

# CS 224d: Assignment #1

**Due date:** 4/19 11:59 PM PST (You are allowed to use three (3) late days maximum for this assignment)

These questions require thought, but do not require long answers. Please be as concise as possible.

We encourage students to discuss in groups for assignments. However, each student must finish the problem set and programming assignment individually, and must turn in her/his assignment. We ask that you abide by the university Honor Code and that of the Computer Science department, and make sure that all of your submitted work is done by yourself

Please review any additional instructions posted on the assignment page at <http://cs224d.stanford.edu/assignments.html>. When you are ready to submit, please follow the instructions on the course website.

## 1 Softmax (10 points)

- (a) (5 points) Prove that softmax is invariant to constant offsets in the input, that is, for any input vector  $\mathbf{x}$  and any constant  $c$ ,

$$\text{softmax}(\mathbf{x}) = \text{softmax}(\mathbf{x} + c)$$

where  $\mathbf{x} + c$  means adding the constant  $c$  to every dimension of  $\mathbf{x}$ . Remember that

$$\text{softmax}(\mathbf{x})_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

Note: In practice, we make use of this property and choose  $c = -\max_i x_i$  when computing softmax probabilities for numerical stability (i.e. subtracting its maximum element from all elements of  $\mathbf{x}$ ). *Handwritten note: This can help compute the softmax with small range differences.*

- (b) (5 points) Given an input matrix of  $N$  rows and  $d$  columns, compute the softmax prediction for each row. Write your implementation in `q1_softmax.py`. You may test by executing `python q1_softmax.py`.

Note: The provided tests are not exhaustive. Later parts of the assignment will reference this code so it is important to have a correct implementation. Your implementation should also be efficient and vectorized whenever possible. A non-vectorized implementation will not receive full credit!

## 2 Neural Network Basics (30 points)

- (a) (3 points) Derive the gradients of the sigmoid function and show that it can be rewritten as a function of the function value (i.e. in some expression where only  $\sigma(x)$ , but not  $x$ , is present). Assume that the input  $x$  is a scalar for this question. Recall, the sigmoid function is

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad \frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x)) = \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{-1}{(1 + e^{-x})^2} \times \frac{\partial e^{-x}}{\partial x} \quad (2)$$

- (b) (3 points) Derive the gradient with regard to the inputs of a softmax function when cross entropy loss is used for evaluation, i.e. find the gradients with respect to the softmax input vector  $\theta$ , when the prediction is made by  $\hat{\mathbf{y}} = \text{softmax}(\theta)$ . Remember the cross entropy function is

$$CE(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_i y_i \log(\hat{y}_i) \quad (3)$$



$$\hat{y} = \text{Saturate}(\theta)$$

$$\mathcal{L}(y, \hat{y}) = -\sum_i y_i \log \hat{y}_i$$

$$y_i = \frac{e^{\theta_i}}{\sum_j e^{\theta_j}}$$

$$\frac{\partial \mathcal{L}(E)}{\partial \theta_k} = -y_i (1 - \hat{y}_k) \quad i=k$$

$$= y_i \hat{y}_k \quad i \neq k$$

$$\frac{\partial \mathcal{L}(E)}{\partial \theta_k} = \frac{\partial \mathcal{L}(E)}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial \theta_k}$$

$$= -\frac{y_i}{\hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial \theta_k}$$

$\hat{y}_i \rightarrow$   
ith  
element  
of  $y$

when  $i=k$

$$\frac{\partial \hat{y}_i}{\partial \theta_i} = \frac{\partial}{\partial \theta_i} \left( \frac{e^{\theta_i}}{\sum_j e^{\theta_j}} \right)$$

$$= \frac{e^{\theta_i}}{\sum_j e^{\theta_j}} - \frac{e^{\theta_i} \cdot e^{\theta_i}}{(\sum_j e^{\theta_j})^2}$$

$$= \hat{y}_i - \hat{y}_i^2$$

$$= \hat{y}_i (1 - \hat{y}_i)$$

when  $i \neq k$

$$\frac{\partial \hat{y}_i}{\partial \theta_k} = \frac{\partial}{\partial \theta_k} \left( \frac{e^{\theta_i}}{\sum_j e^{\theta_j}} \right)$$

$$= \frac{-e^{\theta_i} \cdot e^{\theta_k}}{(\sum_j e^{\theta_j})^2} = -\hat{y}_i \hat{y}_k$$

$i=k \quad \& \quad y_i=1$

$$\frac{\partial \mathcal{L}(E)}{\partial \theta_j} = (\hat{y}_j - 1) \quad i=j$$

$$= \hat{y}_j$$

$$\frac{\partial \mathcal{L}(E)}{\partial \theta} = \hat{y} - y$$

$$\delta_1 = \frac{\partial \mathcal{L}(E)}{\partial \theta}$$

$$= \frac{\partial \mathcal{L}(E)}{\partial z_2} = \hat{y} - y$$

$$\delta_2 = \frac{\partial \mathcal{L}(E)}{\partial h} \cdot \frac{\partial \hat{z}_2}{\partial h}$$

$$= \delta_1 \cdot W_2^T$$

$$\delta_3 = \frac{\partial \mathcal{L}(E)}{\partial z_1} \cdot \frac{\partial h}{\partial h}$$

$$= \delta_2 \cdot \sigma'(z_1)$$

$$\frac{\partial \mathcal{L}(E)}{\partial x} = \frac{\partial \mathcal{L}(E)}{\partial z_1} \frac{\partial z_1}{\partial x}$$

$$= \underline{\underline{\delta_3 W_1^T}}$$



$$CE_i = -y_i \log \hat{y}_i \quad \frac{\partial CE_i}{\partial \hat{y}_i} = -\frac{y_i}{\hat{y}_i} \quad \frac{\partial \hat{y}_i}{\partial z_2} \quad \left. \begin{array}{l} \text{derived above} \\ \frac{\partial CE}{\partial z_2} = \hat{y} - y \end{array} \right\}$$

where  $y$  is the one-hot label vector, and  $\hat{y}$  is the predicted probability vector for all classes. (Hint: you might want to consider the fact many elements of  $y$  are zeros, and assume that only the  $k$ -th dimension of  $y$  is one.)

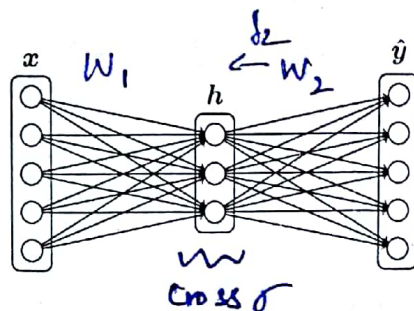
- (c) (6 points) Derive the gradients with respect to the inputs  $x$  to an one-hidden-layer neural network (that is, find  $\frac{\partial J}{\partial x}$  where  $J$  is the cost function for the neural network). The neural network employs sigmoid activation function for the hidden layer, and softmax for the output layer. Assume the one-hot label vector is  $y$ , and cross entropy cost is used. (feel free to use  $\sigma'(x)$  as the shorthand for sigmoid gradient, and feel free to define any variables whenever you see fit)

$$z_1 = xW_1 + b_1$$

$$h_1 = \sigma(z_1)$$

$$z_2 = h_1W_2 + b_2$$

$$\hat{y} = \text{softmax}(z_2)$$



Recall that the forward propagation is as follows

$$\text{grad } z_1 = \frac{\partial CE}{\partial z_1} = \text{grad } h_1 \cdot \text{sigmoid-grad}(h_1)$$

$$\hat{y} = \text{softmax}(hW_2 + b_2)$$

Note that here we're assuming that the input vector (thus the hidden variables and output probabilities) is a row vector to be consistent with the programming assignment. When we apply the sigmoid function to a vector, we are applying it to each of the elements of that vector.  $W_i$  and  $b_i$  ( $i = 1, 2$ ) are the weights and biases, respectively, of the two layers.

$$\text{grad } W_1 = \frac{\partial CE}{\partial W_1} = \frac{\partial CE}{\partial z_1} \cdot \frac{\partial z_1}{\partial W_1}$$

$$\text{grad } b_1 = \sum \text{grad } z_1$$

- (d) (2 points) How many parameters are there in this neural network, assuming the input is  $D_x$ -dimensional, the output is  $D_y$ -dimensional, and there are  $H$  hidden units?  $(D_x + 1)H + (H + 1) \cdot D_y$
- (e) (4 points) Fill in the implementation for the sigmoid activation function and its gradient in `q2.sigmoid.py`. Test your implementation using python `q2.sigmoid.py`. Again, thoroughly test your code as the provided tests may not be exhaustive.
- (f) (4 points) To make debugging easier, we will now implement a gradient checker. Fill in the implementation for `gradcheck_naive` in `q2.gradcheck.py`. Test your code using python `q2.gradcheck.py`.
- (g) (8 points) Now, implement the forward and backward passes for a neural network with one sigmoid hidden layer. Fill in your implementation in `q2.neural.py`. Sanity check your implementation with `q2.neural.py`.

### 3 word2vec (40 points + 5 bonus)

- (a) (3 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

$$\hat{y}_o = p(o | c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)} \quad (4)$$

where  $w$  denotes the  $w$ -th word and  $u_w$  ( $w = 1, \dots, W$ ) 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  $v_c$ .



$$J_i = -y_i \log \hat{y}_i \quad ; \quad \hat{y}_i = \frac{e^{u_i^T v_c}}{\sum_{w=1}^W e^{u_w^T v_c}}$$

$$\frac{\partial J_i}{\partial v_c} = \frac{\partial}{\partial v_c} \left( -u_i^T v_c + \log \sum_{w=1}^W e^{u_w^T v_c} \right)$$

$$= -u_i + \frac{\sum_{w=1}^W e^{u_w^T v_c} \cdot u_w}{\sum_{w=1}^W e^{u_w^T v_c}}$$

$$= -u_i + \sum_{w=1}^W \hat{y}_w u_w$$

$$\boxed{U(\hat{y} - y)}$$

Vectorized?!

$$\frac{\partial J_i}{\partial u_w} = \frac{\partial}{\partial u_w} \left( -u_i^T v_c + \log \sum_{w=1}^W e^{u_w^T v_c} \right)$$

w is not in context of c

$$= \frac{\partial}{\partial u_w} \left( \log \sum_{w=1}^W e^{u_w^T v_c} \right) = \frac{1}{\sum_{w=1}^W e^{u_w^T v_c}} \cdot e^{u_w^T v_c} \cdot v_c = \underline{\underline{\hat{y}_w v_c}}$$

w is in context of c

$$= \frac{\partial}{\partial u_w} \left( -u_i^T v_c + \log \sum_{w=1}^W e^{u_w^T v_c} \right)$$

$$= -v_c + \frac{1}{\sum_{w=1}^W e^{u_w^T v_c}} e^{u_w^T v_c} \cdot v_c = -v_c + \hat{y}_w v_c$$

$$= \underline{\underline{(\hat{y}_w - 1) v_c}}$$

$$\boxed{v_c (\hat{y} - y)^T}$$

*Hint: It will be helpful to use notation from question 2. 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*

$$J_{\text{softmax-CE}}(\mathbf{o}, \mathbf{v}_c, \mathbf{U}) = \text{CE}(\mathbf{y}, \hat{\mathbf{y}}) \quad (5)$$

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) (3 points) As in the previous part, derive gradients for the “output” word vectors  $\mathbf{u}_w$ ’s (including  $\mathbf{u}_o$ ).
- (c) (6 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  $\mathbf{o}$ . Assume that  $K$  negative samples (words) are drawn, and they are  $1, \dots, K$ , respectively for simplicity of notation ( $o \notin \{1, \dots, K\}$ ). Again, for a given word,  $\mathbf{o}$ , denote its output vector as  $\mathbf{u}_o$ . The negative sampling loss function in this case is

$$J_{\text{neg-sample}}(\mathbf{o}, \mathbf{v}_c, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \quad (6)$$

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) (8 points) Derive gradients for all of the word vectors for skip-gram and CBOW 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  $\mathbf{o}$  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

$$J_{\text{skip-gram}}(\text{word}_{c-m \dots c+m}) = \sum_{-m \leq j \leq m, j \neq 0} F(\mathbf{w}_{c+j}, \mathbf{v}_c) \quad (7)$$

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

CBOW is slightly different. Instead of using  $\mathbf{v}_c$  as the predicted vector, we use  $\hat{\mathbf{v}}$  defined below. For (a simpler variant of) CBOW, we sum up the input word vectors in the context

$$\hat{\mathbf{v}} = \sum_{-m \leq j \leq m, j \neq 0} \mathbf{v}_{c+j} \quad (8)$$

then the CBOW cost is

$$J_{\text{CBOW}}(\text{word}_{c-m \dots c+m}) = F(\mathbf{w}_c, \hat{\mathbf{v}}) \quad (9)$$

*Note: To be consistent with the  $\hat{\mathbf{v}}$  notation such as for the code portion, for skip-gram  $\hat{\mathbf{v}} = \mathbf{v}_c$ .*

- (e) (12 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 `q3_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 skip-gram model. When you are done, test your implementation by running `python q3_word2vec.py`.  
*Note: If you choose not to implement CBOW (part h), simply remove the `NotImplementedError` so that your tests will complete.*



$$J = -\log \sigma(u_0^T v_c) - \sum_{k=1}^K \log \sigma(-u_k^T v_c)$$

$$\frac{\partial J}{\partial v_c} = -\frac{1 \cdot \sigma'(u_0^T v_c) \cdot u_0}{\sigma(u_0^T v_c)} - \sum_{k=1}^K \frac{1 \cdot \sigma'(-u_k^T v_c) u_k}{\sigma(-u_k^T v_c)}$$

$$\begin{aligned} \sigma'(x) &= \sigma(x)(1 - \sigma(x)) \\ &= (\sigma(u_0^T v_c) - 1) u_0 - \sum_{k=1}^K (\sigma(-u_k^T v_c) - 1) u_k \end{aligned}$$

$u_0 \rightarrow$  where  $o$  is in the o/p word

$$\frac{\partial J}{\partial u_0} = \frac{\sigma(u_0^T v_c - 1) v_c}{\text{where } k \text{ is other sampling word'}}$$

$$\frac{\partial J}{\partial u_k} = -(\sigma(-u_k^T v_c) - 1) v_c \quad k = 1, 2, 3, \dots, K$$

$$\frac{\partial J_{\text{skip-gram}}(\text{word } c-m \dots c+m)}{\partial u} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(w_{c+j}, v_c)}{\partial u}$$

$$\frac{\partial J_{\text{skip-gram}}(\text{word } c-m \dots c+m)}{\partial v_c} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(w_{c+j}, v_c)}{\partial v_c}$$

$$\frac{\partial J_{\text{skip-gram}}(\text{word } c-m \dots c+m)}{\partial v_j} = 0 \quad \text{for all } j \neq c$$

(f) (4 points) Complete the implementation for your SGD optimizer in `q3_sgd.py`. Test your implementation by running `python q3_sgd.py`.

(g) (4 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 q3_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 `q3_word_vectors.png` in your project directory. **Include the plot in your homework write up.** Briefly explain in at most three sentences what you see in the plot.

(h) Extra credit (5 points) Implement the CBOW model in `q3_word2vec.py`. *Note: This part is optional but the gradient derivations for CBOW in part (d) are not!*

## 4 Sentiment Analysis (20 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 and softmax regression. Fill in the implementation in `q4_softmaxreg.py`. Sanity check your implementation with `python q4_softmaxreg.py`.

(b) (2 points) Explain in fewer than three sentences why we want to introduce regularization when doing classification (in fact, most machine learning tasks).

(c) (4 points) Fill in the hyperparameter selection code in `q4_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 30% accuracy on dev.*

(d) (4 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 `q4_reg_acc.png` in your homework write up.** Briefly explain in at most three sentences what you see in the plot.



$$\frac{\partial J_{\text{bow}}(\text{word}_{c-m \dots c+m})}{\partial U} = \frac{\partial F(w_c, \hat{U})}{\partial U}$$

$$\frac{\partial J_{\text{bow}}(\text{word}_{c-m \dots c+m})}{\partial U_j} = \frac{\partial F(w_c, \hat{U})}{\partial \hat{U}} \quad j \in \{c-m, \dots, c+m\}$$

$$\frac{\partial J_{\text{bow}}(\text{word}_{c-m \dots c+m})}{\partial U_j} = 0 \quad j \notin \{c-m, \dots, c+m\}$$