial o^{(t)}}{\partial V}$$

$$E = \sum_{t=1}^n E^{(t)}$$

$$\frac{\partial E}{\partial V} = \sum_{t=1}^n \frac{\partial E^{(t)}}{\partial o^{(t)}} \cdot \frac{\partial o^{(t)}}{\partial V}$$



8.1 网络记忆能力

8.2 循环神经网络 (RNN)

8.3 随时间反向传播 (BPTT)

8.4 双向循环神经网络



8.4 双向循环神经网络

对于语言模型来说，很多时候光看前面的词是不够的，比如下面这句话：

我的手机坏了，我打算_____一部新手机。

如果我们只看横线前面的词，手机坏了，那么我是打算修一修？换一部新的？还是大哭一场？这些都是无法确定的。

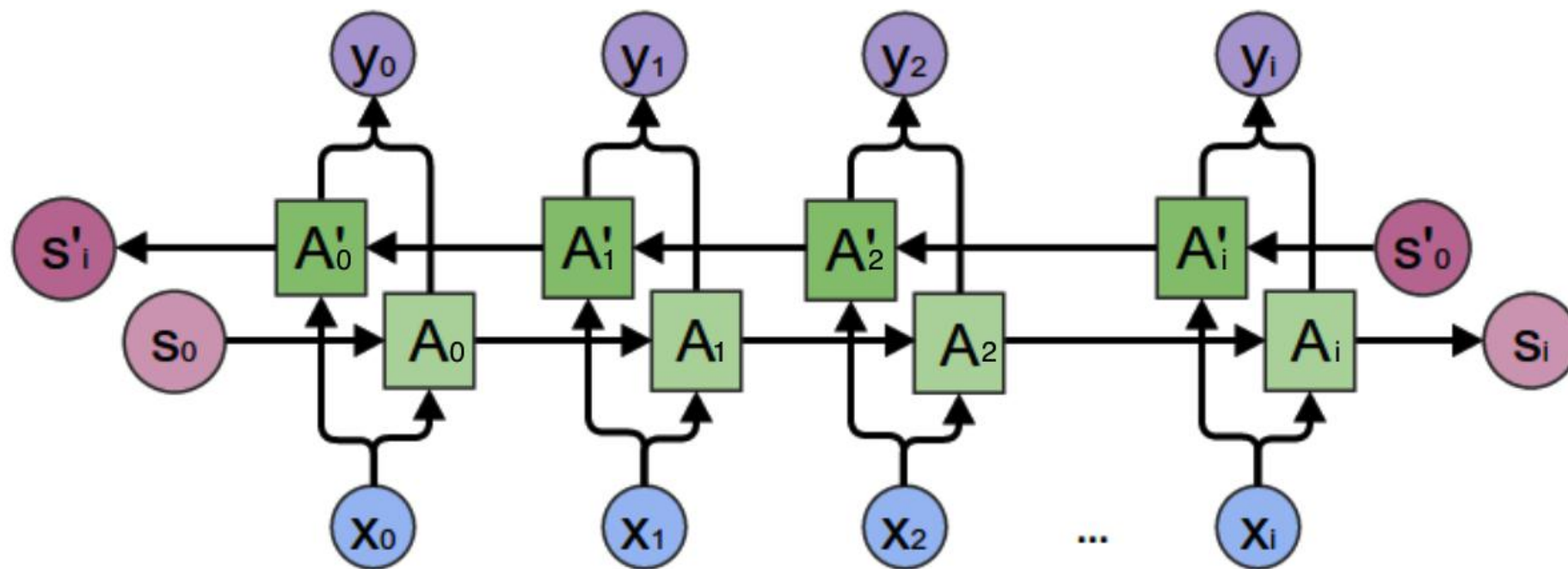
但如果我们也看到了横线后面的词是『一部新手机』，那么，横线上的词填『买』的概率就大得多了。

在此任务中，可以增加一个按照时间的逆序来传递信息的网络层，增强网络的能力。



8.4 双向循环神经网络

基本循环神经网络无法对此进行建模，因此，我们需要双向循环神经网络，如下图所示：



双向循环神经网络由两层循环神经网络组成，它们的输入相同，只是信息传递的方向不同。



8.4 双向循环神经网络

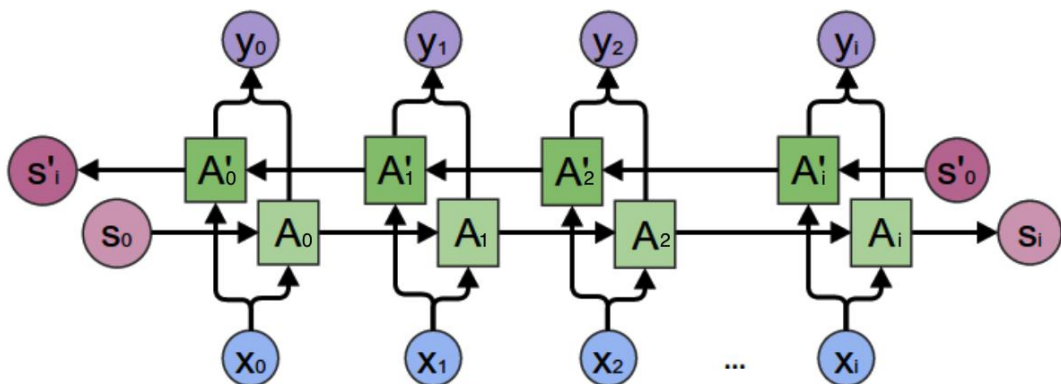
双向循环神经网络的隐藏层要保存两个值，一个值 A 参与正向计算，另一个值 A' 参与反向计算。最终的输出值取决于 A 和 A' 。其计算方法为（以 y_2 为例）：

$$y_2 = g(VA_2 + V'A'_2)$$

A_2 和 A'_2 则分别计算：

$$A_2 = f(WA_1 + Ux_2)$$

$$A'_2 = f(W'A'_3 + U'x_2)$$





8.4 双向循环神经网络

正向计算时，隐藏层的值 s_t 与 s_{t-1} 有关；反向计算时，隐藏层的值 s_t' 与 s_{t+1}' 有关；最终的输出取决于正向和反向计算的**加和**。

$$o_t = g(Vs_t + V's_t')$$

$$s_t = f(Ux_t + Ws_{t-1})$$

$$s_t' = f(U'x_t + W's_{t+1}')$$

从上面三个公式得到，正向计算和反向计算**不共享权重**，也就是说 U 和 U' 、 W 和 W' 、 V 和 V' 都是不同的**权重矩阵**。



主要参考文献

<https://zybuluo.com/hanbingtao/note/541458>

Recurrent Neural Networks Tutorial <http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/>

Pascanu R, Mikolov T, Bengio Y. On the difficulty of training recurrent neural networks[C] //International conference on machine learning. 2013: 1310-1318.

Attention and Augmented Recurrent Neural Networks. <https://distill.pub/2016/augmented-rnns/>

The Unreasonable Effectiveness of Recurrent Neural Networks.
<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

Karpathy A, Fei-Fei L. Deep visual-semantic alignments for generating image descriptions[C]//Proceedings of the IEEE conference on computer vision and pattern recognition. 2015: 3128-3137.



北京交通大学《深度学习》课程组成员

景丽萍: <http://faculty.bjtu.edu.cn/8249/>

桑基韬: <http://faculty.bjtu.edu.cn/9129/>

张淳杰: <http://faculty.bjtu.edu.cn/9371/>

万怀宇: <http://faculty.bjtu.edu.cn/8793/>

滕 竹: <http://faculty.bjtu.edu.cn/8902/>

原继东: <http://faculty.bjtu.edu.cn/9076/>

丛润民: <http://faculty.bjtu.edu.cn/9374/>

夏佳楠: <http://faculty.bjtu.edu.cn/9430/>

许万茹

杨 扩

