

CS224n - Assignment 2

By Dominique Paul

Question 1

(c)

Placeholder Variables:

- Placeholder variables are empty Tensorflow variables which we use to tell Tensorflow that we will replace these values later on. We would not want to use constants in any model to make it useable for other data as well.

Feed Dictionaries:

- Feed Dictionaries are used for training for new variables to be fed in each iteration.

(d)

Automatic Differentiation

- Tensorflows method of building a graph before filling in values allows for a structure that knows how to calculate the backward pass for each mathematical operation as we construct our graph. Therefore, we don't have to take care of calculating the gradients ourselves.

Question 2

(a)

Stack	Buffer	New Dependency	Transition
[ROOT]	[I, parsed, this, sentence, correctly]		Initial configuration
[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, sentence]	[correctly]	Parsed → sentence	RIGHT-ARC

[ROOT, parsed, correctly]	[]		SHIFT
[ROOT, parsed]	[]	Parsed → correctly	RIGHT-ARC
[ROOT]	[]	ROOT → parsed	RIGHT-ARC

(b) A sentence containing n words will be parsed in how many steps (in terms of n)? Briefly explain why.

A sentence of n words will be parsed in $2*n$ steps. Each word has to be shifted at least once, and each word has to be assigned an arc operation to another word at least once.

(f) (2 points, written) We will regularize our network by applying Dropout^[3]. During training this randomly sets units in the hidden layer \mathbf{h} to zero with probability p_{drop} and then multiplies \mathbf{h} by a constant γ (dropping different units each minibatch). We can write this as

$$\mathbf{h}_{drop} = \gamma \mathbf{d} \circ \mathbf{h}$$

where $\mathbf{d} \in \{0, 1\}^{D_h}$ (D_h is the size of \mathbf{h}) is a mask vector where each entry is 0 with probability p_{drop} and 1 with probability $(1 - p_{drop})$. γ is chosen such that the value of \mathbf{h}_{drop} in expectation equals \mathbf{h} :

$$\mathbb{E}_{p_{drop}}[\mathbf{h}_{drop}]_i = \mathbf{h}_i$$

for all $0 < i < D_h$. What must γ equal in terms of p_{drop} ? Briefly justify your answer.

$$\text{Epsilon} = 1 / (1 - p_{drop})$$

It is stated that epsilon is chosen in a manner that the expected value of \mathbf{h}_{drop} equals \mathbf{h} . Therefore if we drop e.g. 90% of the values, then we would have to make the remaining 10% ten times as big to have an expected value of the sum of the values equal to the original sum.

(g) (4 points, written) We will train our model using the Adam^[4] optimizer. Recall that standard SGD uses the update rule

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} J_{minibatch}(\boldsymbol{\theta})$$

where $\boldsymbol{\theta}$ is a vector containing all of the model parameters, J is the loss function, $\nabla_{\boldsymbol{\theta}} J_{minibatch}(\boldsymbol{\theta})$ is the gradient of the loss function with respect to the parameters on a minibatch of data, and α is the learning rate. Adam uses a more sophisticated update rule with two additional steps^[5]

- (i) First, Adam uses a trick called *momentum* by keeping track of \mathbf{m} , a rolling average of the gradients:

$$\begin{aligned}\mathbf{m} &\leftarrow \beta_1 \mathbf{m} + (1 - \beta_1) \nabla_{\boldsymbol{\theta}} J_{\text{minibatch}}(\boldsymbol{\theta}) \\ \boldsymbol{\theta} &\leftarrow \boldsymbol{\theta} - \alpha \mathbf{m}\end{aligned}$$

where β_1 is a hyperparameter between 0 and 1 (often set to 0.9). Briefly explain (you don't need to prove mathematically, just give an intuition) how using \mathbf{m} stops the updates from varying as much. Why might this help with learning?

- (ii) Adam also uses *adaptive learning rates* by keeping track of \mathbf{v} , a rolling average of the magnitudes of the gradients:

$$\begin{aligned}\mathbf{m} &\leftarrow \beta_1 \mathbf{m} + (1 - \beta_1) \nabla_{\boldsymbol{\theta}} J_{\text{minibatch}}(\boldsymbol{\theta}) \\ \mathbf{v} &\leftarrow \beta_2 \mathbf{v} + (1 - \beta_2) (\nabla_{\boldsymbol{\theta}} J_{\text{minibatch}}(\boldsymbol{\theta}) \circ \nabla_{\boldsymbol{\theta}} J_{\text{minibatch}}(\boldsymbol{\theta})) \\ \boldsymbol{\theta} &\leftarrow \boldsymbol{\theta} - \alpha \circ \mathbf{m} / \sqrt{\mathbf{v}}\end{aligned}$$

where \circ and $/$ denote elementwise multiplication and division (so $\mathbf{z} \circ \mathbf{z}$ is elementwise squaring) and β_2 is a hyperparameter between 0 and 1 (often set to 0.99). Since Adam divides the update by $\sqrt{\mathbf{v}}$, which of the model parameters will get larger updates? Why might this help with learning?

(i) The parameter \mathbf{m} holds information on the past weight updates. If the current weight update has been close to zero, but past updates have been quite large, then the new update will nevertheless be of considerable size. The momentum factor \mathbf{m} can help overcome local minima through its consideration of past updates and also ignores some of the noise caused by the batches.

(ii) The factor affects the updates by making those gradients with the on-average larger gradients smaller and those gradients with on-average smaller values bigger. This helps small updates from plateauing and getting stuck.

Question 3

(3c) Cross entropy is defined as

$$CE = - \sum y_u \log(\hat{y}_u)$$

Considering that y is one hot encoded and will only equal 1 once we can also write

$$CE = - \log \hat{y}_i$$

Equally we can write the Perplexity as:

$$PP^{(t)}(y^{(t)}, \hat{y}^{(t)}) = \frac{1}{\hat{y}_i^{(t)}}$$

Using the logarithm rules we can rewrite CE as

$$CE = \log \frac{1}{\hat{y}_i}$$

And therefore we can write CE as using the perplexity

$$CE = \log PP^{(t)}(y^{(t)}, \hat{y}^{(t)})$$

Due to this logarithmic transformation minimizing the arithmetic mean of the cross entropy is equivalent to minimizing the geometric mean of the perplexity

The expected \hat{y}_i value for \hat{y}_i would be $E(\hat{y}_i) = \frac{1}{|V|}$. This corresponds to an expected value of $E(PP) = 1 / \frac{1}{|V|} = |V|$ and thus an expected cross entropy error of $\log(|V|)$, which for $|V| = 10,000$ would equal $\log(10,000) = 9.210$

(313) First we ~~rewrite~~ write some simplifications as follows

$$z = h^{(+1)} H + e^{(+)} I + b_1$$

$$v = h^{(+)} V + b_2$$

Then:

$$\delta_1 = \frac{\partial \mathcal{J}}{\partial v} = y^{(+)} - y^{(+)}$$

$$\delta_2 = \frac{\partial \mathcal{J}}{\partial z} = \delta_1 v^T \circ h \circ (1-h)$$

which we use to answer the original question:

$$\frac{\partial \mathcal{J}}{\partial b_2} = \delta_1$$

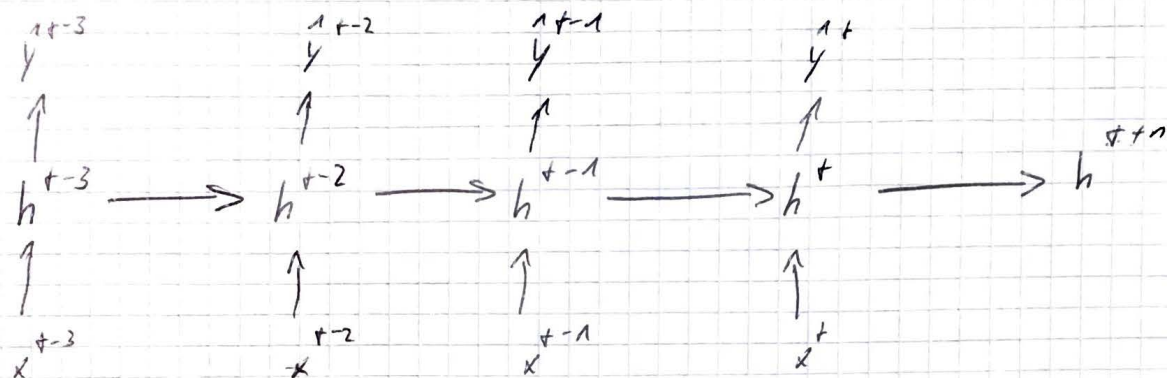
$$\frac{\partial \mathcal{J}}{\partial L_x^{(+)}} = \delta_2^{\bullet} I^T$$

$$\frac{\partial \mathcal{J}}{\partial I} = (e^{(+)})^T \delta_2$$

$$\frac{\partial \mathcal{J}}{\partial H} = h^{(+1)T} \delta_2^{(+)}$$

$$\frac{\partial \mathcal{J}}{\partial h^{+-1}} = \delta_2^{(+)} H^T$$

(3c)



$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}^{t-1}} = \delta^{t-1} \sigma'(\mathbf{v}^{t-1}) \mathbf{I}^T$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{I}} = \mathbf{e}^{t-1T} \delta^{t-1} \sigma'(\mathbf{v}^{t-1})$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}} = \mathbf{h}^{t-2T} \delta^{t-1} \sigma'(\mathbf{v}^{t-1})$$

where $\sigma'(\mathbf{v}^{t-1}) = \text{diag}(\mathbf{h}^{t-1} \circ (1 - \mathbf{h}^{t-1}))$

(3d)

Forward propagation: $O(|V| D_h + d D_h + D_h^2)$

Backward propagation: $O(\tau(|V| D_h + d D_h + D_h^2))$

→ Back-prop takes τ times as long as forward propagation

→ if $|V| \gg D_h$ then $|V| D_h$ takes very long as well