Deep Learning for Natural Language Processing

Stephen Clark et al. University of Cambridge and DeepMind





8. Long Short Term Memory

Felix Hill DeepMind





RNNs: A recap







What is the forward computation like?

 $h_t = tanh(Uh_{t-1} + Vx_t)$





What is the forward computation like?

 $h_t = tanh(Uh_{t-1} + Vx_t)$ $h_{t+1} = tanh(Uh_t + Vx_{t+1})$





What is the forward computation like?

 $h_{t} = tanh(Uh_{t-1} + Vx_{t})$ $h_{t+1} = tanh(Uh_{t} + Vx_{t+1})$ $= tanh(Utanh(Uh_{t-1} + Vx_{t}) + Vx_{t+1})$ = tanh(Utanh(Utanh(Utanh.....))





"Vanishing" gradients

c(f(x), y)

story all about how my life got flipped turned

now this is a

 \mathcal{X}

$$\frac{dC}{dw_1} \propto \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \cdots \times w_n \times \sigma'(z_n) \times \frac{dC}{da_n}$$

where

 $a_i = \sigma(z_i)$



"Vanishing" gradients

y

an

a2

a1

 \overline{W}

upside =

 $\overline{w}_{\underline{n}\underline{-}}$

 w_n

(or exploding)



Long Short Term Memory (LSTM)

Hochreiter and Schmidhuber, 1997







LSTM





Thanks to http://colah.github.io/posts/2015-08-Understanding-LSTMs/



xt



Be careful with the past!



Thanks to http://colah.github.io/posts/2015-08-Understanding-LSTMs/





Sigmoid "gates"

$$f_t = \sigma(W_{f_h}h_{t-1} + W_{f_x}x_t + b_f)$$

$$i_t = \sigma(W_{i_h}h_{t-1} + W_{i_x}x_t + b_i)$$

$$o_t = \sigma(W_{o_h} h_{t-1} + W_{o_x} x_t + b_o)$$









New stuff to consider

$$f_{t} = \sigma(W_{f_{h}}h_{t-1} + W_{f_{x}}x_{t} + b_{f})$$

$$i_{t} = \sigma(W_{i_{h}}h_{t-1} + W_{i_{x}}x_{t} + b_{i})$$

$$o_{t} = \sigma(W_{o_{h}}h_{t-1} + W_{o_{x}}x_{t} + b_{o})$$



$$\tilde{C}_t = \tanh(W_{c_h}h_{t-1} + W_{c_x}x_t + b_c)$$

$$\sum$$
Look familiar?







Carefully update your cell!!

$$f_{t} = \sigma(W_{f_{h}}h_{t-1} + W_{f_{x}}x_{t} + b_{f})$$

$$i_{t} = \sigma(W_{i_{h}}h_{t-1} + W_{i_{x}}x_{t} + b_{i})$$

$$o_{t} = \sigma(W_{o_{h}}h_{t-1} + W_{o_{x}}x_{t} + b_{o})$$

$$\tilde{C}_t = \tanh(W_{c_h}h_{t-1} + W_{c_x}x_t + b_c)$$









Use your cell (if you want!!)

$$f_{t} = \sigma(W_{f_{h}}h_{t-1} + W_{f_{x}}x_{t} + b_{f})$$

$$i_{t} = \sigma(W_{i_{h}}h_{t-1} + W_{i_{x}}x_{t} + b_{i})$$

$$o_{t} = \sigma(W_{o_{h}}h_{t-1} + W_{o_{x}}x_{t} + b_{o})$$

$$\tilde{C}_{t} = \tanh(W_{c_{h}}h_{t-1} + W_{c_{x}}x_{t} + b_{c})$$

$$C_{t} = C_{t-1} * f_{t} + \tilde{C}_{t} * i_{t}$$

 $h_t = o_t * \tanh(C_t)$









Why does it work?







How many weights in an LSTM?

1: 100,000 words in my vocab

2: layers C_t have 500 units





Can we go deeper?









Cell sensitive to position in line:

The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae-pressed forward into boats and into the ice-covered water and did not, surrender.





Cell that turns on inside quotes:

"You mean to imply that I have nothing to eat out of... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."





Cell that robustly activates inside if statements:







A large portion of cells are not easily interpretable. Here is a typical example: /* Unpack a filter field's string representation from user-space * buffer. */ char *audit_unpack_string(void **bufp, size_t *remain, size_t len) { char *str; if (!*bufp || (len == 0) || (len > *remain)) return ERR_PTR(-EINVAL); /* Of the currently implemented string fields, PATH_MAX * defines the longest valid length. */





Why is it called Long Short Term Memory?





NLP for the LSTM generation



"Read" the sentence both ways and concatenate h_t at each time step

"LSTM-based deep learning models for non-factoid answer selection." Ming Tan et al. 2015

References

A problem

"Learning long-term dependencies with gradient descent is difficult". Bengio et al. 1994

A solution

"Long short-term memory". Hochreiter and Schmidhuber, 1997

A(nother) cottage industry "LSTM-based deep learning models for non-factoid answer selection." Ming Tan et al. 2015

A very clear explanation of LSTMs http://colah.github.io/posts/2015-08-Understanding-LSTMs/