

CS224n Assignment 2: Understanding word2vec

David Lee

July 16, 2019

The given equations:

- the softmax function

$$P(O = o|C = c) = \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \quad (1)$$

- the loss function

$$\mathbf{J}_{\text{naive-softmax}}(v_c, o, U) = -\log P(O = o|C = c) \quad (2)$$

- the sigmoid function

$$\sigma(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{x}}} = \frac{e^{\mathbf{x}}}{e^{\mathbf{x}} + 1} \quad (3)$$

(a) Show that the naive-softmax loss given in Equation (2) is the same as the cross-entropy loss between \mathbf{y} and $\hat{\mathbf{y}}$; i.e., show that

$$-\sum_{w \in \text{Vocab}} y_w \log(\hat{y}_w) = -\log(\hat{y}_o) \quad (4)$$

Because the true label y_w is a one-hot vector. And it is 1 when w is that vocabulary, otherwise it is 0.

$$y_w = \begin{cases} 1, & \text{if } w = o \\ 0, & \text{if } w \neq o \end{cases}$$

That is, the term can be reduced and represented with the negative log of its predicted output vector.

(b) Compute the partial derivative of $J_{\text{naive-softmax}}(v_c, o, U)$ with respect to v_c . Please write your answer in terms of y , \hat{y} , and U reference 1¹ : Classification and Loss Evaluation - Softmax and Cross Entropy Loss

reference 2² : Derivatives of log and exp

$$\begin{aligned}
\frac{\partial J_{\text{naive-softmax}}(v_c, o, U)}{\partial v_c} &= \frac{\partial -\log P(O = o | C = c)}{\partial v_c} = \frac{\partial -\log \frac{\exp(u_o^T v_c)}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)}}{\partial v_c} \\
&= -\left(\frac{\partial \log \exp(u_o^T v_c)}{\partial v_c} - \frac{\partial \log \sum_{w \in \text{Vocab}} \exp(u_w^T v_c)}{\partial v_c} \right) \\
&= -\frac{1}{\exp(u_o^T v_c)} \frac{\partial \exp(u_o^T v_c)}{\partial v_c} + \frac{1}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)} \frac{\partial \sum_{w \in \text{Vocab}} \exp(u_w^T v_c)}{\partial v_c} \\
&= -\frac{1}{\exp(u_o^T v_c)} \exp(u_o^T v_c) u_o + \sum_{w \in \text{Vocab}} \frac{\exp(u_w^T v_c)}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)} u_w \\
&= -u_o + \sum_{w \in \text{Vocab}} P(O = w | C = c) u_w \\
&= U^T(\hat{y} - y)
\end{aligned}$$

(b) Compute the partial derivatives of $J_{\text{naive-softmax}}(v_c, o, U)$ with respect to each of the outside word vectors, 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 y , \hat{y} , and v_c . Basically shared the same derivation of the first part thus skip some steps (partial derivative on the log)

¹<https://deepnotes.io/softmax-crossentropy>

²<https://www.themathpage.com/aCalc/exponential.htm>

$$\begin{aligned}
\frac{\partial \mathbf{J}_{\text{naive-softmax}}(v_c, o, U)}{\partial \mathbf{v}_w} &= \frac{\partial -\log P(O = o | C = c)}{\partial \mathbf{v}_w} = \frac{\partial -\log \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)}}{\partial \mathbf{v}_w} \\
&= -\left(\frac{\partial \log \exp(\mathbf{u}_o^T \mathbf{v}_c)}{\partial \mathbf{v}_w} - \frac{\partial \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)}{\partial \mathbf{v}_w} \right)
\end{aligned}$$

If $w = o$:

$$\begin{aligned}
&= -\mathbf{v}_c + \frac{\exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \mathbf{v}_c \\
&= (P(O = o | C = c) - 1) \mathbf{v}_c
\end{aligned}$$

If $w \neq o$:

$$\begin{aligned}
&= 0 + \frac{\exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \mathbf{v}_c \\
&= P(O = o | C = c) \mathbf{v}_c
\end{aligned}$$

Thus, in summary:

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(v_c, o, U)}{\partial \mathbf{v}_w} = (\hat{\mathbf{y}} - \mathbf{y})^T \mathbf{v}_c$$

(d) The sigmoid function is given by Equation 3, Please compute the derivative of $\sigma(x)$ with respect to \mathbf{x} , where \mathbf{x} is a vector

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

(e) 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 $\mathbf{u}_1, \dots, \mathbf{u}_K$. 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:

$$\mathbf{J}_{\text{neg-sample}}(v_c, o, U) = -\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \quad (5)$$

for a sample w_1, \dots, w_K , where $\sigma(\cdot)$ is the sigmoid function. ³

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 a negative sample \mathbf{u}_k . Please write your answers in terms of the vectors \mathbf{v}_c , \mathbf{u}_o , and \mathbf{u}_k , where $k \in [1, K]$. After you've done this, describe with one sentence why this loss function is much more efficient to compute than the naive-softmax loss. Note, you should be able to use your solution to part (d) to help compute the necessary gradients here.

1. $\mathbf{J}_{\text{neg-sample}}(v_c, o, U)$ with respect to \mathbf{v}_c

$$\begin{aligned} \frac{\partial \mathbf{J}_{\text{neg-sample}}(v_c, o, U)}{\partial \mathbf{v}_c} &= \frac{\partial -\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c))}{\partial \mathbf{v}_c} \\ &= -\frac{\sigma(\mathbf{u}_o^T \mathbf{v}_c)(1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c))}{\sigma \mathbf{u}_o^T \mathbf{v}_c} \frac{\partial \mathbf{u}_o^T \mathbf{v}_c}{\partial \mathbf{v}_c} - \sum_{k=1}^K \frac{\partial \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c))}{\partial \mathbf{v}_c} \\ &= -(1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \mathbf{u}_o + \sum_{k=1}^K (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \mathbf{u}_k \end{aligned}$$

³Note: the loss 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 assignment code. Ultimately, this is the same objective function.

2. $\mathbf{J}_{\text{neg-sample}}(v_c, o, U)$ with respect to \mathbf{u}_o

$$\begin{aligned}\frac{\partial \mathbf{J}_{\text{neg-sample}}(v_c, o, U)}{\partial \mathbf{u}_o} &= \frac{\partial -\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c))}{\partial \mathbf{u}_o} \\ &= \frac{\partial (-\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)))}{\partial \mathbf{u}_o} \\ &= -(1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \mathbf{v}_c\end{aligned}$$

3. $\mathbf{J}_{\text{neg-sample}}(v_c, o, U)$ with respect to \mathbf{u}_k

$$\begin{aligned}\frac{\partial \mathbf{J}_{\text{neg-sample}}(v_c, o, U)}{\partial \mathbf{u}_k} &= \frac{\partial -\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c))}{\partial \mathbf{u}_k} \\ &= \frac{\partial (-\log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)))}{\partial \mathbf{u}_k} \\ &= (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \mathbf{v}_c\end{aligned}$$

(f) Suppose the center word is $c = w_t$ and the context window is $[w_{tm}, \dots, w_{t1}, w_t, w_{t+1}, \dots, w_{t+m}]$, where m is the context window size. Recall that for the skip-gram version of word2vec, the total loss for the context window is:

$$\mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U}) \quad (6)$$

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:

1. $\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{U}$
2. $\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_c$
3. $\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_w$ when $w \neq c$

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

1.

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

2.

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

3. when $w \neq c$

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