

# Backpropagation, Word embeddings, Recurrent Neural Networks

Data Mining in Action: Trends

# Plan

1. Backpropagation algorithm for training
2. Word embeddings as a basis for doing Deep NLP
3. RNNs

# Backpropagation

- Generally, it is just a chain rule on the computational graph
- For details we follow these links:
  - <https://goo.gl/LbEHy4>
  - <https://goo.gl/maQ8KS>
  - =)

# Word Embeddings

*based on awesome CS224d Stanford course*

[http://cs224d.stanford.edu/lecture\\_notes/notes1.pdf](http://cs224d.stanford.edu/lecture_notes/notes1.pdf)

<http://cs224d.stanford.edu/lectures/CS224d-Lecture2.pdf>

<http://cs224d.stanford.edu/lectures/CS224d-Lecture3.pdf>

# Word Embeddings

- **Question:** How can we represent word as a vector?

- **Some ideas:**

1. One-hot
2. Construct word-word cooccurrence matrix
3. Variations of (2): word-word  $\rightarrow$  word-doc, word-window

- **Problems**

- The (1) method does not preserve similarity properties between words.

motel [0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0] AND  
hotel [0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0] = 0

- For (2) and (3) we need to compute & store VERY large matrices (usually you have 10-100k of words) and the VECTOR DIMENSIONALITY is increases with increasing the vocabulary.

# Word Embeddings

- Idea: store “most” of the important information in a fixed, small number of dimensions: a dense vector.
- How to do this?
- **Idea:**
  - Perform SVD of word-word cooccurrence matrix.
- **Problem:**
  - Quadratic cost to perform SVD + all problems with word-word matrix.

# Word Embeddings

- Word2Vec (CBOW, Continuous Bag of Words)
- Predict surrounding words in a window of length  $m$  of every word.
- Objective function: Maximize the log probability of any context word given the current center word:

- $$J(\theta) = \frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log p(w_{t+j} | w_t)$$

- Where  $\theta$  represents all variables we optimize

# Word Embeddings

- Predict surrounding words in a window of length  $m$  of every word
- For  $p(w_{t+j}|w_t)$  the simplest first formulation is

$$p(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)}$$

- where  $o$  is the outside (or output) word id,  $c$  is the center word id,  $u$  and  $v$  are “center” and “outside” vectors of  $o$  and  $c$
- Every word has two vectors!



# Word Embeddings

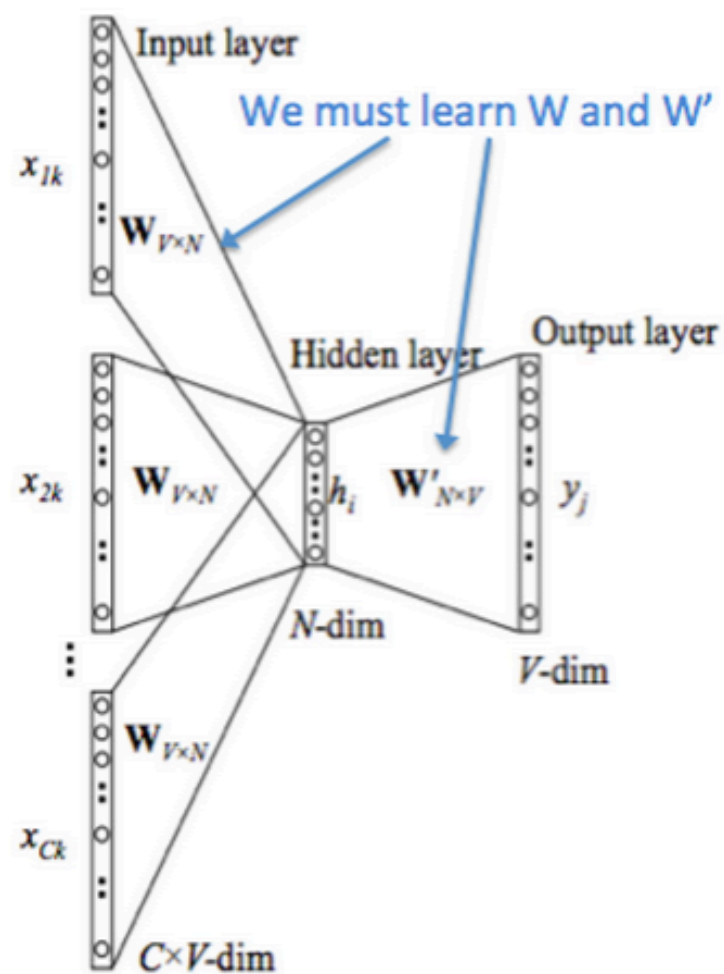


Figure 1: This image demonstrates how CBOW works and how we must learn the transfer matrices

# Word Embeddings

These representations are *very good* at encoding dimensions of similarity!

- Analogies testing dimensions of similarity can be solved quite well just by doing vector subtraction in the embedding space

Syntactically

- $x_{apple} - x_{apples} \approx x_{car} - x_{cars} \approx x_{family} - x_{families}$

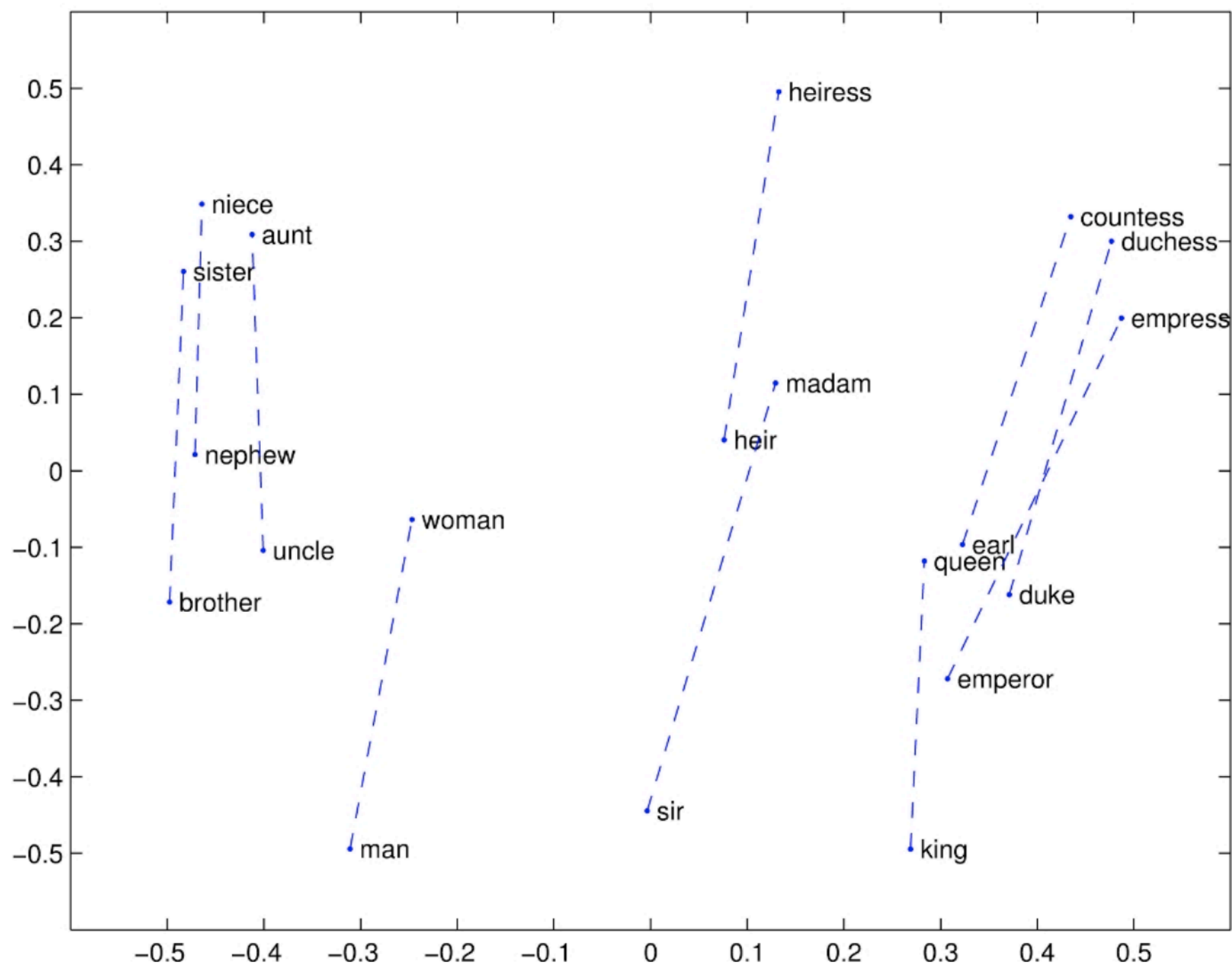
- Similarly for verb and adjective morphological forms

Semantically (Semeval 2012 task 2)

- $x_{shirt} - x_{clothing} \approx x_{chair} - x_{furniture}$

- $x_{king} - x_{man} \approx x_{queen} - x_{woman}$

# Word Embeddings



# Word Embeddings

- Other methods with the same idea:
  - word2vec (skip-gram), doc2vec, GloVe, AdaGram, Ida2vec...
- Implementations:
  - gensim (has implementations of w2v, d2v, glove)  
<https://radimrehurek.com/gensim/>
  - AdaGram (<https://github.com/sbos/AdaGram.jl>)

# Recurrent Neural Networks

*links:*

<https://karpathy.github.io/2015/05/21/rnn-effectiveness/>  
[www.wildml.com/2015/10/recurrent-neural-networks-tutorial-part-3-backpropagation-through-time-and-vanishing-gradients/](http://www.wildml.com/2015/10/recurrent-neural-networks-tutorial-part-3-backpropagation-through-time-and-vanishing-gradients/)  
[cs224d.stanford.edu/lecture\\_notes/notes4.pdf](http://cs224d.stanford.edu/lecture_notes/notes4.pdf)

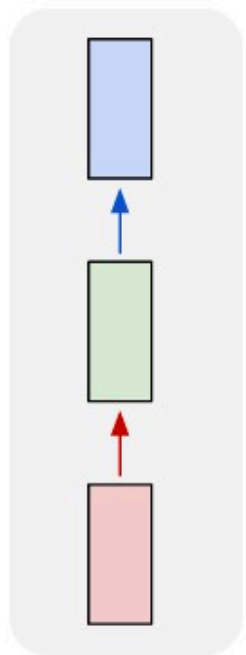
# Plan

- What can you do with them
- Vanilla RNN
- Vanishing/exploding gradient problem
- LSTM, GRU
- Bidirectional RNN
- Some applications of RNNs

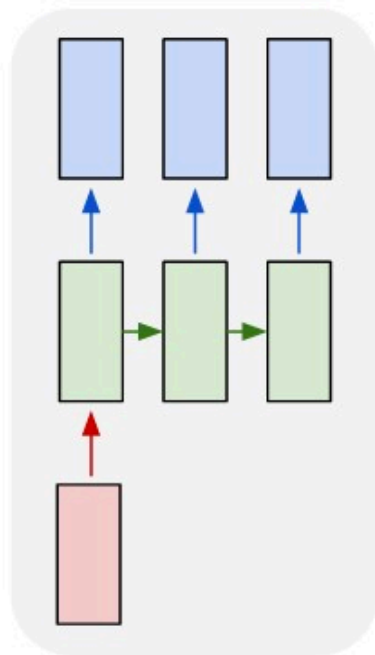
# What is RNN and its purpose?

- A neural network architecture for working with sequence data, e.g. texts, stock prices, ...

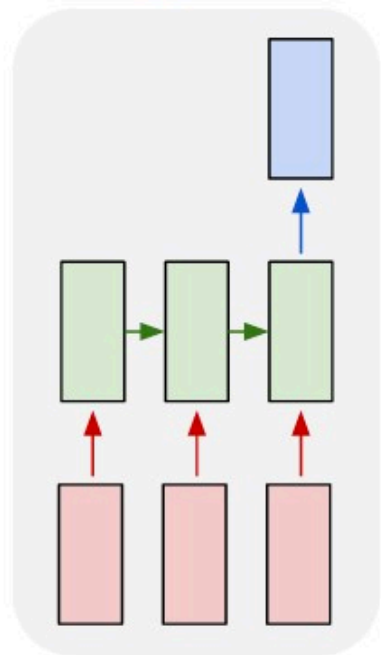
one to one



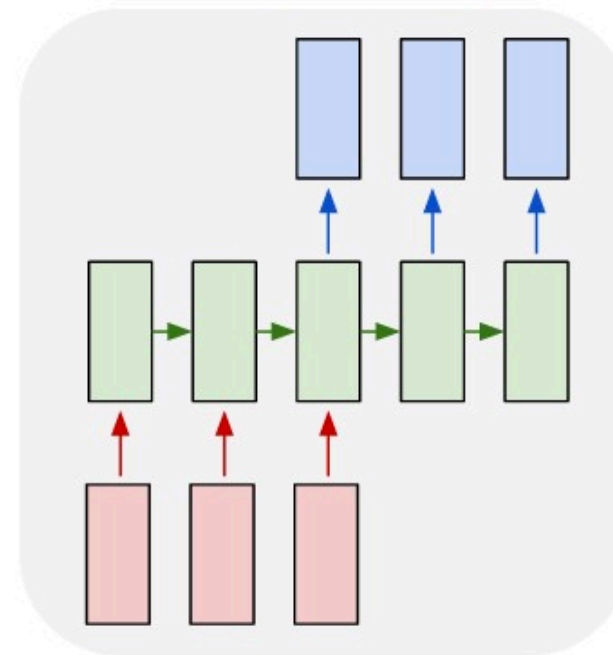
one to many



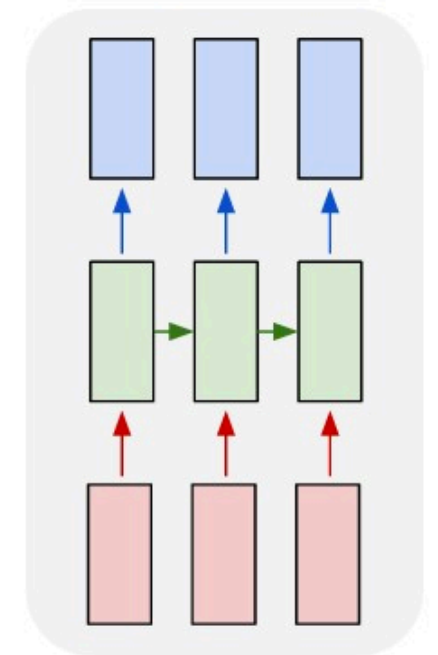
many to one



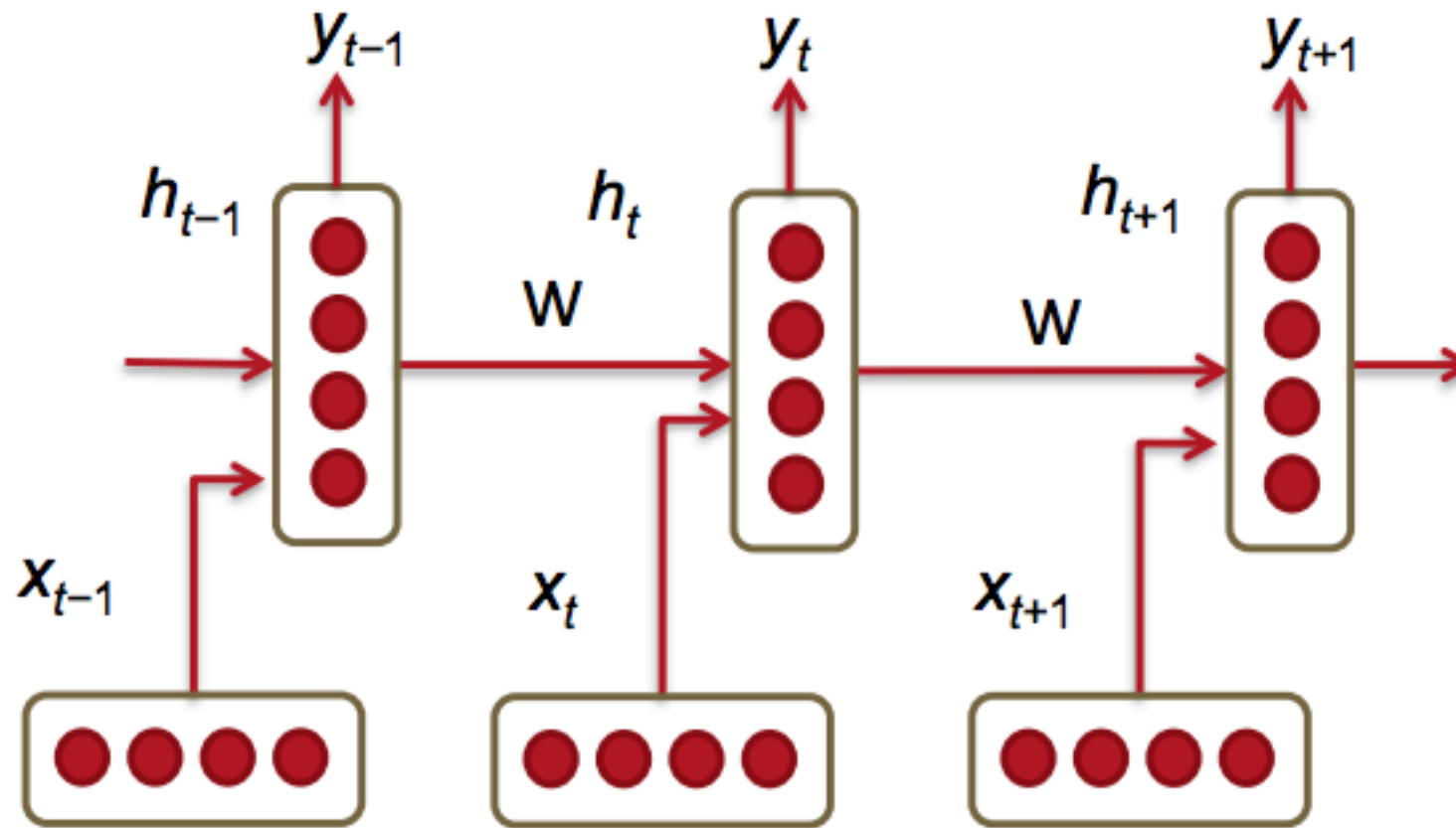
many to many



many to many



# Vanilla RNN



$$h_t = \sigma(W^{(hh)}h_{t-1} + W^{(hx)}x_t)$$

$$\hat{y}_t = \text{softmax}(W^{(S)}h_t)$$



# Vanishing/exploding gradient

$$J^{(t)}(\theta) = - \sum_{j=1}^{|V|} y_{t,j} \times \log(\hat{y}_{t,j})$$

$$\frac{\partial E_t}{\partial W} = \sum_{k=1}^t \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial W}$$

$$E_t == J^{(t)}(\theta)$$

$$\frac{\partial h_t}{\partial h_k} = \prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}} = \prod_{j=k+1}^t W^T \times \text{diag}[f'(h_{j-1})]$$

$$\left\| \frac{\partial h_j}{\partial h_{j-1}} \right\| \leq \| W^T \| \left\| \text{diag}[f'(h_{j-1})] \right\| \leq \beta_W \beta_h$$

$$\left\| \frac{\partial h_t}{\partial h_k} \right\| = \left\| \prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}} \right\| \leq (\beta_W \beta_h)^{t-k}$$

# Exploding grad solution

---

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

---

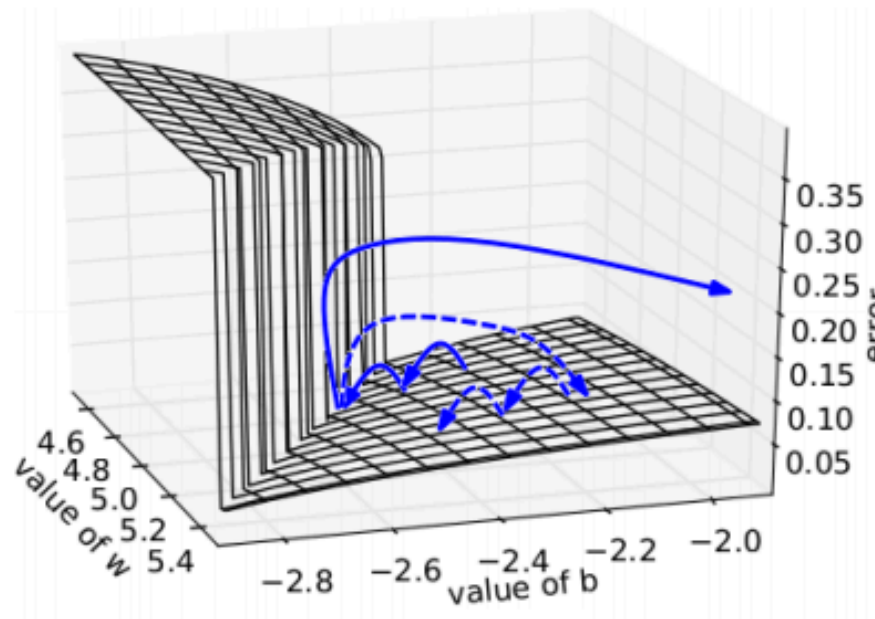


Figure 5: Gradient explosion clipping visualization

# Vanishing grad solution

- Initialize  $W$  with identity matrix
- Use ReLU activation
- Or skip to the next slides with LSTM & GRU

# LSTM

- Vanilla RNN cells can not capture long-term dependencies. Let's develop another cell!

$$i_t = \sigma(W^{(i)}x_t + U^{(i)}h_{t-1}) \quad (\text{Input gate})$$

$$f_t = \sigma(W^{(f)}x_t + U^{(f)}h_{t-1}) \quad (\text{Forget gate})$$

$$o_t = \sigma(W^{(o)}x_t + U^{(o)}h_{t-1}) \quad (\text{Output/Exposure gate})$$

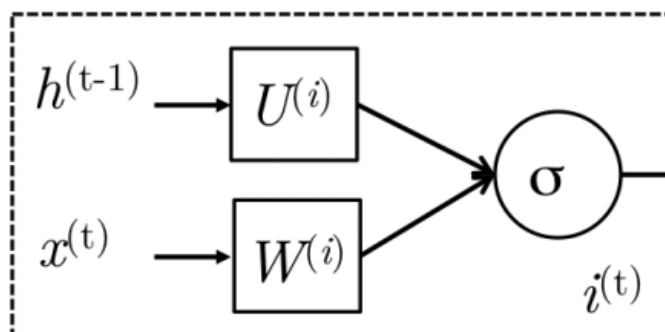
$$\tilde{c}_t = \tanh(W^{(c)}x_t + U^{(c)}h_{t-1}) \quad (\text{New memory cell})$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \quad (\text{Final memory cell})$$

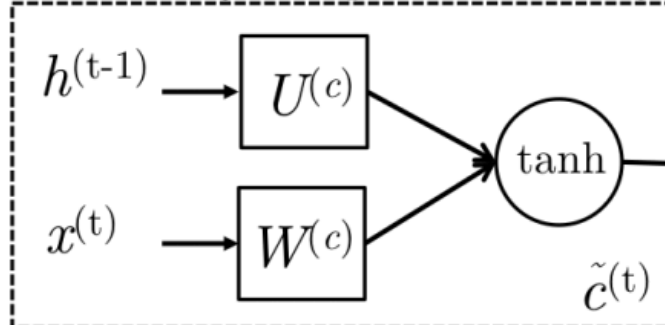
$$h_t = o_t \circ \tanh(c_t)$$

# LSTM

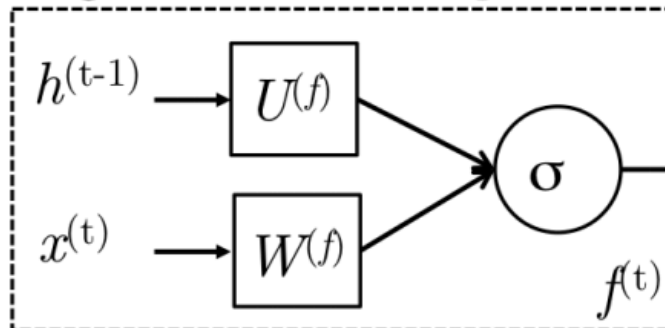
Input: Does  $x^{(t)}$  matter?



New memory: Compute new memory

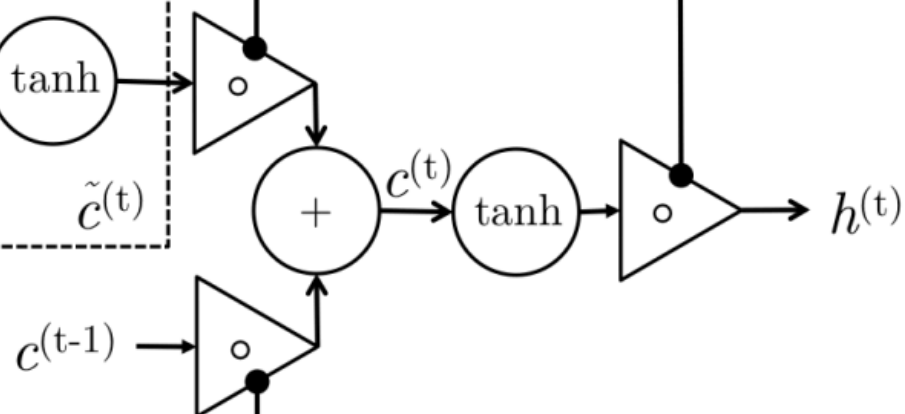
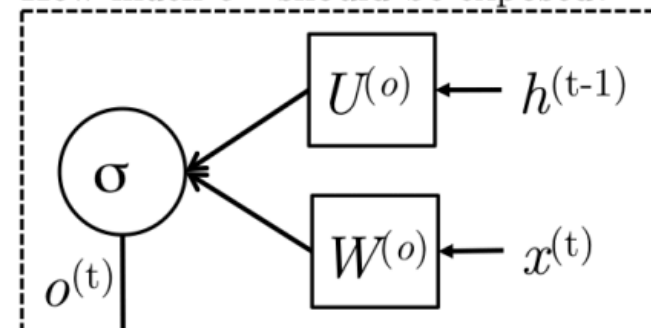


Forget: Should  $c^{(t-1)}$  be forgotten?



Output/Exposure:

How much  $c^{(t)}$  should be exposed?



# GRU

- Another cell architecture but idea is the same: capture long-dependencies + eliminate vanishing gradient problem.

$$z_t = \sigma(W^{(z)}x_t + U^{(z)}h_{t-1}) \quad \text{(Update gate)}$$

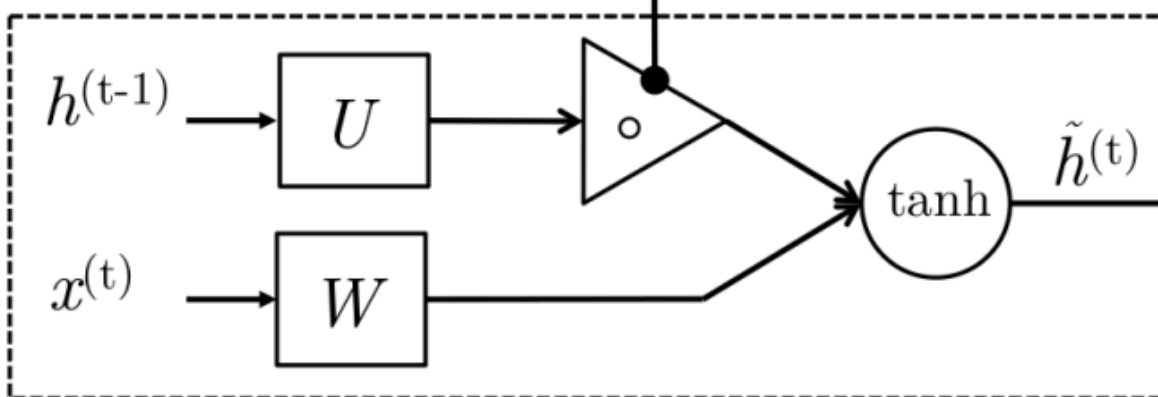
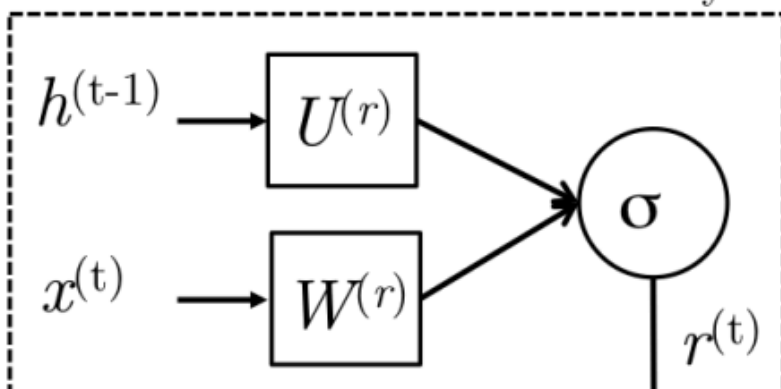
$$r_t = \sigma(W^{(r)}x_t + U^{(r)}h_{t-1}) \quad \text{(Reset gate)}$$

$$\tilde{h}_t = \tanh(r_t \circ U h_{t-1} + W x_t) \quad \text{(New memory)}$$

$$h_t = (1 - z_t) \circ \tilde{h}_t + z_t \circ h_{t-1} \quad \text{(Hidden state)}$$

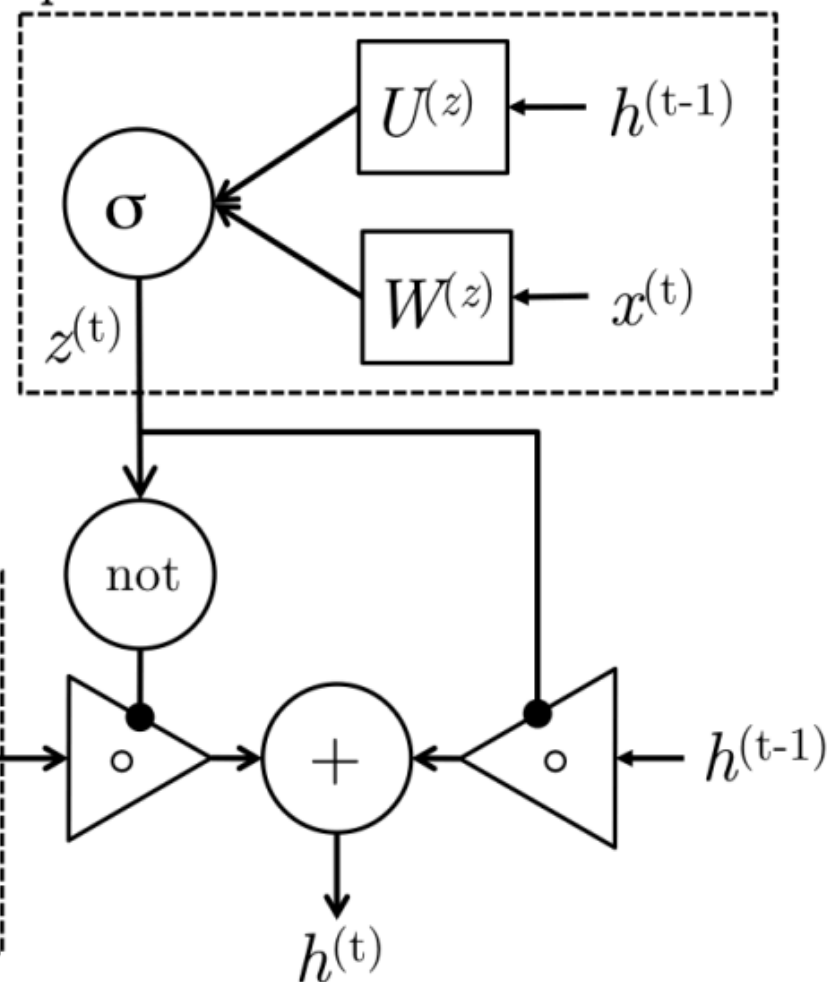
# GRU

Reset: Include  $h^{(t-1)}$  in new memory?

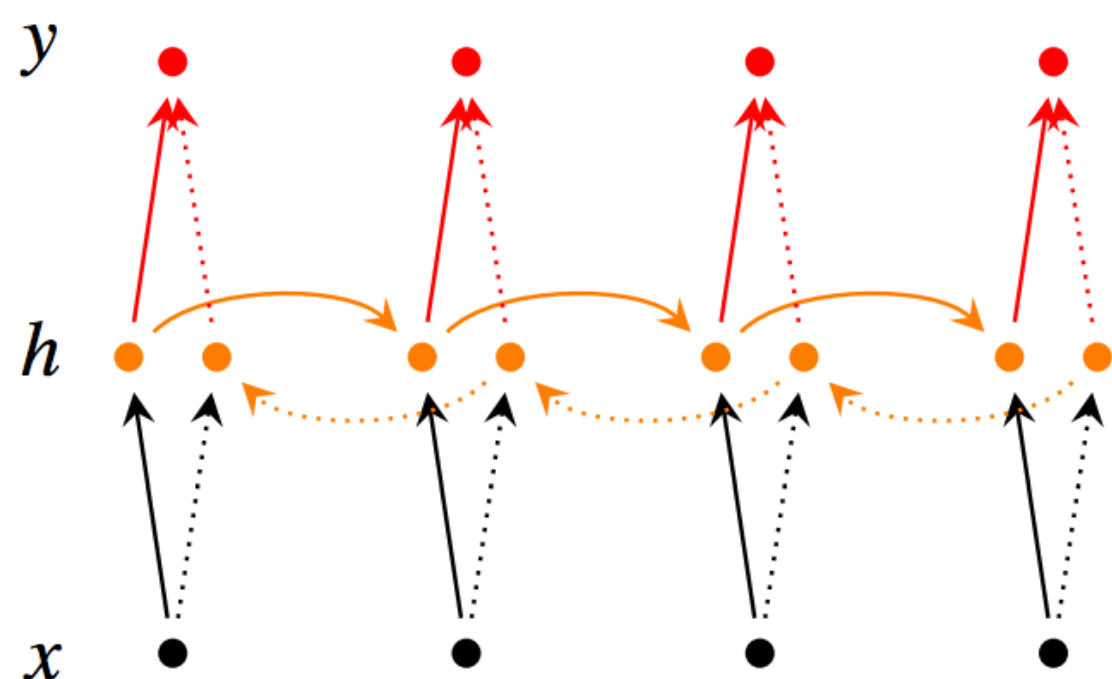


New memory: Compute new memory based on current word input  $x^{(t)}$  and potentially  $h^{(t-1)}$

Update: How much  $h^{(t-1)}$  in next state?



# Bidirectional RNN



$$\vec{h}_t = f(\vec{W}x_t + \vec{V}\vec{h}_{t-1} + \vec{b})$$

$$\overleftarrow{h}_t = f(\overleftarrow{W}x_t + \overleftarrow{V}\overleftarrow{h}_{t+1} + \overleftarrow{b})$$

$$\hat{y}_t = g(Uh_t + c) = g(U[\vec{h}_t; \overleftarrow{h}_t] + c)$$

Figure 6: A bi-directional RNN model



# Applications

- Machine Translation, e.g. Google Neural Machine Translation  
<https://arxiv.org/pdf/1609.08144v2.pdf>
- Conversational models & dialog systems  
<https://arxiv.org/pdf/1506.05869.pdf>
- Image captioning  
<https://cs.stanford.edu/people/karpathy/deepimagesent/>
- Speech recognition  
<http://www.jmlr.org/proceedings/papers/v32/graves14.pdf>
- And many more...