Recurrent Neural Networks and Encoder-Decoder Models

Diego Marcheggiani

University of Amsterdam ILLC

Unsupervised Language Learning 2016

Outline

Word Embeddings Recap

Recurrent Neural Networks

Encoder Decoder Approach

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @

Outline

Word Embeddings Recap

Recurrent Neural Networks

Encoder Decoder Approach

◆□▶ ◆□▶ ◆ ≧▶ ◆ ≧▶ ○ ≧ ○ � � �

▶ In traditional NLP words are represented as a one-hot vector

- In traditional NLP words are represented as a one-hot vector
- Syntactically and semantically correlated words appear as completely different symbols

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへぐ

- In traditional NLP words are represented as a one-hot vector
- Syntactically and semantically correlated words appear as completely different symbols
- Add general features to the word, pos, prefixes, etc. to overcome the problem of data sparsity

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

 machine learning algorithms "understand" that dog is similar to dogs

- In traditional NLP words are represented as a one-hot vector
- Syntactically and semantically correlated words appear as completely different symbols
- Add general features to the word, pos, prefixes, etc. to overcome the problem of data sparsity

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

 machine learning algorithms "understand" that dog is similar to dogs

but what about cat and dog ? or good and amazing ?

"You shall know a word by the company it keeps!" (J. R. Firth, 1957)

(ロ)、(型)、(E)、(E)、 E) の(()

"You shall know a word by the company it keeps!" (J. R. Firth, 1957)

the meaning of a word is given by the context in which the word appears

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへぐ

"You shall know a word by the company it keeps!" (J. R. Firth, 1957)

- the meaning of a word is given by the context in which the word appears
- if two words appear in similar contexts (company), they have, to some extent, similar meanings.

- ロ ト - 4 回 ト - 4 □

"You shall know a word by the company it keeps!" (J. R. Firth, 1957)

- the meaning of a word is given by the context in which the word appears
- if two words appear in similar contexts (company), they have, to some extent, similar meanings.

How can we formalize this intuition?

Co-occurence matrix

Corpus sentences		Co-occurre	nce counts	Vector
He also found five fish swimming in murky water in an old <u>bathtub</u> .]	the	12	$\begin{pmatrix} 12\\ 9 \end{pmatrix}$
We do abhor dust and dirt, and stains on the <u>bathtub</u> , and any kind of filth.		a of	7	7
Above At the far end of the garden room a <u>bathtub</u> has been planted with herbs for the winter.		and	6	6
They had been drinking Cisco, a fruity, wine-based fluid that smells and tastes like a mixture of cough syrup and <u>bathtub</u> gin.	→	in 	5	
Science finds that a surface tension on the water can draw the boats together, like toy boats in a <u>bathtub</u> .		like water	2	2
In fact, the godfather of gloom comes up with a plot that takes in Windsor Davies (the ghost of sitcoms past), a		boat	2	2
"I'll tell him,' said the Dean from the bathroom above the sound of bathwater falling from a great height into the ample Edwardian <u>bathtub</u> .		from	2	$\begin{array}{c} 2\\ 1\end{array}$
		toy	1	1
		god- father	1	1
		Cisco	1	

◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□▶

Word similarity



Vanilla count-based approaches results in huge sparse matrices



Vanilla count-based approaches results in huge sparse matrices

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへぐ

They can be approximated to dense matrices using SVD.

- Vanilla count-based approaches results in huge sparse matrices
- They can be approximated to dense matrices using SVD.
- Bengio et al. 2006 introduced a neural language model for learning word representations.

- Vanilla count-based approaches results in huge sparse matrices
- They can be approximated to dense matrices using SVD.
- Bengio et al. 2006 introduced a neural language model for learning word representations. expensive

- Most popular model nowadays is Google's skip-gram word2vec [Mikolov et al. 2013].
- The idea is to predict context words of a given word.

Skip-gram word2vec

- $W \in \mathbb{R}^{|V| imes d}$ is the word embedding matrix
- $C \in \mathbb{R}^{|V| \times d}$ is the context embedding matrix
- V is the words dictionary
- d is the vector size of the word representation
- the size of the context window is given as hyperparameter

For each context, word pair $(c, w) \in D$, we want to maximize

$$p(c|w) = \frac{\exp(C(c) \cdot W(w))}{\sum_{c'} \exp(C(c') \cdot W(w))},$$
(1)

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Skip-gram word2vec

For each context, word pair $(c, w) \in \mathcal{D}$, we want to maximize

$$p(c|w) = \frac{\exp(C(c) \cdot W(w))}{\sum_{c'} \exp(C(c') \cdot W(w))},$$
(2)

Cons

• $\sum_{c'}$ is intractable, but negative sampling does the job.

Skip-gram word2vec

For each context, word pair $(c, w) \in D$, we want to maximize

$$p(c|w) = \frac{\exp(C(c) \cdot W(w))}{\sum_{c'} \exp(C(c') \cdot W(w))},$$
(2)

▲□▶▲□▶▲≡▶▲≡▶ ≡ めぬぐ

Cons

• $\sum_{c'}$ is intractable, but negative sampling does the job. Pros

Log-linear, easy to train!

Ok, well done! Now what?



Ok, well done! Now what?

- word embeddings as input to neural networks
- since deep learning models seek to solve non-convex optimization problems, starting from a good point in the parameters space usually helps.

Outline

Word Embeddings Recap

Recurrent Neural Networks

Encoder Decoder Approach

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @



Input layer (x)

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @



◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @

• W_1 is a parameter matrix and b_1 is the bias



▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

- W_1 is a parameter matrix and b_1 is the bias
- f is a non-linear activation function
- h is a hidden layer of the network



▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

- W_1 is a parameter matrix and b_1 is the bias
- f is a non-linear activation function
- h is a hidden layer of the network
- W_2 , b_2 more parameters

Multi-layer perceptron



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

- W_1 is a parameter matrix and b_1 is the bias
- f is a non-linear activation function
- h is a hidden layer of the network
- W_2 , b_2 more parameters

Multi-layer perceptron



- W_1 is a parameter matrix and b_1 is the bias
- f is a non-linear activation function
- h is a hidden layer of the network
- W_2 , b_2 more parameters



▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

Standard feedforward neural networks cannot deal with variable length input.

Language modeling

Compute the probability of a sentence, or given a sequence of words predicting the word that comes next.

• trigram language model $p(x_t|x_{t-1}, x_{t-2})$



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

- cannot "see" the subject so cannot predict the right verb
- increase the size of the model.

Language modeling

Compute the probability of a sentence, or given a sequence of words predicting the word that comes next.

• trigram language model $p(x_t|x_{t-1}, x_{t-2})$



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

- cannot "see" the subject so cannot predict the right verb
- increase the size of the model.
- what about longer dependencies ?



▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = のへで



▲□▶ ▲□▶ ▲ 三▶ ▲ 三 ● ● ●



▲口 ▶ ▲圖 ▶ ▲ 臣 ▶ ▲ 臣 ▶ 二 臣



・ロト ・四ト ・ヨト ・ヨト

æ



ヘロト 人間 とくほとくほとう

æ



・ロト ・ 一下・ ・ ヨト ・

æ



æ

<ロ> (日) (日) (日) (日) (日)



▲ロト ▲暦 ト ▲ 臣 ト ▲ 臣 ト ● 回 ● の Q (2)



In each step we use information coming from all the previous steps.

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

Language modeling (reprise)



▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

Long dependencies are captured

Language modeling (reprise)



Long dependencies are captured (at least in theory)

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

Training RNNs

- Total error is the sum of errors at each time step $\sum_{t} E(y_t^*, y_t)$
- Error is calculated as cross entropy loss $E(y_t^*, y_t) = -y_t^* \log y_t$
- Calculating the gradients for V depends only on the current time step t.
- ▶ For the recurrent parameters *W* is a bit different.



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Since we sum errors, we also sum gradients at each time step.

$$\frac{\partial E}{\partial W} = \sum_{t} \frac{\partial E_{t}}{\partial W}$$

Following standard backprop derivation, for each time step t we end-up having:





For t = 3 we have:

$$\frac{\partial E_3}{\partial W} = \sum_{k=1}^2 \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial s_k} \frac{\partial s_k}{\partial W}$$

that is unfolded as:





< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

$$\frac{\partial E_3}{\partial W} = \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial W} + \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial W}$$



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

$$\frac{\partial E_3}{\partial W} = \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial W} + \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial W}$$



$$\frac{\partial E_3}{\partial W} = \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial W} + \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial W}$$



$$\frac{\partial E_3}{\partial W} = \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial W} + \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial W}$$



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─ 臣 ─ のへで

$$\frac{\partial E_3}{\partial W} = \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial W} + \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial W}$$



◆ロト ◆母 ト ◆臣 ト ◆臣 ト ○臣 ○ のへで

$$\frac{\partial E_3}{\partial W} = \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \frac{\partial s_1}{\partial W} + \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial W}$$



◆ロ ▶ ◆母 ▶ ◆臣 ▶ ◆臣 ▶ ○ 臣 ○ のへで

RNN learning issues

Vanishing (exploding) gradient for long dependencies.

- Exploding gradient solution:
 - Gradient clipping clipping the norm of the gradient if exceeds a certain threshold.
- Vanishing gradient solution:
 - Long short-term memory networks (LSTM) [Hochreiter and Schmidhuber 1997]
 - ▶ Gated recurrent unit networks (GRU) [Cho et al. 2014]



ヨト くヨト 二

-

Outline

Word Embeddings Recap

Recurrent Neural Networks

Encoder Decoder Approach

◆□ ▶ ◆□ ▶ ◆ 臣 ▶ ◆ 臣 ▶ ○ 臣 ○ の Q @

- What if we want an output of variable length with respect to the input?
- Can we use RNN for tasks where we have a sequence as input and we want another sequence as output?

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- What if we want an output of variable length with respect to the input?
- Can we use RNN for tasks where we have a sequence as input and we want another sequence as output?

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

machine translation: Mary eats apples. ->

- What if we want an output of variable length with respect to the input?
- Can we use RNN for tasks where we have a sequence as input and we want another sequence as output?
- machine translation: Mary eats apples. -> Marie mange des pommes.
- output question answering: Tim is playing in his room. ||Where is Tim? ->

- What if we want an output of variable length with respect to the input?
- Can we use RNN for tasks where we have a sequence as input and we want another sequence as output?
- machine translation: Mary eats apples. -> Marie mange des pommes.
- question answering: Tim is playing in his room. ||Where is Tim? -> Tim is in his room.

- What if we want an output of variable length with respect to the input?
- Can we use RNN for tasks where we have a sequence as input and we want another sequence as output?
- machine translation: Mary eats apples. -> Marie mange des pommes.
- question answering: Tim is playing in his room. ||Where is Tim? -> Tim is in his room.

Encoder decoder a.k.a. sequence to(2) sequence learning. [Sutskever et al. 2014]

Sequence to sequence learning - Encoding

 $p(z_1,...,z_{T'}|x_1,...,x_T)$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

Sequence to sequence learning - Encoding

 $p(z_1,...,z_{T'}|x_1,...,x_T)$

Encode the entire input sequence in a single vector

$$s_t = f(U \cdot x_t + W \cdot s_{t-1})$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで



Sequence to sequence learning - Decoding

 $p(z_1,...,z_{T'}|x_1,...,x_T)$

proceed like in a standard RNN starting from the encoding vector

 $y_1 = softmax(V \cdot s_{T_e})$



Sequence to sequence learning - Decoding

 $p(z_1,...,z_{T'}|x_1,...,x_T)$

proceed like in a standard RNN starting from the encoding vector

$$y_1 = softmax(V \cdot s_{T_e})$$

and then continue decoding the second word in the sequence with

$$s_1^d = f(U \cdot z_1 + W \cdot s_{\mathcal{T}_e}), \quad y_2 = softmax(V \cdot s_1^s)$$



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Sequence to sequence learning - Decoding

 $p(z_1,...,z_{T'}|x_1,...,x_T)$

Begin to decode the output sequence conditioning on this vector

$$y_1 = softmax(V \cdot s_{T_e})$$

continue decoding the second word in the sequence with

$$s_1^d = f(U \cdot z_1 + W \cdot s_{T_e}), \quad y_2 = softmax(V \cdot s_1^d)$$

keep decoding until an <eos> token is predicted.



Encoder decoder





イロト イヨト イヨト

э

Conclusions

- Vanishing gradient is even more problematic in this setting.
- LSTMs or GRUs are crucial to capture long term dependencies.
- Very appealing in task like machine translation.



Lecture recap

- Word embeddings recap
- Neural networks recap
- Recurrent neural networks
- Sequence to sequence learning with neural networks

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

That's it!

Thanks for your attention! Questions?

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで