

CS224N Natural Language Processing with Deep Learning Assignment #2 Solution

Roy C.K. Chan

In this assignment, let $N = |Vocab|$ be the size of the vocabulary, and D be the dimension of word vectors. There is a subtle difference in defining matrices \mathbf{U} and \mathbf{V} between the written and coding parts. In the written part, the word vectors are column vectors such that $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{D \times N}$; whereas in the coding part, the word vectors are row vectors so that $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{N \times D}$.

Problem 1(a).

$$\begin{aligned} - \sum_{w \in Vocab} y_w \log(\hat{y}_w) &= - \sum_{w \in Vocab} 1\{w = o\} \log(\hat{y}_w) \\ &= - \log(\hat{y}_o) \end{aligned}$$

Problem 1(b).

$$\begin{aligned} \mathbf{J} &= \mathbf{J}_{naive-softmax}(\mathbf{v}_c, o, \mathbf{U}) \\ &= - \log(\hat{y}_o) \\ &= - \log \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \\ &= -\mathbf{u}_o^T \mathbf{v}_c + \log \sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c) \\ \implies \nabla_{\mathbf{v}_c} \mathbf{J} &= -\mathbf{u}_o + \frac{\sum_{w' \in Vocab} \exp(\mathbf{u}_{w'}^T \mathbf{v}_c) \cdot \mathbf{u}_{w'}}{\sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \\ &= -\mathbf{u}_o + \sum_{w' \in Vocab} \frac{\exp(\mathbf{u}_{w'}^T \mathbf{v}_c)}{\sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \cdot \mathbf{u}_{w'} \\ &= -\mathbf{u}_o + \sum_{w' \in Vocab} \hat{y}_{w'} \cdot \mathbf{u}_{w'} \\ &= -\mathbf{U} \mathbf{y} + \mathbf{U} \hat{\mathbf{y}} \\ &= \mathbf{U}(\hat{\mathbf{y}} - \mathbf{y}). \end{aligned}$$

Problem 1(c).

$$\begin{aligned}\mathbf{J} &= \mathbf{J}_{naive-softmax}(\mathbf{v}_c, o, \mathbf{U}) \\ &= -\mathbf{u}_o^T \mathbf{v}_c + \log \sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c)\end{aligned}$$

When $w = o$,

$$\begin{aligned}\nabla_{\mathbf{u}_w} \mathbf{J} = \nabla_{\mathbf{u}_o} \mathbf{J} &= -\mathbf{v}_c + \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c) \cdot \mathbf{v}_c}{\sum_{w' \in Vocab} \exp(\mathbf{u}_{w'}^T \mathbf{v}_c)} \\ &= -\mathbf{v}_c + \hat{y}_o \cdot \mathbf{v}_c \\ &= \mathbf{v}_c(\hat{y}_o - 1) \\ &= \mathbf{v}_c(\hat{y}_w - 1).\end{aligned}$$

When $w \neq o$,

$$\begin{aligned}\nabla_{\mathbf{u}_w} \mathbf{J} &= \mathbf{0} + \frac{\exp(\mathbf{u}_w^T \mathbf{v}_c) \cdot \mathbf{v}_c}{\sum_{w' \in Vocab} \exp(\mathbf{u}_{w'}^T \mathbf{v}_c)} \\ &= \mathbf{v}_c \cdot \hat{y}_w.\end{aligned}$$

Combining the two cases, we have

$$\begin{aligned}\nabla_{\mathbf{u}_w} \mathbf{J} &= \mathbf{v}_c(\hat{y}_w - 1\{w = o\}) \\ &= \mathbf{v}_c(\hat{y}_w - y_w),\end{aligned}$$

for all w . Hence,

$$\nabla_U \mathbf{J} = \mathbf{v}_c(\hat{\mathbf{y}} - \mathbf{y})^T.$$

Problem 1(d). Consider the case when $\mathbf{x} \in \mathbb{R}$ is a scalar, we have

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

When \mathbf{x} is a vector, the sigmoid function is simply applied to it elementwise, so

$$\nabla_{\mathbf{x}} \sigma(\mathbf{x}) = \text{diag}\left(\sigma(\mathbf{x}) \odot (1 - \sigma(\mathbf{x}))\right),$$

where \odot is the elementwise product.

Problem 1(e).

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

In the following, we will use the identity $\sigma(x) + \sigma(-x) = 1$ to simplify the expressions. Firstly,

$$\begin{aligned}
&\nabla_{\mathbf{v}_c} \mathbf{J} \\
&= -\frac{1}{\sigma(\mathbf{u}_o^T \mathbf{v}_c)} \cdot \sigma(\mathbf{u}_o^T \mathbf{v}_c) \cdot (1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \cdot \mathbf{u}_o - \sum_{k=1}^K \frac{1}{\sigma(-\mathbf{u}_k^T \mathbf{v}_c)} \cdot \sigma(-\mathbf{u}_k^T \mathbf{v}_c) \cdot (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c))(-\mathbf{u}_k) \\
&= -(1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \cdot \mathbf{u}_o + \sum_{k=1}^K (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \cdot \mathbf{u}_k \\
&= -\sigma(-\mathbf{u}_o^T \mathbf{v}_c) \cdot \mathbf{u}_o + \sum_{k=1}^K \sigma(\mathbf{u}_k^T \mathbf{v}_c) \cdot \mathbf{u}_k.
\end{aligned}$$

Secondly,

$$o \notin \{w_1, \dots, w_K\} \implies \mathbf{u}_o \neq \mathbf{u}_k \quad \forall k = 1, \dots, K,$$

so we have

$$\begin{aligned}
\nabla_{\mathbf{u}_o} \mathbf{J} &= -\frac{1}{\sigma(\mathbf{u}_o^T \mathbf{v}_c)} \cdot \sigma(\mathbf{u}_o^T \mathbf{v}_c) \cdot (1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \cdot \mathbf{v}_c - \mathbf{0} \\
&= -\sigma(-\mathbf{u}_o^T \mathbf{v}_c) \cdot \mathbf{v}_c.
\end{aligned}$$

Lastly, consider the multiset $\{\mathbf{u}_1, \dots, \mathbf{u}_K\}$, where the element \mathbf{u}_k has multiplicity n_k . Therefore,

$$\begin{aligned}
\nabla_{\mathbf{u}_k} \mathbf{J} &= \mathbf{0} - \sum_{\substack{1 \leq k' \leq K \\ \mathbf{u}_{k'} = \mathbf{u}_k}} \frac{1}{\sigma(-\mathbf{u}_{k'}^T \mathbf{v}_c)} \cdot \sigma(-\mathbf{u}_{k'}^T \mathbf{v}_c) \cdot (1 - \sigma(-\mathbf{u}_{k'}^T \mathbf{v}_c))(-\mathbf{v}_c) \\
&= n_k \cdot \sigma(\mathbf{u}_k^T \mathbf{v}_c) \cdot \mathbf{v}_c.
\end{aligned}$$

Gradient computation in naive-softmax requires summing over $\mathcal{O}(|Vocab|)$ terms, whereas that of negative sampling requires $\mathcal{O}(K)$ terms only, where $K \ll |Vocab|$.

Problem 1(f)(i).

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

Problem 1(f)(ii).

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

Problem 1(f)(iii).

$$\frac{\partial \mathbf{J}_{skip-gram}}{\partial \mathbf{v}_w} = \mathbf{0},$$

when $w \neq c$, since \mathbf{v}_c is the only column vector in \mathbf{V} related to $\mathbf{J}_{skip-gram}$.

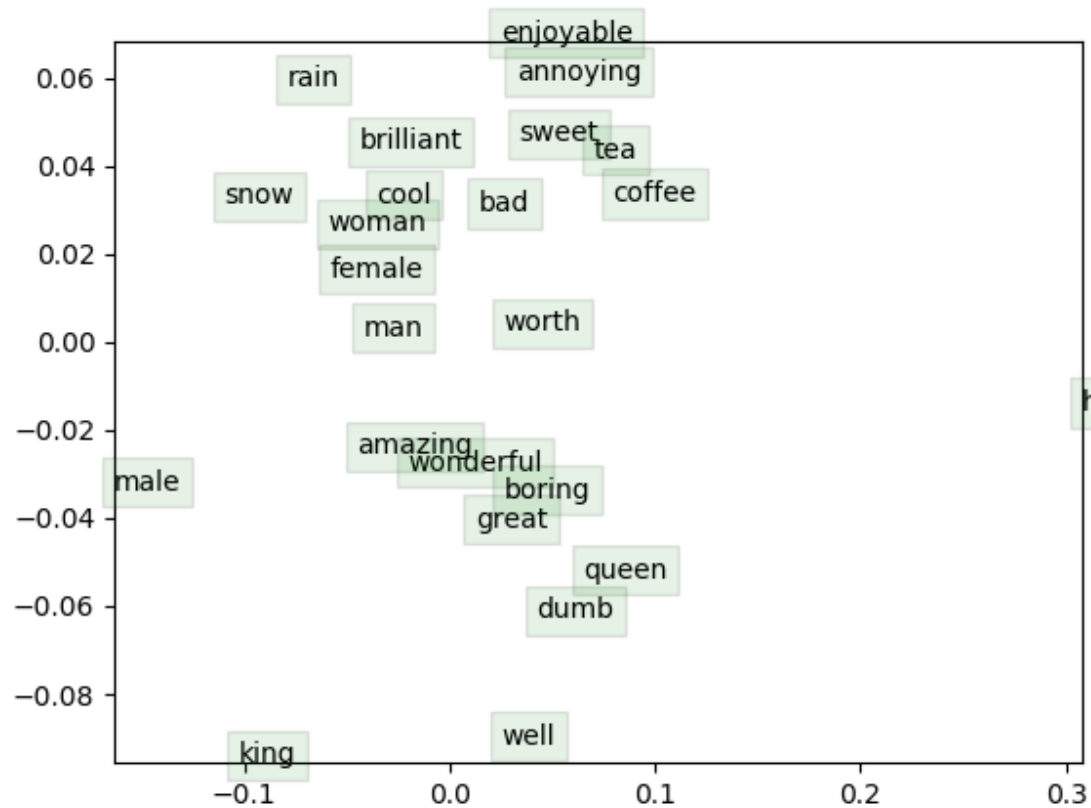


Figure 1: Visualization of Word Vectors

Problem 2(c). Some synonyms cluster together, such as {“amazing”, “wonderful”, “great”} and {“woman”, “female”}. However, antonyms may also be close to such clusters, for example, “boring” is next to {“amazing”, “wonderful”, “great”} and “man” is next to {“woman”, “female”}. Analogies, such as “man : king :: woman : queen”, may not hold in this 2D plot.