

Solution to Assignment 1, Problem 1(a)

Akash Rana

Softmax Prove that *softmax* is invariant to constant offsets in the input, that is, for any input vector x and any constant c ,

$$\text{softmax}(\mathbf{x}) = \text{softmax}(\mathbf{x} + c), \quad (1)$$

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

$$\text{softmax}(\mathbf{x})_j = \frac{e^{x_j}}{\sum_{k=1} e^{x_k}} \quad (2)$$

Solution:

$$\text{softmax}(\mathbf{x} + c)_j = \frac{e^{(x_j+c)}}{\sum_{k=1} e^{(x_k+c)}} \quad (3)$$

$$= \frac{e^c}{e^c} \frac{e^{(x_j)}}{\sum_{k=1} e^{(x_k)}} \quad (4)$$

$$= \text{softmax}(\mathbf{x})_j \quad (5)$$

Solution to Assignment 1, Problem 1(b)

Akash Rana

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!

```
import numpy as np

def softmax(x):
    #x.shape by default gives column value
    if len(x.shape) > 1:
        # Matrix
        ### YOUR CODE HERE
        c = np.max(x, axis=1).reshape(-1,1) #-1 here means, internally
        #numpy is just calculating, to get the missing dimension.
        x = np.exp(x-c) / np.sum(np.exp(x-c), axis=1).reshape(-1,1)
        print (np.shape(x))
        ### END YOUR CODE
    else:
        # Vector
        ### YOUR CODE HERE
        c = np.max(x)
        x = np.exp(x-c) / np.sum(np.exp(x-c))
        ### END YOUR CODE

    assert x.shape == orig_shape
    return x
```

Solution to Assignment 1, Problem 2(a)*Akash Rana*

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.

Denote the sigmoid function as $\sigma(z)$,

$$\sigma(z) = \frac{1}{1 + e^{-z}}, \tag{6}$$

Using chain rule,

$$\begin{aligned} \sigma'(z) &= \frac{-1}{(1 + e^{-z})^2} \times (-e^{-z}) \\ &= \frac{1}{1 + e^{-z}} \left(\frac{e^{-z}}{1 + e^{-z}} \right) \\ &= \left(\frac{1}{1 + e^{-z}} \right) \left(\frac{1 + e^{-z}}{1 + e^{-z}} - \frac{1}{1 + e^{-z}} \right) \\ &= \sigma(z)(1 - \sigma(z)) \end{aligned}$$

Solution to Assignment 1, Problem 2(b)

Akash Rana

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 θ , when the prediction is made by $\hat{y} = \text{softmax}(\theta)$. Remember the cross entropy function is

$$\text{CE}(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_i y_i \log \hat{y}_i \quad (7)$$

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.

Cross entropy error function for multi-class output,

$$\text{CE}(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_i y_i \log \hat{y}_i \quad (8)$$

Computing the gradient yields,

$$\frac{\partial(\text{CE})}{\partial \hat{y}_i} = - \frac{y_j}{\hat{y}_i} \quad (9)$$

$$\frac{\partial(\text{CE})}{\partial \theta_k} = \frac{\partial(\text{CE})}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial \theta_k} \quad (10)$$

$$= - \frac{y_j}{\hat{y}_i} \frac{\partial \hat{y}_i}{\partial \theta_k} \quad (11)$$

Calculating the partial derivative of \hat{y}_i (derivation using the quotient rule):

$$\text{if } i = k : \frac{\partial y_i}{\partial \theta_i} = \frac{\partial e^{\theta_i}}{\partial \theta_i} \quad (12)$$

$$= \frac{e^{\theta_i} \Sigma_\theta - e^{\theta_i} e^{\theta_i}}{\Sigma_\theta^2} \quad (13)$$

$$= \frac{e^{\theta_i}}{\Sigma_\theta} \frac{\Sigma_\theta - e^{\theta_i}}{\Sigma_\theta} \quad (14)$$

$$= \frac{e^{\theta_i}}{\Sigma_\theta} \left(1 - \frac{e^{\theta_i}}{\Sigma_\theta}\right) = \hat{y}_i(1 - \hat{y}_i) \quad (15)$$

$$\text{if } i \neq k : \frac{\partial y_i}{\partial \theta_j} = \frac{\partial e^{\frac{\theta_i}{\Sigma_\theta}}}{\partial \theta_j} \quad (16)$$

$$= \frac{0 - e^{\theta_i} e^{\theta_j}}{\Sigma_\theta^2} \quad (17)$$

$$= - \frac{e^{\theta_i} e^{\theta_j}}{\Sigma_\theta \Sigma_\theta} \quad (18)$$

$$= - \hat{y}_i y_k \quad (19)$$

Combining Equations 9, 15, 19, yields

$$\frac{\partial(\text{CE})}{\partial \theta_k} = \begin{cases} -y_j(1 - \hat{y}_k) & \text{for } i = k \\ y_j \hat{y}_k & \text{for } i \neq k \end{cases} \quad (20)$$

Requiring y_j to be non-zero, imposes that the auxiliary condition, $k = j$ and $y_j = 1$, hence it follows immediately,

$$\frac{\partial(\text{CE})}{\partial \theta_j} = \begin{cases} (\hat{y}_j - 1) & \text{for } i = j \\ \hat{y}_j & \text{for } i \neq j \end{cases} \quad (21)$$

Which is equivalent to

$$\frac{\partial(\text{CE})}{\partial \boldsymbol{\theta}} = \hat{\mathbf{y}} - \mathbf{y} \quad (22)$$

Solution to Assignment 1, Problem 2(c)

Akash Rana

Derive the gradients with respect to the inputs \mathbf{x} to an one-hidden-layer neural network (that is, find $\partial J / \partial \mathbf{x}$ 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 \mathbf{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).

In order to simplify the notation used to solve the problem, define the following terms:

$$\mathbf{x}^{(2)} \equiv \mathbf{h} \quad (23)$$

$$\mathbf{z}_i \equiv \mathbf{x}^{(i)} \mathbf{W}_i + \mathbf{b}_i \quad (24)$$

Now, to calculate $\partial J / \partial \mathbf{x}^1$, one can use the back propagation algorithm. Starting with the chain rule and then results from Question 2(b):

$$\frac{\partial J}{\partial \mathbf{x}^{(2)}} = \left(\left(\frac{\partial J}{\partial \mathbf{z}_2} \frac{\partial \mathbf{z}_2}{\partial \mathbf{x}^{(2)}} \right) \cdot \frac{\partial \mathbf{x}^{(2)}}{\partial \mathbf{z}_1} \right) * \frac{\partial \mathbf{z}_1}{\partial \mathbf{x}^{(1)}} \quad (25)$$

$$\frac{\partial J}{\partial \mathbf{z}_2} = \hat{\mathbf{y}} - \mathbf{y} \quad (26)$$

and

$$\frac{\partial \mathbf{z}_i}{\partial \mathbf{x}^{(i)}} = \mathbf{W}_i^\top \quad (27)$$

Sigmoid (σ) derivative can be found in Question 2(a), but we define:

$$\frac{\partial \mathbf{x}^{(2)}}{\partial \mathbf{z}_1} \equiv \sigma'(\mathbf{z}_1) \quad (28)$$

Combining these, and using \cdot to denote element-wise product:

$$\frac{\partial J}{\partial \mathbf{z}_1} = (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{W}_2^\top \cdot \sigma'(\mathbf{z}_1) \quad (29)$$

Finally, using the results from Equation 27 (but for the first layer):

$$\frac{\partial J}{\partial \mathbf{x}^{(1)}} = (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{W}_2^\top \cdot \sigma'(\mathbf{z}_1) \cdot \mathbf{W}_1^\top \quad (30)$$

Solution to Assignment 1, Problem 2(d)*Akash Rana*

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?

W_1 must have dimensions: $D_x \times H$. The bias (\mathbf{b}_1) for the first layer must have dimensions H . Adding these two together, yields $(D_x + 1) \times H$. Proceeding to the second layer, there must be $H \times D_y$ parameters associated with the weight matrix W_2 . The bias (\mathbf{b}_2) for the second layer must have dimensions D_y elements. This yields,

$$(D_x + 1) \times H + D_y \times (H + 1) \tag{31}$$

weights, for each input vector of dimensions D_x .

Solution to Assignment 1, Problem 3(a)

Akash Rana

Assume you are given a predicted word vector \mathbf{v}_c corresponding to the center word c for `skipgram`, and word prediction is made with the `softmax` function found in `word2vec` models

$$\hat{\mathbf{y}}_o = p(\mathbf{o} \mid \mathbf{c}) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{j=1}^{|W|} \exp(\mathbf{u}_j^\top \mathbf{v}_c)} \quad (32)$$

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 (v'_w in the lecture notes). Assume the 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 notation from question 2. For instance, letting $\hat{\mathbf{y}}$ be the vector of the softmax prediction for every word, \mathbf{y} as the expected word vector, and the loss function

$$J_{\text{softmax-CE}}(\mathbf{o}, \mathbf{v}_c, \mathbf{U}) = CE(\mathbf{y}, \hat{\mathbf{y}}) \quad (33)$$

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.*

Applying cross-entropy cost to the above softmax probability defined above:

$$J = -\log p = -\mathbf{u}_o^\top \mathbf{v}_c + \log \sum_{j=1}^{|V|} \exp(\mathbf{u}_j^\top \mathbf{v}_c) \quad (34)$$

Let $z_j = \mathbf{u}_j^\top \mathbf{v}_c$, and δ_j^i be the indicator function, then

$$\frac{\partial J}{\partial z_k} = -\delta_k^i + \frac{\exp(\mathbf{u}_i^\top \mathbf{v}_c)}{\sum_{j=1}^{|V|} \exp(\mathbf{u}_j^\top \mathbf{v}_c)} \quad (35)$$

Now, using the chain rule, we can calculate,

$$\frac{\partial J}{\partial \mathbf{v}_c} = \frac{\partial J}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{v}_c} \quad (36)$$

$$= \sum_{j=1}^{|V|} \mathbf{u}_j^\top \left(\frac{e^{z_j}}{\sum_{k=1}^{|V|} e^{z_k}} - 1 \right) \quad (37)$$

$$= \sum_{k=1}^{|V|} P(\mathbf{u}_j | \mathbf{v}_c) \mathbf{u}_j - \mathbf{u}_j \quad (38)$$

Solution to Assignment 1, Problem 3(b)*Akash Rana*

As in the previous problem, derive gradients for the ‘output’ word vectors \mathbf{u}_w ’s (including \mathbf{u}_o)

This follows immediately from the chain rule:

$$\frac{\partial J}{\partial \mathbf{u}_j} = \frac{\partial J}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{u}_j} \tag{39}$$

$$= \mathbf{v}_c \left(\frac{\exp(\mathbf{u}_0^\top \mathbf{v}_c)}{\sum_{j=1}^{|V|} \exp(\mathbf{u}_j^\top \mathbf{v}_c)} - \delta_j^0 \right) \tag{40}$$

Solution to Assignment 1, Problem 3(c)

Akash Rana

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 $\mathbf{o}(\mathbf{u}_o)$. Assume that K negative samples (words) are drawn, and they are $\mathbf{u}_1, \dots, \mathbf{u}_k$, respectively for simplicity of notation ($k \in \{1, \dots, K\}$ and $o \notin \{1, \dots, K\}$). Again for a given word, \mathbf{o} , denote its output vector as \mathbf{u}_o . The negative sampling loss function in this case is,

$$J(\mathbf{u}_o, \mathbf{v}_c, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \quad (41)$$

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.

$$\frac{\partial J}{\partial \mathbf{v}_c} = \frac{\partial J}{\partial \mathbf{p}_i} \frac{\partial \mathbf{p}_i}{\partial \mathbf{x}_j} \frac{\partial \mathbf{x}_i}{\partial \mathbf{v}_c} + \sum_{k=1}^K \frac{\partial J}{\partial \mathbf{p}_i} \frac{\partial \mathbf{p}_i}{\partial \mathbf{z}_k} \frac{\partial \mathbf{z}_k}{\partial \mathbf{v}_c} \quad (42)$$

$$= -(\mathbf{1} - \mathbf{p}_i) \mathbf{U}_i + \sum_{k=1}^K (\mathbf{1} - \mathbf{q}_k) \mathbf{U}_k \quad (43)$$

$$\frac{\partial J}{\partial \mathbf{u}_j} = \frac{\partial J}{\partial \mathbf{p}_i} \frac{\partial \mathbf{p}_i}{\partial \mathbf{x}_j} \frac{\partial \mathbf{x}_i}{\partial \mathbf{u}_j} + \sum_{k=1}^K \frac{\partial J}{\partial \mathbf{p}_i} \frac{\partial \mathbf{p}_i}{\partial \mathbf{z}_k} \frac{\partial \mathbf{z}_k}{\partial \mathbf{u}_j} \quad (44)$$

$$= -(\mathbf{1} - \mathbf{p}_i) \delta_{ij} \mathbf{V}_c + \sum_{k=1}^K (\mathbf{1} - \mathbf{q}_k) \delta_{ik} \mathbf{V}_c \quad (45)$$

where, δ_{xy} = kronecker delta
 when $x = y = 1$
 0, otherwise.

Solution to Assignment 1, Problem 3(d)

Akash Rana

Derive the gradients for all of the word vectors for **skip-gram** and **CBOW** (optional) 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. You can denote the ‘input’ and ‘output’ word vectors for word_k as \mathbf{v}_k and \mathbf{u}_k respectively for convenience.

Hint: feel free to use $F(\mathbf{o}, \mathbf{v}_c)$ (where \mathbf{o} is the expected word) as a placeholder for $J_{\text{softmax-CE}}(\mathbf{o}, \mathbf{v}_c, \dots)$ or $J_{\text{neg-sampling}}(\mathbf{o}, \mathbf{v}_c, \dots)$ 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 $\frac{\partial F_i(\mathbf{o}, \mathbf{v}_c)}{\partial \lambda}$.

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

$$J_{\text{skip-gram}} = \sum_{-m \leq j \leq m, j \neq 0} F(\mathbf{w}_{c+j}, \mathbf{v}_c) \quad (46)$$

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 prediction vector, we use $\hat{\mathbf{v}}$ defined below. For (a simpler variant of) **CBOW**, we sum up the input word vectors in the context:

$$\hat{\mathbf{v}} = \sum_{-m \leq j \leq m, j \neq 0} \mathbf{v}_{c+j} \quad (47)$$

then the **CBOW** cost is

$$J_{\text{CBOW}}(\text{word}_{c-m}, \dots, \text{word}_{c+m}) = F(\mathbf{w}_c, \hat{\mathbf{v}}) \quad (48)$$

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

For **skip-gram**:

$$\frac{\partial J_{\text{skip-gram}}}{\partial U} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(\mathbf{w}_{c+j}, \mathbf{v}_c)}{\partial U} \quad (49)$$

$$\frac{\partial J_{\text{skip-gram}}}{\partial v_c} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(\mathbf{w}_{c+j}, \mathbf{v}_c)}{\partial v_c} \quad (50)$$

$$\frac{\partial J_{\text{skip-gram}}}{\partial v_i} = 0, \text{ For } i \neq c \quad (51)$$

For **CBOW**:

$$\frac{\partial J_{\text{CBOW}}}{\partial U} = \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial U} \quad (52)$$

$$\frac{\partial J_{\text{CBOW}}}{\partial v_j} = \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial v} \quad (53)$$

$$\frac{\partial J_{\text{CBOW}}}{\partial v_j} = 0, \text{ For } i \neq c \quad (54)$$