

Assignment 3: Dependency Parsing

1. a) Momentum: $m \leftarrow \beta_1 m + (1 - \beta_1) \nabla_{\theta} J(\theta)$
i) $\theta \leftarrow \theta - \alpha m$

Consider the taco shape



If this is a loss landscape, the y direction could have high but oscillating derivatives, but m will converge to the average of the oscillations, which will be much smaller and more stable. This helps to speed up convergence.

ii) Adaptive learning rates:

$$m \leftarrow \beta_1 m + (1 - \beta_1) \nabla_{\theta} J(\theta)$$
$$v \leftarrow \beta_2 v + (1 - \beta_2) (\nabla_{\theta} J(\theta) \odot \nabla_{\theta} J(\theta))$$
$$\theta \leftarrow \theta - \alpha m / \sqrt{v}$$

v is a running average of gradient magnitude per parameter, so when dividing by \sqrt{v} , parameters with historically smaller gradients will get larger updates. This helps to accelerate convergence by preventing params w/ large gradients from exploding and prevents params w/ small gradients from converging very slowly.

b) i) Since each $J_i \sim \text{Bernoulli}(1 - p_{\text{drop}})$,

$$E[h_{\text{drop}}] = \gamma k (1 - p_{\text{drop}}) = k$$

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

ii) Dropout introduces randomness, which is good for training a more generalizable network, but the very spirit of validation/testing is to have a deterministic evaluation of the network. Besides, dropout is a sort of bagging, and the ensemble learners perform better together.

2. a) ROOT $\xrightarrow{\text{I}}$ $\xrightarrow{\text{parsed}}$ $\xrightarrow{\text{this}}$ $\xrightarrow{\text{sentence}}$ $\xrightarrow{\text{correctly}}$

Stack	Buffer	New Dependency	Transition
[ROOT]	[I, parsed, this, sentence, correctly]		Initial Config
[ROOT, I]	[parsed, this, sentence, correctly]		SHIFT
[ROOT, I, parsed]	[this, sentence, correctly]		SHIFT
[ROOT, I, parsed]	[this, sentence, correctly]	parsed \rightarrow I	LEFT-ARC
[ROOT, parsed, this]	[sentence, correctly]	"	SHIFT-ARC
[ROOT, parsed, this, sentence]	[correctly]	"	SHIFT
[ROOT, parsed, sentence]	[correctly]	sentence \rightarrow this	LEFT-ARC
[ROOT, parsed]	[correctly]	parsed \rightarrow sentence	RIGHT-ARC
[ROOT, parsed, correctly]	[]		SHIFT
[ROOT, parsed]	[]	parsed \rightarrow correctly	RIGHT-ARC
[ROOT]	[]	ROOT \rightarrow parsed	RIGHT-ARC

b) Since every word must go from the buffer to the stack $\Rightarrow n$ SHIFTS
and since every word must be a dependent and removed $\Rightarrow n$ ARCS
 $\Rightarrow 2n$ steps.

4) 17. $\log_2(17) = 4.087$

4) i) I disembarked and was heading to a wedding fearing my death.

Error type: Verb phrase attachment error

Incorrect dependency: wedding \rightarrow fearing

Correct dependency: heading \rightarrow fearing

Since "fearing my death" doesn't modify "wedding", the verb phrase modifies the verb "heading" (modifies how I am heading to a wedding).

ii) It makes me want to rush out and rescue people from dilemmas of their own making.

Error type: Modifier attachment error

Incorrect dependency: making \rightarrow their

Correct dependency: own \rightarrow their