

# CS 224N Assignment # 3

1. (a)

$$(i) \quad m \leftarrow \beta_1 m + (1 - \beta_1) \nabla_{\theta} J_{\text{mini-batch}}(\theta)$$

$$\theta \leftarrow \theta - \alpha m$$

Using  $\beta_1$  tracks the history of the gradients by the rolling average. Controlling  $\beta_1$  on the previous state  $m_{i-1}$ , we can let the recent changes have a bit more importance than the current gradient, thus reduce the variance

(ii) Small  $\nu$  values will get larger updates.

This can be helpful for learning because it can get recent parameters moving more efficiently along the axes and thus expediate the convergence.

$$(b) (i) \quad \gamma = \frac{1}{1 - P_{\text{drop}}} = \frac{1}{P_{\text{keep}}}$$

$$E_{\text{drop}}[h_{\text{drop}}]_i = P_{\text{drop}} \cdot 0 + (1 - P_{\text{drop}}) \gamma \cdot h_i = h_i$$

$$\gamma = \frac{1}{1 - P_{\text{drop}}}$$

(ii) dropout is used for preventing overfitting.

2. (a)

Stack	Buffer	Arc	Transition
[ ROOT ]	[ I, parsed, this, sentence, correctly ]		Initial
[ ROOT, I ]	[ parsed, this, sentence, correctly ]		Shift
[ ROOT, I, parsed ]	[ this, sentence, correctly ]		Shift
[ ROOT, parsed ]	[ this, sentence, correctly ]	parsed $\rightarrow$ I	Left
[ ROOT, parsed, this ]	[ sentence, correctly ]		Shift
[ ROOT, parsed, this, sentence ]	[ correctly ]		Shift
[ ROOT, parsed, sentence ]	[ correctly ]	sentence $\rightarrow$ this	Left
[ ROOT, parsed ]	[ correctly ]	parsed $\rightarrow$ sentence	Right
[ ROOT, parsed, correctly ]	[ ]		
[ ROOT, parsed ]	[ ]	parsed $\rightarrow$ correctly	Shift Right
[ ROOT ]	[ ]		
		ROOT $\rightarrow$ parsed	Right

(b)  $O(n)$

The worst case is  $2 \cdot n$ , which is linear.

(e) Final model performance

dev UAS	test UAS
88.32	88.47