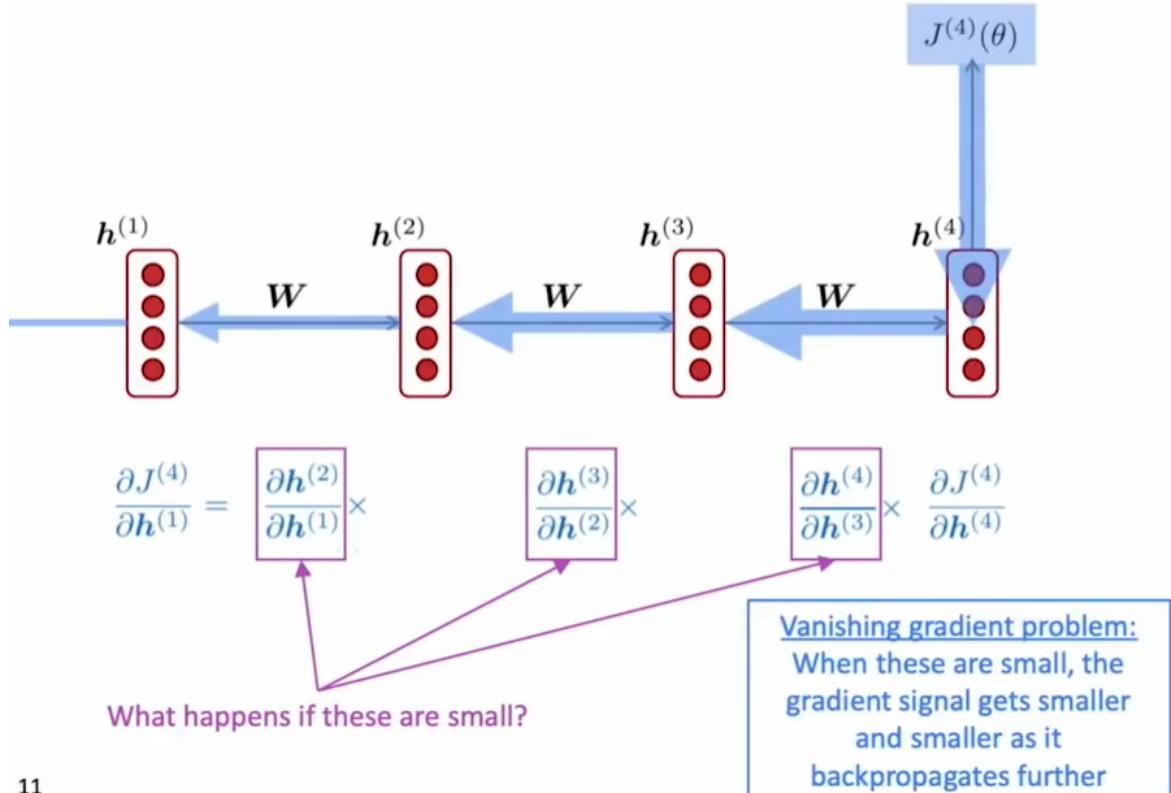


Vanishing gradient problem

链式法则可能导致梯度消失。



11

数学细节推导，造成梯度消失的是 W_h 的幂次：

Vanishing gradient proof sketch

- Recall: $\mathbf{h}^{(t)} = \sigma(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1)$
- Therefore: $\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} = \text{diag}(\sigma'(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1)) \mathbf{W}_h$ (chain rule)
- Consider the gradient of the loss $J^{(i)}(\theta)$ on step i , with respect to the hidden state $\mathbf{h}^{(j)}$ on some previous step j .

$$\begin{aligned} \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(j)}} &= \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \prod_{j < t \leq i} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \\ &= \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \boxed{\mathbf{W}_h^{(i-j)}} \prod_{j < t \leq i} \text{diag}(\sigma'(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1)) \end{aligned}$$

↑
(value of $\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}}$)

If W_h is small, then this term gets vanishingly small as i and j get further apart

当 W_h 的最大特征值大于1，梯度就会以指数速度消失;大于1，梯度以指数速度爆炸;

Vanishing gradient proof sketch

- Consider matrix L2 norms:

$$\left\| \frac{\partial J^{(i)}(\theta)}{\partial h^{(j)}} \right\| \leq \left\| \frac{\partial J^{(i)}(\theta)}{\partial h^{(i)}} \right\| \|W_h\|^{(i-j)} \prod_{j < t \leq i} \left\| \text{diag} \left(\sigma' \left(W_h h^{(t-1)} + W_x x^{(t)} + b_1 \right) \right) \right\|$$

- Pascanu et al showed that if the largest eigenvalue of W_h is less than 1, then the gradient $\left\| \frac{\partial J^{(i)}(\theta)}{\partial h^{(j)}} \right\|$ will shrink exponentially
 - Here the bound is 1 because we have sigmoid nonlinearity
- There's a similar proof relating a largest eigenvalue >1 to exploding gradients

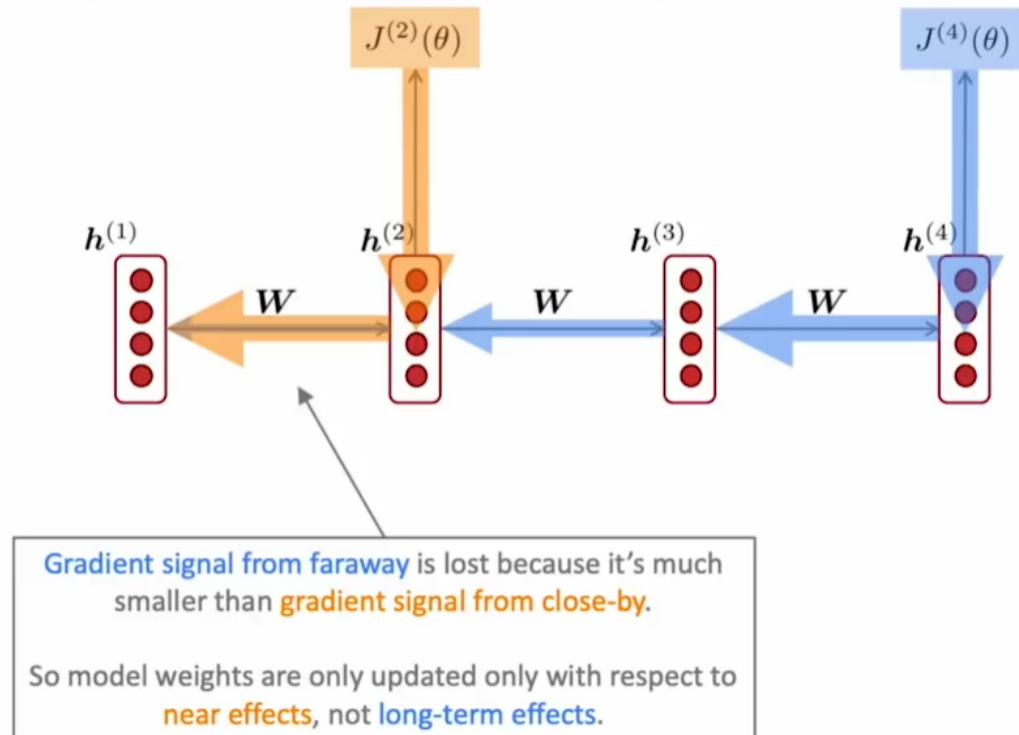
13

Source: "On the difficulty of training recurrent neural networks", Pascanu et al, 2013. <http://proceedings.mlr.press/v28/pascanu13.pdf>

梯度消失有什么问题?

会导致 $h^{(i)}$ 受 $J^{(i)}$ 的影响远大于 J^{i+j} , 即 $h^{(i)}$ 对 long-term 的影响在梯度下降时考虑很少。

Why is vanishing gradient a problem?



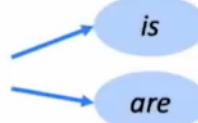
14

几个梯度消失的例子

Effect of vanishing gradient on RNN-LM

- LM task: When she tried to print her tickets, she found that the printer was out of toner. She went to the stationery store to buy more toner. It was very overpriced. After installing the toner into the printer, she finally printed her _____
- To learn from this training example, the RNN-LM needs to model the dependency between “tickets” on the 7th step and the target word “tickets” at the end.
- But if gradient is small, the model can't learn this dependency
 - So the model is unable to predict similar long-distance dependencies at test time

Effect of vanishing gradient on RNN-LM

- LM task: The writer of the books __ 
- Correct answer: The writer of the books is planning a sequel
- Syntactic recency: The writer of the books is (correct)
- Sequential recency: The writer of the books are (incorrect)
- Due to vanishing gradient, RNN-LMs are better at learning from sequential recency than syntactic recency, so they make this type of error more often than we'd like [Linzen et al 2016]

梯度爆炸有什么问题？

每次梯度下降幅度过大，甚至可能出现NaN和Inf。

解决方法

- gradient clipping

Algorithm 1 Pseudo-code for norm clipping

```
 $\hat{g} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$ 
if  $\|\hat{g}\| \geq threshold$  then
     $\hat{g} \leftarrow \frac{threshold}{\|\hat{g}\|} \hat{g}$ 
end if
```

LSTM and GRU

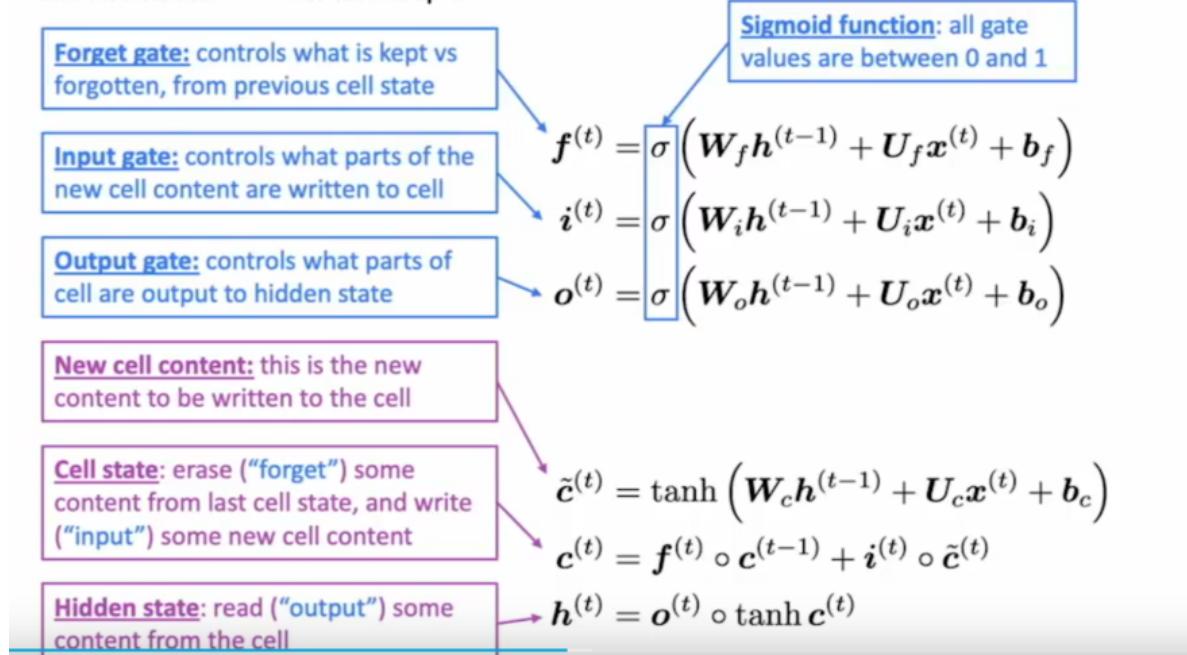
Long Short-Term Memory (LSTM)

Long Short-Term Memory (LSTM)

- A type of RNN proposed by Hochreiter and Schmidhuber in 1997 as a solution to the vanishing gradients problem.
- On step t , there is a hidden state $h^{(t)}$ and a cell state $c^{(t)}$
 - Both are vectors length n
 - The cell stores long-term information
 - The LSTM can erase, write and read information from the cell
- The selection of which information is erased/written/read is controlled by three corresponding gates
 - The gates are also vectors length n
 - On each timestep, each element of the gates can be open (1), closed (0), or somewhere in-between.
 - The gates are dynamic. Their output is based on the current context

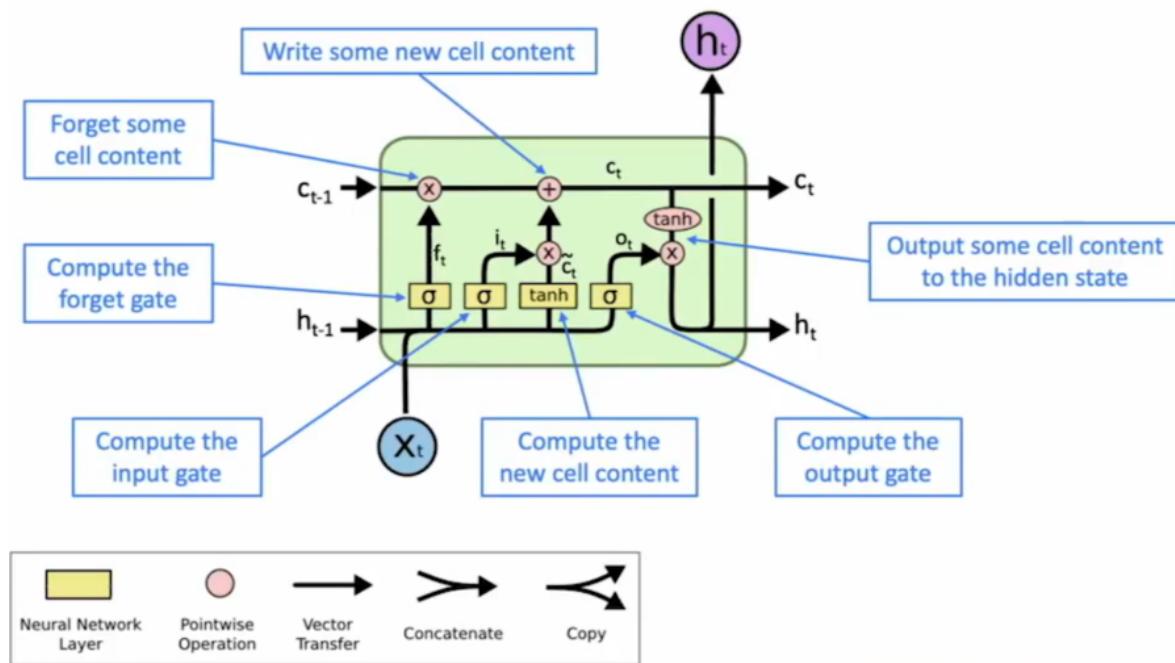
Long Short-Term Memory (LSTM)

We have a sequence of inputs $x^{(t)}$, and we will compute a sequence of hidden states $h^{(t)}$ and cell states $c^{(t)}$. On timestep t :



Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:



注:

- LSTM中通常把 h 视作输出，类似RNN的 y ，因为 c 是对外不可见的。
- 其实一切只和 c 有关， h 只是 c 经过一些操作得到的。
- 用tanh只是加上激活函数，类似cnn的线性层后面会加上一个relu。

LSTM如何解决vanishing gradients?

RNN的vanishing gradients问题会导致“长期的影响丢失”，即远处传来的梯度很小。

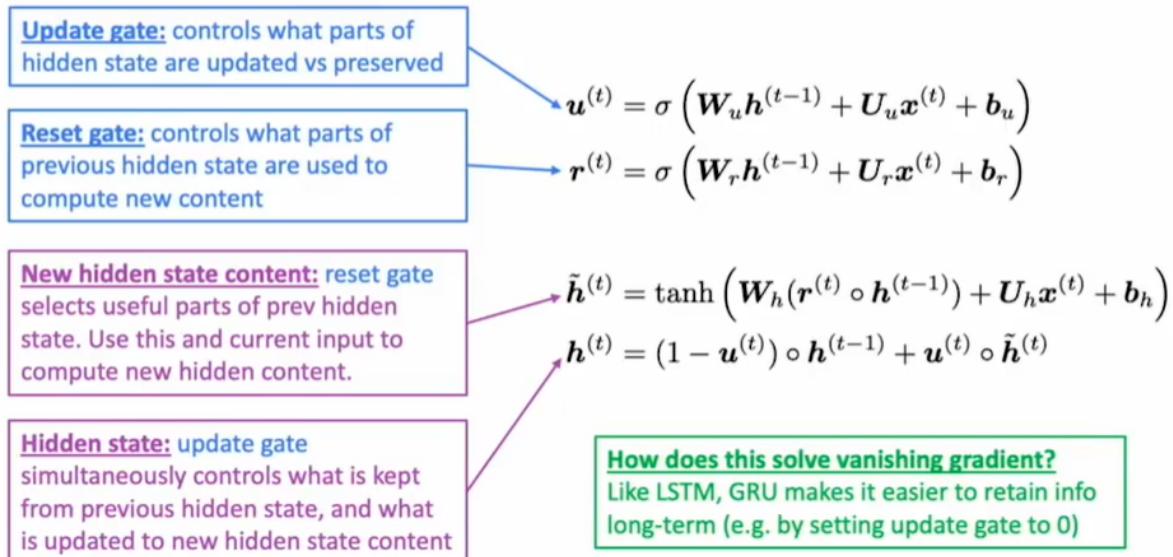
LSTM的解决方式是通过f(遗忘门)来控制保留以前的信息，如果f全1，就会全部保留。换句话说，LSTM并不是从数学上(反向传播)上来避免vanishing gradients，而是减轻了它的后果，更加robust。

Gated Recurrent Units (GRU)

GRU的提出是为了简化LSTM，同时保留LSTM的优点。

Gated Recurrent Units (GRU)

- Proposed by Cho et al. in 2014 as a simpler alternative to the LSTM.
- On each timestep t we have input $x^{(t)}$ and hidden state $h^{(t)}$ (no cell state).



LSTM vs GRU

- Researchers have proposed many gated RNN variants, but LSTM and GRU are the most widely-used
- The biggest difference is that GRU is quicker to compute and has fewer parameters
- There is no conclusive evidence that one consistently performs better than the other
- LSTM is a good default choice (especially if your data has particularly long dependencies, or you have lots of training data)
- Rule of thumb: start with LSTM, but switch to GRU if you want something more efficient

其他网络解决vanishing gradients

ResNet

aka skip-connections。有效是因为identity connection可以至少保持恒等映射。

DenseNet

所有层有加上skip-connections。

HighwayNet

借鉴了LSTM的gate

总结：梯度消失不只是RNN的问题，但在RNN里尤其严重，因为重复使用相同的矩阵。

Other Variants

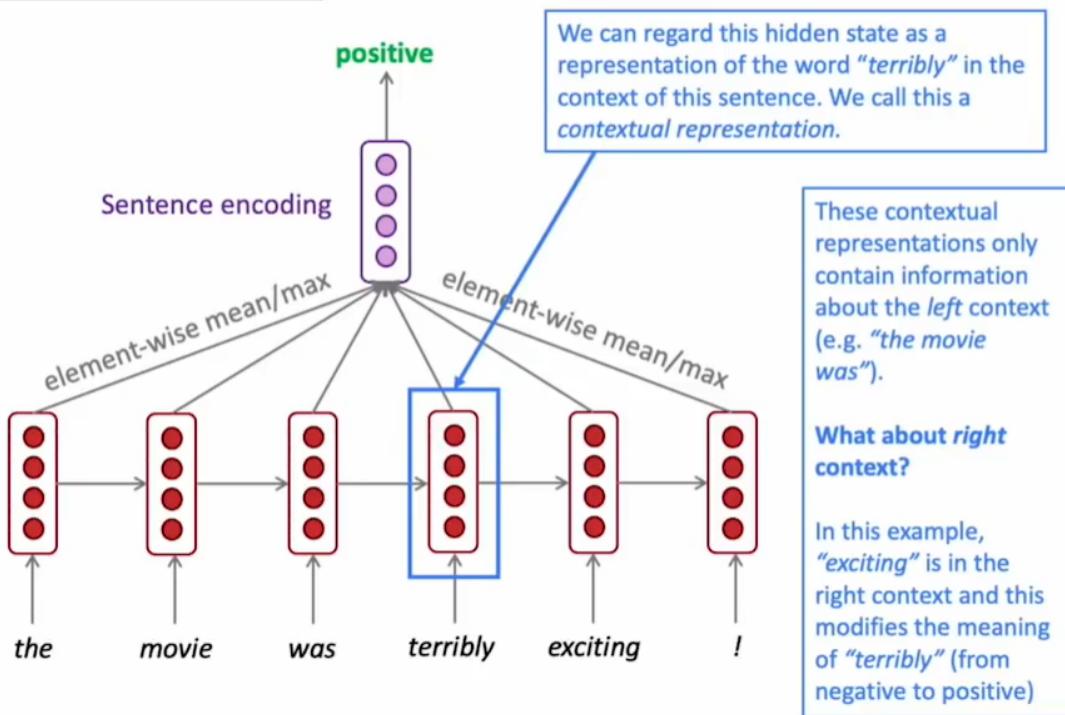
Bidirectional RNNs

在sentiment classification任务中，每个t的h都作为当前t的特征。用RNN只能包含left context。

在下面的例子中和，希望能包含right context，因为right context正确地解释了terribly。

Bidirectional RNNs: motivation

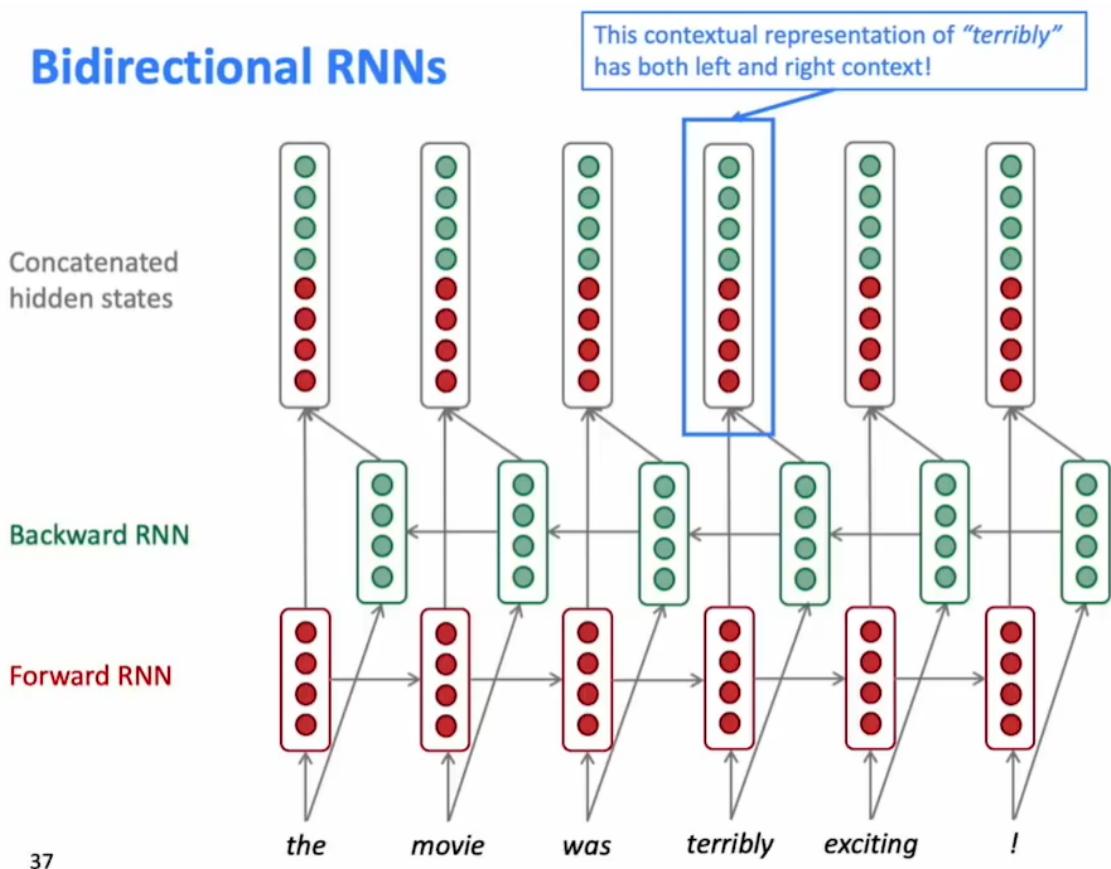
Task: Sentiment Classification



36

双向RNN的结构如下：

Bidirectional RNNs



37

Bidirectional RNNs

On timestep t :

This is a general notation to mean “compute one forward step of the RNN” – it could be a vanilla, LSTM or GRU computation.

$$\begin{aligned} \text{Forward RNN } \vec{h}^{(t)} &= \text{RNN}_{\text{FW}}(\vec{h}^{(t-1)}, \mathbf{x}^{(t)}) \\ \text{Backward RNN } \overleftarrow{h}^{(t)} &= \text{RNN}_{\text{BW}}(\overleftarrow{h}^{(t+1)}, \mathbf{x}^{(t)}) \end{aligned} \quad \left. \begin{array}{l} \text{Generally, these} \\ \text{two RNNs have} \\ \text{separate weights} \end{array} \right\}$$

Concatenated hidden states $\mathbf{h}^{(t)} = [\vec{h}^{(t)}; \overleftarrow{h}^{(t)}]$

We regard this as “the hidden state” of a bidirectional RNN.
This is what we pass on to the next parts of the network.

38

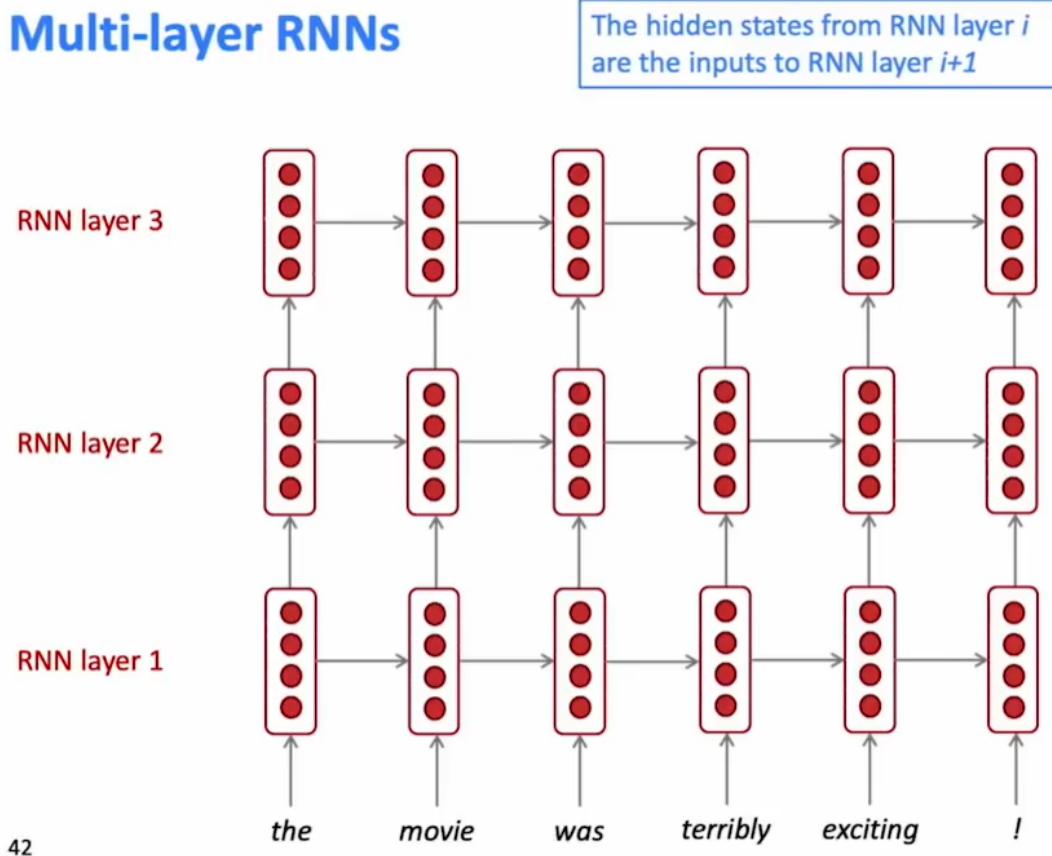
注:

- 双向RNN不能用于LM，只能用于已知完整的sentence。
- 在有完整sentence的情况下，bidirection很有用！

Multi-layer RNNs

在RNN的宽度上扩展，又称stacked RNN。

多层RNN的结构如下：



注：

- Pytorch中会先算第一层，再算第二层，以此类推。但是逻辑上可以先算第一个单词，再算第二个。
- High-performing RNN通常都是多层的，但一般2-4层即可，加上skip-connections可能8层(BERT有24层)。
- 没有几百层的原因是RNN计算是串行的，代价过高。
- 当RNN层数很多时，通常需要加上skip-connections。