

$$P(O = o \mid C = c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (2) \quad J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\log P(O = o \mid C = c). \quad (1)$$

- (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 $\mathbf{y}, \hat{\mathbf{y}}$ are vectors and \hat{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.

$$(2) \sum_{w \in \text{Vocab}} y_w \log(y_w) = \sum_{w \in \text{Vocab}} y_w \log(\hat{y}_w) = \sum_{w \in \text{Vocab}} \log y_w = \log \hat{y}_o$$

$y_w = 1 \quad w \in O$
 $y_w = 0 \quad w \notin O$
 \mathbf{y} is one-hot encoder

$\hat{y}_w = 1$
 $w \in O$

(b) (7 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).

$$(b)(i) \frac{\partial J}{\partial \mathbf{v}_c} = \frac{\partial}{\partial \mathbf{v}_c} -\log P(O=o \mid C=c) = \frac{\partial}{\partial \mathbf{v}_c} -\log \frac{e^{\mathbf{u}_o^\top \mathbf{v}_c}}{\sum_{w \in V} e^{\mathbf{u}_w^\top \mathbf{v}_c}} =$$

$$\frac{\partial}{\partial \mathbf{v}_c} \left(-\log e^{\mathbf{u}_o^\top \mathbf{v}_c} - \log \sum_{w \in V} e^{\mathbf{u}_w^\top \mathbf{v}_c} \right) = \frac{\partial}{\partial \mathbf{v}_c} \left(\log \sum_{w \in V} e^{\mathbf{u}_w^\top \mathbf{v}_c} - \mathbf{u}_o^\top \mathbf{v}_c \right) =$$

$$\frac{\sum_{w \in V} \mathbf{u}_w \cdot e^{\mathbf{u}_w^\top \mathbf{v}_c}}{\sum_{w \in V} e^{\mathbf{u}_w^\top \mathbf{v}_c}} - \mathbf{u}_o^\top = \sum_{w \in V} \mathbf{u}_w P(O=w \mid C=c) - \mathbf{u}_o^\top =$$

$$\mathbf{U}^\top \hat{\mathbf{y}} - \mathbf{u}_o^\top = \mathbf{U}^\top (\hat{\mathbf{y}} - \mathbf{y})$$

$y_o = 0 \quad w \notin O$
 $y_o = 1 \quad w \in O$

(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 two terms. Provide an interpretation of how each of these terms improves the word vector when this gradient is subtracted from the word vector v_c .

(iv) 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}). Now, suppose you would like to classify phrases as being positive or negative. 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 .

(ii) when there are no words that are embedded or when it classifies perfectly the outside words and there are no other words that are predicted for the center word.

When $\hat{y} = y$

(iii) the first term is only the outside words that are related to the center word, adding it will help the gradient to keep the relevant words with high score. The second term is the probability for each word in the vocabulary will be an outside word for the center word, subtracting it will push away the non outside words.

(iv) $U_x = \alpha U_y$ When L2 normalization equals $U_x = U_y$

$$\frac{\alpha U_x}{\sqrt{\sum \alpha^2 U_x^2}} = \frac{\alpha U_x}{\alpha \sqrt{\sum U_x^2}} = \frac{U_x}{\sqrt{\sum U_x^2}}$$

(c) (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.

$$(c) \frac{\partial J}{\partial \mathbf{u}_w} = \frac{\partial}{\partial \mathbf{u}_w} \left(\log \sum_{w \in V} e^{\mathbf{u}_w^T \mathbf{v}_c} - \mathbf{u}_o^T \mathbf{v}_c \right) =$$

$$\frac{\sum_{w \in V} \mathbf{v}_c \cdot e^{\mathbf{u}_w^T \mathbf{v}_c}}{\sum_{w \in V} e^{\mathbf{u}_w^T \mathbf{v}_c}} \cdot \frac{\partial \mathbf{u}_w^T \mathbf{v}_c}{\partial \mathbf{u}_w} = \sum_{w \in V} \mathbf{v}_c \cdot P(o=w|c=c) - \frac{\partial \mathbf{u}_o^T \mathbf{v}_c}{\partial \mathbf{u}_w} =$$

$$V_c \hat{y}_w - V_c \frac{\partial U_0}{\partial U_w} = V_c \left(\hat{y}_w - \frac{\partial U_0}{\partial U_w} \right) =$$

$$V_c (\hat{y}_w - 1) \quad \text{when } w=0$$

$$V_c \hat{y}_w \quad \text{when } w \neq 0$$

$$= V_c (\hat{y}_w - y_w)$$

y is the hot encoder

- (d) (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.

$$\frac{\partial J}{\partial U_w} = V_c (\hat{y}_w - y_w) \quad \text{so} \quad \frac{\partial J}{\partial \mathbf{U}} = V_c (\hat{\mathbf{y}} - \mathbf{y}) =$$

$$\left[\frac{\partial J}{\partial U_1}, \frac{\partial J}{\partial U_2}, \dots, \frac{\partial J}{\partial U_w} \right]$$

- (e) (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 0.⁵

$$\frac{\partial \max(\alpha x, x)}{\partial x} = \frac{\partial \begin{cases} x & \text{when } x > 0 \\ \alpha x & \text{when } x < 0 \end{cases}}{\partial x} = \begin{cases} 1 & \text{when } x > 0 \\ \alpha & \text{when } x < 0 \end{cases}$$

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

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.

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

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

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

- (g) (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:

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

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 $J_{\text{neg-sample}}$ 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 (f) to help compute the necessary gradients here.

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{v}_c} &= \frac{\partial -\log \sigma(\mathbf{u}_o^\top \mathbf{v}_c)}{\partial \mathbf{v}_c} - \frac{\partial \sum_{s=1}^K \log \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)}{\partial \mathbf{v}_c} \\ &= \frac{\sigma(\mathbf{u}_o^\top \mathbf{v}_c) \cdot (1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c))}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} - \sum_{s=1}^K \frac{\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c) (1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c))}{\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)} \\ &= (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1) \mathbf{u}_o + \sum_{s=1}^K [1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)] \mathbf{u}_{w_s} \end{aligned}$$

$$\frac{\partial J}{\partial U_o} = \frac{\partial -\log \sigma(U_o^T V_c)}{\partial U_o} - \frac{\partial \sum_{s=1}^K \log \sigma(-U_{w_s}^T V_c)}{\partial U_o} =$$

$$-\frac{\sigma(U_o^T V_c)(1 - \sigma(U_o^T V_c)) \cdot V_c}{\sigma(U_o^T V_c)} = V_c (\sigma(U_o^T V_c) - 1)$$

$$\frac{\partial J}{\partial U_{w_s}} = \frac{\partial -\log \sigma(U_o^T V_c)}{\partial U_{w_s}} - \frac{\partial \sum_{s=1}^K \log \sigma(-U_{w_s}^T V_c)}{\partial U_{w_s}} =$$

$$\sum_{s=1}^K \frac{\sigma(-U_{w_s} V_c) \cdot (1 - \sigma(-U_{w_s} V_c)) \cdot V_c}{\sigma(-U_{w_s} V_c)} =$$

$$\sum_{s=1}^K V_c \cdot [1 - \sigma(-U_{w_s} V_c)]$$

(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 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.

$$\sigma(U_o^T V_c) - 1, \quad 1 - \sigma(-U_{w_s} V_c)$$

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

(iii) it takes small number of k instead all of the vocabulary.

- (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 $s \in [1, K]$. 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.

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

$$= \frac{\partial}{\partial \mathbf{u}_{w_s}} \sum_{w_j = w_s} \log \sigma(-\mathbf{u}_{w_j}^\top \mathbf{v}_c) + \sum_{w_j \neq w_s} \log \sigma(-\mathbf{u}_{w_j}^\top \mathbf{v}_c) =$$

$$= \sum_{w_j = w_s} [\sigma(-\mathbf{u}_{w_j}^\top \mathbf{v}_c) - 1] \mathbf{v}_c$$

- (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!

$$(i) \frac{\partial \mathbf{J}_{\text{sg}}}{\partial \mathbf{U}} = \sum_{-m \leq j \leq m} \frac{\partial \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{U}}$$

$$(ii) \frac{\partial \mathbf{J}_{\text{sg}}}{\partial \mathbf{v}_c} = \sum_{-m \leq j \leq m} \frac{\partial \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{v}_c}$$

$$(iii) \frac{\partial \mathbf{J}_{\text{sg}}}{\partial \mathbf{v}_w} = 0 \quad \text{when } w \neq c$$

(c) (2 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. You will need to fetch the datasets first. To do this, run `sh get_datasets.sh`. 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 and the compute power of your machine (an efficient implementation takes one to two hours). Plan accordingly!

After 40,000 iterations, the script will finish and 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.** In at most three sentences, briefly explain what you see in the plot. This may include, but is not limited to, observations on clusters and words that you expect to cluster but do not.

We can easily see that familiar words are clustered together female, woman..
we can also see that there are similar vectors between similar compntext
king queen and male female.



