

1 ASSIGNMENT 2

- (a) (3 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 $\mathbf{y}, \hat{\mathbf{y}}$ are vectors and $\hat{\mathbf{y}}_o$ is a scalar):

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

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

As we know that \mathbf{y} is a one hot vector, so

$$- \sum_{w \in \text{Vocab}} \mathbf{y}_w \log(\hat{\mathbf{y}}_w) = -1 \times \mathbf{y}_o \log(\hat{\mathbf{y}}_o) = -\log(\hat{\mathbf{y}}_o)$$

- (b) (5 points) Compute the partial derivative of $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ with respect to \mathbf{v}_c . Please write your answer in terms of \mathbf{y} , $\hat{\mathbf{y}}$, and \mathbf{U} . Additionally, answer the following two questions with one sentence each: (1) When is the gradient zero? (2) Why does subtracting this gradient, in the general case when it is nonzero, make \mathbf{v}_c a more desirable vector (namely, a vector closer to outside word vectors in its window)?

- **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 $\mathbf{y}_1, \mathbf{y}_2, \dots$).

(1)

$$\begin{aligned} J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) &= -\log P(o|c) \\ &= -\log \frac{\exp(u_o^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(u_w^T \mathbf{v}_c)} \\ &= \log \sum \exp(u_w^T \mathbf{v}_c) - u_o^T \mathbf{v}_c \\ &\quad \text{Therefore} \\ \frac{\partial J}{\partial \mathbf{v}_c} &= \frac{1}{\sum \exp(u_w^T \mathbf{v}_c)} \frac{\partial \sum \exp(u_w^T \mathbf{v}_c)}{\partial \mathbf{v}_c} - u_o \\ &= \frac{1}{\sum \exp(u_w^T \mathbf{v}_c)} \sum \exp(u_w^T \mathbf{v}_c) \mathbf{u}_w - u_o \\ &= \sum P(w|c) \mathbf{u}_w - u_o \\ &\quad \text{When } \partial J \text{ is 0, we have} \\ u_o &= \sum P(w|c) \mathbf{u}_w \end{aligned}$$

(2)

That is because when we are doing this gradient descent, $P(o|c)$ tends to go up, which means outside words are closer to the center word

- (c) (5 points) Compute the partial derivatives of $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ 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 $\mathbf{y}_1, \mathbf{y}_2, \dots$). Note that \mathbf{u}_w is a vector while $\mathbf{y}_1, \mathbf{y}_2, \dots$ are scalars.

$$\begin{aligned}\frac{\partial J}{\partial u_w} &= \frac{1}{\sum \exp(u_w^T v_c)} \frac{\partial \sum \exp(u_w^T v_c)}{\partial u_w} - \frac{\partial u_o^T v_c}{\partial u_w} \\ &= \frac{1}{\sum \exp(u_w^T v_c)} \exp(u_w^T v_c) v_c - \frac{\partial u_o^T v_c}{\partial u_w} \\ &= P(w|c) v_c - \frac{\partial u_o^T v_c}{\partial u_w}\end{aligned}$$

$$\text{When } w = o, \text{ then } y_o = 1, \frac{\partial u_o^T v_c}{\partial u_w} = v_c$$

$$\frac{\partial J}{\partial u_o} = (\hat{y}_o - 1) v_c = (\hat{y}_o - y_o) v_c$$

$$\text{When } w \neq o, \text{ then } y_w = 0, \frac{\partial u_o^T v_c}{\partial u_w} = 0$$

$$\frac{\partial J}{\partial u_w} = (\hat{y}_w - 0) v_c = (\hat{y}_w - y_w) v_c$$

- (d) (1 point) Write down the partial derivative of $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ with respect to \mathbf{U} . Please break down your answer in terms of $\frac{\partial \mathbf{J}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_1}$, $\frac{\partial \mathbf{J}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_2}$, \dots , $\frac{\partial \mathbf{J}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{|\text{Vocab}|}}$. The solution should be one or two lines long.

$$\begin{aligned}\frac{\partial J}{\partial u_i} &= (\hat{y}_i - y_i) v_c \\ i &= 1, 2, \dots, |\text{Vocab}|\end{aligned}$$

- (e) (2 points) The ReLU (Rectified Linear Unit) activation function is given by Equation 4:

$$f(x) = \max(0, x) \quad (4)$$

Please compute the derivative of $f(x)$ with respect to x , where x is a scalar. You may ignore the case that the derivative is not defined at 0.⁵

$$\text{When } x < 0$$

$$\frac{df}{dx} = 0$$

$$\text{When } x > 0$$

$$\frac{df}{dx} = 1$$

- (f) (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. Hint: you may want to write your answer in terms of $\sigma(x)$.

$$\begin{aligned}
\frac{d\sigma}{dx} &= \frac{d \frac{e^x}{e^x+1}}{dx} \\
&= \frac{e^x(e^x+1) - e^x \times e^x}{(e^x+1)^2} \\
&= \frac{e^x}{(e^x+1)^2} \\
&= \sigma(x)(1-\sigma(x))
\end{aligned}$$

$w_i \neq w_j$ 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:

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

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

- (i) Please repeat parts (b) and (c), computing the partial derivatives of $\mathbf{J}_{\text{neg-sample}}$ 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 $s \in [1, K]$. **Note:** you should be able to use your solution to part (f) to help compute the necessary gradients here.
- (ii) In lecture, we learned that an efficient implementation of backpropagation leverages the re-use of previously-computed partial derivatives. Which quantity could you reuse between the three partial derivatives to minimize duplicate computation? Write your answer in terms of $\mathbf{U}_{o, \{w_1, \dots, w_K\}} = [\mathbf{u}_o, -\mathbf{u}_{w_1}, \dots, -\mathbf{u}_{w_K}]$, a matrix with the outside vectors stacked as columns, and $\mathbf{1}$, a $(K+1) \times 1$ vector of 1's.⁸
- (iii) Describe with one sentence why this loss function is much more efficient to compute than the naive-softmax loss.

(i)

$$\begin{aligned}
\frac{\partial J}{\partial v_c} &= -\frac{1}{\sigma(u_o^T v_c)} \sigma(u_o^T v_c)(1 - \sigma(u_o^T v_c))u_o - \sum_{s=1}^K \frac{1}{\sigma(-u_{w_s}^T v_c)} \sigma(-u_{w_s}^T v_c)(1 - \sigma(-u_{w_s}^T v_c))(-u_{w_s}) \\
&= -(1 - \sigma(u_o^T v_c))u_o + \sum_{s=1}^K (1 - \sigma(-u_{w_s}^T v_c))u_{w_s} \\
\frac{\partial J}{\partial u_o} &= -(1 - \sigma(u_o^T v_c))v_c \\
\frac{\partial J}{\partial u_{w_s}} &= (1 - \sigma(-u_{w_s}^T v_c))v_c
\end{aligned}$$

(ii)

$$(1 - \sigma(u_o^T v_c))$$

(iii)

This loss function does not calculate all the losses from every word in vocabulary.

And it also eliminates the use of logarithm

- (h) (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 $\mathbf{u}_{w_1}, \dots, \mathbf{u}_{w_K}$. 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:

$$\mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \quad (7)$$

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

Compute the partial derivative of $\mathbf{J}_{\text{neg-sample}}$ with respect to a negative sample \mathbf{u}_{w_s} . Please write your answers in terms of the vectors \mathbf{v}_c and \mathbf{u}_{w_s} , where $s \in [1, K]$. 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.

$$\frac{\partial J}{\partial \mathbf{u}_{w_s}} = \sum_{i \wedge (w_s = w_i)}^K (1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \mathbf{v}_c$$

- (i) (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:

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

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

Write down three partial derivatives:

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

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

Once you’re done: Given that you computed the derivatives of $\mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ with respect to all the model parameters \mathbf{U} and \mathbf{V} in parts (a) to (c), you have now computed the derivatives of the full loss function $\mathbf{J}_{\text{skip-gram}}$ with respect to all parameters. You’re ready to implement word2vec!

(i)

$$\sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{U}}$$

(ii)

$$\sum_{-m \leq j \leq m, j \neq 0} (1 - \sigma(\mathbf{u}_{w_j}^\top \mathbf{v}_{w_0})) \mathbf{u}_{w_0} + \sum_{s=1}^K (1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_{w_0})) \mathbf{u}_{w_s}$$

(iii)

