

# CS 224n Assignment 3.

Abhishek Goswami.

January 28, 2019

## 1. Machine Learning & Neural Networks

### (a) Adam Optimizer.

- i.  $\mathbf{m}$  is rolling average of the gradients.  $\theta_1$  is a hyper parameter between 0 and 1.

$$\mathbf{m} \leftarrow \beta_1 \mathbf{m} + (1 - \beta_1) \nabla_{\theta} \mathbf{J}_{\text{minibatch}}(\theta). \quad (1)$$

The fact that  $\beta_1$  is set to a high value (0.9 often) suggests that  $\mathbf{m}$  is changing very little at each step. This in turn ensures that the model parameters do not bounce around when moving towards the local optimum.

- ii.  $\mathbf{v}$  is the rolling average of the magnitude of the gradients.

$$\theta \leftarrow \theta - \alpha \odot \mathbf{m} / \sqrt{\mathbf{v}}. \quad (2)$$

Since we are dividing by  $\sqrt{\mathbf{v}}$  it means that the larger gradients will get smaller updates. Conversely, the smaller gradients will get larger updates. This will help the learning algorithm to move off flat areas.

### (b) Dropout.

- i.  $\gamma$  in terms of  $p_{\text{drop}}$

$$\begin{aligned} h_i &= \mathbb{E}_{p_{\text{drop}}}[\mathbf{h}_{\text{drop}}]_i \\ &= \mathbb{E}_{p_{\text{drop}}}[\gamma d_i h_i] \\ &= p_{\text{drop}}(0) + (1 - p_{\text{drop}})\gamma h_i \\ &= (1 - p_{\text{drop}})\gamma h_i \end{aligned} \quad (3)$$

So,

$$\gamma = \frac{1}{1 - p_{\text{drop}}} \quad (4)$$

$$(5)$$

- ii. Dropout is a regularization technique that we want to use during training to prevent overfitting on the training data.

At test time, we should not do dropout, else it would result in randomness in predictions. One thing we should do during test time is to scale the predictions appropriately to account for the expected drop probability.

Alternatively, we can use **inverted dropout** so we do the scaling at training time itself, and leave the code untouched during test time.

Stack	Buffer	New dependency	Transition
[ROOT]	[I, parsed, this, sentence, correctly]		start
[ROOT, I]	[parsed, this, sentence, correctly]		SHIFT
[ROOT, I, parsed]	[this, sentence, correctly]		SHIFT
[ROOT, parsed]	[this, sentence, correctly]	parsed->I	LEFT-ARC
[ROOT, parsed, this]	[sentence, correctly]		SHIFT
[ROOT, parsed, this, sentence]	[correctly]		SHIFT
[ROOT, parsed, sentence]	[correctly]	sentence->this	LEFT-ARC
[ROOT, parsed]	[correctly]	parsed->sentence	RIGHT-ARC
[ROOT, parsed, correctly]	[]		SHIFT
[ROOT, parsed]	[]	parsed->correctly	RIGHT-ARC
[ROOT]	[]	ROOT->parsed	RIGHT-ARC

Table 1: Transitions for sentence “*I parsed this sentence correctly*”.

## 2. Neural Transition-Based Dependency Parsing

- (a) Table 1 shows the sequence of transitions needed for parsing the sentence : “*I parsed this sentence correctly*”
- (b) A sentence containing  $n$  words will be parsed in  $2n$  steps. It is either (a) SHIFT or (b) one of  $\{LEFT - ARC | RIGHT - ARC\}$
- (c) Coding exercise
- (d) Coding exercise
- (e) Coding exercise
- (f) Getting wrong.