

# Assignment #4

Due Time: Mid-night at December 27<sup>th</sup> 2017

## 1 word2vec (70 points)

- (a) (5 points) Assume you are given a predicted word vector  $v_c$  corresponding to the center word  $c$  for skipgram, and word prediction is made with the softmax function found in word2vec models

$$\hat{y}_o = p(o|c) = \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w=1}^W \exp(\mathbf{u}_w^T \mathbf{v}_c)} \quad (1)$$

where  $w$  denotes the  $w$ -th word and  $\mathbf{u}_w (w = 1, \dots, W)$  are the “output” word vectors for all words in the vocabulary. Assume cross entropy cost is applied to this prediction and word  $o$  is the expected word (the  $o$ -th element of the one-hot label vector is one), derive the gradients with respect to  $\mathbf{v}_c$ .

*Hint: It will be helpful to use simple notation. For instance, letting  $\hat{\mathbf{y}}$  be the vector of softmax predictions for every word,  $\mathbf{y}$  as the expected word vector, and the loss function*

$$J_{\text{softmax-CE}}(\mathbf{o}, \mathbf{v}_c, \mathbf{U}) = \text{CE}(\mathbf{y}, \hat{\mathbf{y}}) \quad (2)$$

Where  $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_W]$  is the matrix of all the output vectors. *Make sure you state the orientation of your vectors and matrices.*

- (b) (5 points) As in the previous part, derive gradients for the “output” word vectors  $\mathbf{u}_w$ ’s (including  $\mathbf{u}_o$ ).
- (c) (10 points) Repeat part (a) and (b) assuming we are using the negative sampling loss for the predicted vector  $\mathbf{v}_c$ , and the expected output word is  $o$ . Assume that  $K$  negative samples (words) are drawn, and they are  $\mathbf{1}, \dots, \mathbf{K}$ , respectively for simplicity of notation ( $o \notin \{1, \dots, K\}$ ). Again, for a given word,  $o$ , denote its output vector as  $\mathbf{u}_o$ . The negative sampling loss function in this case is

$$J_{\text{neg-sample}}(\mathbf{o}, \mathbf{v}_c, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \quad (3)$$

where  $\sigma(\cdot)$  is the sigmoid function.

After you’ve done this, describe with one sentence why this cost function is much more efficient to compute than the softmax-CE loss (you could provide a speed-up ratio, i.e. the runtime of the softmax-CE loss divided by the runtime of the negative sampling loss).

*Note: the cost function here is the negative of what Mikolov et al had in their original paper, because we are doing a minimization instead of maximization in our code.*

- (d) (10 points) Derive gradients for all of the word vectors for skip-gram given the previous parts and given a set of context words  $[\text{word}_{c-m}, \dots, \text{word}_{c-1}, \text{word}_c, \text{word}_{c+1}, \dots, \text{word}_{c+m}]$ , where  $m$  is the context size. Denote the “input” and “output” word vectors for  $\text{word}_k$  as  $\mathbf{v}_k$  and  $\mathbf{u}_k$  respectively.

*Hint: feel free to use  $F(\mathbf{o}, \mathbf{v}_c)$  (where  $o$  is the expected word) as a placeholder for the  $J_{\text{softmax-CE}}(\mathbf{o}, \mathbf{v}_c, \dots)$  or  $J_{\text{neg-sample}}(\mathbf{o}, \mathbf{v}_c, \dots)$  cost functions in this part — you’ll see that this is a useful abstraction for the coding part. That is, your solution may contain terms of the form  $\frac{\partial F(\mathbf{o}, \mathbf{v}_c)}{\partial \dots}$ .*

Recall that for skip-gram, the cost for a context centered around  $c$  is

$$J_{\text{skip-gram}}(\text{word}_{c-m \dots c+m}) = \sum_{-m \leq j \leq m, j \neq 0} F(\mathbf{w}_{c+j}, \mathbf{v}_c) \quad (4)$$

where  $\mathbf{w}_{c+j}$  refers to the word at the  $j$ -th index from the center.

(e) (30 points) In this part you will implement the word2vec models and train your own word vectors with stochastic gradient descent (SGD). First, write a helper function to normalize rows of a matrix in `word2vec.py`. In the same file, fill in the implementation for the softmax and negative sampling cost and gradient functions. Then, fill in the implementation of the cost and gradient functions for the skipgram model. When you are done, test your implementation by running `python word2vec.py`.

(f) (10 points) Show time! Now we are going to load some real data and train word vectors with everything you just implemented! We are going to use the Stanford Sentiment Treebank (SST) dataset to train word vectors, and later apply them to a simple sentiment analysis task. There is no additional code to write for this part; just run `python run.py`.

*Note: The training process may take a long time depending on the efficiency of your implementation (**an efficient implementation takes approximately an hour**). Plan accordingly!*

When the script finishes, a visualization for your word vectors will appear. It will also be saved as `word_vectors.png` in your project directory. **Include the plot in your homework write up.**

## 2 Sentiment Analysis (30 points)

Now, with the word vectors you trained, we are going to perform a simple sentiment analysis. For each sentence in the Stanford Sentiment Treebank dataset, we are going to use the average of all the word vectors in that sentence as its feature, and try to predict the sentiment level of the said sentence. The sentiment level of the phrases are represented as real values in the original dataset, here we'll just use five classes:

“very negative”, “negative”, “neutral”, “positive”, “very positive”

which are represented by 0 to 4 in the code, respectively. For this part, you will learn to train a softmax regressor with SGD, and perform train/dev validation to improve generalization of your regressor.

(a) (10 points) Implement a sentence featurizer. Fill in the implementation in `softmaxreg.py`. Sanity check your implementation with `python softmaxreg.py`.

(b) (15 points) Fill in the hyperparameter selection code in `sentiment.py` to search for the “optimal” regularization parameter. **What value did you select? Report your train, dev, and test accuracies. Justify your hyperparameter search methodology in at most one sentence.**

*Note: you should be able to attain at least 25% accuracy on dev.*

(c) (5 points) Plot the classification accuracy on the train and dev set with respect to the regularization value, using a logarithmic scale on the x-axis. This should have been done automatically. **Include `reg_acc.png` in your homework write up.**