

# CS 378 Lecture 16

Today

- RNNs
- LSTMs (the type of RNN you will be using)
- Implementation

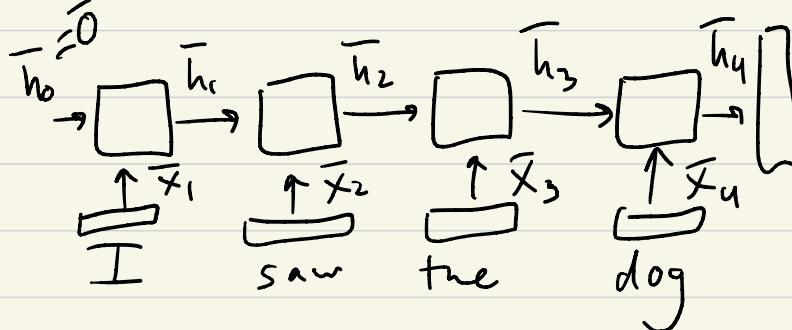
Announcements

- AY
- Midterm back soon

LM:  $P(\bar{w})$  or

$P(w_i | w_1, \dots, w_{i-1})$   
"predict the next word"

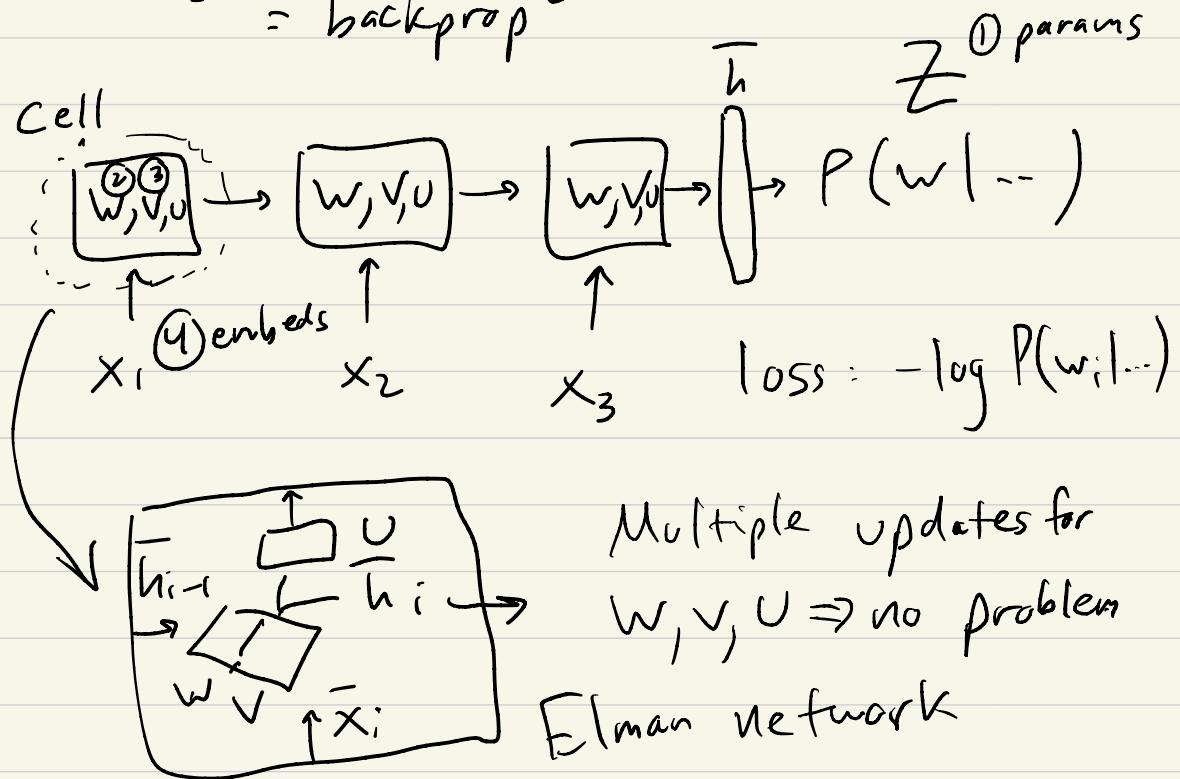
Recap RNNs + language modeling



$Z : |V| \times d$   
parameter matrix

$$P(w | I \text{ saw the dog}) = \text{softmax}(Z^T \bar{h})$$

Training "Backpropagation through time"  
 = backprop



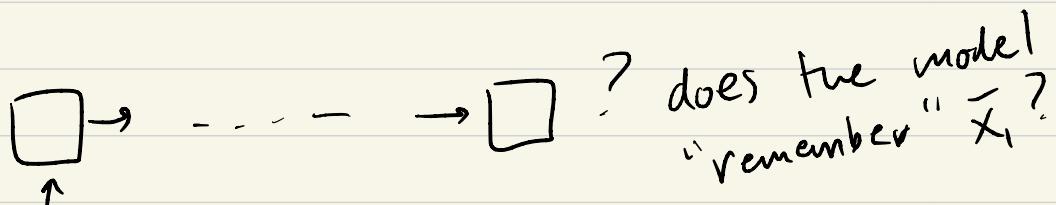
## Long Short-term memory networks (LSTMs)

Many types of RNNs



# LSTMs (1998)

short-term memory: what the model can remember in its state



LONG short-term memory  
(remember for longer)

Problem w/ Elman networks

vanishing/exploding gradients

$$\begin{array}{c} \square \rightarrow \square \rightarrow \square \rightarrow h_3 \\ \uparrow \quad \uparrow \quad \uparrow \\ \bar{x}_1 \quad \bar{x}_2 \quad \bar{x}_3 \end{array} \quad \begin{aligned} \bar{h}_i &= \tanh(W\bar{x}_i + V\bar{h}_{i-1}) \\ \bar{h}_3 &= \tanh(W\bar{x}_3 + V \cdot \\ &\quad \tanh(W\bar{x}_2 + V \cdot \\ &\quad \tanh(W\bar{x}_1 + \bar{0}))) \end{aligned}$$

Assume  $\tanh$  is the identity function

$$\bar{h}_3 = W\bar{x}_3 + VVW\bar{x}_2 + V^2W\bar{x}_1$$

after  $n$  steps  $\Rightarrow V^{n-1}\bar{x}_1$

### LSTM gates

~~Elman:~~  $\bar{h}_i = \tanh(W\bar{x}_i + V\bar{h}_{i-1})$

input  
gate

Gated:  $\bar{h}_i = \bar{h}_{i-1} \odot \bar{f} + \text{function}(\bar{x}_i, \bar{h}_{i-1}) \odot i$

prev state      elementwise  $\times$

$\bar{f}$ : forget gate, values in  $[0, 1]$

$$\bar{h}_{i-1} \quad \boxed{\phantom{0}} \quad \odot \quad \boxed{\begin{array}{c} \text{X} \\ \text{✓} \\ \text{✓} \\ \text{✓} \end{array}} = \boxed{\begin{array}{c} \text{X} \\ \text{✓} \\ \text{✓} \\ \text{✓} \end{array}}$$

If  $\bar{f} = 1$ :  
 $\bar{h}_{i-1}$  is  
totally  
preserved

(added for  
polf exercise)  
bias

Where do  $f$ ,  $i$  come from?

$$f = \text{sigmoid} \left( w^{(1)} \bar{x}_i + w^{(2)} \bar{h}_{i-1} + b_{\text{forget}} \right)$$

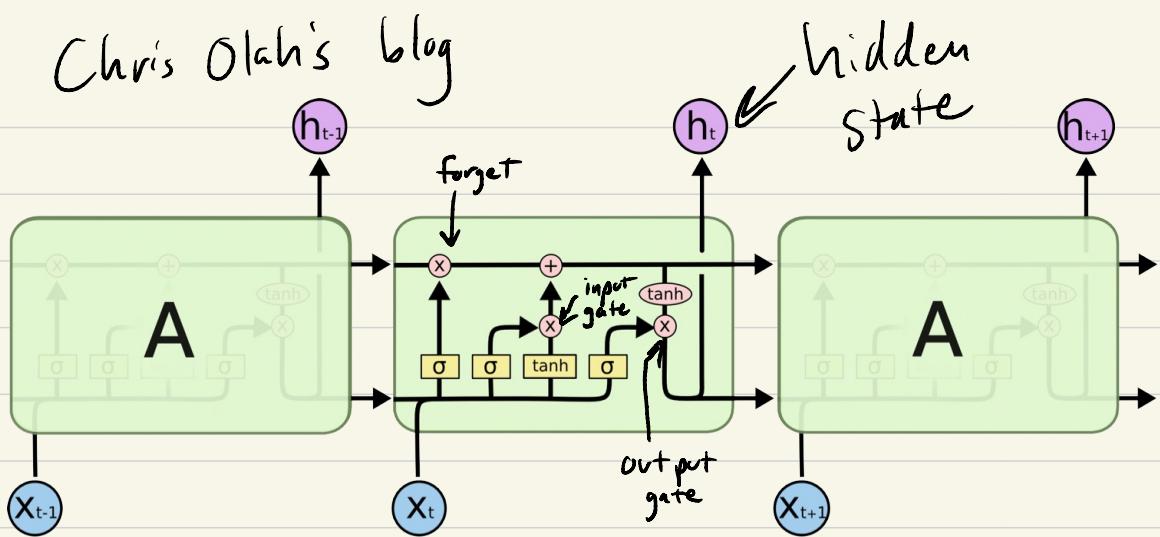
$$i = \text{sigmoid} \left( w^{(3)} \bar{x}_i + w^{(4)} \bar{h}_{i-1} + b_{\text{input}} \right)$$

sigmoid

$$\sigma = \frac{e^x}{1+e^x}$$

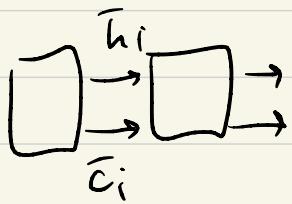
$$\bar{h}_{i-1} \cdot \bar{f}^{\text{forget}} \rightarrow \bar{h}_i + \bar{x}_i \cdot \bar{i} \rightarrow \bar{h}_i^{\text{update}}$$

Chris Olah's blog



LSTM: 8 weight matrices

hidden state  $\overline{h}$   
cell state  $\overline{c}$  ] tuple of the  
LSTM state



Part 1:

discussed lstm-lecture.py

