

Vanishing Gradients and Fancy RNNs

Shubham Gupta

February 3, 2020

1 Introduction

- Learn about problems with RNN
- More RNN variants

2 Vanishing gradient problem

- Occurs in RNN
- Small gradient in each step reduces the overall gradient signal as it back-propogates further.
- Why is it a problem?
 - Gradient signal from faraway is lost because it's much smaller than gradient signal from closeby.
 - Model weights are only updated with respect to *near effects*, not long term effects.
 - *TLDR*: Model will not learn the parameters well and hence will have weak predictability.
 - Gradient is the effect of the past on the future.
 - If it doesnt learn the parameters, then either
 - * No dependency at t and $t+1$
 - * Or it learns wrong parameters to capture true dep between t and $t + 1$
 - Syntactic recency: Pays attention to syntax of sentence i.e longer language dependency
 - Sequential recency: Pays attention to things that only happen recently
 - Due to vanishing gradient problem, RNN learns sequency recency more.

3 Exploding gradients

- Gradient too big \implies SGD update too big $\theta^{new} = \theta^{old} - \alpha \delta_\theta J(\theta)$
- Solution: **Gradient clipping**
- If norm of gradient is above threshold, normalize gradient before applying SGD update
- Normalize gradient by setting max and min thresholds. This will prevent gradient from changing drastically, thereby avoiding exploding gradients problem.

3.1 Fix vanishing gradients problem

- Separate memory for longer dependencies
- Solution: **LSTM**
- At step t , there is hidden state h^t and cell state c^t
- Can erase, read and write cell state
- Gates control whether they will write, read, etc.
- Gates are also vectors
- Gates are dynamic. Diff on each step t .
- Gates are as follows:
 - Forget gate: $\sigma(W_f h^{t-1} + U_f x^t + b_f)$
 - Input gate: $\sigma(W_i h^{t-1} + U_i x^t + b_i)$
 - Output gate: $\sigma(W_o h^{t-1} + U_o x^t + b_o)$
- New cell content: $c^t = \tanh(W_c h^{t-1} + U_c x^t + b_c)$
- Forget some info using the forget gate
- Hidden state read output from some cell
- Solves vanishing gradient problem