

CS 224n Assignment 2: word2vec

Lukashevich Ilya

November 28, 2020

1 Written part: Understanding word2vec

Let's write some formulas, which were introduced in the recap section of *word2vec* algorithm, that will be used in solutions. The goal of the skip-gram *word2vec* algorithm is to accurately learn the probability distribution $P(O | C)$. The formula for the probability that the word o falls within the contextual window of c is:

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

For a single pair of words c and o the loss is given by:

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

(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

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

Solution. We know that \mathbf{y} shows the true empirical distribution and it is a one-hot vector with 1 for the true outside word, and 0 everywhere else. So, the sum contains only one non-zero item:

$$-\sum_{w \in \text{Vocab}} y_w \log(\hat{y}_w) = -(0 + \dots + 0 + 1 \cdot \log(\hat{y}_o) + 0 + \dots + 0) = -\log(\hat{y}_o)$$

(b) 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} .

Solution.

$$\begin{aligned} \frac{\partial \mathbf{J}_{\text{naive-softmax}}}{\partial \mathbf{v}_c} &= -\frac{\partial}{\partial \mathbf{v}_c} \log \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} = -\frac{\partial}{\partial \mathbf{v}_c} \log \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \\ &= -\frac{\partial}{\partial \mathbf{v}_c} \log \exp(\mathbf{u}_o^T \mathbf{v}_c) + \frac{\partial}{\partial \mathbf{v}_c} \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c) \\ &= -\mathbf{u}_o + \frac{1}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \frac{\partial}{\partial \mathbf{v}_c} \sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c) \\ &= -\mathbf{u}_o + \frac{1}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c) \frac{\partial}{\partial \mathbf{v}_c} \mathbf{u}_x^T \mathbf{v}_c \\ &= -\mathbf{u}_o + \frac{1}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c) \mathbf{u}_x \\ &= -\mathbf{u}_o + \sum_{x \in \text{Vocab}} \frac{\exp(\mathbf{u}_x^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \mathbf{u}_x \\ &= -\mathbf{u}_o + \sum_{x \in \text{Vocab}} P(O = x | C = c) \mathbf{u}_x = -\mathbf{U} \mathbf{y} + \mathbf{U} \hat{\mathbf{y}} = \mathbf{U}(\hat{\mathbf{y}} - \mathbf{y}) \end{aligned}$$

(c) 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' vector, and $w \neq o$, for all other words. Please write down your answer in terms of \mathbf{y} , $\hat{\mathbf{y}}$, and \mathbf{v}_c .

Solution.

$$\begin{aligned}\frac{\partial \mathbf{J}_{\text{naive-softmax}}}{\partial \mathbf{u}_w} &= -\frac{\partial}{\partial \mathbf{u}_w} \log \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c)} \\ &= -\frac{\partial}{\partial \mathbf{u}_w} \log \exp(\mathbf{u}_o^T \mathbf{v}_c) + \frac{\partial}{\partial \mathbf{u}_w} \log \sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c)\end{aligned}$$

Now let's consider two different cases:

1. $w = o$:

$$\begin{aligned}\frac{\partial \mathbf{J}_{\text{naive-softmax}}}{\partial \mathbf{u}_o} &= -\mathbf{v}_c + \frac{1}{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c)} \sum_{z \in \text{Vocab}} \frac{\partial}{\partial \mathbf{u}_o} \exp(\mathbf{u}_z^T \mathbf{v}_c) \\ &= -\mathbf{v}_c + \frac{1}{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c)} \frac{\partial}{\partial \mathbf{u}_o} \exp(\mathbf{u}_o^T \mathbf{v}_c) \\ &= -\mathbf{v}_c + \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c)} \mathbf{v}_c \\ &= -\mathbf{v}_c + P(O = o \mid C = c) \mathbf{v}_c \\ &= (P(O = w = o \mid C = c) - 1) \mathbf{v}_c\end{aligned}$$

2. $w \neq o$:

$$\begin{aligned}\frac{\partial \mathbf{J}_{\text{naive-softmax}}}{\partial \mathbf{u}_w} &= \frac{\partial}{\partial \mathbf{u}_w} \log \sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c) \\ &= \frac{\exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c)} \mathbf{v}_c \\ &= P(O = w \mid C = c) \mathbf{v}_c \\ &= (P(O = w \mid C = c) - 0) \mathbf{v}_c\end{aligned}$$

To conclude, we have

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}}{\partial \mathbf{u}_w} = (\hat{y}_w - y_w) \mathbf{v}_c.$$

(d) The sigmoid function is given by Equation (4):

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} \quad (4)$$

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 $\sigma(x)$.

Solution.

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

(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:

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

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

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{u}_o , \mathbf{v}_c , 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.

Solution.

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

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

This loss function is much more efficient to compute than the naive-softmax loss because at each step we do not go through all the words in the vocabulary, which is quite expensive in terms on speed costs.

(f) 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:

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

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.

Solution.

$$1. \quad \partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{U} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \partial \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U}) / \partial \mathbf{U}$$

$$2. \quad \partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_c = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \partial \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U}) / \partial \mathbf{v}_c$$

$$3. \quad \partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_w = 0 \text{ when } w \neq c$$