

自然语言处理

Natural Language Processing

权小军 教授

中山大学数据科学与计算机学院

quanxj3@mail.sysu.edu.cn

课程回顾

基本概念

语句 $s = w_1 w_2 \dots w_m$ 的先验概率：

$$\begin{aligned} p(s) &= p(w_1) \times p(w_2/w_1) \times p(w_3/w_1w_2) \times \dots \\ &\quad \times p(w_m/w_1\dots w_{m-1}) \\ &= \prod_{i=1}^m p(w_i | w_1 \dots w_{i-1}) \end{aligned}$$

当 $i=1$ 时, $p(w_1|w_0) = p(w_1)$ 。

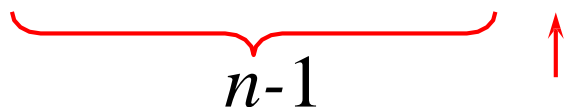
语言模型！！

基本概念

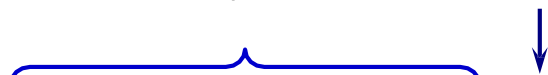
□ 如何划分等价类

- 将两个历史映射到同一个等价类，当且仅当这两个历史中的最近 $n-1$ 个基元相同，即：

$$H_1: w_1 w_2 \dots \dots w_{i-n+1} w_{i-n+2} \dots w_{i-1} w_i \dots \dots$$


 $n-1$

$$H_2: v_1 v_2 \dots \dots v_{k-n+1} v_{k-n+2} \dots v_{k-1} v_k \dots \dots$$



$$S(w_1, w_2, \dots, w_i) = S(v_1, v_2, \dots, v_k)$$

$$\text{iff } H_1: (w_{i-n+1}, \dots, w_i) = H_2: (v_{k-n+1}, \dots, v_k)$$

基本概念

□ 这种情况下的语言模型称为 n 元文法(n -gram)模型

□ 通常地,

- 当 $n=1$ 时, 即出现在第 i 位上的基元 w_i 独立于历史。
一元文法也被写为 uni-gram 或 monogram;
- 当 $n=2$ 时, 2-gram (bi-gram) 被称为1阶马尔可夫链;
- 当 $n=3$ 时, 3-gram(tri-gram)被称为2阶马尔可夫链,
依次类推。

Lecture 8: 语言模型（中）

语言模型

- 1. 统计语言模型
- 2. 神经语言模型

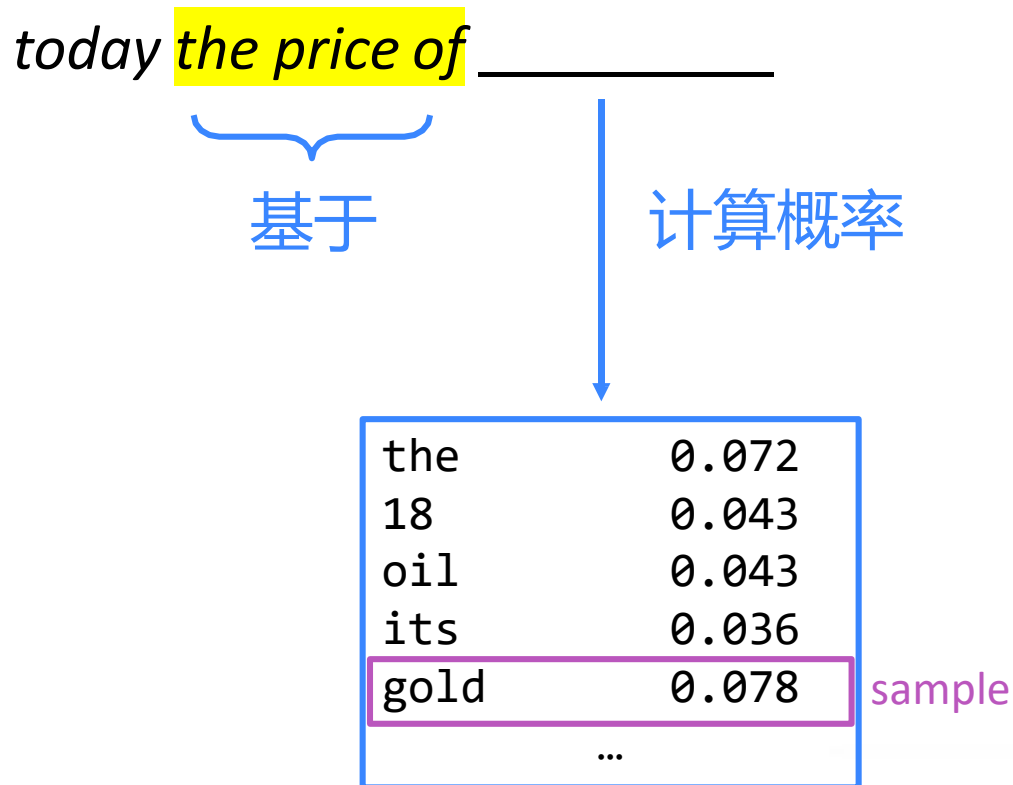
基于语言模型的文本生成

- You can also use a Language Model to generate words.

例句: *today the price of* ____

基于语言模型的文本生成

- You can also use a Language Model to generate words.



基于语言模型的文本生成

- You can also use a Language Model to generate text.

today the price of gold per ton , while production of shoe lasts and shoe industry , the bank intervened just after it considered and rejected an imf demand to rebuild depleted european stocks , sept 30 end primary 76 cts a share .

基于语言模型的文本生成

- You can also use a Language Model to generate text.

today the price of gold per ton , while production of shoe lasts and shoe industry , the bank intervened just after it considered and rejected an imf demand to rebuild depleted european stocks , sept 30 end primary 76 cts a share .

Surprisingly grammatical!

...but incoherent.

Neural Language Model

- Window-based neural model

as the proctor started the clock.

The students opened their_____

Neural Language Model

- Window-based neural model?

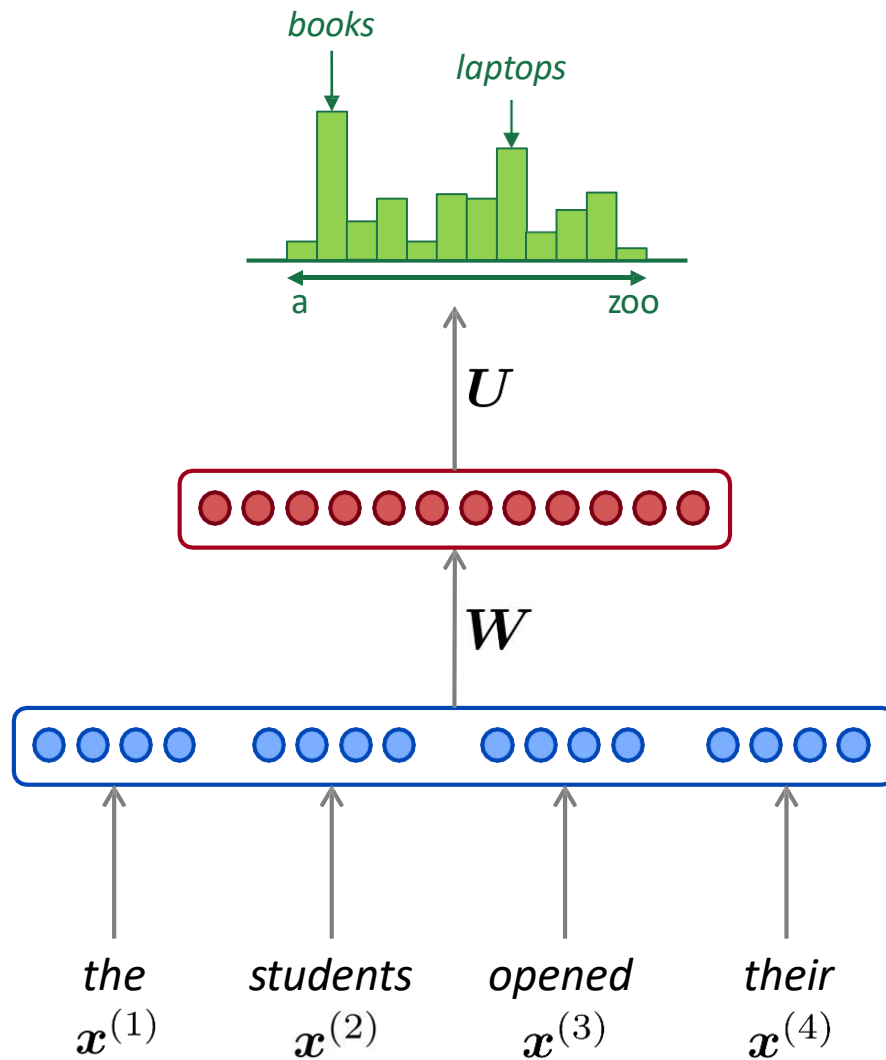
~~as the proctor started the clock.~~

丢弃

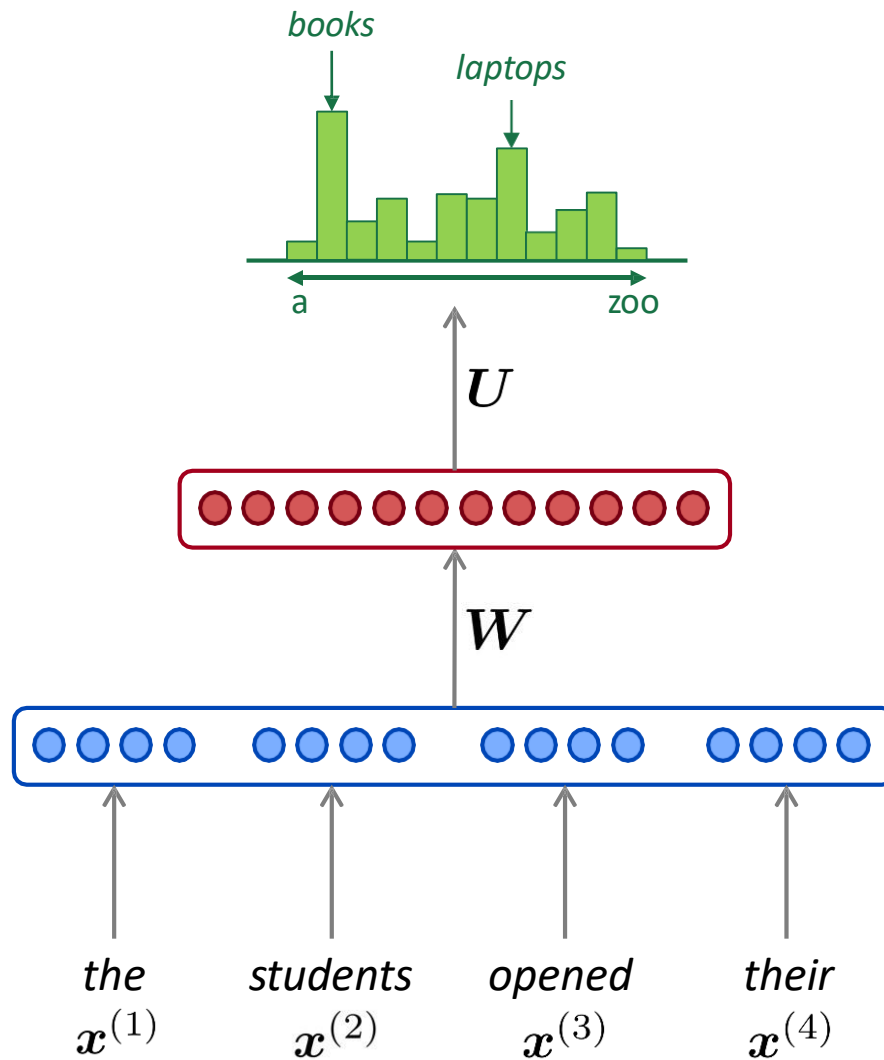
The students opened their_____

固定窗口

A fixed-window neural Language Model



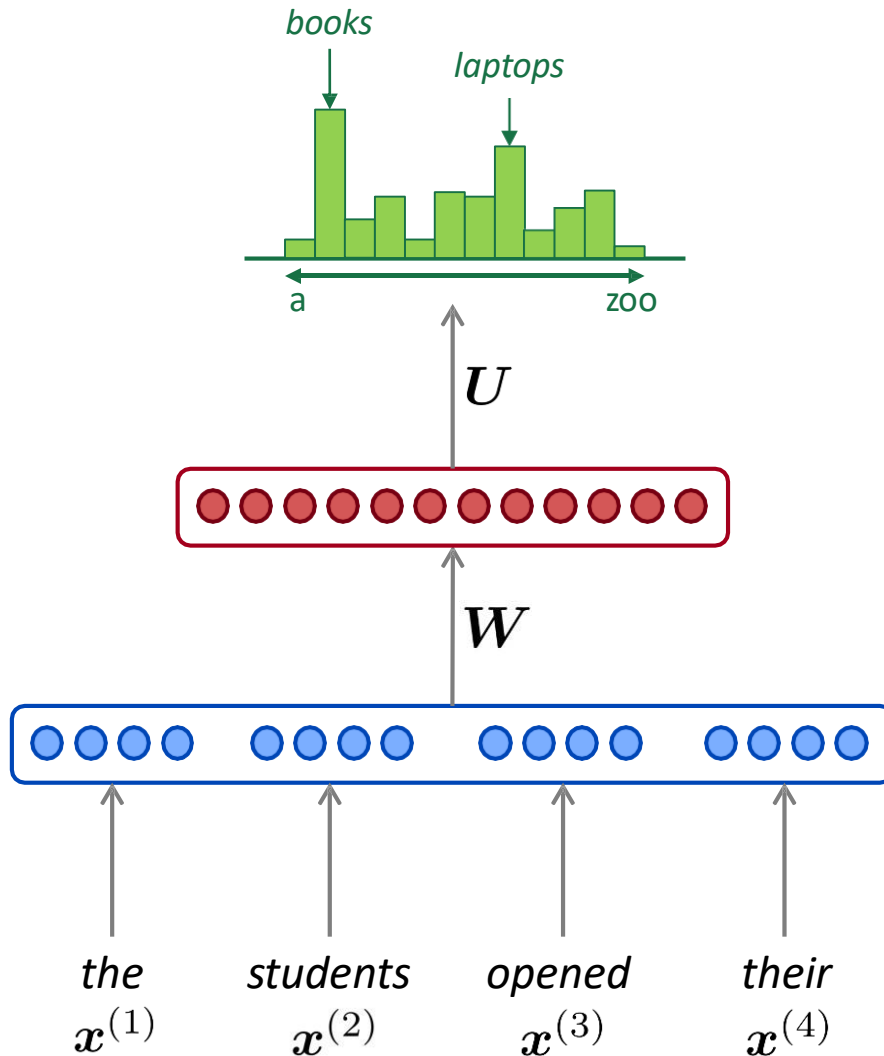
A fixed-window neural Language Model



Level 1: words / one-hot vectors

$x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$

A fixed-window neural Language Model



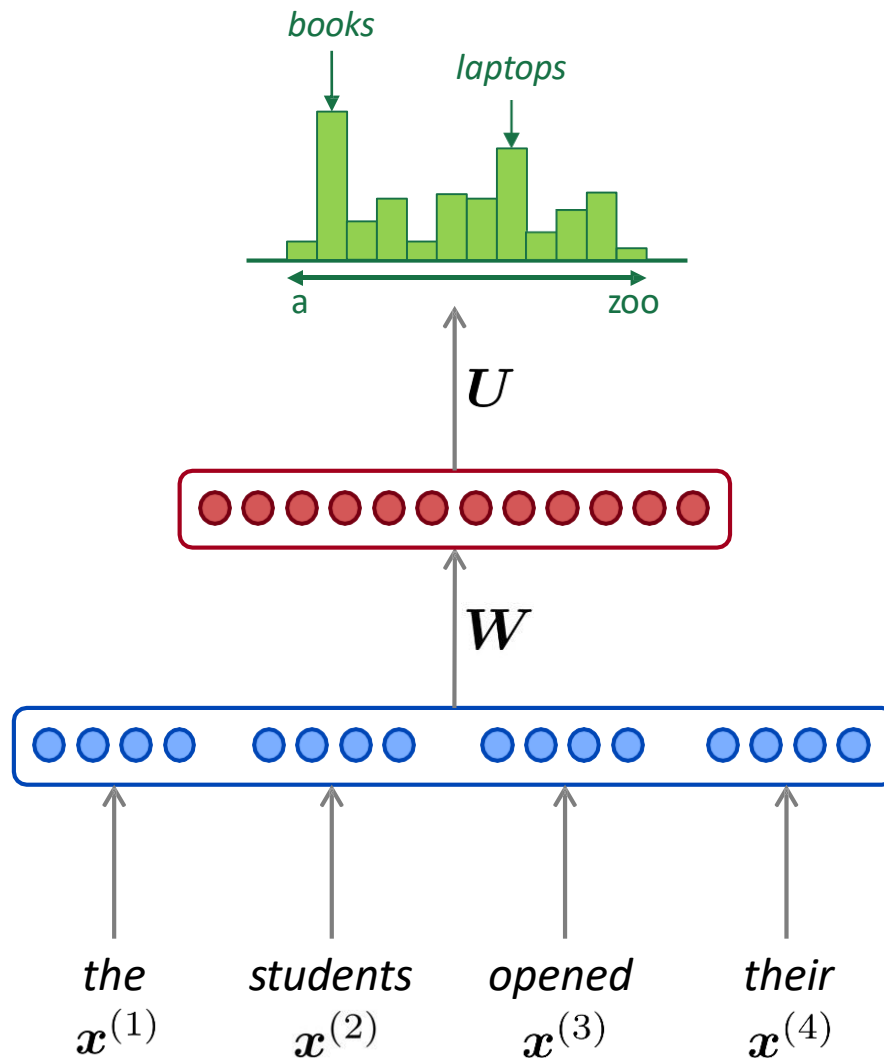
Level 2: concatenated word embeddings

$$e = [e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}]$$

Level 1: words / one-hot vectors

$$x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$$

A fixed-window neural Language Model



Level 3: hidden layer

$$h = f(We + b_1)$$

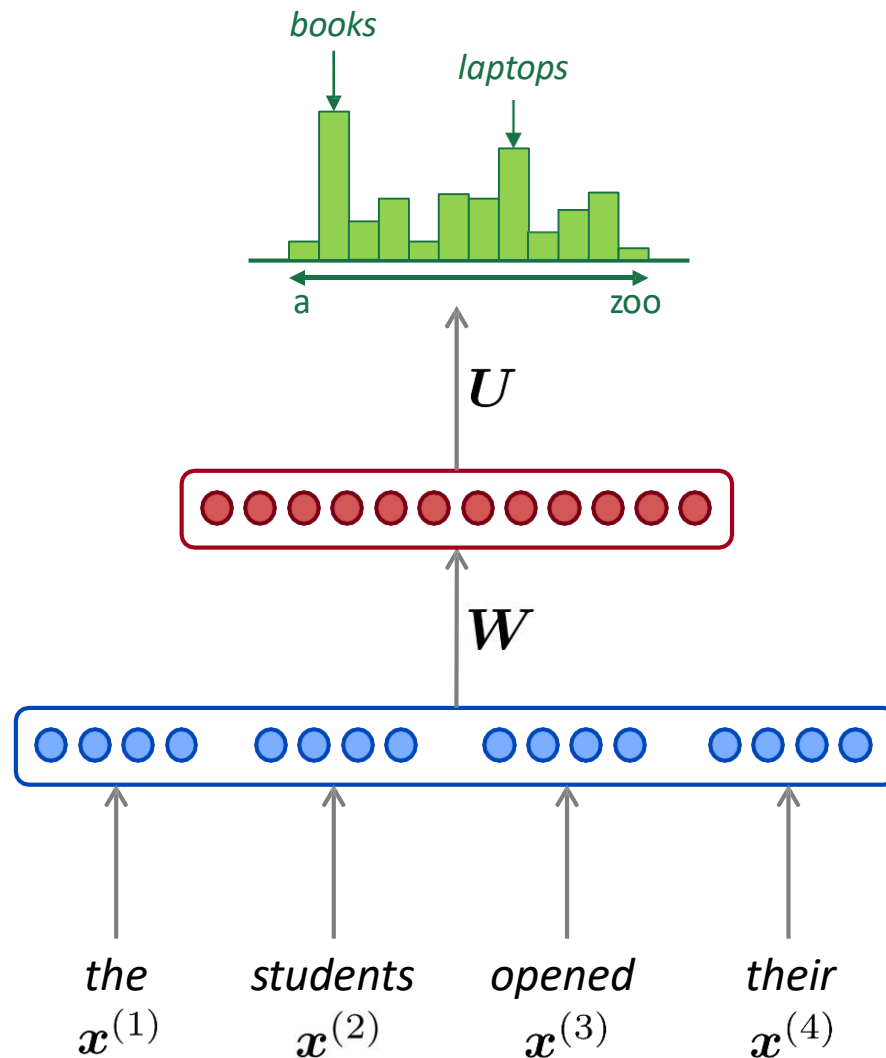
Level 2: concatenated word embeddings

$$e = [e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}]$$

Level 1: words / one-hot vectors

$$x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$$

A fixed-window neural Language Model



Level 4: output distribution

$$\hat{y} = \text{softmax}(Uh + b_2) \in \mathbb{R}^{|V|}$$

Level 3: hidden layer

$$h = f(We + b_1)$$

Level 2: concatenated word embeddings

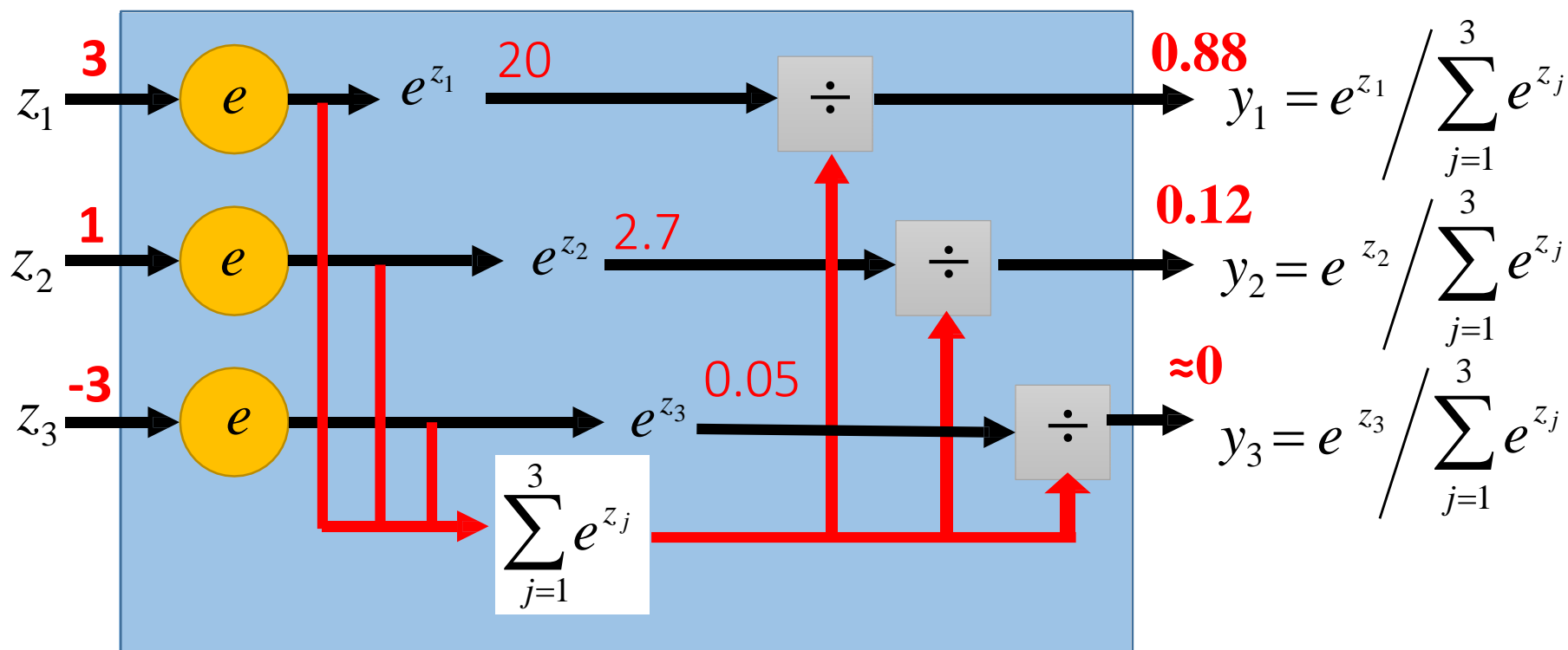
$$e = [e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}]$$

Level 1: words / one-hot vectors

$$x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$$

Softmax

- Softmax 层作为一个输出层



典型应用： Word2vec

Word2vec

“A word is known by the company it keeps”



Word2vec

Word Representations

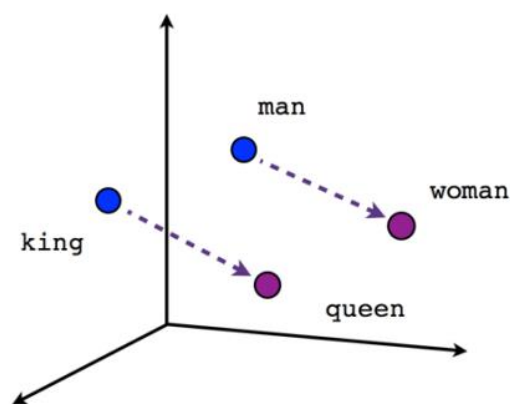
Traditional Method - Bag of Words Model	Word Embeddings
<ul style="list-style-type: none">• Uses one hot encoding• Each word in the vocabulary is represented by one bit position in a HUGE vector.• For example, if we have a vocabulary of 10000 words, and “Hello” is the 4th word in the dictionary, it would be represented by: 0 0 0 1 0 0 0 0 0 0• Context information is not utilized	<ul style="list-style-type: none">• Stores each word in as a point in space, where it is represented by a vector of fixed number of dimensions (generally 300)• Unsupervised, built just by reading huge corpus• For example, “Hello” might be represented as : [0.4, -0.11, 0.55, 0.3 . . . 0.1, 0.02]• Dimensions are basically projections along different axes, more of a mathematical concept.

Word2vec

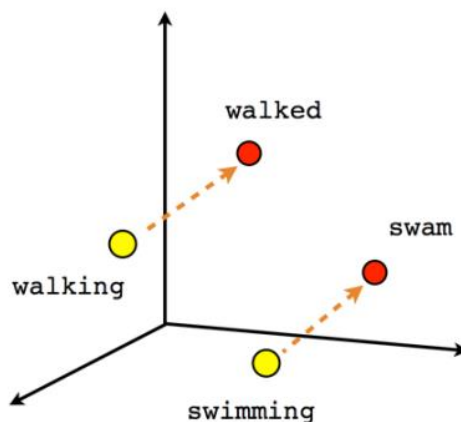
- 词向量给NLP问题提供一个全新的视角
- Word2vec通过一种无监督的方式获取词向量

Word2vec

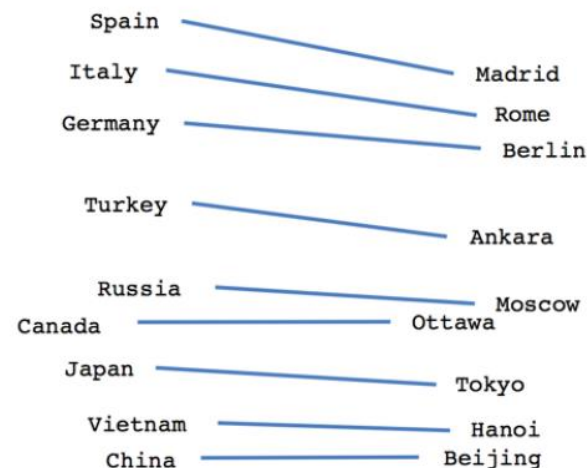
典型的例子:



Male-Female



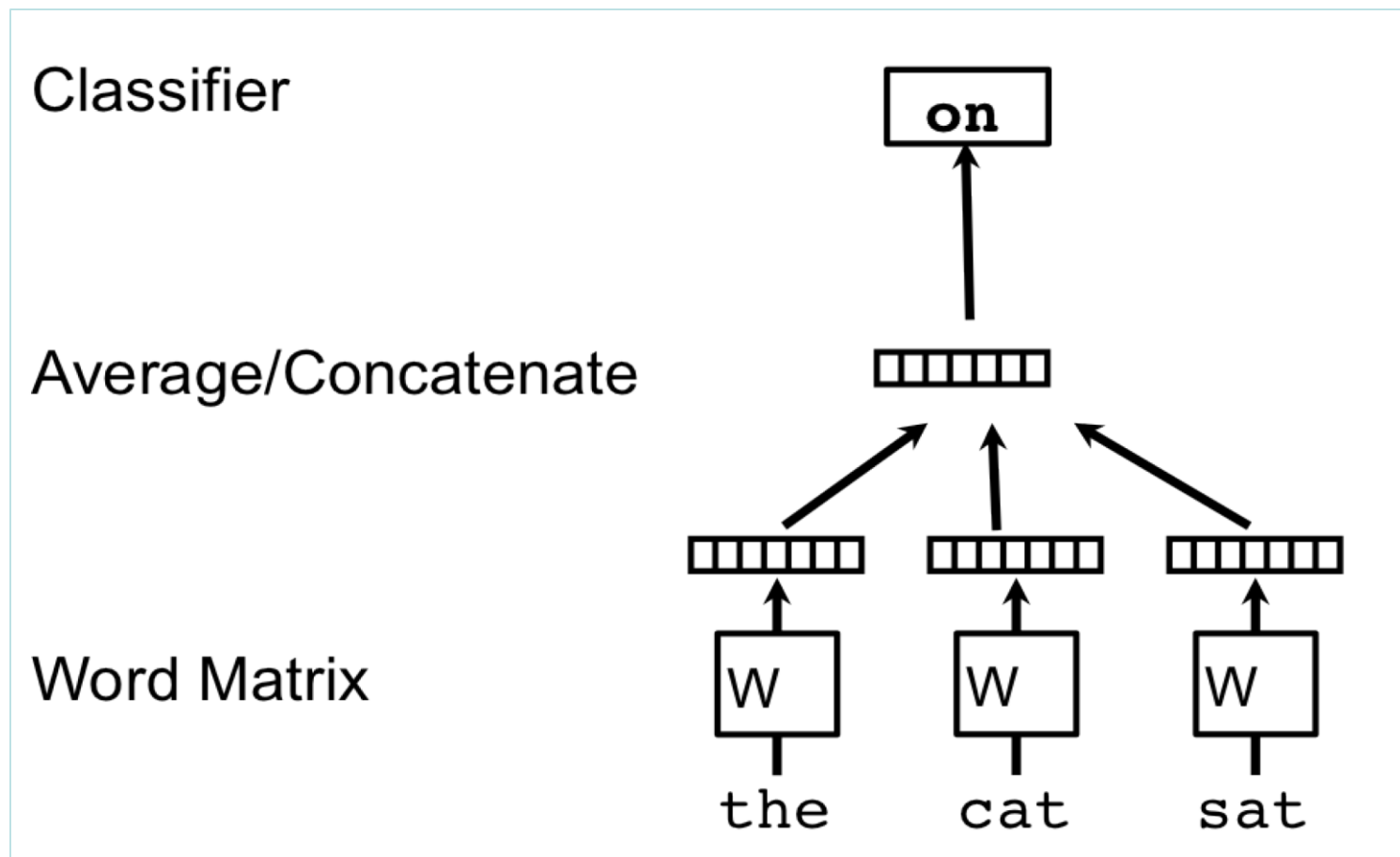
Verb tense



Country-Capital

$$\text{vector}[\text{Queen}] = \text{vector}[\text{King}] - \text{vector}[\text{Man}] + \text{vector}[\text{Woman}]$$

Word2vec原理



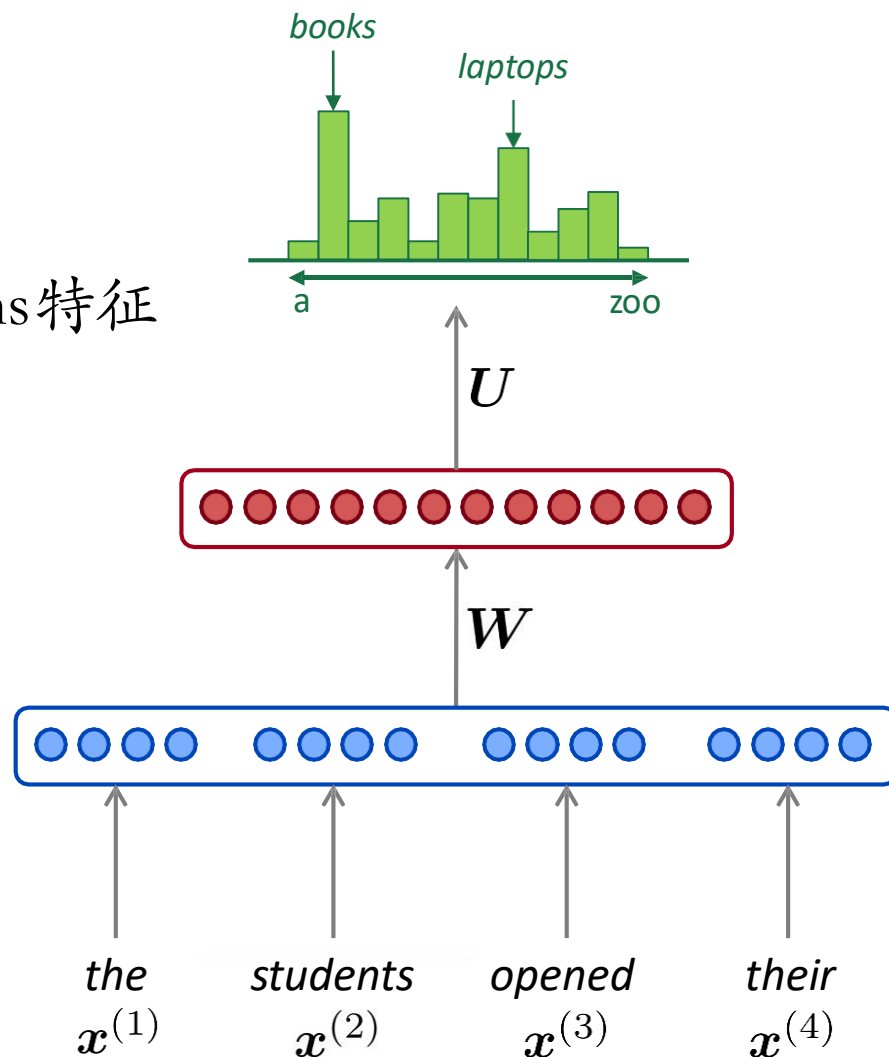
A fixed-window neural Language Model

和n-gram对比分析

A fixed-window neural Language Model

和 n -gram 语言模型相比:

- 不存在稀疏性问题
- 不用存储所有已知的 n -grams 特征



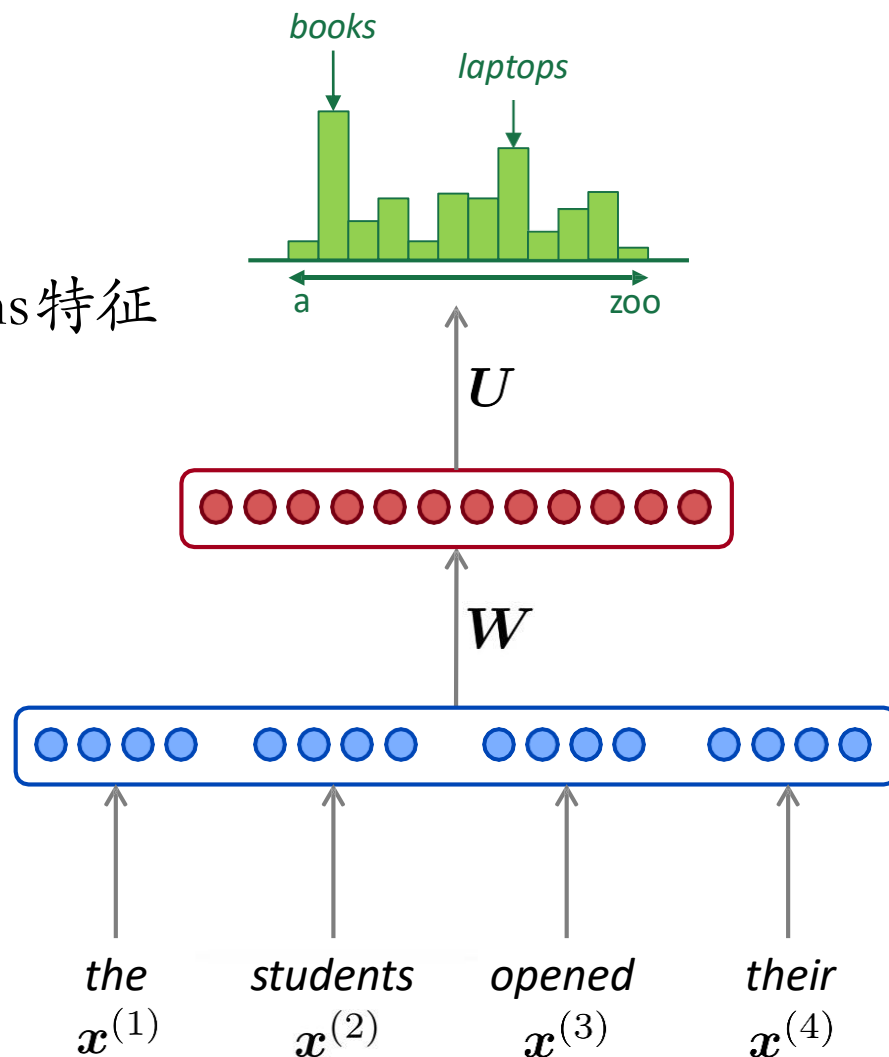
A fixed-window neural Language Model

和 n -gram语言模型相比:

- 不存在稀疏性问题
- 不用存储所有已知的 n -grams特征

存在的问题:

- 窗口太小, 效果有限
- 窗口太大, W 也会变大



A fixed-window neural Language Model

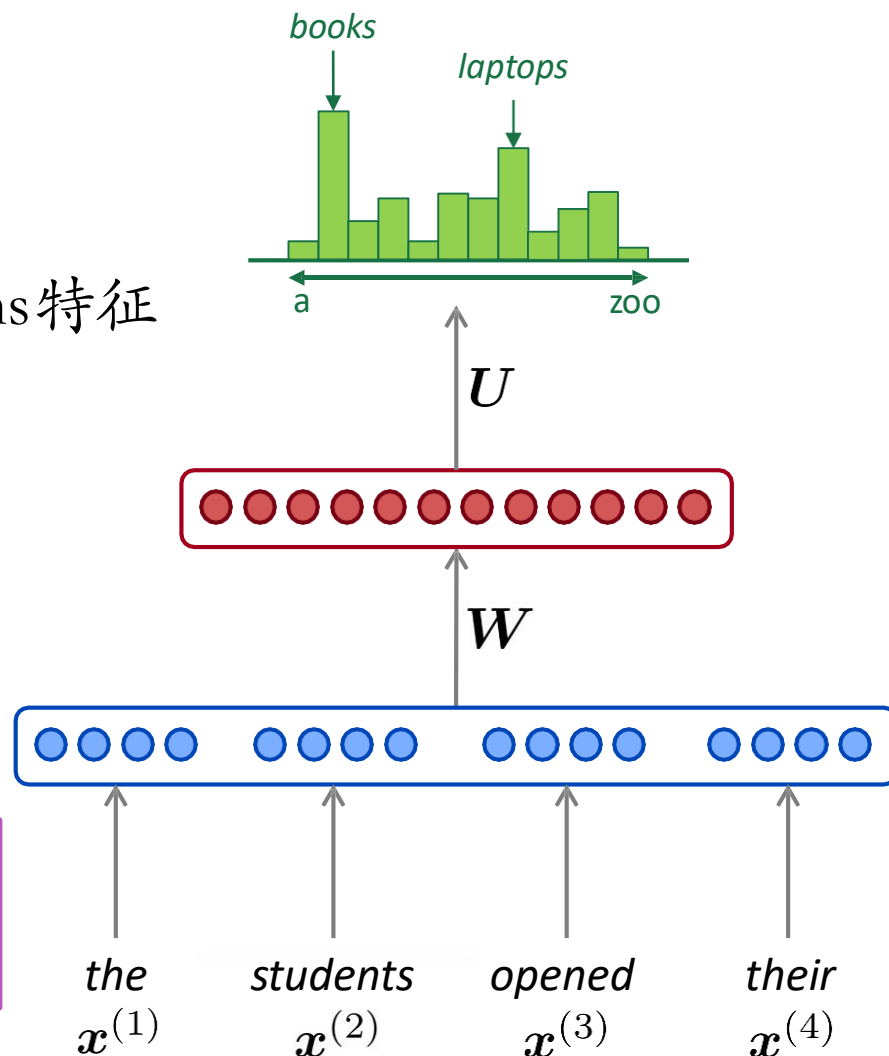
和 n -gram语言模型相比:

- 不存在稀疏性问题
- 不用存储所有已知的 n -grams特征

存在的问题:

- 窗口太小
- 窗口太大, W 也会变大

需要一种能处理任意
长度输入的架构

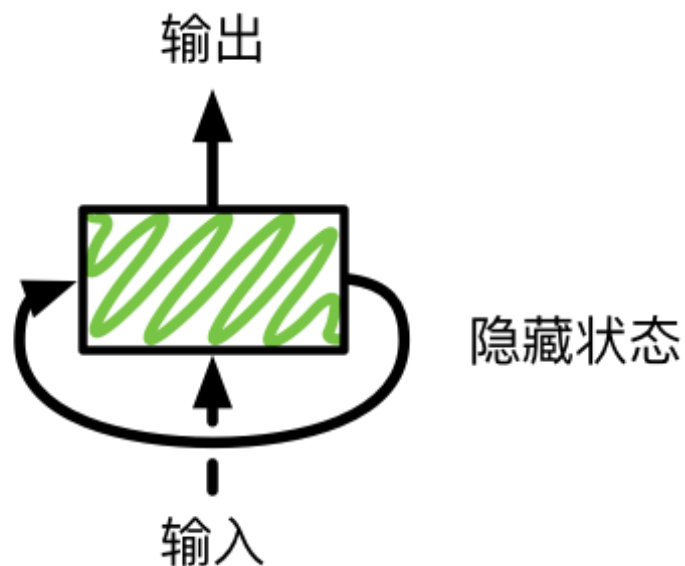


循环神经网络

Recurrent Neural Networks

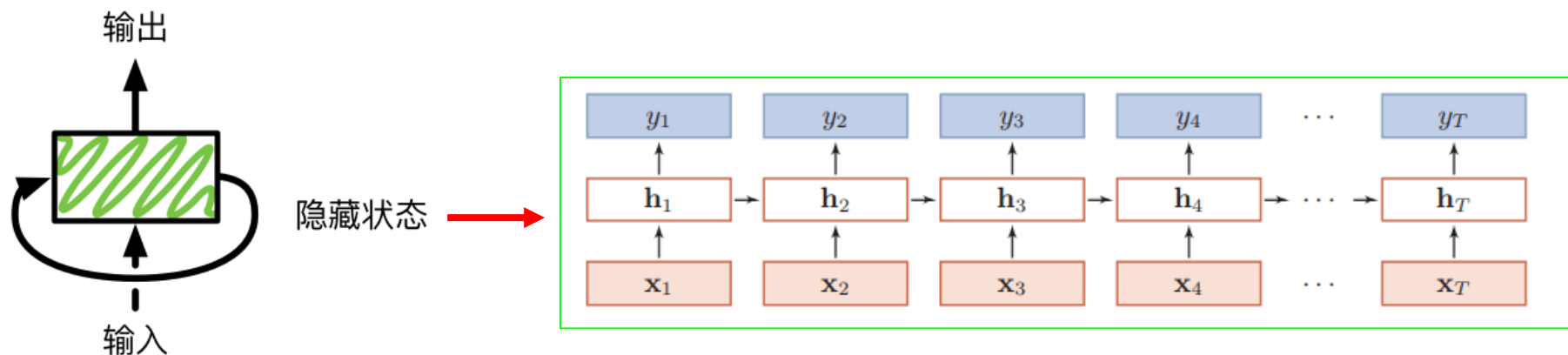
循环神经网络

- 循环神经网络主要用于处理（变长）序列数据



循环神经网络

- 循环神经网络主要用于处理（变长）序列数据

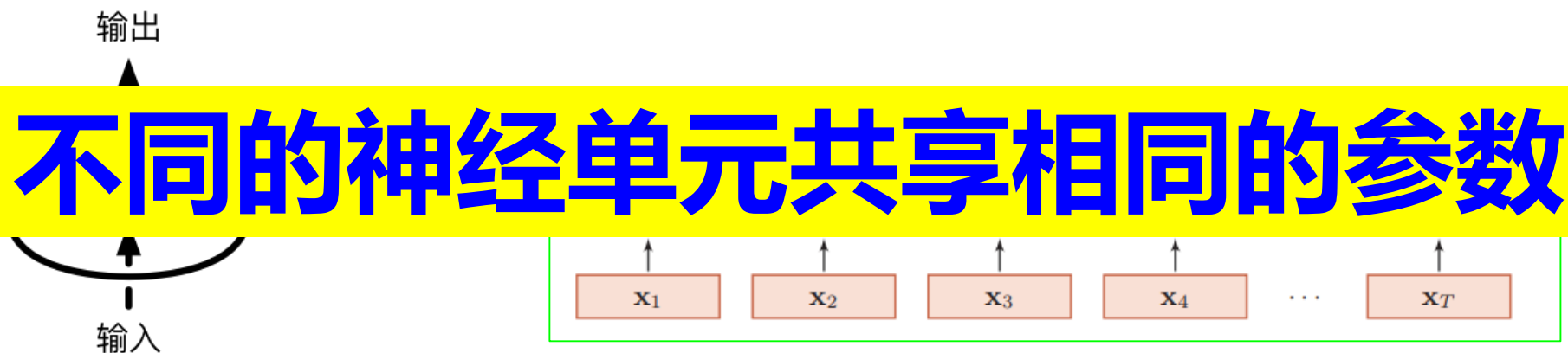


$$\mathbf{h}_t = f(U\mathbf{h}_{t-1} + W\mathbf{x}_t + \mathbf{b}),$$

$$\mathbf{y}_t = V\mathbf{h}_t,$$

循环神经网络

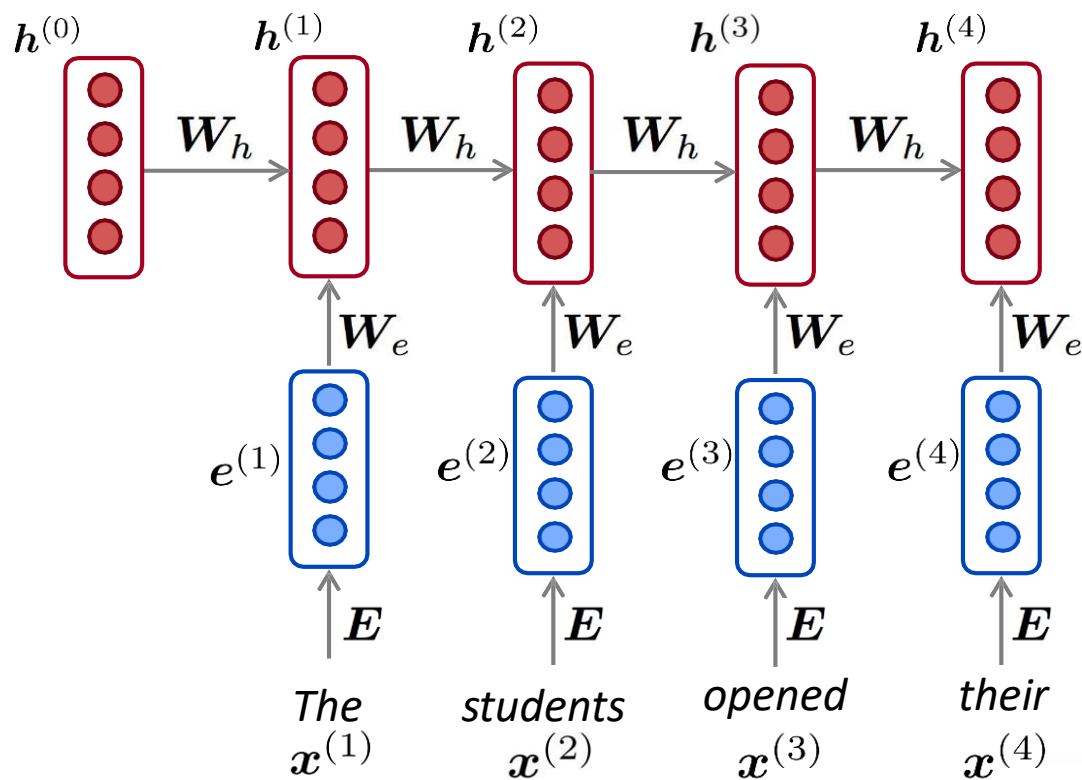
- 循环神经网络主要用于处理（变长）序列数据



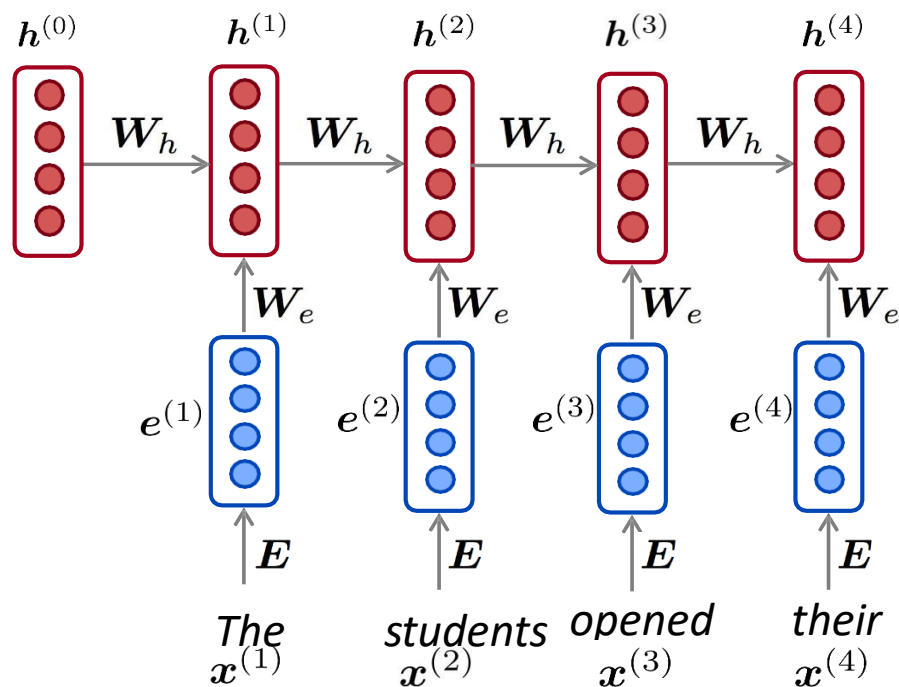
$$\mathbf{h}_t = f(U\mathbf{h}_{t-1} + W\mathbf{x}_t + \mathbf{b}),$$

$$\mathbf{y}_t = V\mathbf{h}_t,$$

RNN语言模型



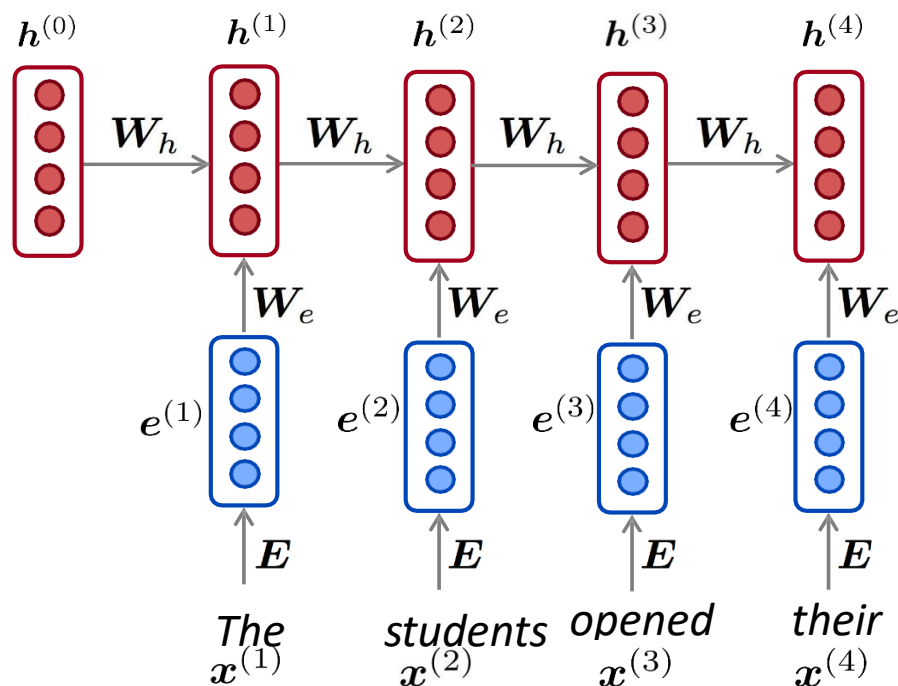
RNN语言模型



Level 1: words / one-hot vectors

$$\mathbf{x}^{(t)} \in \mathbb{R}^{|V|}$$

RNN语言模型



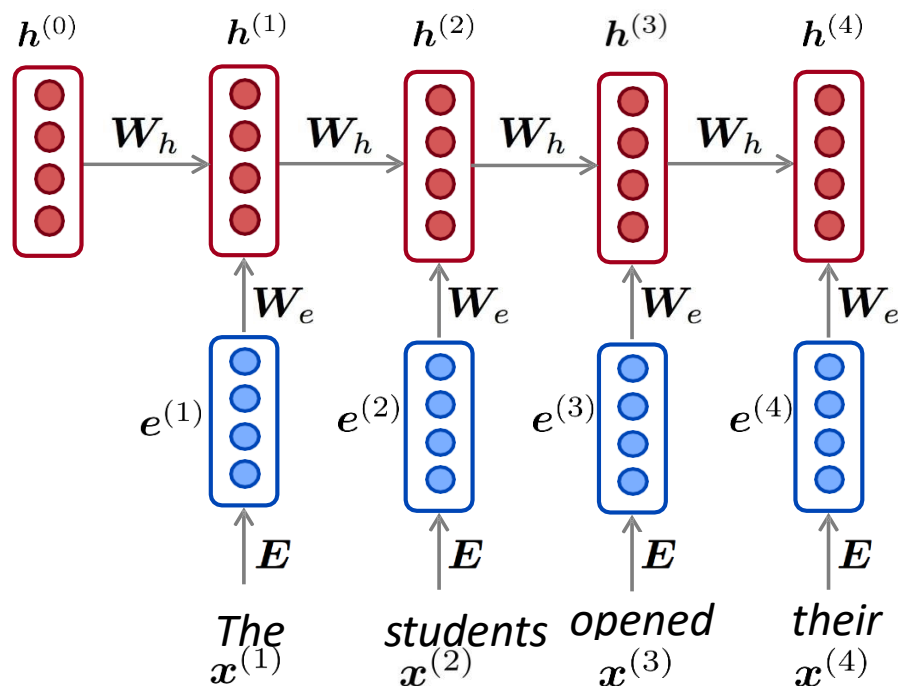
Level 2: word embeddings

$$e^{(t)} = E x^{(t)}$$

Level 1: words / one-hot vectors

$$x^{(t)} \in \mathbb{R}^{|V|}$$

RNN语言模型



Level 3: hidden states

$$h^{(t)} = \sigma \left(W_h h^{(t-1)} + W_e e^{(t)} + b_1 \right)$$

$h^{(0)}$ is the initial hidden state

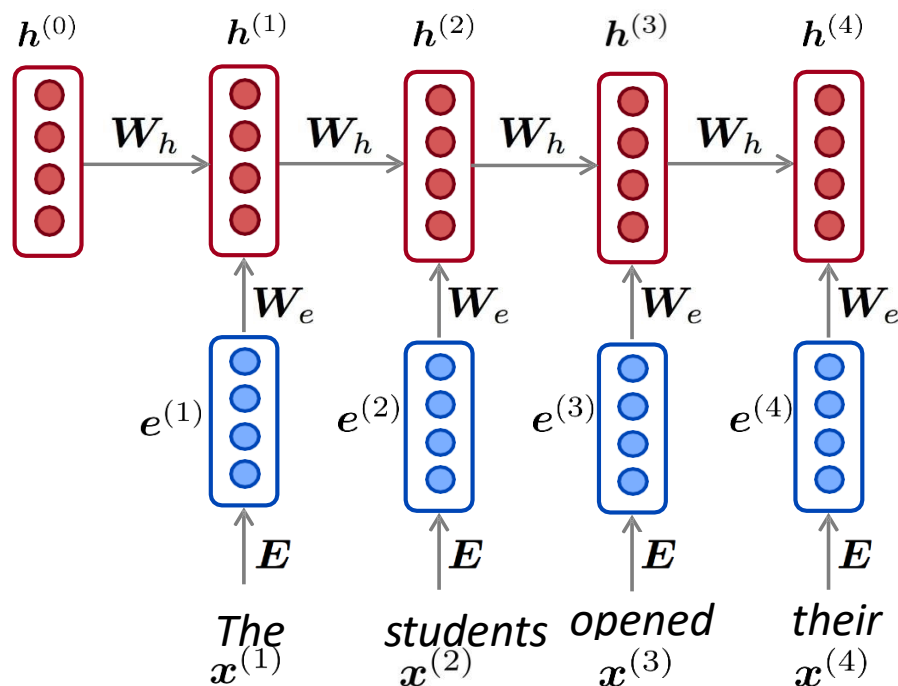
Level 2: word embeddings

$$e^{(t)} = E x^{(t)}$$

Level 1: words / one-hot vectors

$$x^{(t)} \in \mathbb{R}^{|V|}$$

RNN语言模型



Level 4: output distribution

$$\hat{y}^{(t)} = \text{softmax} \left(U h^{(t)} + b_2 \right) \in \mathbb{R}^{|V|}$$

Level 3: hidden states

$$h^{(t)} = \sigma \left(W_h h^{(t-1)} + W_e e^{(t)} + b_1 \right)$$

$h^{(0)}$ is the initial hidden state

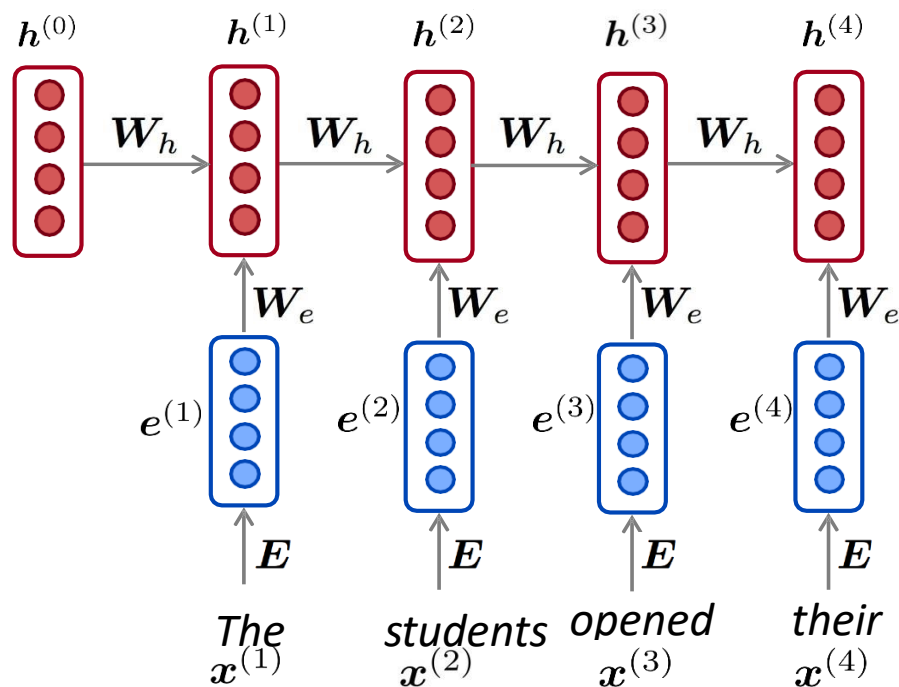
Level 2: word embeddings

$$e^{(t)} = E x^{(t)}$$

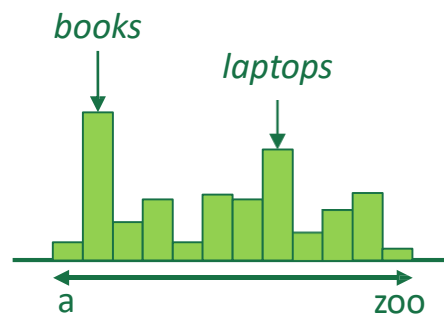
Level 1: words / one-hot vectors

$$x^{(t)} \in \mathbb{R}^{|V|}$$

RNN语言模型



$$\hat{y}^{(4)} = P(x^{(5)} | \text{the students opened their})$$



RNN Language Model

RNN 的优点:

- 能够处理任意长度的输入;

RNN Language Model

RNN 的优点:

- 能够处理任意长度的输入;
- t 时刻可以访问之前任意时刻的信息;

RNN Language Model

RNN 的优点:

- 能够处理任意长度的输入;
- t 时刻可以访问之前任意时刻的信息;
- 对于较长的输入, 模型的大小不变;

RNN Language Model

RNN 的优点:

- 能够处理任意长度的输入;
- t 时刻可以访问之前任意时刻的信息;
- 对于较长的输入, 模型的大小不变;
- 权重共享

RNN Language Model

RNN 的不足:

- 循环计算过程较慢;

RNN Language Model

RNN 的不足:

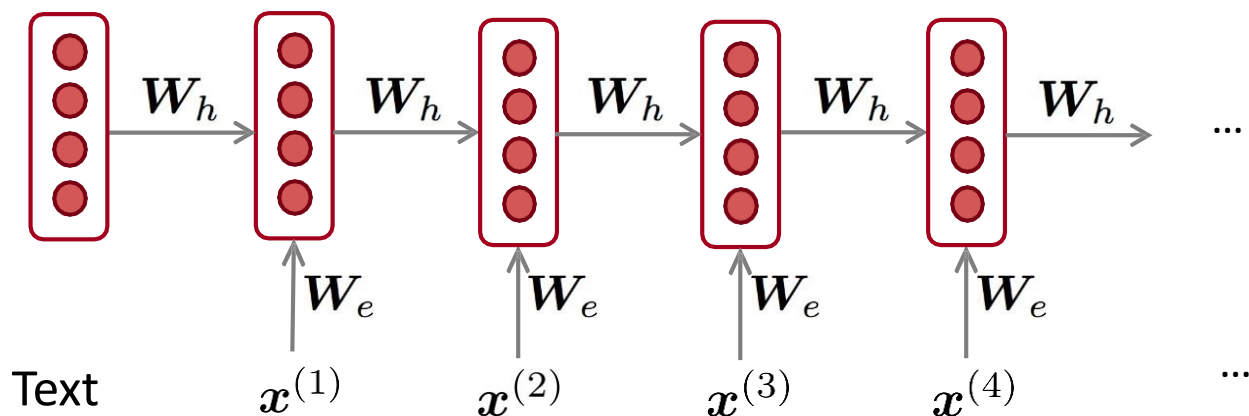
- 循环计算过程较慢;
- 实际运用中, 很难访问距当前时刻较远的信息;

RNN语言模型的训练

- Step 1: 输入文本数据集，每个文本表示为： $x^{(1)}, \dots, x^{(T)}$ ；

RNN语言模型的训练

- Step 1: 输入文本数据集，每个文本表示为： $x^{(1)}, \dots, x^{(T)}$ ；

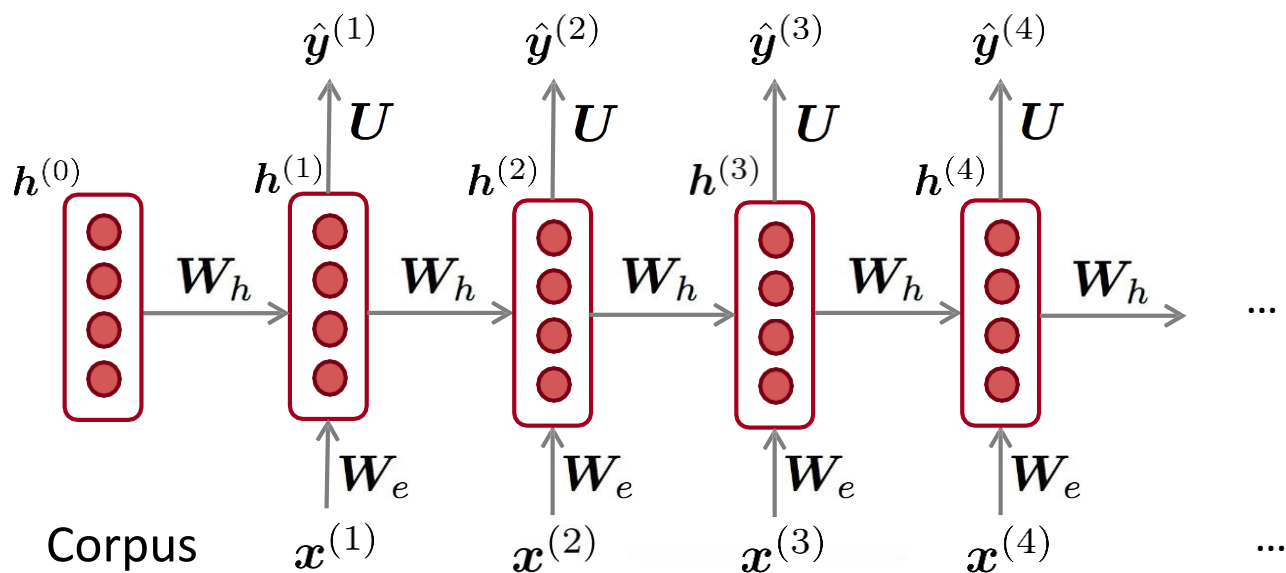


RNN语言模型的训练

- Step 1: 输入文本数据集，每个文本表示为： $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$ ；
- Step 2: 逐词把每条文本输入给RNN语言模型，计算每一时刻的输出概率 $\hat{y}^{(t)}$ ；

RNN语言模型的训练

- Step 1: 输入文本数据集，每个文本表示为： $x^{(1)}, \dots, x^{(T)}$ ；
- Step 2: 逐词把每条文本输入给RNN语言模型，计算每一时刻的输出概率 $\hat{y}^{(t)}$ ；



RNN语言模型的训练

- Step 1: 输入文本数据集，每个文本表示为： $x^{(1)}, \dots, x^{(T)}$ ；
- Step 2: 逐词把每条文本输入给RNN语言模型，计算每一时刻的输出概率 $\hat{y}^{(t)}$ ；
- Step 3: 用交叉熵(**cross-entropy**)计算每一时刻的损失(loss)

RNN语言模型的训练

- Step 1: 输入文本数据集，每个文本表示为： $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$ ；
- Step 2: 逐词把每条文本输入给RNN语言模型，计算每一时刻的输出概率 $\hat{\mathbf{y}}^{(t)}$ ；
- Step 3: 用交叉熵(**cross-entropy**)计算每一时刻的损失(loss)
 - 方法：使用时刻 t 的预测分布 $\hat{\mathbf{y}}^{(t)}$ 和下一个真实词 $\mathbf{x}^{(t+1)}$ 的表示 $\mathbf{y}^{(t)}$ ：

$$J^{(t)}(\theta) = CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = - \sum_{w \in V} \mathbf{y}_w^{(t)} \log \hat{\mathbf{y}}_w^{(t)} = - \log \hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)}$$

RNN语言模型的训练

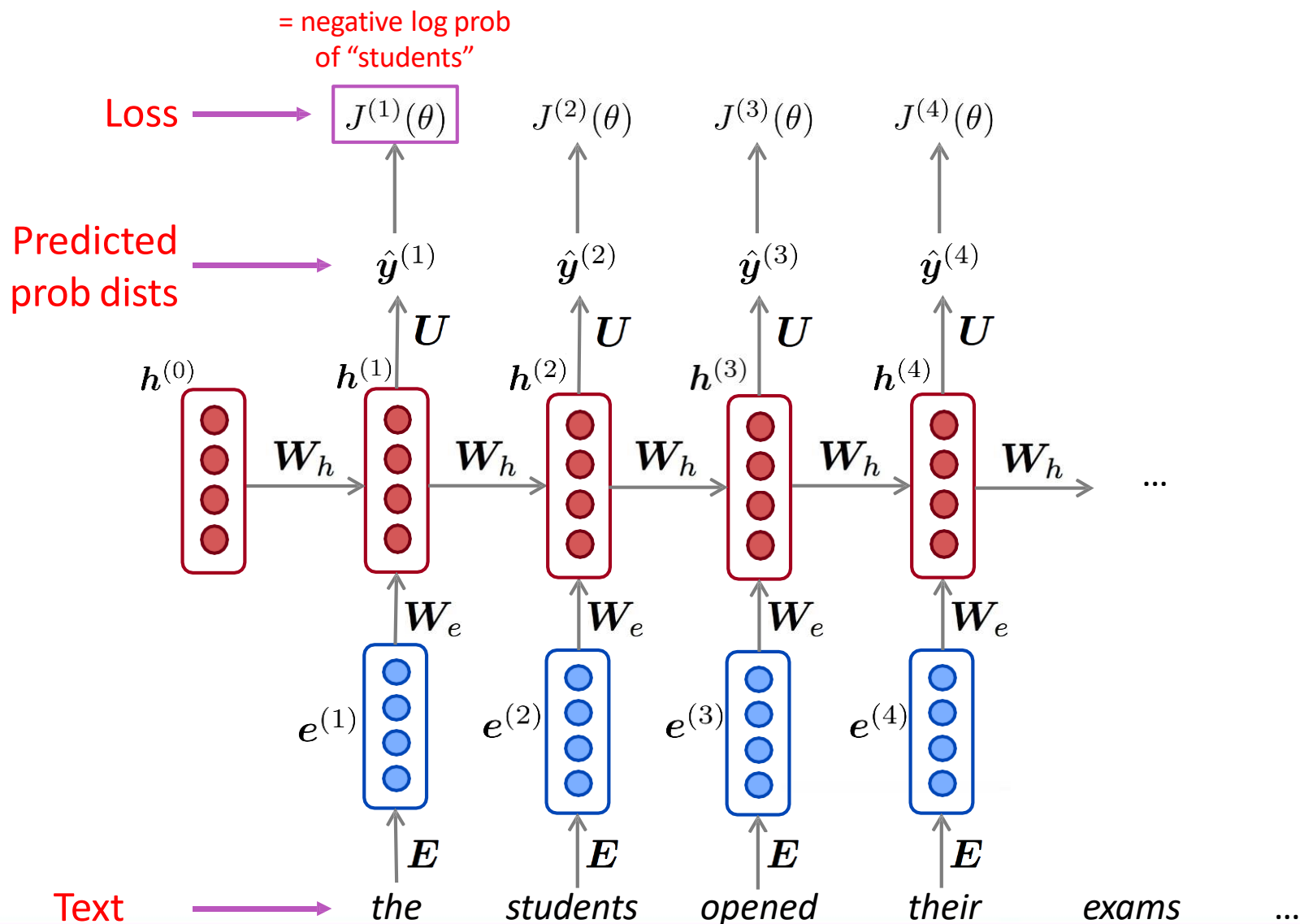
- Step 1: 输入文本数据集，每个文本表示为： $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$ ；
- Step 2: 逐词把每条文本输入给RNN语言模型，计算每一时刻的输出概率 $\hat{\mathbf{y}}^{(t)}$ ；
- Step 3: 用交叉熵(**cross-entropy**)计算每一时刻的损失(loss)
 - 方法：使用时刻 t 的预测分布 $\hat{\mathbf{y}}^{(t)}$ 和下一个真实词 $\mathbf{x}^{(t+1)}$ 的表示 $\mathbf{y}^{(t)}$ ：

$$J^{(t)}(\theta) = CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = - \sum_{w \in V} \mathbf{y}_w^{(t)} \log \hat{\mathbf{y}}_w^{(t)} = - \log \hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)}$$

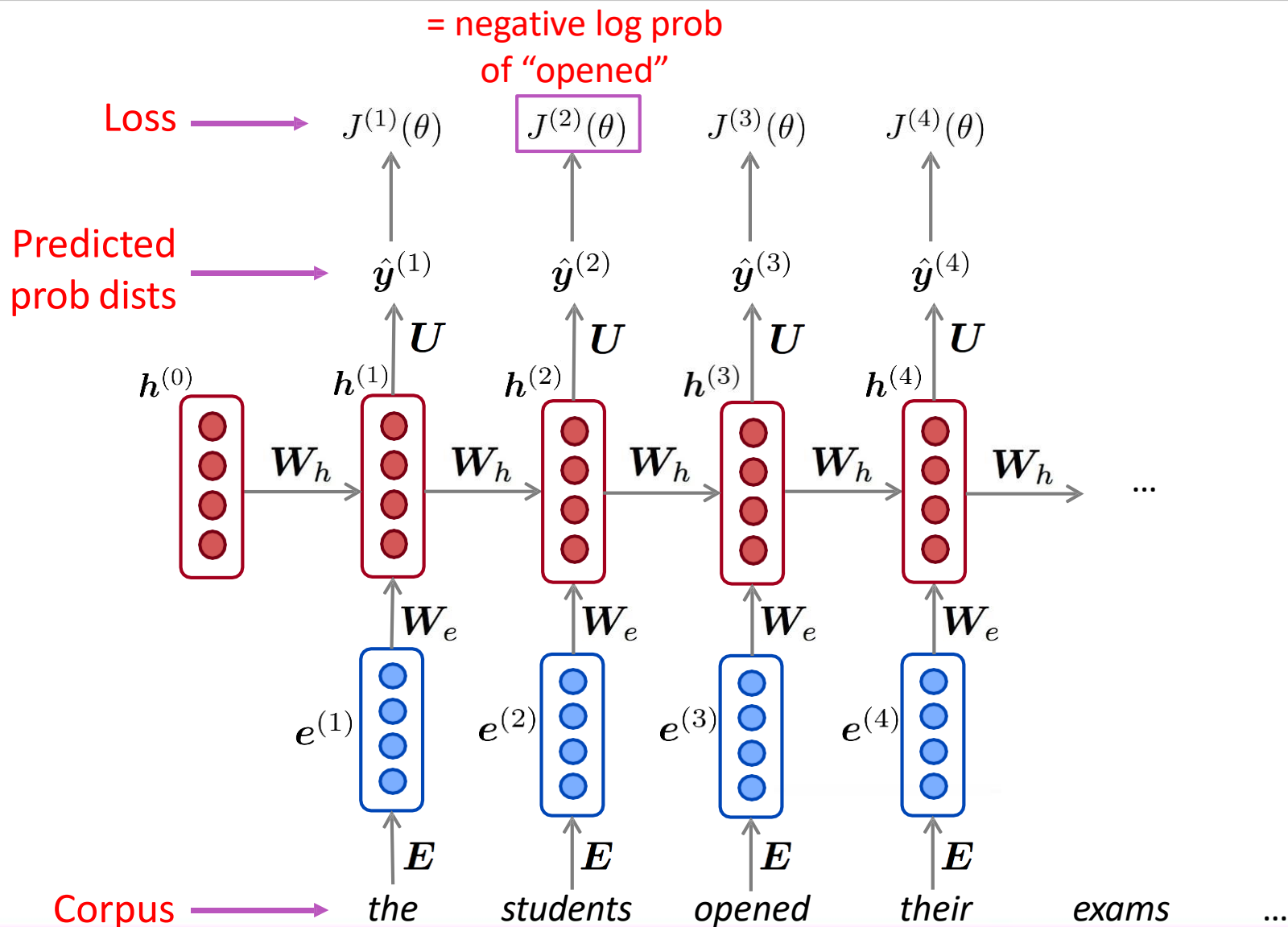
- Step 4: 计算所有文本的**平均损失**：

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^T - \log \hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)}$$

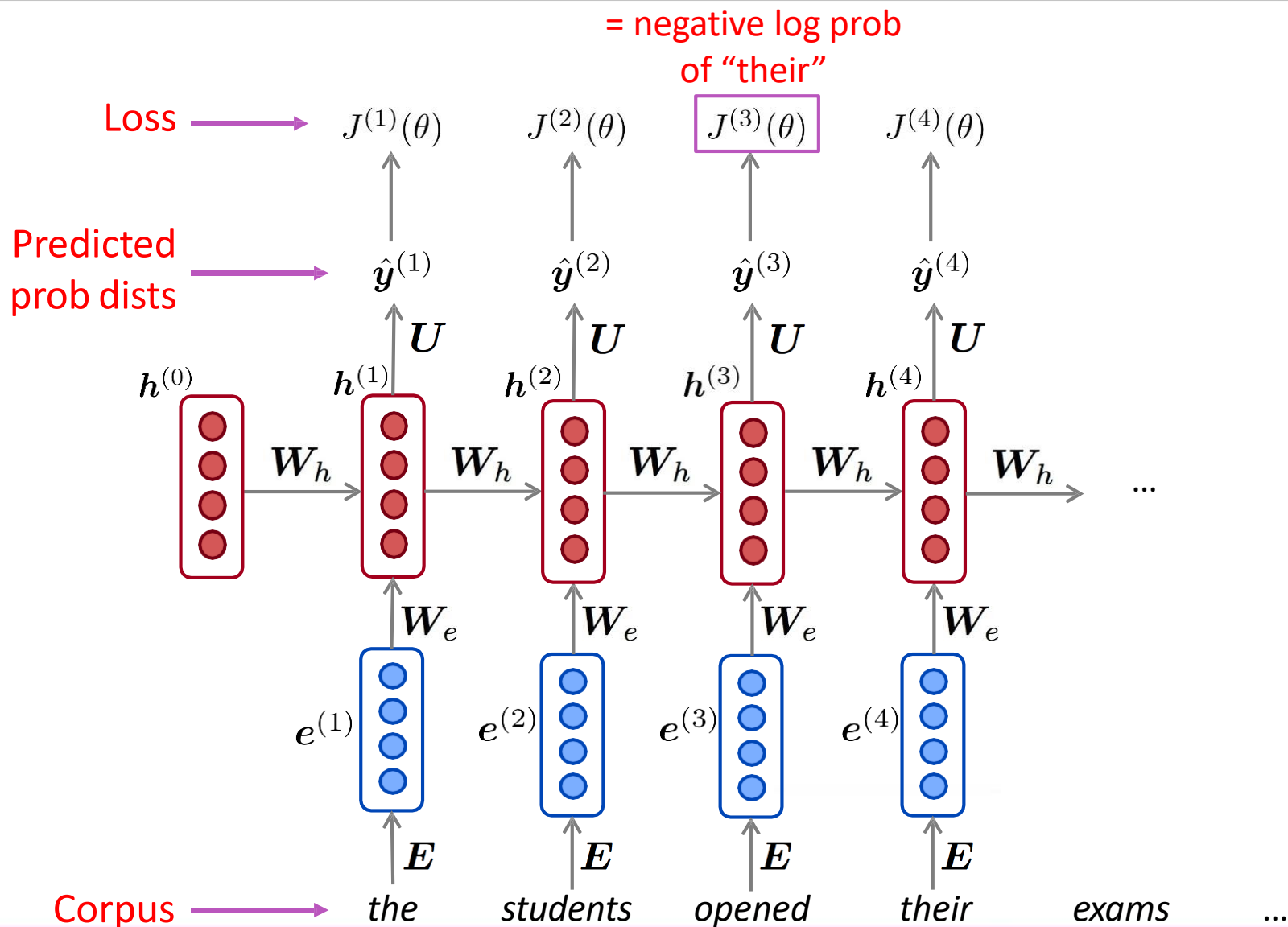
RNN语言模型的训练



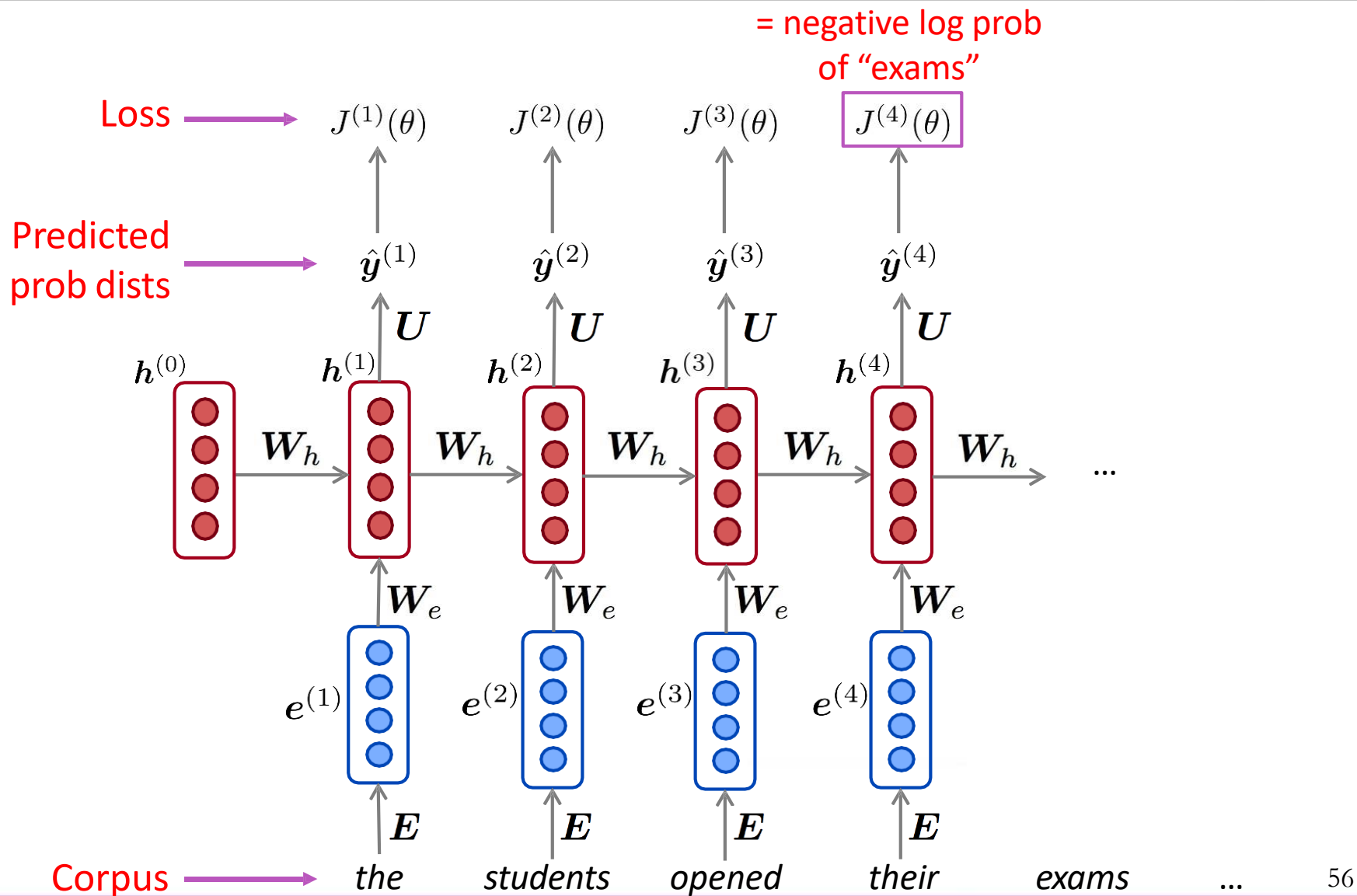
RNN语言模型的训练



RNN语言模型的训练

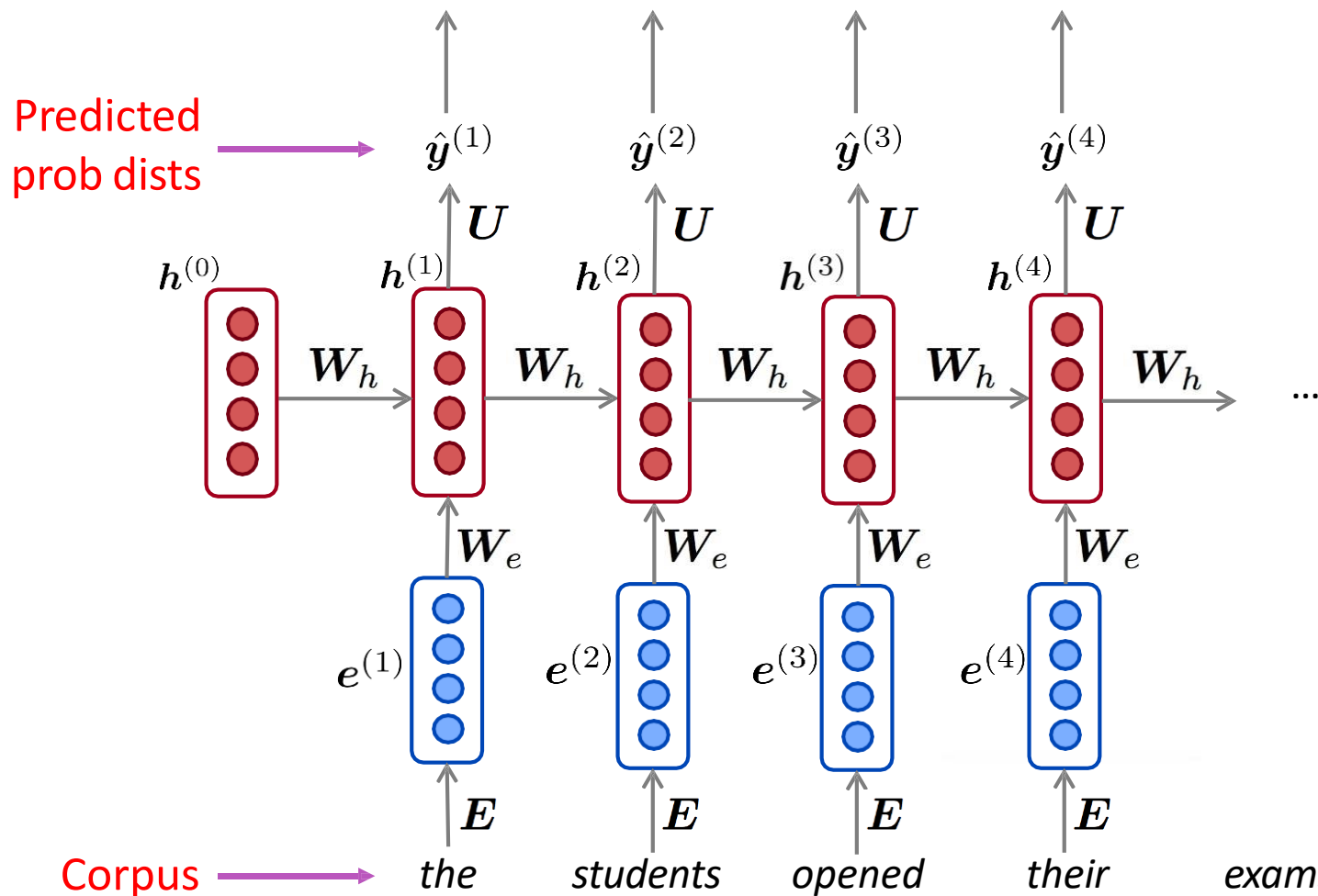


RNN语言模型的训练



RNN语言模型的训练

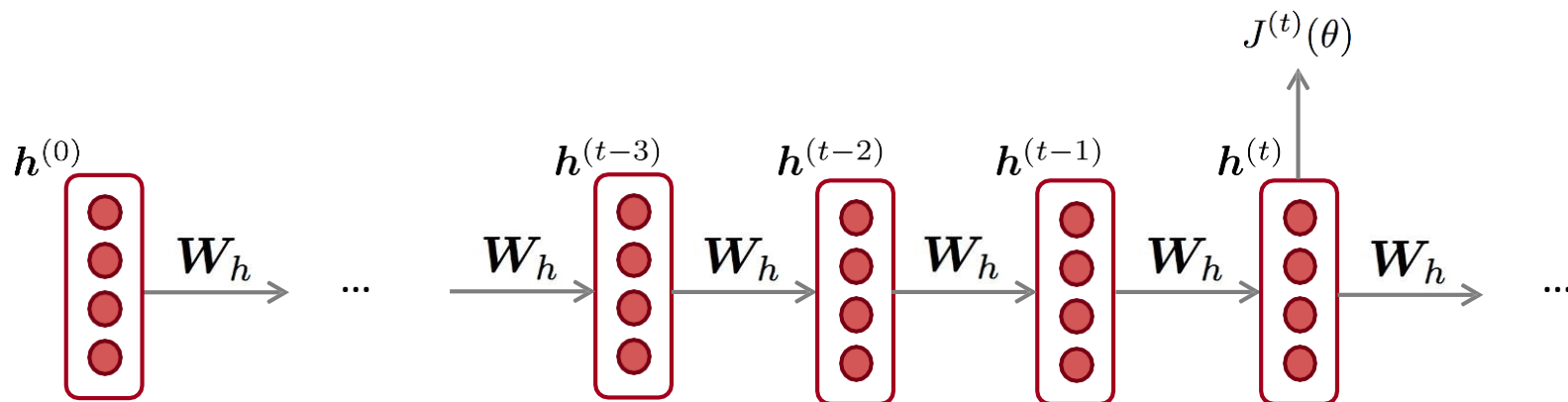
Loss $\longrightarrow J^{(1)}(\theta) + J^{(2)}(\theta) + J^{(3)}(\theta) + J^{(4)}(\theta) + \dots = J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta)$



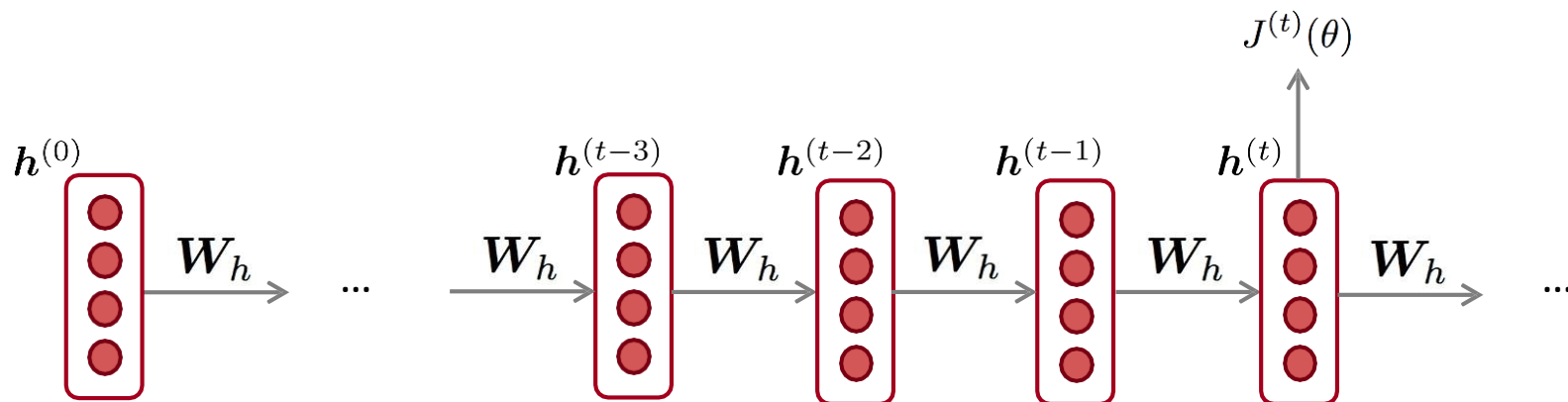
RNN语言模型的训练

- 通常，使用随机梯度下降(Stochastic Gradient Descent) 计算一小部分随机数据(mini-batch)的损失；
- 根据得到的损失，计算梯度并更新参数；

RNN的反向传播

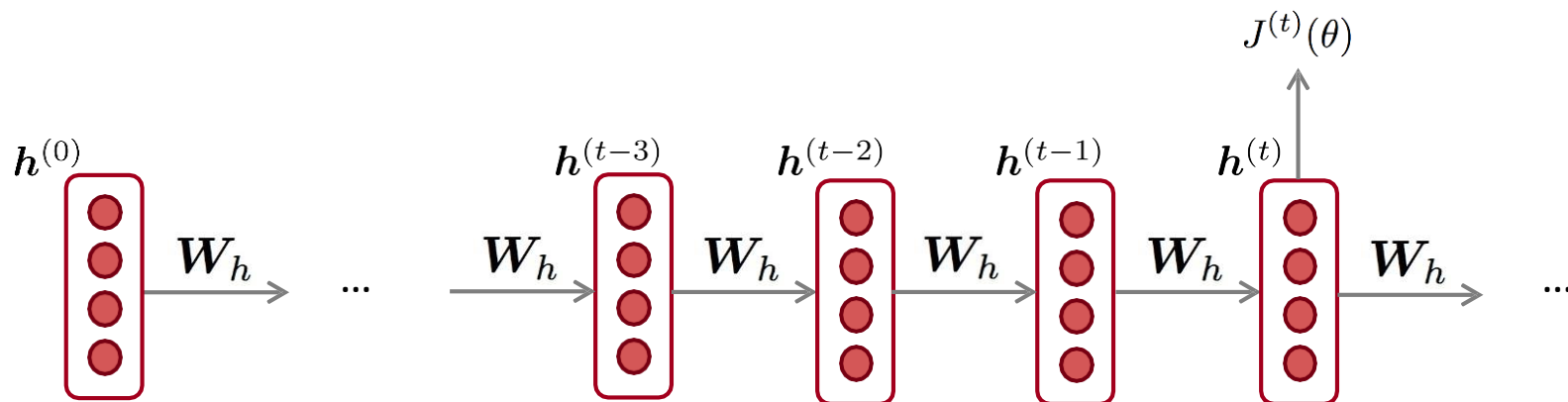


RNN的反向传播



问题: $J^{(t)}(\theta)$ 对权重矩阵 W_h 的偏导是什么?

RNN的反向传播

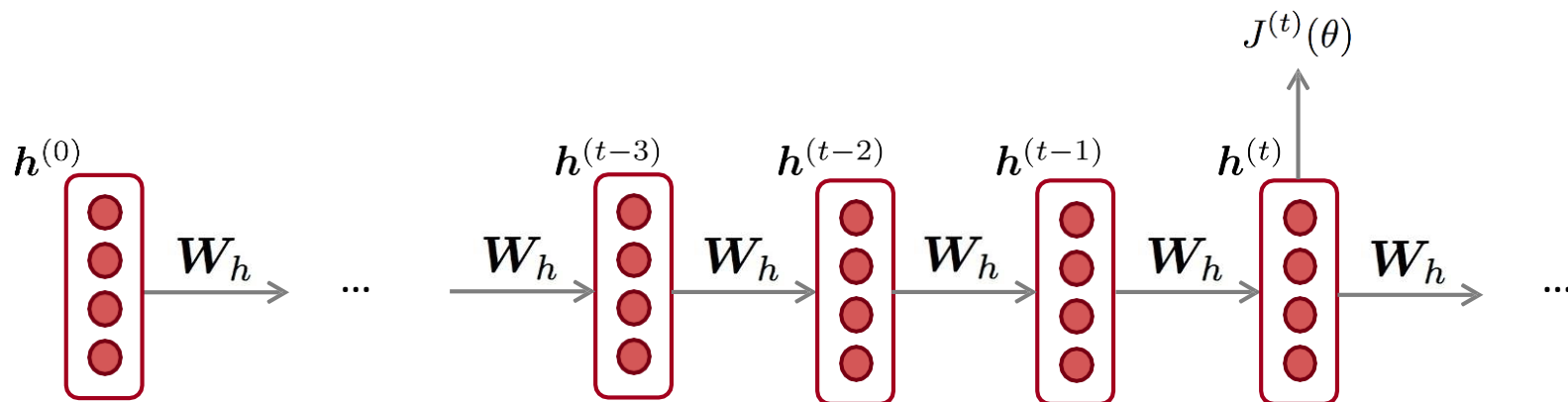


问题: $J^{(t)}(\theta)$ 对权重矩阵 W_h 的偏导是什么?

答案:
$$\frac{\partial J^{(t)}}{\partial W_h} = \sum_{i=1}^t \frac{\partial J^{(i)}}{\partial W_h} \Big|_{(i)}$$

对于重复出现的变量，它的梯度是每次出现的梯度之和。

RNN的反向传播



问题: $J^{(t)}(\theta)$ 对权重矩阵 W_h 的偏导是什么?

答案:
$$\frac{\partial J^{(t)}}{\partial W_h} = \sum_{i=1}^t \frac{\partial J^{(t)}}{\partial W_h} \Big|_{(i)}$$

对于重复出现的变量，它的梯度是每次出现的梯度之和。

WHY?

RNN的反向传播

多变量链式法则：

- Given a multivariable function $f(x, y)$, and two single variable functions $x(t)$ and $y(t)$, here's what the multivariable chain rule says:

$$\underbrace{\frac{d}{dt} f(\mathbf{x}(t), \mathbf{y}(t))}_{\text{Derivative of composition function}} = \frac{\partial f}{\partial \mathbf{x}} \frac{d\mathbf{x}}{dt} + \frac{\partial f}{\partial \mathbf{y}} \frac{d\mathbf{y}}{dt}$$

Derivative of composition function

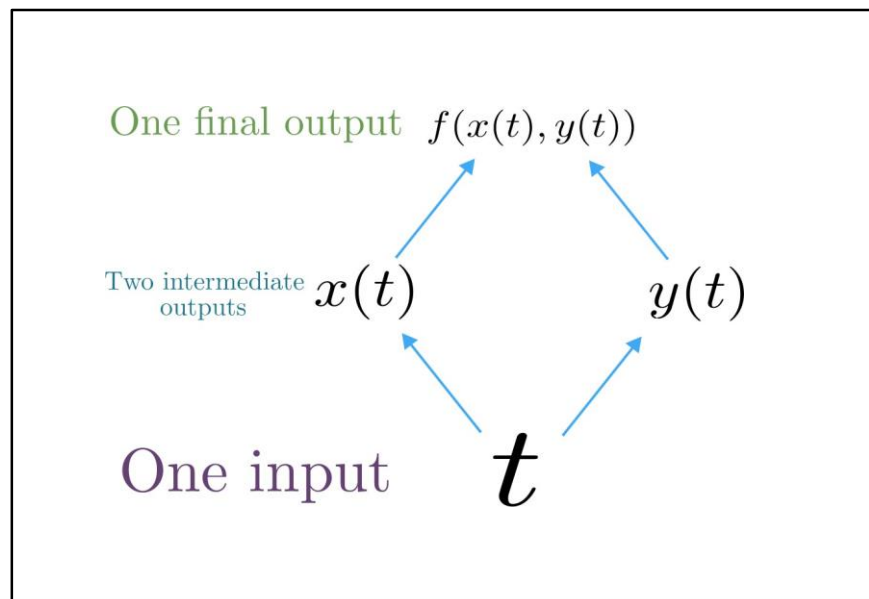
RNN的反向传播

多变量链式法则：

- Given a multivariable function $f(x, y)$, and two single variable functions $x(t)$ and $y(t)$, here's what the multivariable chain rule says:

$$\underbrace{\frac{d}{dt} f(\mathbf{x}(t), \mathbf{y}(t))}_{\text{Derivative of composition function}} = \frac{\partial f}{\partial \mathbf{x}} \frac{d\mathbf{x}}{dt} + \frac{\partial f}{\partial \mathbf{y}} \frac{d\mathbf{y}}{dt}$$

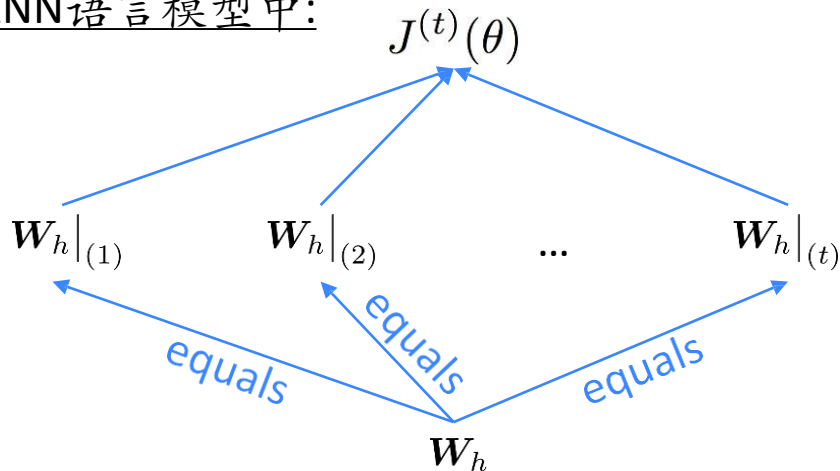
Derivative of composition function



RNN的反向传播

多变量链式法则：

RNN语言模型中：



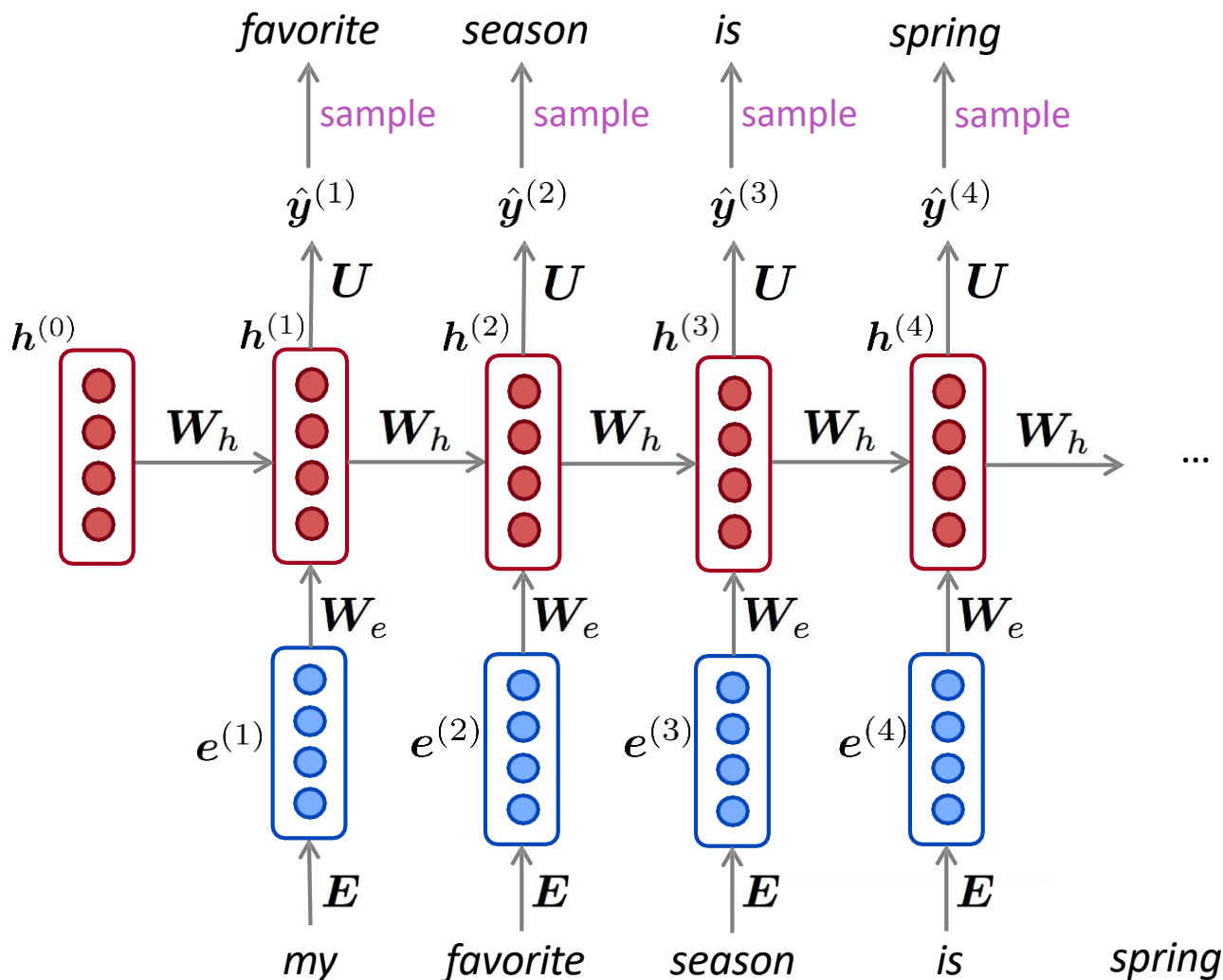
运用多变量链式法则：

$$\begin{aligned}\frac{\partial J^{(t)}}{\partial W_h} &= \sum_{i=1}^t \frac{\partial J^{(t)}}{\partial W_h} \bigg|_{(i)} \boxed{\frac{\partial W_h|_{(i)}}{\partial W_h}} = 1 \\ &= \sum_{i=1}^t \frac{\partial J^{(t)}}{\partial W_h} \bigg|_{(i)}\end{aligned}$$

RNN语言模型用于文本生成

- 和n-gram 语言模型一样，RNN语言模型可以通过循环采样生成文本：

RNN语言模型用于文本生成



RNN语言模型用于文本生成

- 可以使用RNN语言模型在各类文本数据上训练，然后生成相应风格的文本；
- 例如，在奥巴马的演说文本上：



The United States will step up to the cost of a new challenges of the American people that will share the fact that we created the problem. They were attacked and so that they have to say that all the task of the final days of war that I will not be able to get this done. The promise of the men and women who were still going

RNN语言模型用于文本生成

- 可以使用RNN语言模型在各类文本数据上训练，然后生成相应风格的文本；
- 例如，在哈利波特小说上：



“Sorry,” Harry shouted, panicking—“I’ll leave those brooms in London, are they?”

“No idea,” said Nearly Headless Nick, casting low close by Cedric, carrying the last bit of treacle Charms, from Harry’s shoulder, and to answer him the common room perched upon it, four arms held a shining knob from when the spider hadn’t felt it seemed. He reached the teams too.

RNN语言模型用于文本生成

- 可以使用RNN语言模型在各类文本数据上训练，然后生成相应风格的文本；
- 例如，在菜谱上：

Title: CHOCOLATE RANCH BARBECUE

Categories: Game, Casseroles, Cookies, Cookies

Yield: 6 Servings

2 tb Parmesan cheese -- chopped

1 c Coconut milk

3 Eggs, beaten

Place each pasta over layers of lumps. Shape mixture into the moderate oven and simmer until firm. Serve hot in bodied fresh, mustard, orange and cheese.

Combine the cheese and salt together the dough in a large skillet; add the ingredients and stir in the chocolate and pepper.



语言模型的评估

□ 评估语言模型的标准方法是计算困惑度(perplexity)

$$\begin{aligned} PP(S) &= P(w_1 w_2 \dots w_N)^{-\frac{1}{N}} \\ &= \sqrt[N]{\frac{1}{p(w_1 w_2 \dots w_N)}} \\ &= \sqrt[N]{\prod_{i=1}^N \frac{1}{p(w_i | w_1 w_2 \dots w_{i-1})}} \end{aligned}$$

语言模型的评估

□ 测试集上的困惑度(perplexity):

$$\text{perplexity}(D_{\text{test}}) = \exp \left\{ -\frac{\sum_{d=1}^M \log p(\mathbf{w}_d)}{\sum_{d=1}^M N_d} \right\}^{[1]}$$

[1] Blei, David M., Andrew Y. Ng, and Michael I. Jordan. "Latent dirichlet allocation." Journal of machine Learning research 3.Jan (2003): 993-1022.

语言模型的评估

RNN能够显著提升模型的困惑度：

n-gram model →

Increasingly complex RNNs ↓

Model	Perplexity
Interpolated Kneser-Ney 5-gram (Chelba et al., 2013)	67.6
RNN-1024 + MaxEnt 9-gram (Chelba et al., 2013)	51.3
RNN-2048 + BlackOut sampling (Ji et al., 2015)	68.3
Sparse Non-negative Matrix factorization (Shazeer et al., 2015)	52.9
LSTM-2048 (Jozefowicz et al., 2016)	43.7
2-layer LSTM-8192 (Jozefowicz et al., 2016)	30
Ours small (LSTM-2048)	43.9
Ours large (2-layer LSTM-2048)	39.8

Perplexity improves (lower is better)

Source: <https://research.fb.com/building-an-efficient-neural-language-model-over-a-billion-words/>

注意!

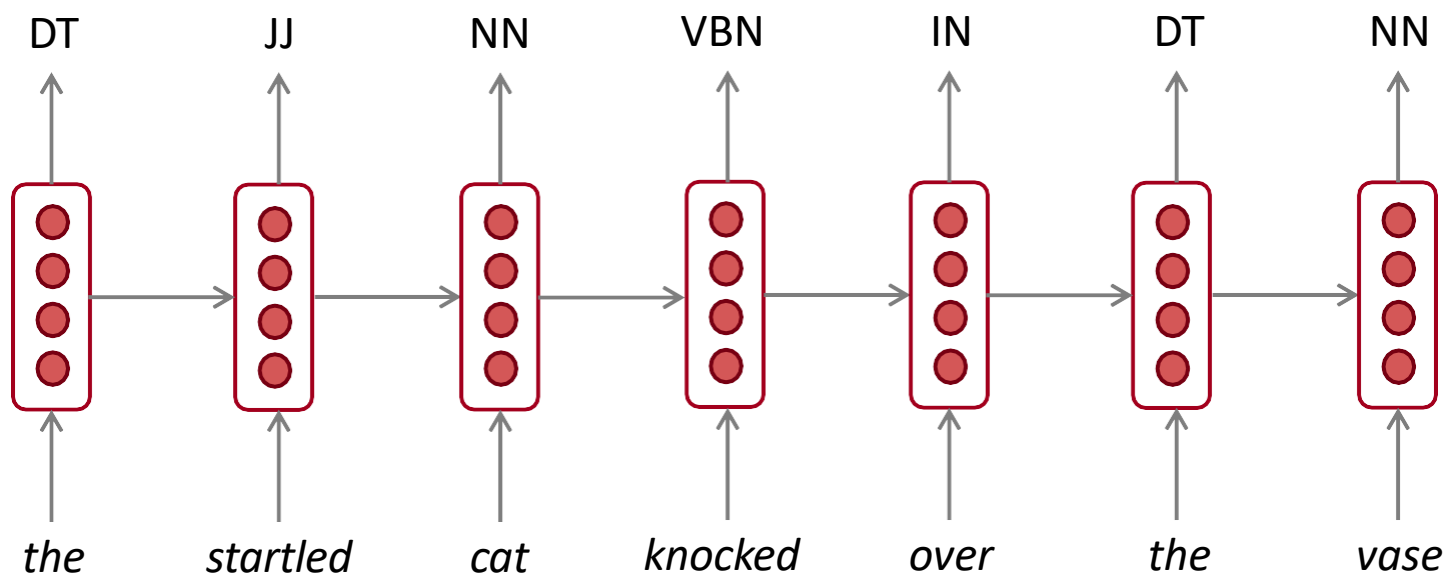
- **语言模型**: 预测下一个词的系统
- **RNN**: 一种神经网络:
 - 输入是任意长度的文本;
 - 共享权重;
 - 每一步都可以产生输出;
- Recurrent Neural Network \neq Language Model
- RNN可以用来构建语言模型, 但RNN还有更多用途;

注意!

- RNN可以用于做序列标注，如词性标注

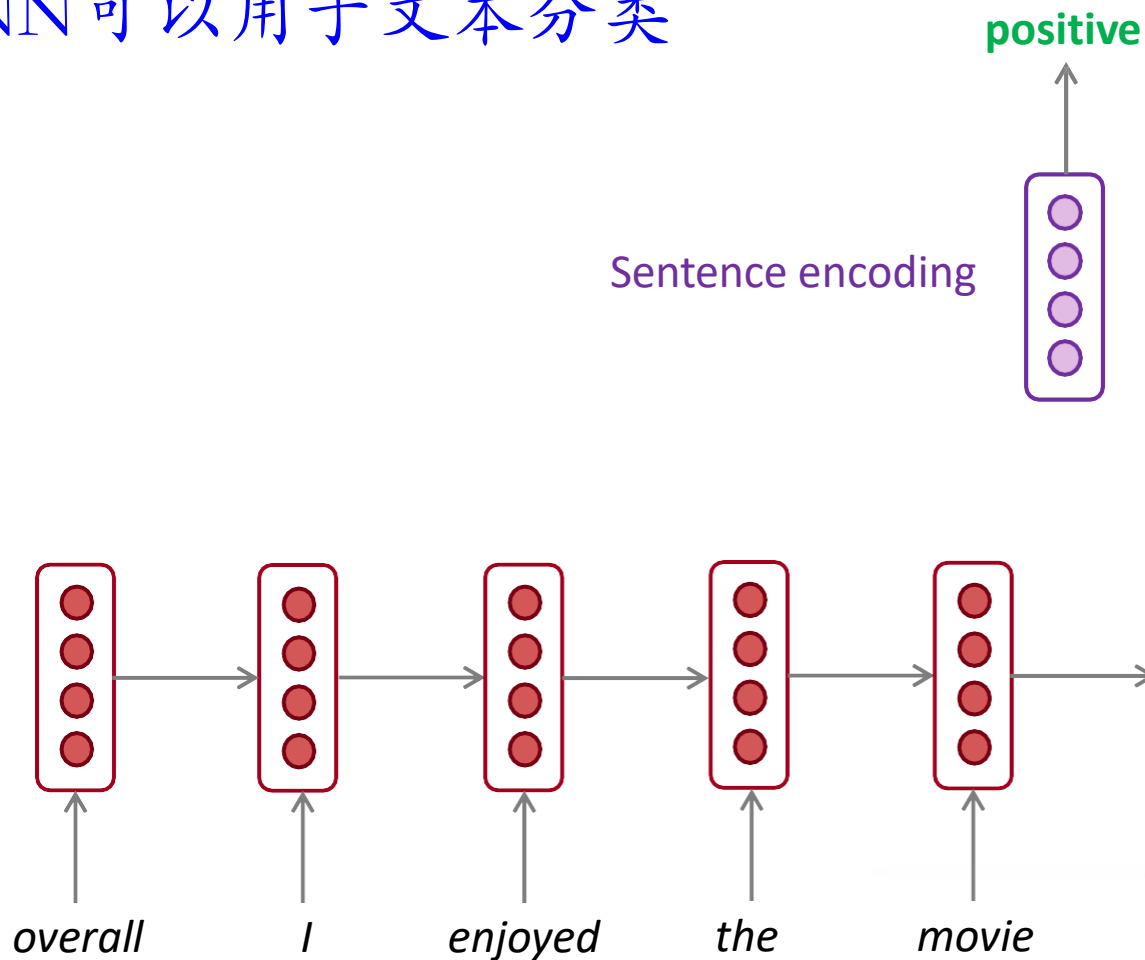
注意!

- RNN可以用于做序列标注，如词性标注



注意!

- RNN可以用于文本分类

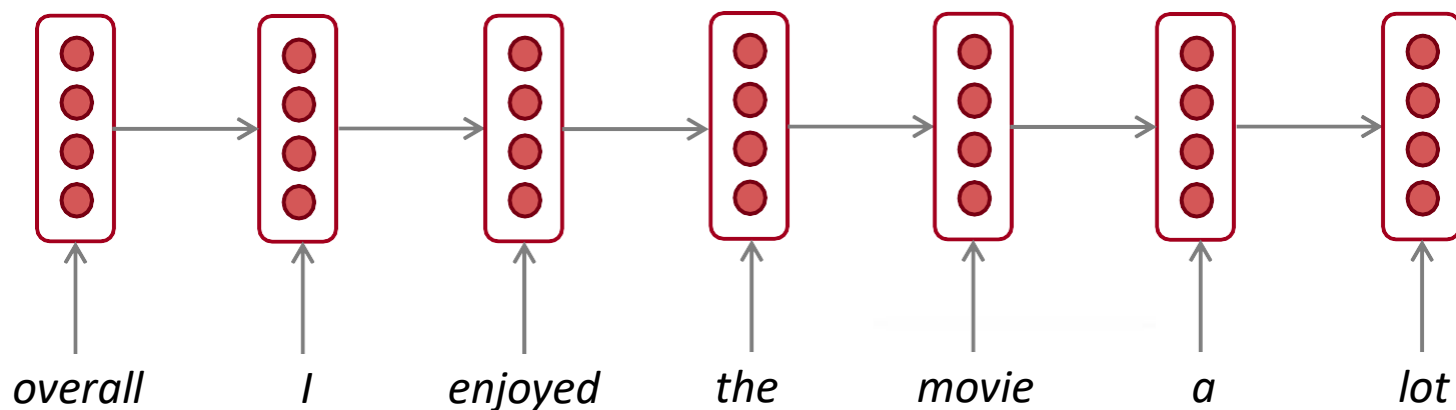


注意!

- RNN可以用于文本分类

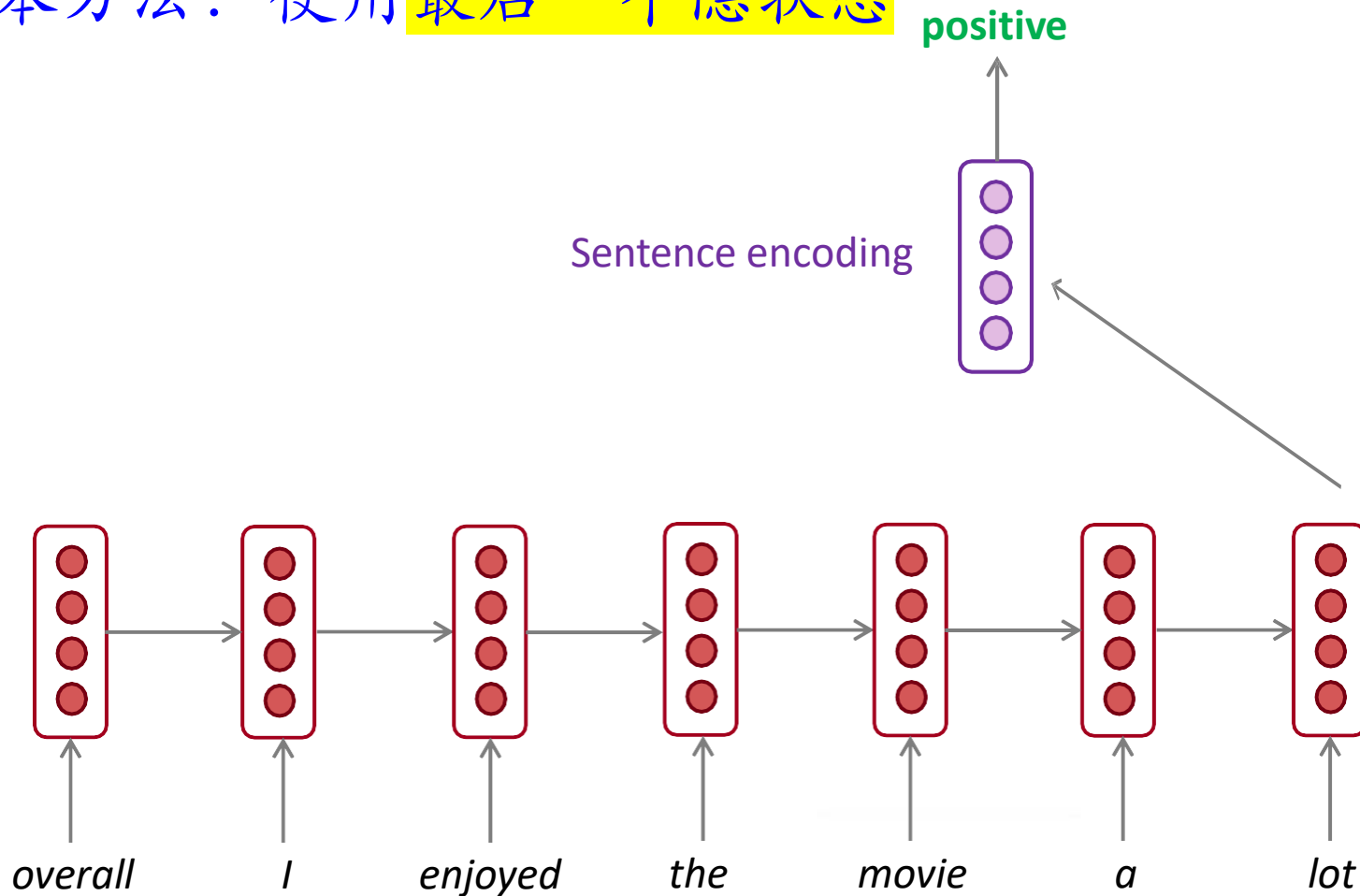
positive

如何计算句子的编码（表示）？



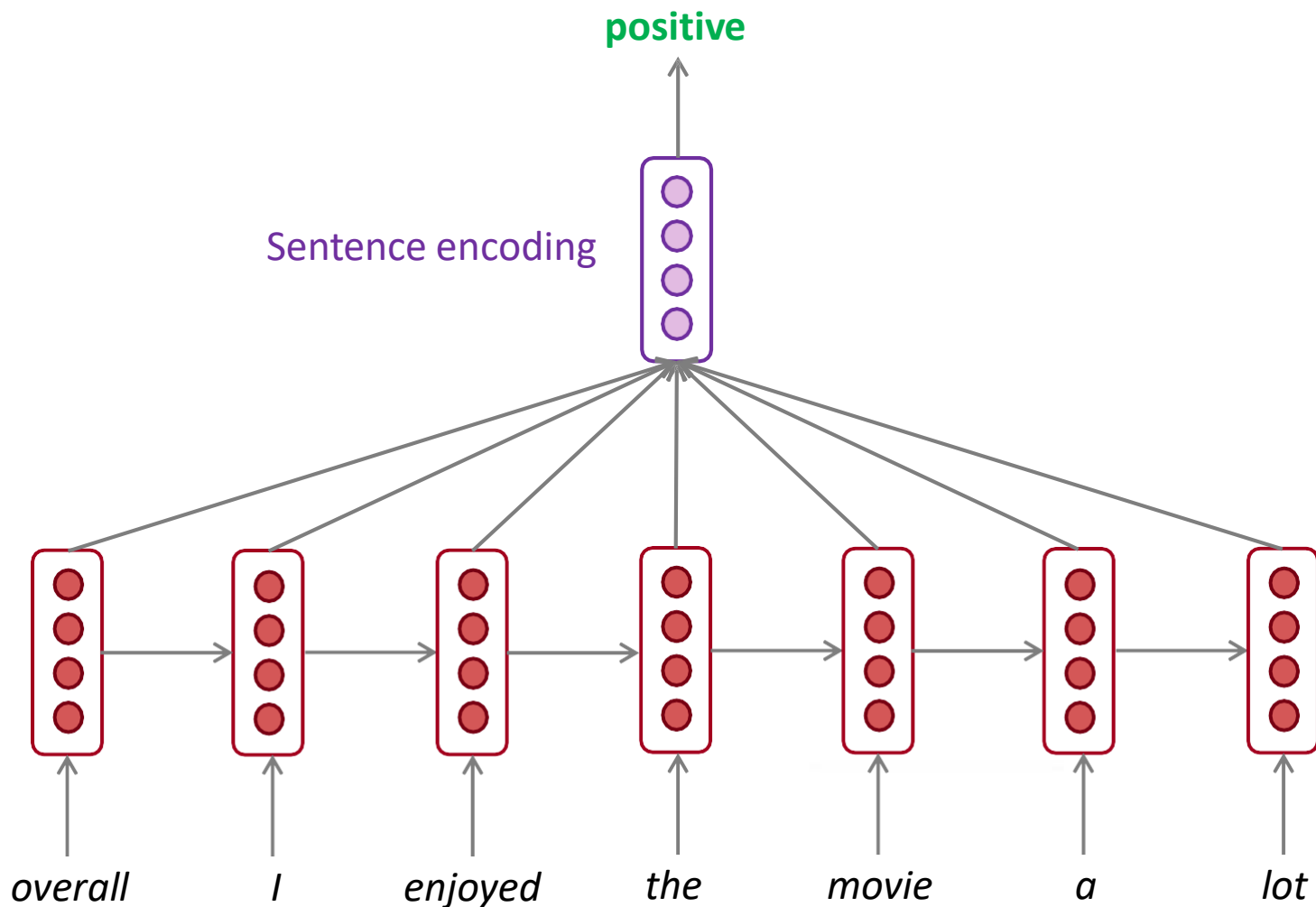
注意!

- 基本方法：使用最后一个隐状态



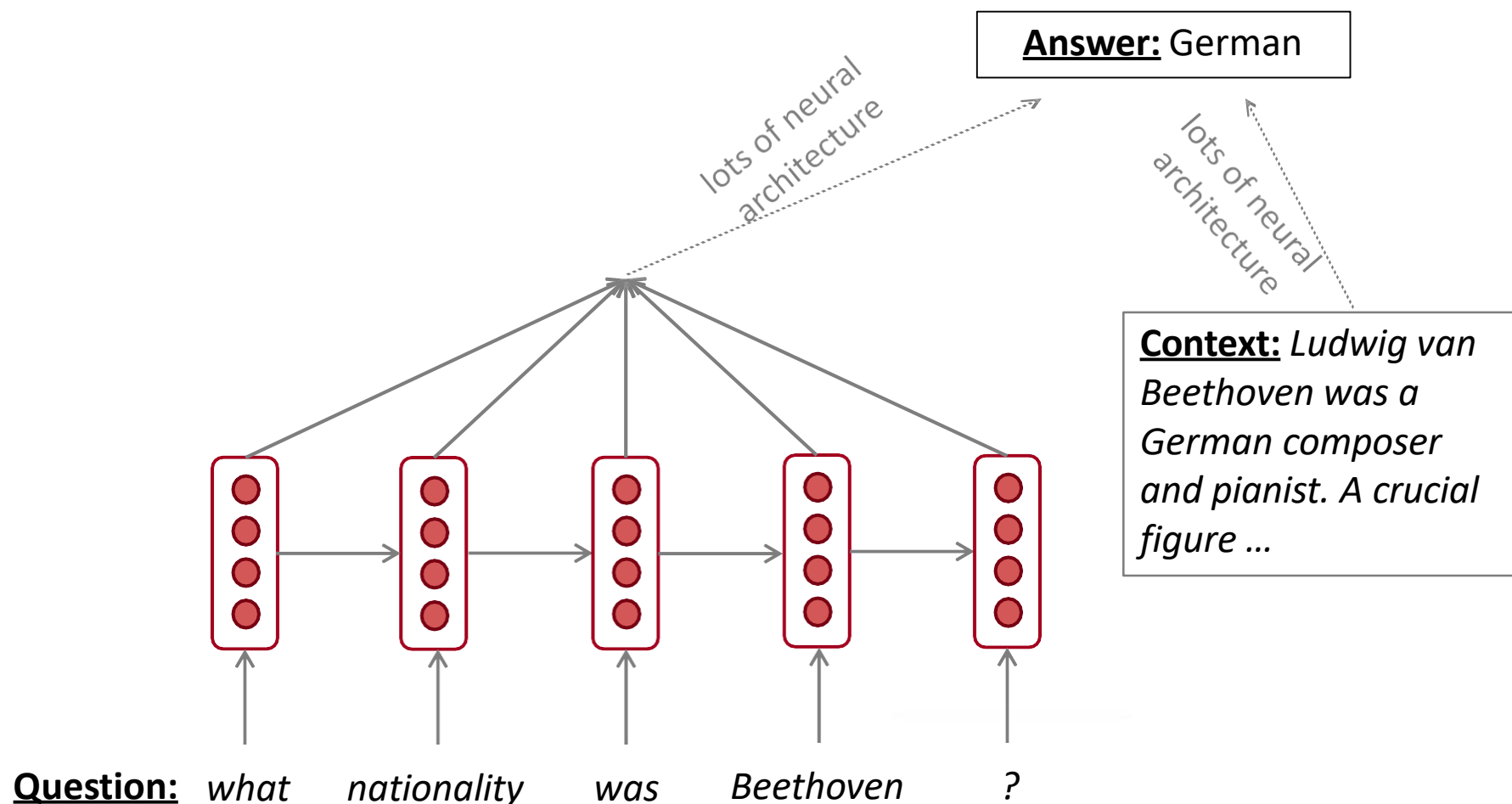
注意!

- 更好方法：取所有隐状态的最大值或者平均值



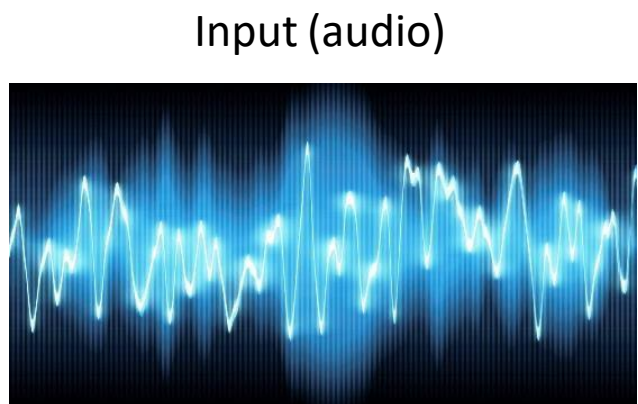
注意!

- RNN还可以用于问答系统和机器翻译等任务

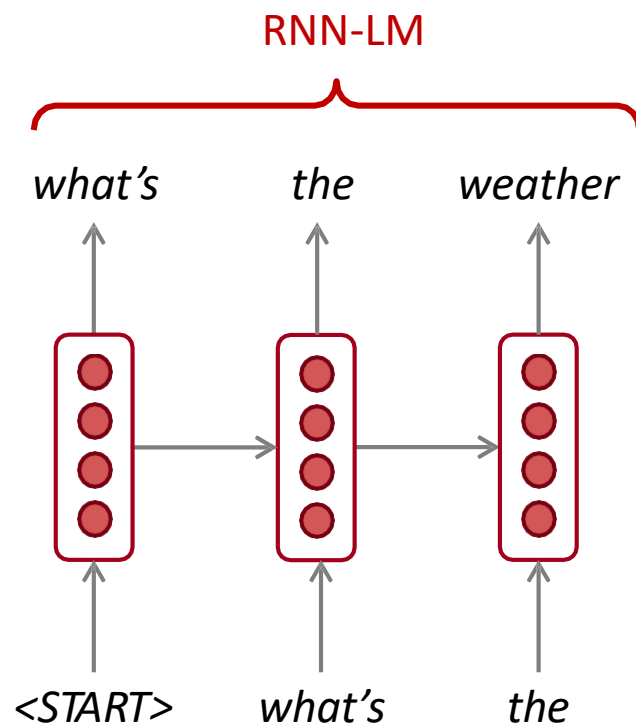


注意!

- RNN还可以用于语音识别



conditioning
.....>



注意!

- 常见的RNN模型
 - LSTM
 - GRU
 - Bidirectional
 - Multi-layer

Thank you!

权小军 中山大学数据科学与计算机学院