

Assignment 2 Written Solutions

- (a) (2 points) Prove that the naive-softmax loss (Equation 2) is the same as the cross-entropy loss between \mathbf{y} and $\hat{\mathbf{y}}$, i.e. (note that \mathbf{y} (true distribution), $\hat{\mathbf{y}}$ (predicted distribution) are vectors and \hat{y}_o is a scalar):

$$-\sum_{w \in \text{Vocab}} \mathbf{y}_w \log(\hat{\mathbf{y}}_w) = -\log(\hat{y}_o). \quad (3)$$

Your answer should be one line. You may describe your answer in words.

$$\text{Cross-Entropy}(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{i=1}^{\text{Vocab}} y_i \log(\hat{y}_i) = -(y_1 \log \hat{y}_1 + \dots + y_o \log \hat{y}_o + \dots + y_{\text{vocab}} \log \hat{y}_{\text{vocab}})$$

But $y_1 = y_2 = \dots = y_{\text{vocab}} = 0 \neq y_o = 1$. So,

$$\text{Cross-Entropy}(\mathbf{y}, \hat{\mathbf{y}}) = -\log \hat{y}_o = -\log(P(O = o | C = c)) = J_{\text{Naive-Softmax}}(\mathbf{v}_c, o, \mathbf{U})$$

- (b) (6 points)

- (i) Compute the partial derivative of $J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ with respect to \mathbf{v}_c . Please write your answer in terms of \mathbf{y} , $\hat{\mathbf{y}}$, \mathbf{U} , and show your work to receive full credit.

- **Note:** Your final answers for the partial derivative should follow the shape convention: the partial derivative of any function $f(x)$ with respect to x should have the **same shape** as x .⁴
- Please provide your answers for the partial derivative in vectorized form. For example, when we ask you to write your answers in terms of \mathbf{y} , $\hat{\mathbf{y}}$, and \mathbf{U} , you may not refer to specific elements of these terms in your final answer (such as y_1, y_2, \dots).

- (ii) When is the gradient you computed equal to zero?

Hint: You may wish to review and use some introductory linear algebra concepts.

- (iii) The gradient you found is the difference between the two terms. Provide an interpretation of how each of these terms improves the word vector when this gradient is subtracted from the word vector \mathbf{v}_c .

$$\begin{aligned} \text{(i)} \quad J_{\text{Naive-Softmax}}(\mathbf{v}_c, o, \mathbf{U}) &= -\log(P(O = o | C = c)) = -\log\left(\frac{e^{u_o^T \mathbf{v}_c}}{\sum_{w \in \text{vocab}} e^{u_w^T \mathbf{v}_c}}\right) \\ &= -\log(e^{u_o^T \mathbf{v}_c}) + \log(\sum_{w \in \text{vocab}} e^{u_w^T \mathbf{v}_c}) = -u_o^T \mathbf{v}_c + \log(\sum_{w \in \text{vocab}} e^{u_w^T \mathbf{v}_c}) \\ \frac{\partial J_{\text{Naive-Softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} &= -u_o^T + \frac{\sum_{w \in \text{vocab}} e^{u_w^T \mathbf{v}_c} u_w^T}{\sum_{w \in \text{vocab}} e^{u_w^T \mathbf{v}_c}} \\ \text{But } \hat{\mathbf{y}}_w &= \log(P(O = o | C = c)) = \frac{e^{u_w^T \mathbf{v}_c}}{\sum_{w \in \text{vocab}} e^{u_w^T \mathbf{v}_c}}. \text{ So,} \\ \frac{\partial J_{\text{Naive-Softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} &= -u_o^T + \sum_{w \in \text{vocab}} \hat{\mathbf{y}}_w u_w^T \end{aligned}$$

We can rewrite u_o^T as $u_o^T = y_1 u_1 + y_2 u_2 + \dots + y_o u_o + \dots y_{vocab} u_{vocab} = \sum_{w \in vocab} y_w u_w^T$

$$\text{as } y = \begin{bmatrix} y_1 = 0 \\ y_2 = 0 \\ y_o = 1 \\ \vdots \\ y_{vocab} = 0 \end{bmatrix}. \text{ In conclusion,}$$

$$\frac{\partial J_{Naive-Softmax}(v_c, o, U)}{\partial v_c} = -\sum_{w \in vocab} y_w u_w^T + \sum_{w \in vocab} \hat{y}_w u_w^T = \sum_{w \in vocab} (-y_w + \hat{y}_w) u_w^T$$

Vectorizing this result yields:

$$\frac{\partial J_{Naive-Softmax}(v_c, o, U)}{\partial v_c} = U^T (\hat{y} - y)$$

(ii) When $\hat{y} = y$, which means that the predicted distribution is equal to the true distribution.

$$(iii) \quad v_c^{New} = v_c^{old} - \alpha \frac{\partial J_{Naive-Softmax}(v_c, o, U)}{\partial v_c} = v_c^{old} - \alpha U(y - \hat{y}) = v_c^{old} - \alpha U y + \alpha U \hat{y}$$

v_c^{old} heads towards the direction of the true distribution and away from the direction of the predicted distribution.

(c) (1 point) In many downstream applications using word embeddings, L2 normalized vectors (e.g. $\mathbf{u}/\|\mathbf{u}\|_2$ where $\|\mathbf{u}\|_2 = \sqrt{\sum_i u_i^2}$) are used instead of their raw forms (e.g. \mathbf{u}). Let's consider a hypothetical downstream task of binary classification of phrases as being positive or negative, where you decide the sign based on the sum of individual embeddings of the words. When would L2 normalization take away useful information for the downstream task? When would it not?

Hint: Consider the case where $\mathbf{u}_x = \alpha \mathbf{u}_y$ for some words $x \neq y$ and some scalar α . Give examples of words x and y which satisfy the given relation and why would they be affected/not affected due to the normalization?

Mostly, normalization is useful and doesn't take away any useful information. However, when we consider a case like the given one, e.g., $u_x = \alpha u_y$, there will be a problem in normalizing the two vectors. They will have the same orientation but different scales.

An example of that is the words "chair" and "chairs". While they may have similar direction and different scales, the task could be to distinguish the plural nouns from the singular ones. Normalization will only keep orientation information while removing the scale information which will cause the 2 words to be the same and the task will fail.

(d) (5 points) Compute the partial derivatives of $J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ with respect to each of the ‘outside’ word vectors, \mathbf{u}_w ’s. There will be two cases: when $w = o$, the true ‘outside’ word vector, and $w \neq o$, for all other words. Please write your answer in terms of \mathbf{y} , $\hat{\mathbf{y}}$, and \mathbf{v}_c . In this subpart, you may use specific elements within these terms as well (such as y_1, y_2, \dots). Note that \mathbf{u}_w is a vector while y_1, y_2, \dots are scalars. Show your work to receive full credit.

$$J_{\text{Naive-Softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\mathbf{u}_o^T \mathbf{v}_c + \log\left(\sum_{w \in \text{vocab}} e^{\mathbf{u}_w^T \mathbf{v}_c}\right)$$

$$\frac{\partial J_{\text{Naive-Softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_w} = -\frac{\partial}{\partial \mathbf{u}_w} (\mathbf{u}_o^T \mathbf{v}_c) + \frac{e^{\mathbf{u}_w^T \mathbf{v}_c} \mathbf{v}_c}{\sum_{w \in \text{vocab}} e^{\mathbf{u}_w^T \mathbf{v}_c}} = -\frac{\partial}{\partial \mathbf{u}_w} (\mathbf{u}_o^T \mathbf{v}_c) + \hat{\mathbf{y}}_w \mathbf{v}_c$$

$$\text{When } u_o = u_w: \frac{\partial}{\partial u_w} (\mathbf{u}_o^T \mathbf{v}_c) = \frac{\partial}{\partial u_w} (\mathbf{u}_w^T \mathbf{v}_c) = \mathbf{v}_c$$

$$\text{When } u_o \neq u_w: \frac{\partial}{\partial u_w} (\mathbf{u}_o^T \mathbf{v}_c) = 0$$

We can say in general that: $\frac{\partial}{\partial u_w} (\mathbf{u}_o^T \mathbf{v}_c) = y_w \mathbf{v}_c$ where $y_w = \begin{cases} 0, & w \neq o \\ 1, & w = o \end{cases}$

$$\frac{\partial J_{\text{Naive-Softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_w} = -y_w \mathbf{v}_c + \hat{\mathbf{y}}_w \mathbf{v}_c = (\hat{\mathbf{y}}_w - y_w) \mathbf{v}_c$$

(e) (1 point) Write down the partial derivative of $J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ with respect to \mathbf{U} . Please break down your answer in terms of the column vectors $\frac{\partial J(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_1}, \frac{\partial J(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_2}, \dots, \frac{\partial J(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{|\text{Vocab}|}}$. No derivations are necessary, just an answer in the form of a matrix.

$$\frac{\partial J_{\text{Naive-Softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{U}} = \left[\frac{\partial J_{\text{Naive-Softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_1}, \dots, \frac{\partial J_{\text{Naive-Softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{\text{vocab}}} \right]$$

$$\text{But from (d): } \frac{\partial J_{\text{Naive-Softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_w} = (\hat{\mathbf{y}}_w - y_w) \mathbf{v}_c$$

$$\frac{\partial J_{\text{Naive-Softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{U}} = [(\hat{\mathbf{y}}_1 - y_1) \mathbf{v}_c, \dots, (\hat{\mathbf{y}}_{\text{vocab}} - y_{\text{vocab}}) \mathbf{v}_c] = (\hat{\mathbf{y}} - \mathbf{y})^T \mathbf{v}_c$$

- (f) (2 points) The Leaky ReLU (Leaky Rectified Linear Unit) activation function is given by Equation 4 and Figure 2:

$$f(x) = \max(\alpha x, x) \quad (4)$$

Where x is a scalar and $0 < \alpha < 1$, please compute the derivative of $f(x)$ with respect to x . You may ignore the case where the derivative is not defined at 0.

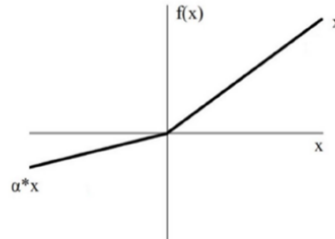


Figure 2: Leaky ReLU

$$f(x) = \begin{cases} \alpha x, & x < 0 \\ x, & x \geq 0 \end{cases}$$

$$\frac{\partial f(x)}{\partial x} = \begin{cases} \alpha, & x < 0 \\ 1, & x \geq 0 \end{cases}$$

- (g) (3 points) The sigmoid function is given by Equation 5:

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} \quad (5)$$

Please compute the derivative of $\sigma(x)$ with respect to x , where x is a scalar. Please write your answer in terms of $\sigma(x)$. Show your work to receive full credit.

$$\sigma(x) = \frac{1}{1 + e^{-x}} = (1 + e^{-x})^{-1}$$

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

- (h) (6 points) Now we shall consider the Negative Sampling loss, which is an alternative to the Naive Softmax loss. Assume that K negative samples (words) are drawn from the vocabulary. For simplicity of notation we shall refer to them as w_1, w_2, \dots, w_K , and their outside vectors as $\mathbf{u}_{w_1}, \mathbf{u}_{w_2}, \dots, \mathbf{u}_{w_K}$.⁶ For this question, assume that the K negative samples are distinct. In other words, $i \neq j$ implies $w_i \neq w_j$ for $i, j \in \{1, \dots, K\}$. Note that $o \notin \{w_1, \dots, w_K\}$. For a center word c and an outside word o , the negative sampling loss function is given by:

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

for a sample w_1, \dots, w_K , where $\sigma(\cdot)$ is the sigmoid function.⁷

- (i) Please repeat parts (b) and (d), computing the partial derivatives of $J_{\text{neg-sample}}$ with respect to \mathbf{v}_c , with respect to \mathbf{u}_o , and with respect to the s^{th} negative sample \mathbf{u}_{w_s} . Please write your answers in terms of the vectors \mathbf{v}_c , \mathbf{u}_o , and \mathbf{u}_{w_s} , where $s \in [1, K]$. Show your work to receive full credit. **Note:** you should be able to use your solution to part (g) to help compute the necessary gradients here.
- (ii) In the lecture, we learned that an efficient implementation of backpropagation leverages the reuse of previously computed partial derivatives. Which quantity could you reuse (we store the common terms in a single matrix/vector) amongst the three partial derivatives calculated above to minimize duplicate computation?

Write your answer in terms of

$U_{o, \{w_1, \dots, w_K\}} = [\mathbf{u}_o, -\mathbf{u}_{w_1}, \dots, -\mathbf{u}_{w_K}]$, a matrix with the outside vectors stacked as columns, and $\mathbf{1}$, a $(K+1) \times 1$ vector of 1's.⁸ Additional terms and functions (other than $U_{o, \{w_1, \dots, w_K\}}$ and $\mathbf{1}$) can be used in your solution.

- (iii) Describe with one sentence why this loss function is much more efficient to compute than the naive-softmax loss.

(i)

$$J_{\text{neg-sample}}(\mathbf{v}_c, o, U) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c))$$

$$\frac{\partial J_{\text{neg-sample}}(\mathbf{v}_c, o, U)}{\partial \mathbf{v}_c} = -\frac{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c))}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} \mathbf{u}_o - \sum_{s=1}^K \frac{\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)(1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c))}{\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)} (-\mathbf{u}_{w_s})$$

$$= -(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)) \mathbf{u}_o + \sum_{s=1}^K (1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \mathbf{u}_{w_s}$$

$$\frac{\partial J_{\text{neg-sample}}(\mathbf{v}_c, o, U)}{\partial \mathbf{u}_o} = -\frac{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c))}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} \mathbf{v}_c - 0 = (1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)) \mathbf{v}_c$$

$$\frac{\partial J_{\text{neg-sample}}(\mathbf{v}_c, o, U)}{\partial \mathbf{u}_{w_s}} = 0 - \frac{\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)(1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c))}{\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)} (-\mathbf{v}_c) = (1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \mathbf{v}_c$$

(ii) We can store $\begin{bmatrix} 1 - \sigma(u_o^T v_c) \\ 1 - \sigma(u_1^T v_c) \\ 1 - \sigma(u_2^T v_c) \\ \vdots \\ 1 - \sigma(u_{w_K}^T v_c) \end{bmatrix}$ as these quantities are calculated in each derivative.

(iii) It computes the gradients using only K samples instead of using the entire vocabulary. It saves both computation time and memory.

(i) (2 points) Now we will repeat the previous exercise, but without the assumption that the K sampled words are distinct. Assume that K negative samples (words) are drawn from the vocabulary. For simplicity of notation we shall refer to them as w_1, w_2, \dots, w_K and their outside vectors as u_{w_1}, \dots, u_{w_K} . In this question, you may not assume that the words are distinct. In other words, $w_i = w_j$ may be true when $i \neq j$ is true. Note that $o \notin \{w_1, \dots, w_K\}$. For a center word c and an outside word o , the negative sampling loss function is given by:

$$J_{\text{neg-sample}}(v_c, o, U) = -\log(\sigma(u_o^T v_c)) - \sum_{s=1}^K \log(\sigma(-u_{w_s}^T v_c)) \quad (7)$$

for a sample w_1, \dots, w_K , where $\sigma(\cdot)$ is the sigmoid function.

Compute the partial derivative of $J_{\text{neg-sample}}$ with respect to a negative sample u_{w_s} . Please write your answers in terms of the vectors v_c and u_{w_s} , where $s \in [1, K]$. Show your work to receive full credit. Hint: break up the sum in the loss function into two sums: a sum over all sampled words equal to w_s and a sum over all sampled words not equal to w_s . Notation-wise, you may write ‘equal’ and ‘not equal’ conditions below the summation symbols, such as in Equation 8.

$$\begin{aligned} J_{\text{neg-sample}}(v_c, o, U) &= -\log(\sigma(u_o^T v_c)) - \sum_{s=1}^K \log(\sigma(-u_{w_s}^T v_c)) \\ &= -\log(\sigma(u_o^T v_c)) - \sum_{\substack{i=1 \\ w_i=w_s}}^{K_1} \log(\sigma(-u_{w_i}^T v_c)) - \sum_{\substack{j=1 \\ w_j \neq w_s}}^{K-K_1} \log(\sigma(-u_{w_j}^T v_c)) \end{aligned}$$

$$\begin{aligned} \frac{\partial J_{\text{neg-sample}}(v_c, o, U)}{\partial u_{w_s}} &= 0 - \sum_{\substack{i=1 \\ w_i=w_s}}^{K_1} \frac{\sigma(-u_{w_s}^T v_c)(1 - \sigma(-u_{w_s}^T v_c))}{\sigma(-u_{w_s}^T v_c)} (-v_c) - 0 \\ &= \sum_{\substack{i=1 \\ w_i=w_s}}^{K_1} (1 - \sigma(-u_{w_s}^T v_c)) v_c = K_1 (1 - \sigma(-u_{w_s}^T v_c)) v_c \end{aligned}$$

- (j) (3 points) Suppose the center word is $c = w_t$ and the context window is $[w_{t-m}, \dots, w_{t-1}, w_t, w_{t+1}, \dots, w_{t+m}]$, where m is the context window size. Recall that for the skip-gram version of word2vec, the total loss for the context window is:

$$J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} J(\mathbf{v}_c, w_{t+j}, \mathbf{U}) \quad (8)$$

Here, $J(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ represents an arbitrary loss term for the center word $c = w_t$ and outside word w_{t+j} . $J(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ could be $J_{\text{naive-softmax}}(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ or $J_{\text{neg-sample}}(\mathbf{v}_c, w_{t+j}, \mathbf{U})$, depending on your implementation.

Write down three partial derivatives:

- (i) $\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{U}}$
- (ii) $\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{v}_c}$
- (iii) $\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{v}_w}$ when $w \neq c$

Write your answers in terms of $\frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{U}}$ and $\frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{v}_c}$. This is very simple – each solution should be one line.

(i)

$$\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{U}} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{U}}$$

(ii)

$$\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{v}_c} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{v}_c}$$

(iii)

$$\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{v}_w} = 0$$