

MIPT Data Visualization Course

Deep Learning for Data with Sequence Structure

Ashuha Arseniy^{1,2}

Bayesian Research Group¹, MIPT²



ars-ashuha.ru/slides

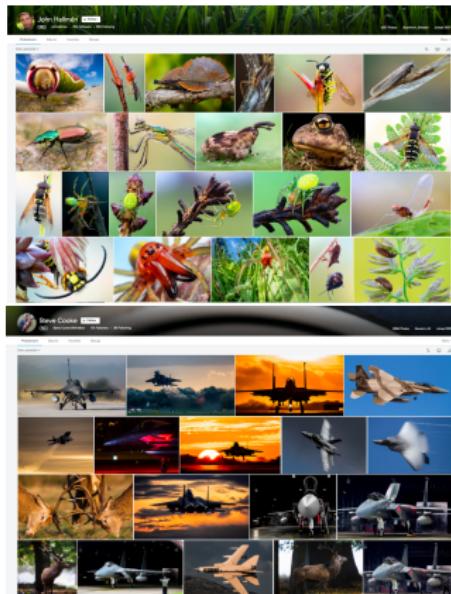
December 20, 2016

Independent and Identically-Distributed

- ▶ Image classification



- ▶ Flicker user's image classification

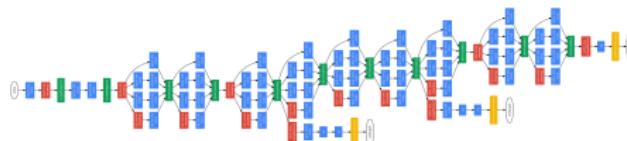


Flicker user's image classification

- Now we need to classify objects jointly, so we had sequence on input



- and predict sequence
 - [‘deer’, ‘plane’, ‘plane’, ‘deer’, ‘plane’]
- What should we do? Do you have any ideas?
- Lets apply CNN for each image to get descriptor.



- We've transform sequence of pictures so sequence of vectors
- Ok data became much simply, what's next?

Recurrent Neural Net Concept

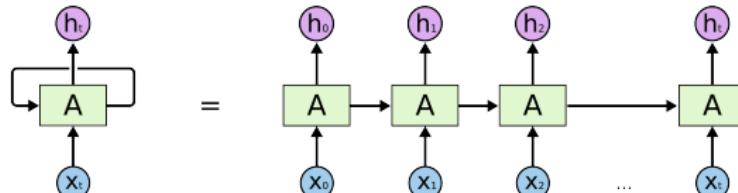
- ▶ Let's introduce memory cell

```
memory = np.ones(num_class)
for pic_vec in ['pic1_vec', 'pic2_vec', 'pic3_vec']:
    pic_class = g(W.dot(pic_vec) * f(memory))
    memory[pic_class] += 1
```

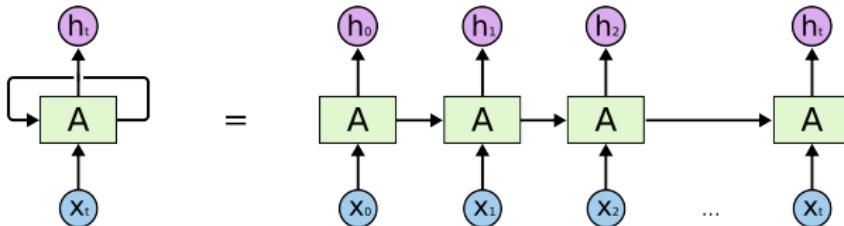
- ▶ Hmm, we will try control memory automatic

```
memory = np.ones(100)
for pic_vec in ['pic1_vec', 'pic2_vec', 'pic3_vec']:
    pic_class = g(Wout.dot([pic_vec, memory]))
    memory += Win.dot(pic_vec) + sigmoid(Wmem.dot(memory))
```

- ▶ Let's plot (memory update it's a little bit different but still it's same)



RNN and Training



- ▶ So we can compute inputs-outputs using this equations

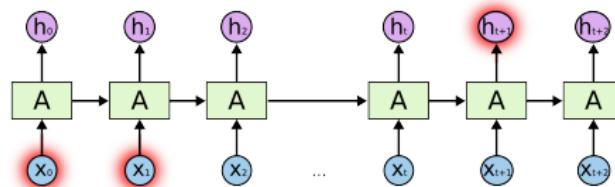
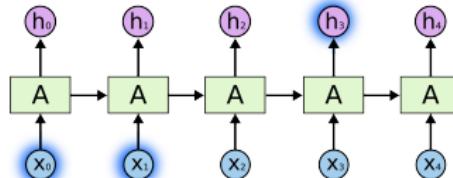
$$h_n = W_x x_n + W_h \sigma(h_{n-1}) \quad \hat{y}_n = \sigma(W_y h_n) \quad \theta = (W_x, W_h, W_y)$$

- ▶ You get several inputs and outputs, how would you train RNNs?
- ▶ Let's sum error for each output

$$L = \sum_t L_t(\hat{y}_t(\theta), y_t), \quad \frac{\partial L}{\partial \theta} = \sum_{1 < t < T} \frac{\partial L_T(\hat{y}_t(\theta), y_t)}{\partial \theta}$$

- ▶ Next we should use ordinary optimization methods

RNN and Training



- Gradient is like this

$$\frac{\partial E_t}{\partial \theta} = \sum_{1 \leq k \leq t} \frac{\partial E_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial \theta} \quad \frac{\partial h_t}{\partial h_k} = \prod_{t \geq i \geq k} \frac{\partial h_i}{\partial h_{i-1}}$$

- Multiply on matrix with small norm means no long deps
- Multiply on matrix with big norm means Gradient Boom

Dealing with gradients Boom

- ▶ What is the simplest solution?
- ▶ Gradient clipping :)

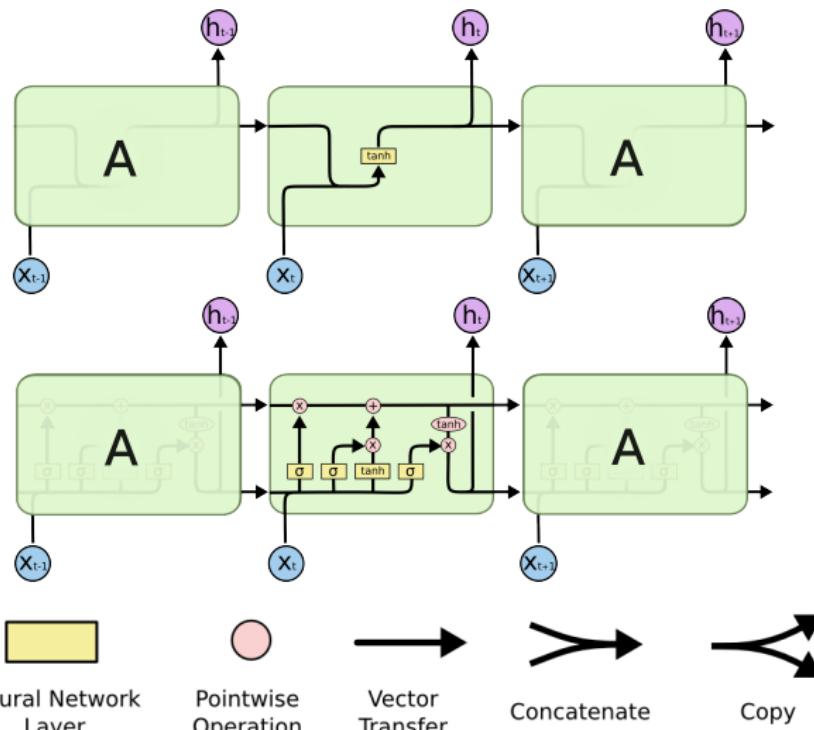
Algorithm 1 Pseudo-code for norm clipping the gradients whenever they explode

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

- ▶ It really works while gradient variance is low

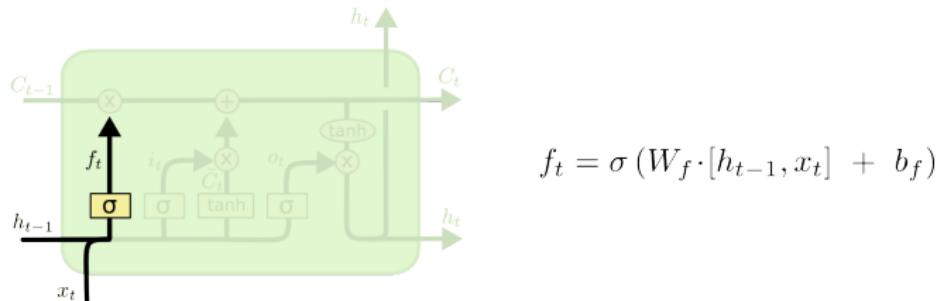
Dealing with long term dependences means LSTM

- ▶ We can't forget
- ▶ We can't understand importance of input information
- ▶ Let's add it

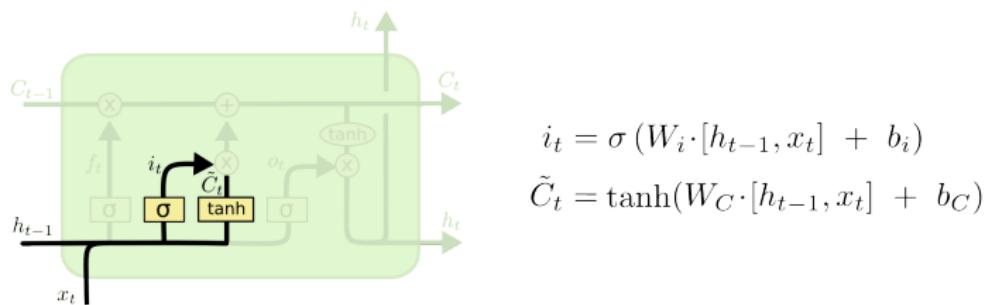


Dealing with long term dependences means LSTM

► Forget Gate

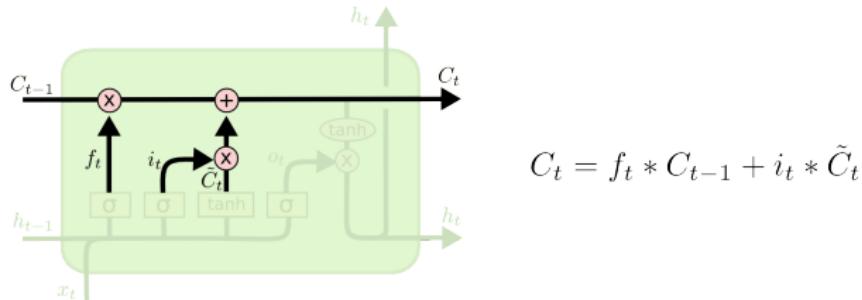


► Input Gate

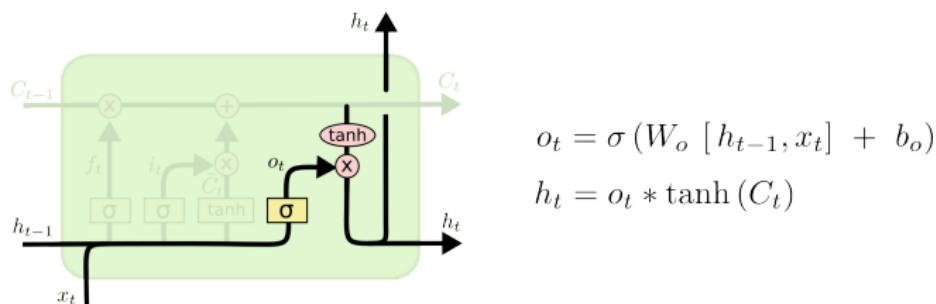


Dealing with long term dependences means LSTM

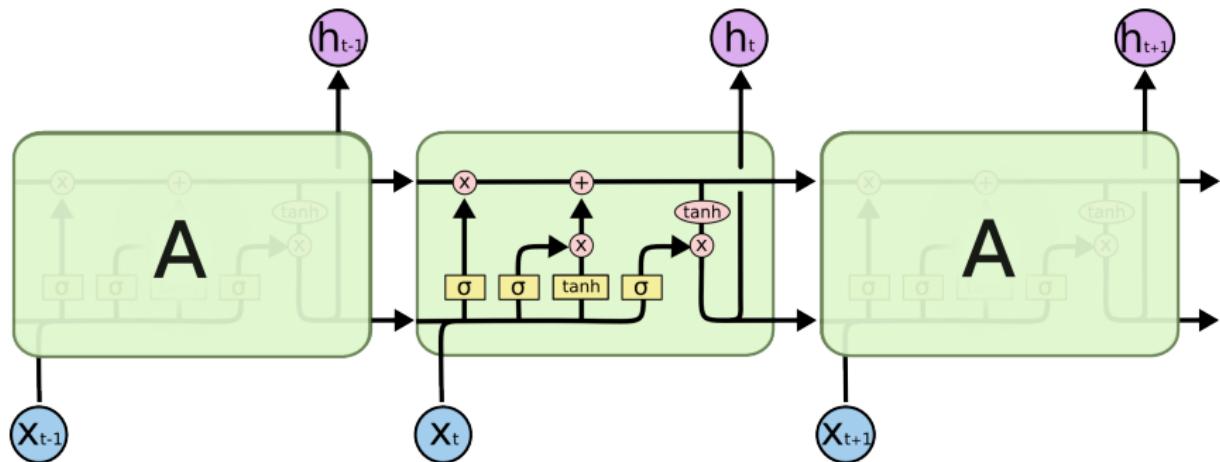
- ▶ Cell update



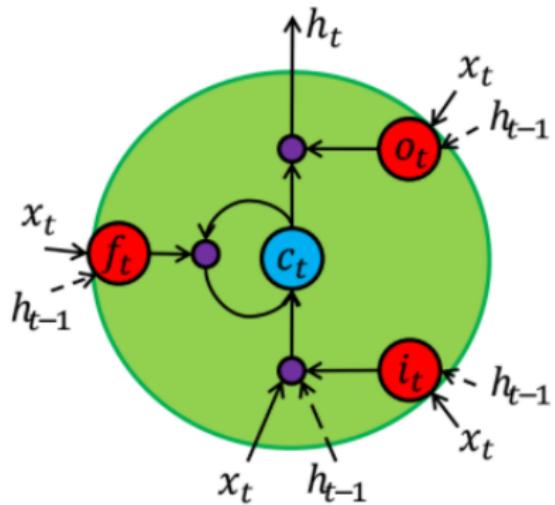
- ▶ Output Gate



Dealing with long term dependences means LSTM



Gradient Vanishing in LSTM

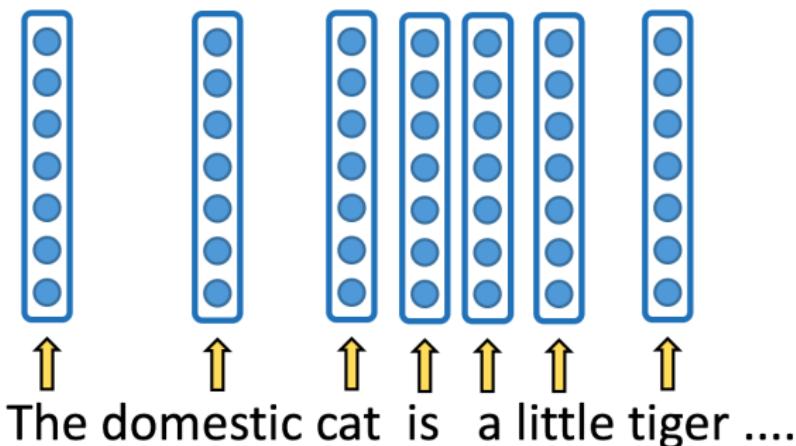


$$c_t = f_t \cdot c_{t-1} + i_t \cdot g(W_x x_t + W_c h_{t-1} + b_c), \quad h_t = o_t \cdot (c_t)$$

$$\frac{\partial c_{t+1}}{\partial c_t} = f_{k+1}$$

Embedding Layer

- ▶ What is the main difference between words and pictures?
- ▶ Number of all words is much less than in picture
- ▶ So we can map identity vector for each word

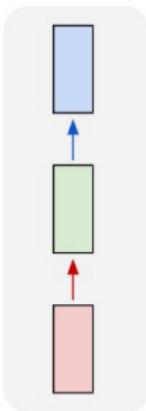


- ▶ Let's initialize with random noise
- ▶ So, we can evaluate gradient by word vectors
- ▶ Word vectors are just model parameters

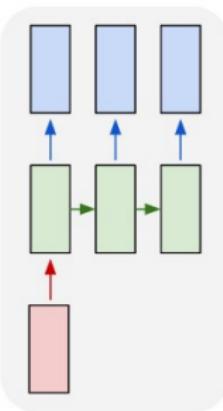
Practical Cases

- Different tasks type

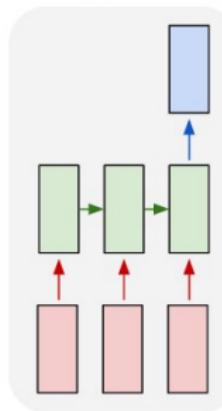
one to one



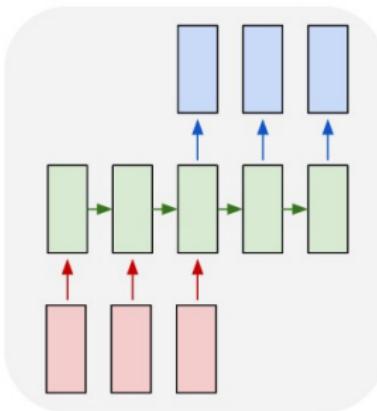
one to many



many to one



many to many

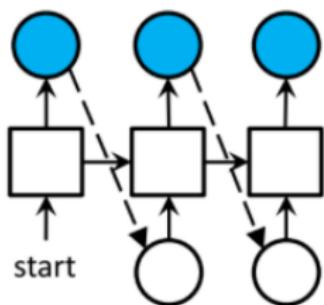


- How to code it

```
decoder = lasagne.layers.LSTMLayer(  
    np.oxes(batch, seq, vectors), num_units=100,  
    cell_init=smth, grad_clipping=5)
```

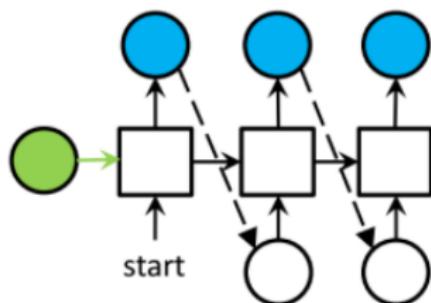
Text/Handwritten Generation

Next symbol/word
Current symbol/word



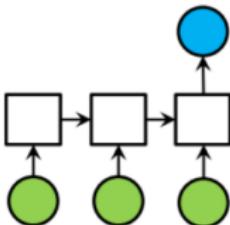
- ▶ <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>
- ▶ <http://www.cs.toronto.edu/~graves/handwriting.html>

Sequence generation



- ▶ <http://www.cs.toronto.edu/~nitish/nips2014demo/index.html>
- ▶ cs.stanford.edu/people/karpathy/deepimagesent/generationdemo/

Sequence classification



Category: Транспорт
SubCategory: Автомобили с пробегом > BMW > X1

BMW X1, 2012

Цена: 1 200 000 руб.

Продавец: Владимир

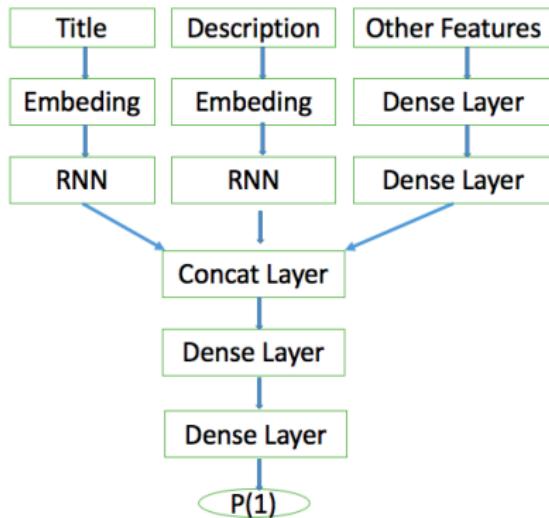
Город: Москва

Описание: Красивый автомобиль в отличном состоянии.

Технические характеристики:

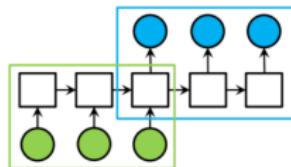
- Фото 30 000 - 34 999 км, 2.0 АТ, бензин, полный привод, кроссовер, левый руль, цвет черный
- 2 литра, 184 л. с. Полный привод, 8 ступ. автомат, климат-контроль, подогрев сидений, звуковые колонки и т.д.
- Не бита, не гнилая.
- Пробег 32 000 км.
- Услуга обслуживания только у официального дилера

Номер объявления: 335309783

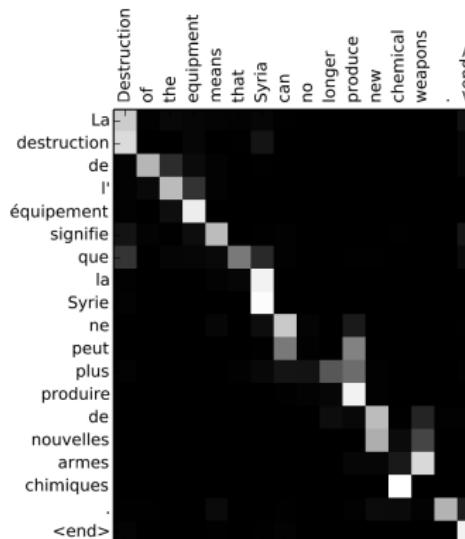


Sequence2Sequence Machine Translation

- ▶ Coding input sentence and generate output



- ▶ Add attention mechanize



Recap

- ▶ Recurrent Deep Neural nets are great
- ▶ RNN has achieved many state of the art awards
- ▶ Input and output may have any length
- ▶ It can be combined with convolution DNN
- ▶ It's Easy to code it by yourself with any DL framework
- ▶ Link's
 - ▶ <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>
 - ▶ <http://www.machinelearning.ru/wiki/index.php?title=DL>
 - ▶ www.deeplearningbook.org