



Recurrent Neural Networks

Phạm Quang Nhật Minh

Aimesoft JSC

minhpham0902@gmail.com



Lecture outline

2

- Recurrent neural networks
- Multi-layer RNNs (Stacked RNNs)
- Bidirectional RNNs
- Two types of RNNs: LSTM and GRU



NLP and Sequential Data

3

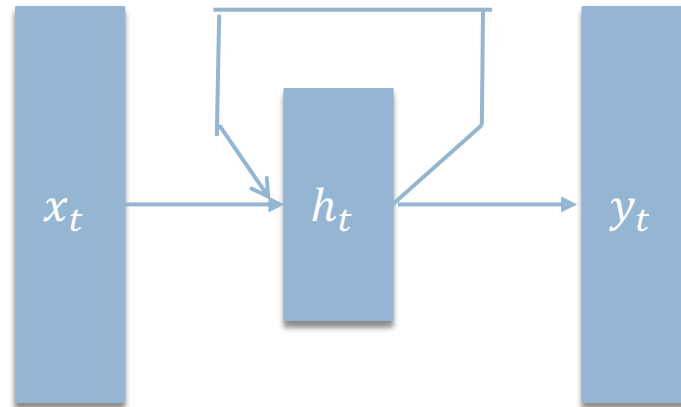
- NLP is full of sequential data
 - Words in sentences
 - Characters in words
 - Sentences in discourse
 - ...
- Sequences have variable lengths
- Long-distance dependencies in languages
 - **He** does not have very much confidence in **himself**.
 - **She** does not have very much confidence in **herself**.



Basic ideas of Recurrent neural networks (RNN)

4

- Use recurrent link
 - Output of $t-1$ step is used as input for the t step
- Used same weights among steps

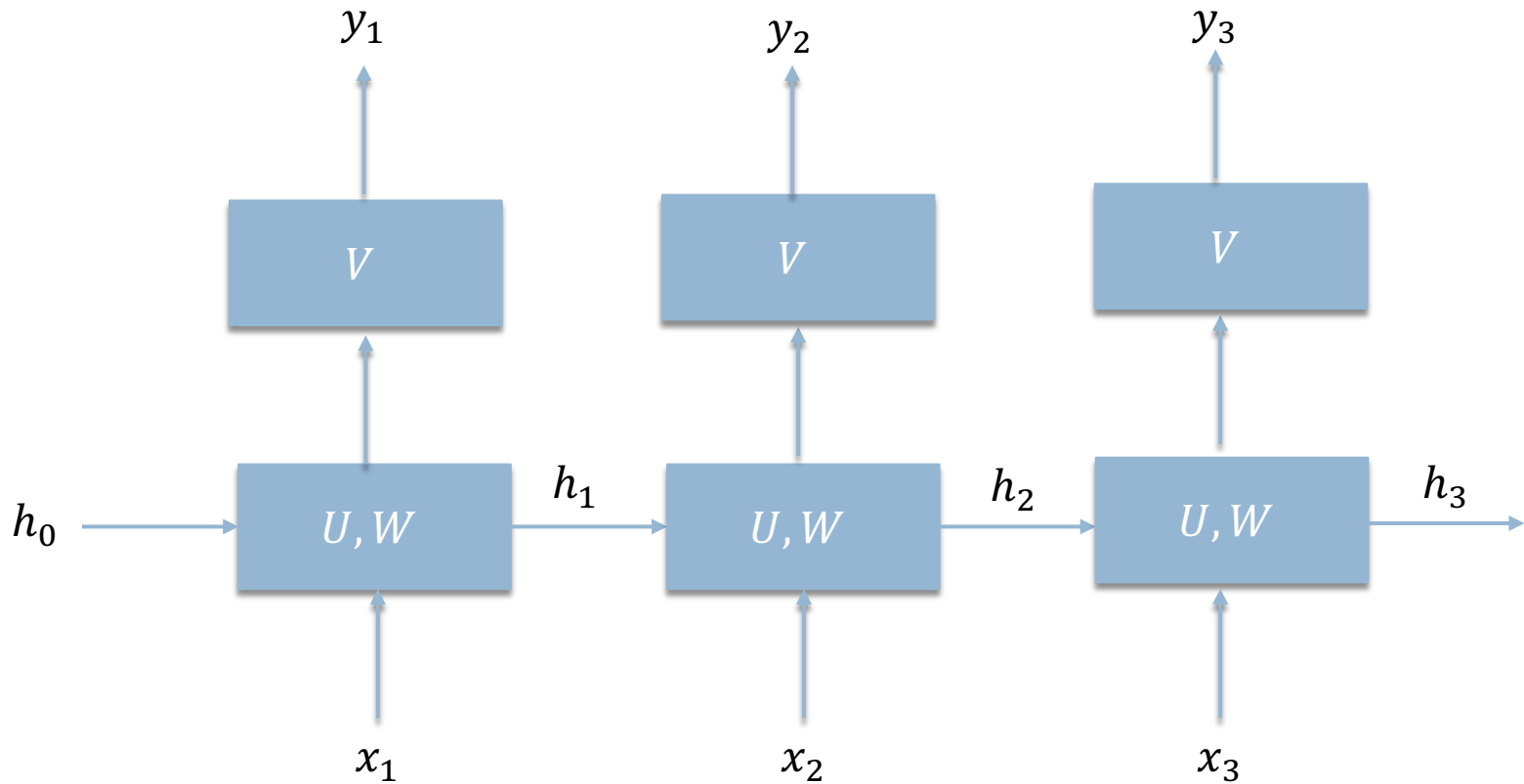




Unrolled representation of RNNs

5

x_i are vectors



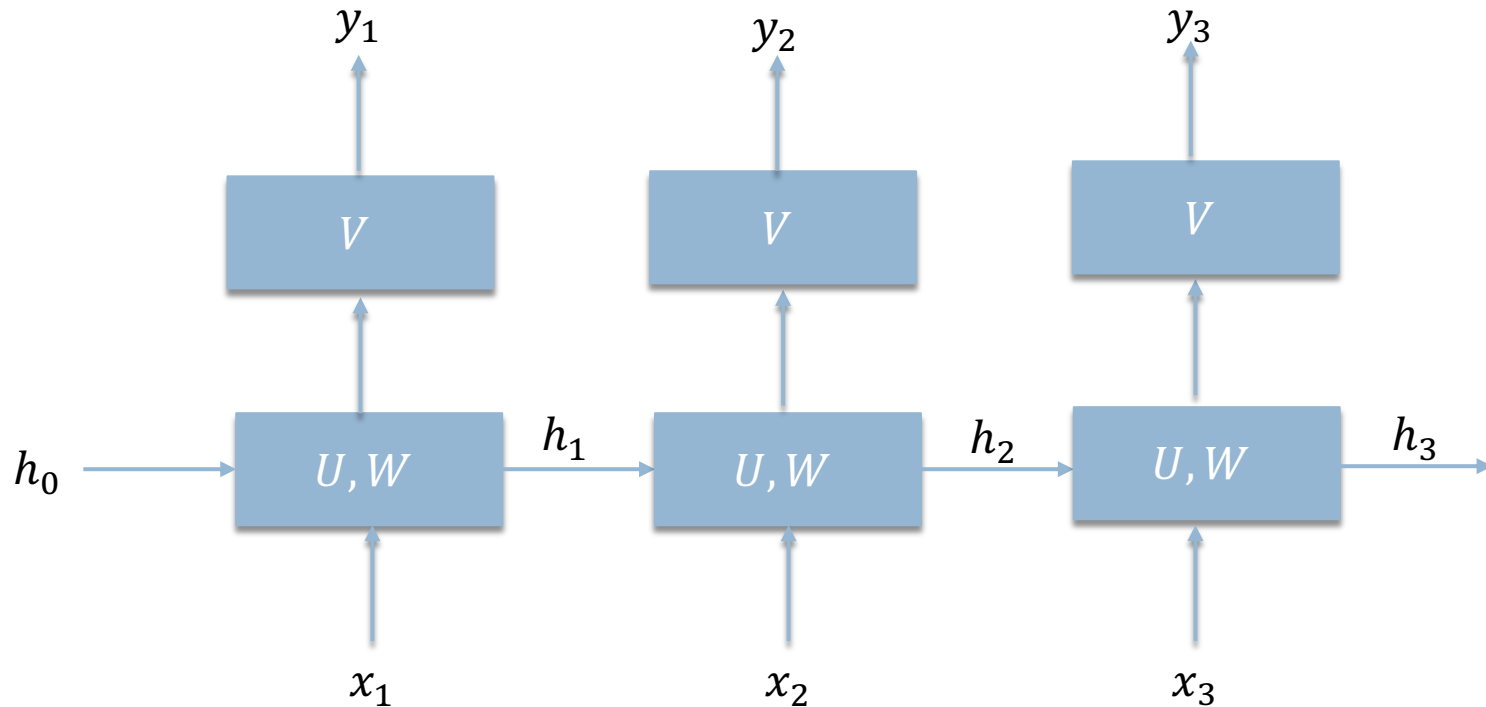


Equations in RNN

6

$$h_t = g(Uh_{t-1} + Wx_t)$$
$$y_t = f(Vh_t)$$
$$W \in \mathbb{R}^{d_h \times d_{in}}, U \in \mathbb{R}^{d_h \times d_h}, V \in \mathbb{R}^{d_{out} \times d_h}$$

d_{in}, d_h, d_{out} are input, hidden, and output layer dimensions

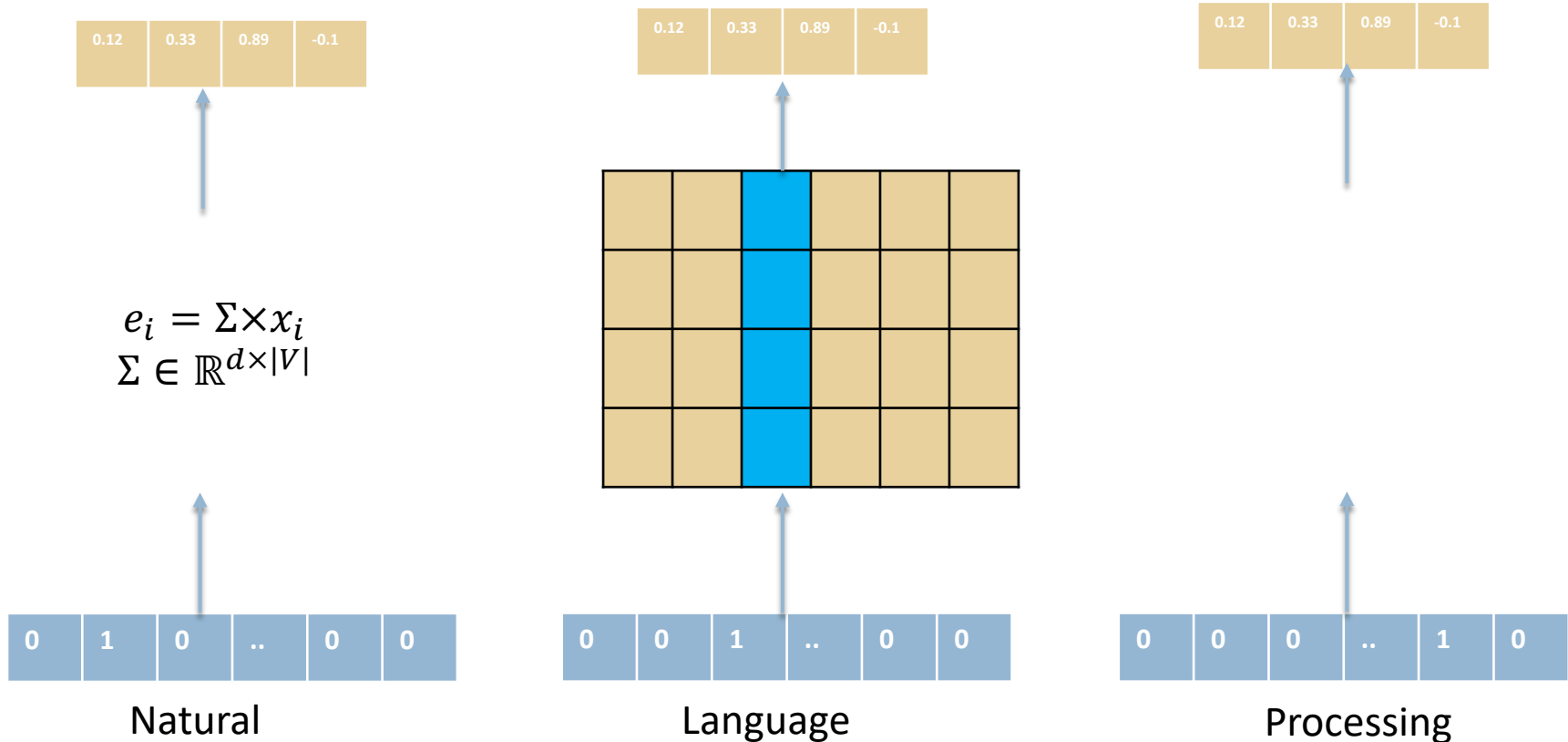




Input Layers of RNN in NLP tasks

7

- Each word is a vector
 - Dense vector obtained by an embedding layer (lookup matrix)





Forward inference in RNN

8

```
function FORWARDRNN( $\mathbf{x}$ , network) returns output sequence  $\mathbf{y}$ 
```

```
   $\mathbf{h}_0 \leftarrow 0$ 
```

```
  for  $i \leftarrow 1$  to LENGTH( $\mathbf{x}$ ) do
```

```
     $\mathbf{h}_i \leftarrow g(\mathbf{U}\mathbf{h}_{i-1} + \mathbf{W}\mathbf{x}_i)$ 
```

```
     $\mathbf{y}_i \leftarrow f(\mathbf{V}\mathbf{h}_i)$ 
```

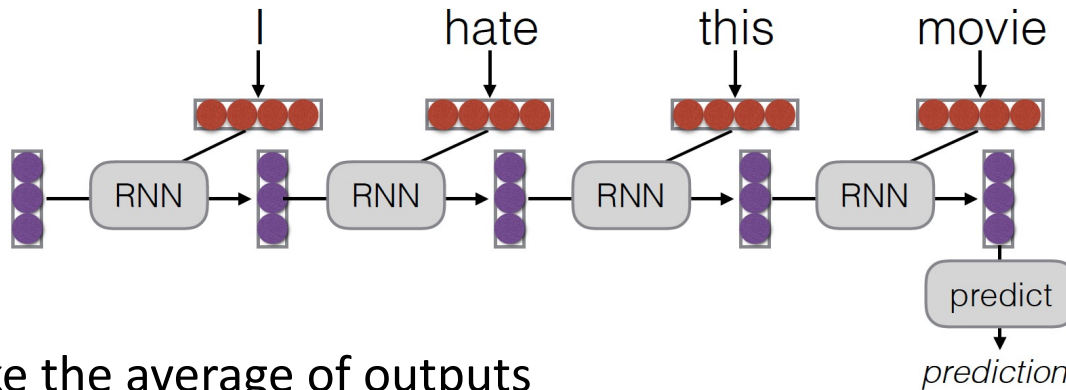
```
  return  $\mathbf{y}$ 
```

The matrices \mathbf{U} , \mathbf{V} and \mathbf{W} are shared across time, while new values for \mathbf{h} and \mathbf{y} are calculated with each time step



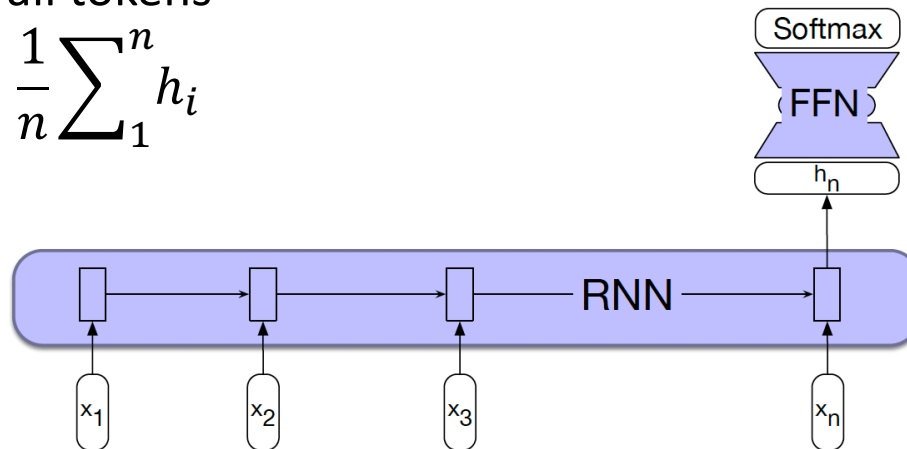
RNN for Sequence Classification

- Use the output of the hidden layer at of the last token as the representation of the entire sequence



We can take the average of outputs of hidden layer for all tokens

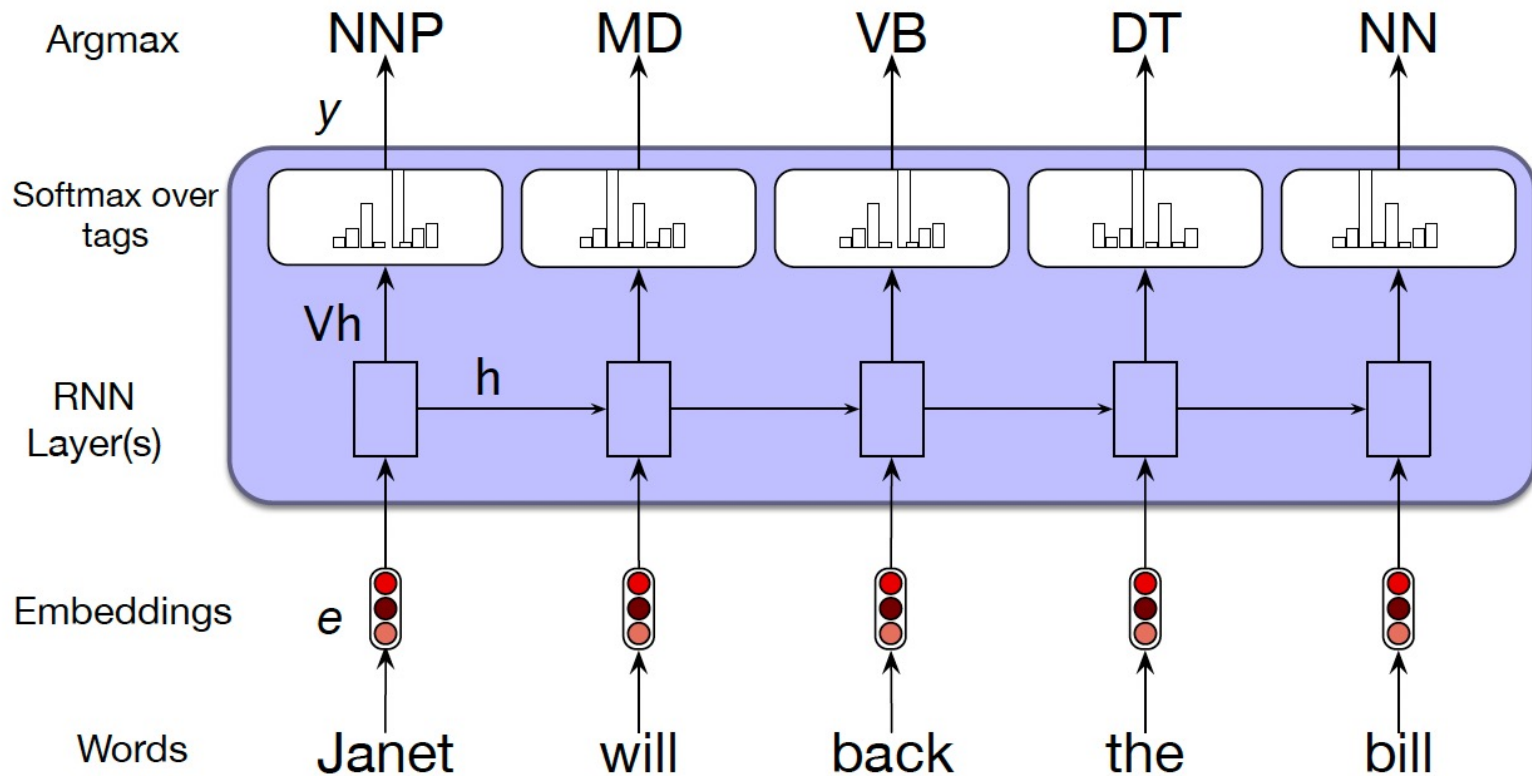
$$h_{mean} = \frac{1}{n} \sum_{i=1}^n h_i$$





RNN for Sequence Labeling

10



Part-of-speech tagging as sequence labeling with a simple RNN. Pre-trained word embeddings serve as inputs and a softmax layer provides a probability distribution over the part-of-speech tags as output at each time step.



Training RNNs

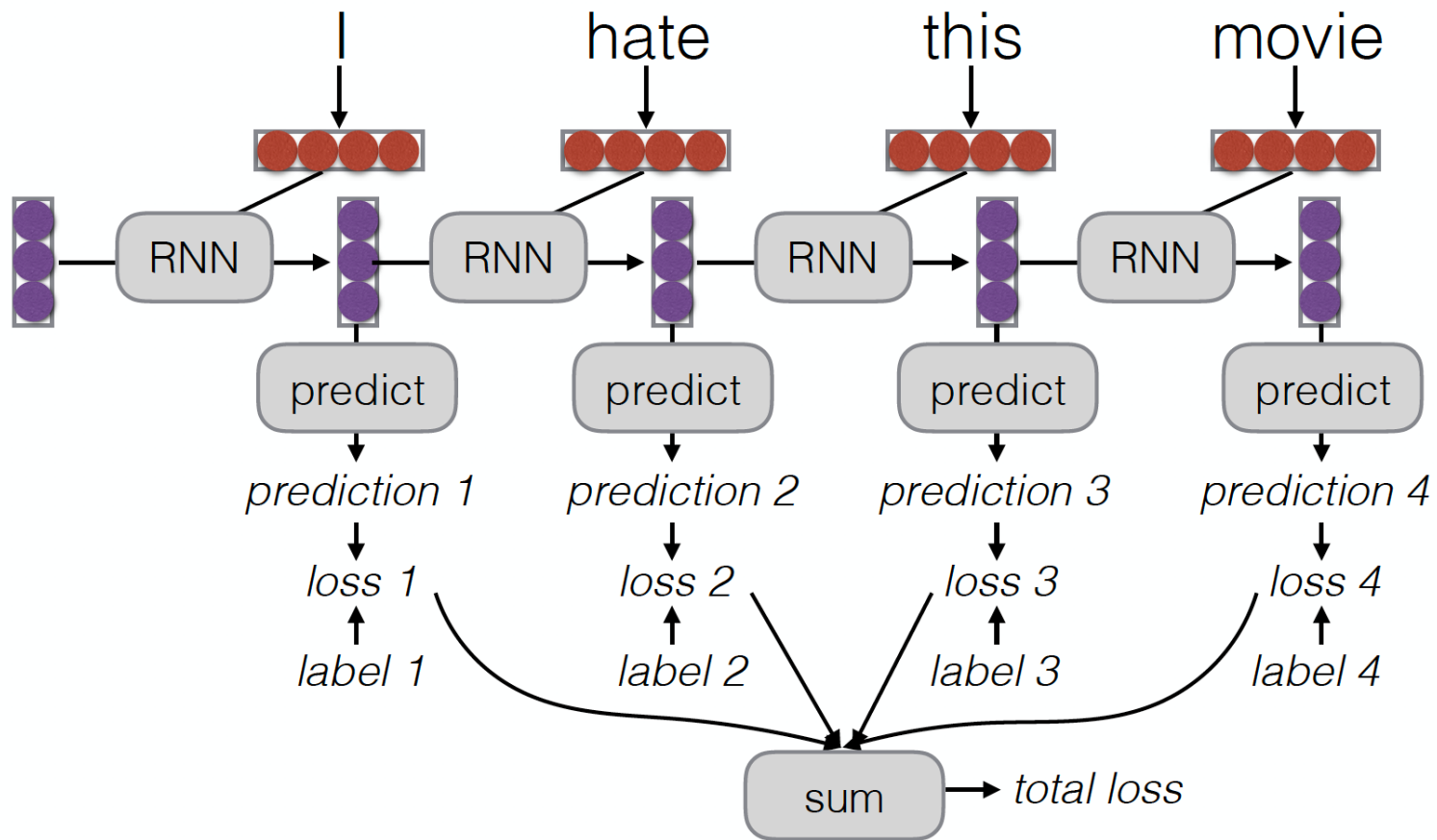
11

- Training RNN steps:
 - Unroll recurrent neural networks
 - Apply backpropagation to calculate gradients
- Algorithm of training RNN is called backpropagation through time (BPTT) (Werbos, 1990)



Training RNNs in tagging problems

12





Lecture outline

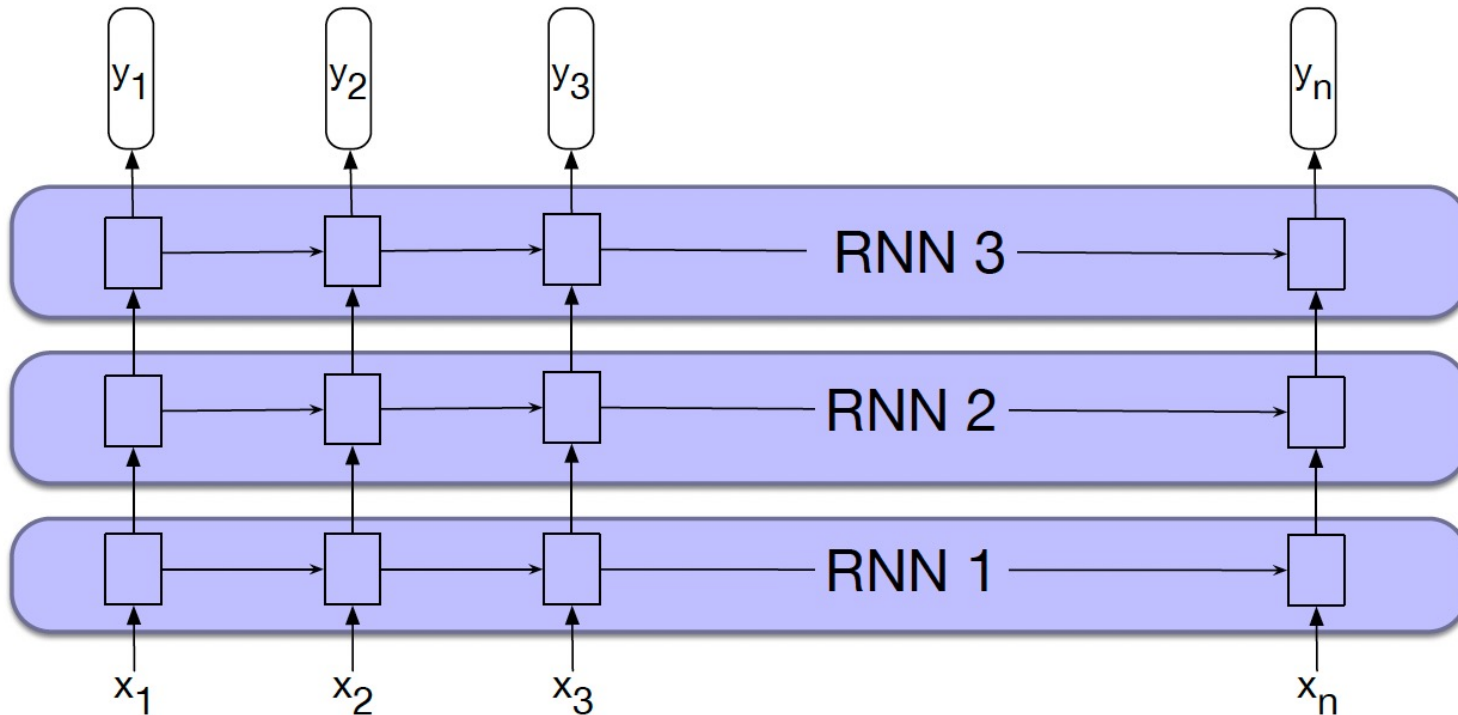
13

- Recurrent neural networks
- **Multi-layer RNNs (Stacked RNNs)**
- Bidirectional RNNs
- Two types of RNNs: LSTM and GRU



Stacked RNNs

14



The output of a lower level serves as the input to higher levels with the output of the last network serving as the final output.

Stacked RNNs generally outperform single-layer networks but the training cost is higher than single-layer RNN networks



Lecture outline

15

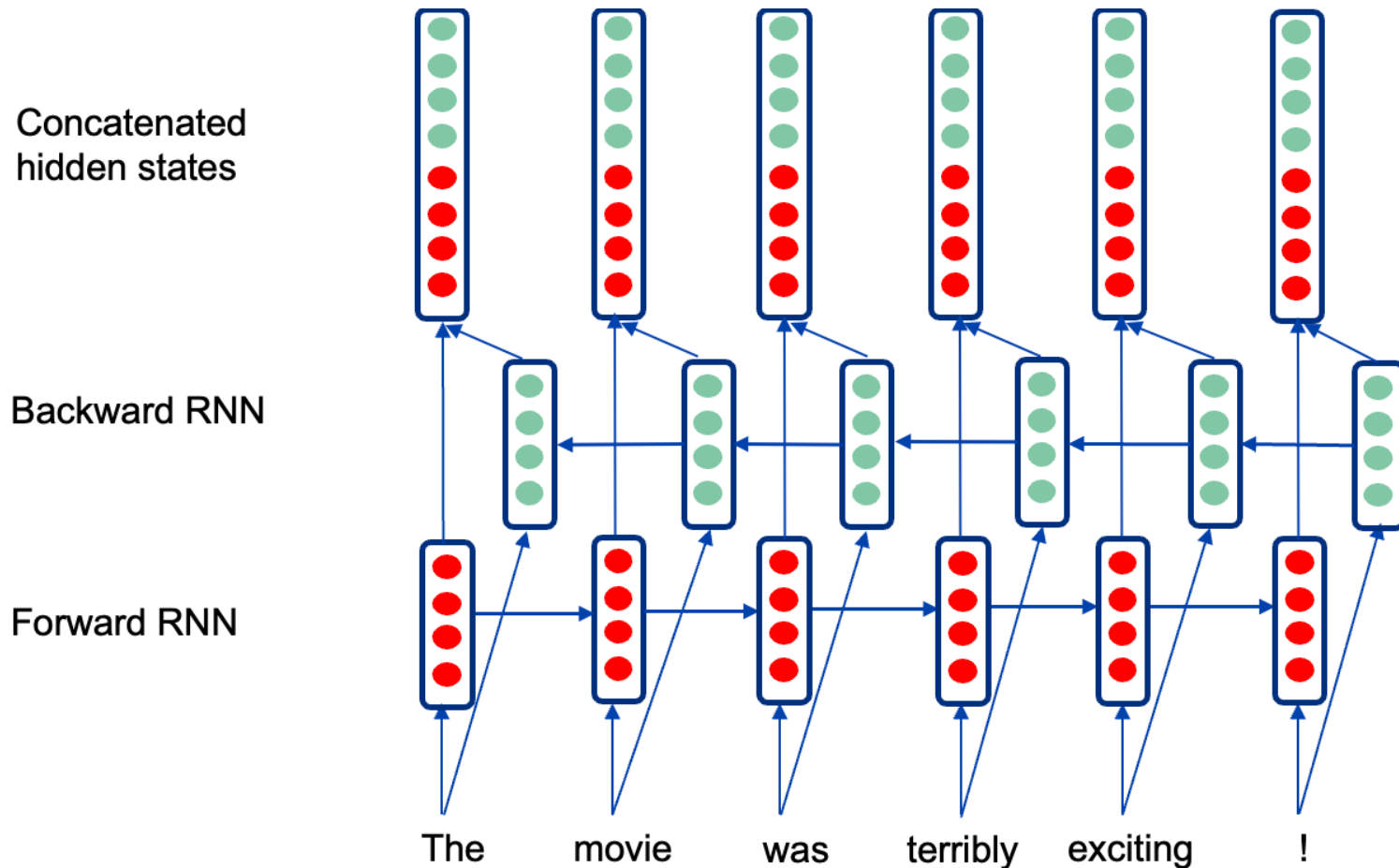
- Recurrent neural networks
- Multi-layer RNNs (Stacked RNNs)
- **Bidirectional RNNs**
- Two types of RNNs: LSTM and GRU



Bidirectional RNNs

16

Context in both directions (left and right) is useful in some task (such as POS Tagging)





Bidirectional RNN: Formulas

17

On time step t :

Forward RNN $\vec{h}^{(t)} = \text{RNN}_{\text{FW}}(x_1, \dots, x_t)$

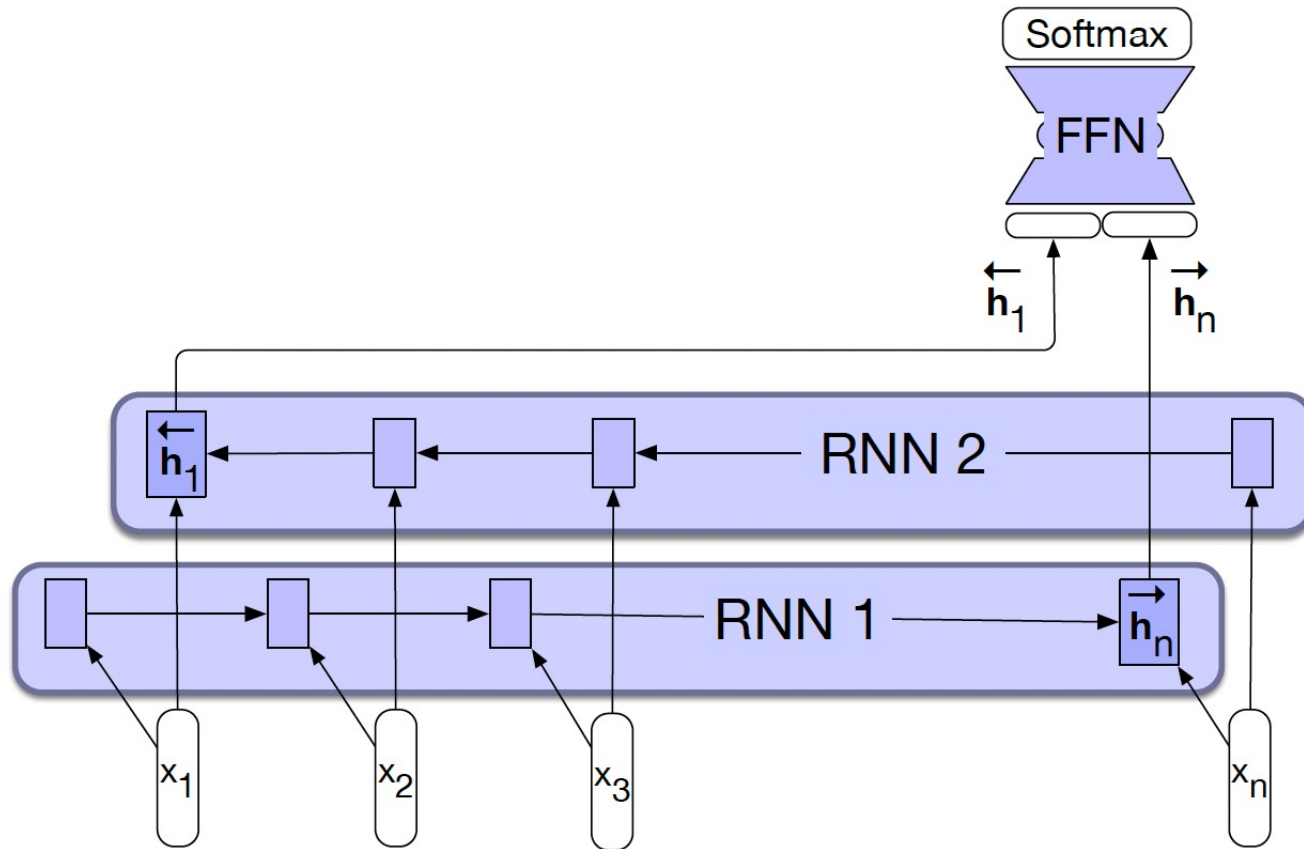
Backward RNN $\overleftarrow{h}^{(t)} = \text{RNN}_{\text{BW}}(x_t, \dots, x_n)$

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



Bidirectional RNN for Sequence Classification

18



The final hidden units from the forward and backward passes are combined to represent the entire sequence. This combined representation serves as input to the subsequent classifier.



Lecture outline

19

- Recurrent neural networks
- Multi-layer RNNs (Stacked RNNs)
- Bidirectional RNNs
- Two types of RNNs: LSTM and GRU



Why Long Short-term Memory (LSTM)?

20

- Standard RNNs could not handle long-term dependencies well
 - “I grew up in France... I speak fluent *French*.”



Why Long Short-term Memory (LSTM)?

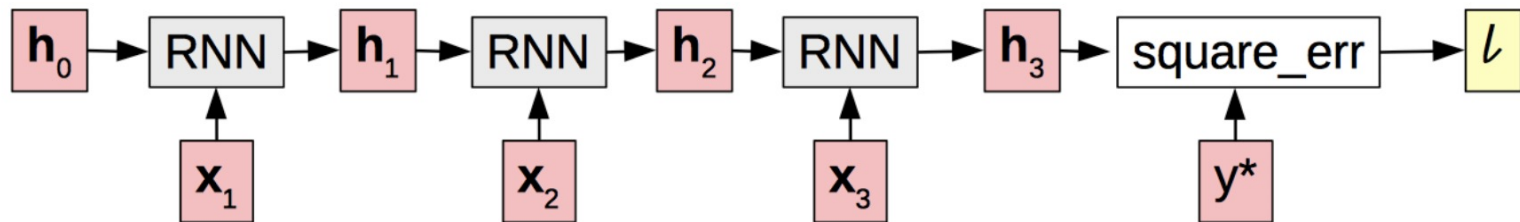
21

■ Vanishing Gradient Descent problem

- Gradients decrease as they get pushed back

$$\frac{dl}{dh_0} = \frac{dh_1}{dh_0} \times \frac{dh_2}{dh_1} \times \frac{dh_3}{dh_2} \times \frac{dl}{dh_3}$$

- $\frac{dl}{dh_0} = \text{tiny}$, $\frac{dl}{dh_1} = \text{small}$, $\frac{dl}{dh_3} = \text{med}$, $\frac{dl}{dh_4} = \text{large}$





Basic ideas of LSTM networks

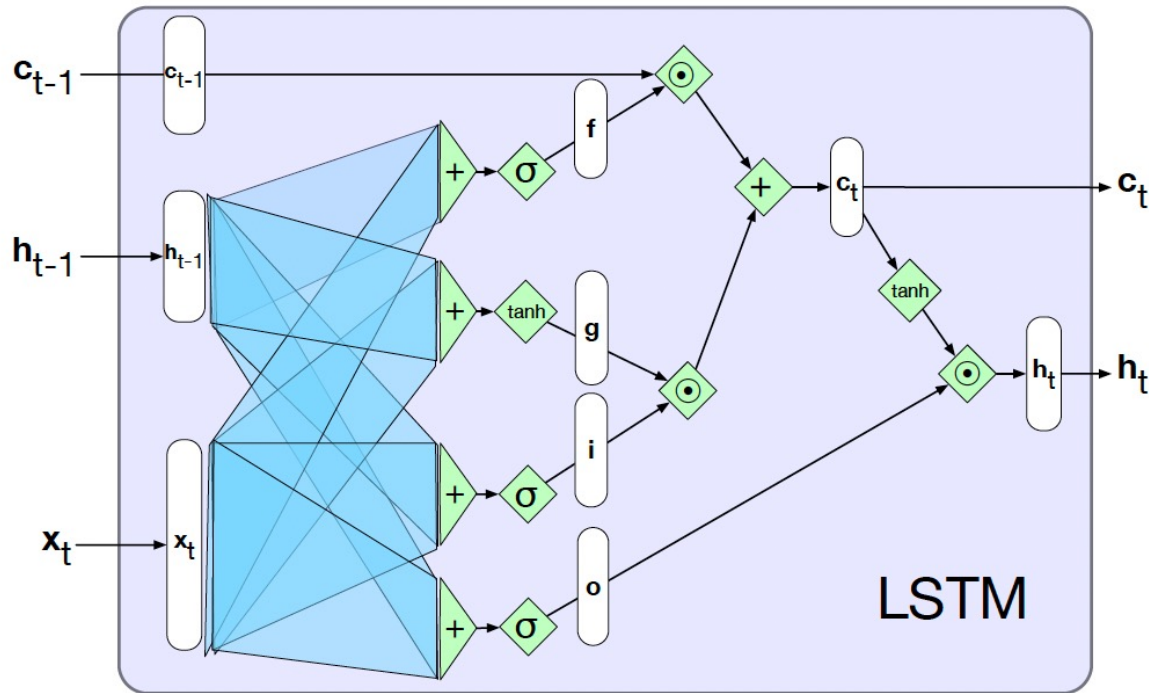
22

- On step t , there is a hidden state *and* a cell state
 - 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 gates



LSTM Network architecture

23



A single LSTM unit displayed as a computation graph. The inputs to each unit consists of the current input, x , the previous hidden state h_{t-1} , , and the previous context c_{t-1} , . The outputs are a new hidden state, h_t and an updated context, c_t .



LSTM Network Equations

24

- **Forget gate:** to delete information from the context that is no needed

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

$$\mathbf{k}_t = \mathbf{c}_{t-1} \odot \mathbf{f}_t$$

- **Input gate:** to select information to add to the current context

$$\mathbf{g}_t = \tanh(\mathbf{U}_g \mathbf{h}_{t-1} + \mathbf{W}_g \mathbf{x}_t)$$

$$\mathbf{i}_t = \sigma(\mathbf{U}_i \mathbf{h}_{t-1} + \mathbf{W}_i \mathbf{x}_t)$$

$$\mathbf{j}_t = \mathbf{g}_t \odot \mathbf{i}_t$$

$$\mathbf{c}_t = \mathbf{j}_t + \mathbf{k}_t$$

- **Output gate:** to decide what information is required for the current hidden state

$$\mathbf{o}_t = \sigma(\mathbf{U}_o \mathbf{h}_{t-1} + \mathbf{W}_o \mathbf{x}_t)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$$



How does LSTM solve vanishing gradients?

25

- The LSTM architecture makes it easier for the RNN to preserve information over many timesteps
- LSTM doesn't *guarantee* that there is no vanishing/exploding gradient, but it does provide an easier way for the model to learn long-distance dependencies