# Assignment 2

## 1 Tensorflow Softmax

**(a)** See file q1_softmax.py

**(b)** See file q1_softmax.py

**(c)** The purpose of placeholder is to hold our input data, and the purpose of feed dictionaries is to feed input values to placeholders. See implementation in file q1_classifier.py

**(d)** See file q1_classifier.py

**(e)** See file q1_classifier.py. After the model's *train_op* is called, the prediction *y_hat* is computed in the forward propagation and the gradient of loss with respect to $W$ and $b$ is computed during the back propagation. The variable $W$ and $b$ will be changed.

## 2 Neural Transition-Based Dependency Parsing

**(a)** The parsing procedure is as follows:

| 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] | [correctly] | parsed→sentence | RIGHT-ARC |
| [ROOT, parsed, correctly] | [] | | SHIFT |
| [ROOT, parsed] | [] | parsed→correctly | RIGHT-ARC |
| [ROOT] | [] | ROOT→parsed | RIGHT-ARC |

Table. 1: Parsing Procedure

**(b)** A sentence containing $n$ words will be parsed in $2 \times n$ steps. This is true because every word will be shifted into the stack once, and will be removed from the stack once, therefore, the total number of steps is $2 \times n$.

**(c)** See file q2_parser_transitions.py

**(d)** See file q2_parser_transitions.py

**(e)** See file q2_initialization.py

**(f)** The constant $\gamma$ can be expressed as:

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

This is true because the expected value of $h_{drop}$ is $(1 - p_{drop})h$.

**(g)**

(i) By using $m$, the amount of steps we take at each update will now become the running average of gradients of all times. And the current gradient calculated only contribute a little to the $m$, therefore it is more stable and can stop the updates from varying too much.

(ii) The parameters that have smaller gradient will get larger updates. This helps the learning because it will smooth the updates we make and try to update all parameters equally.

**(h)** See file q2_parser_model.py