

CS 224n Assignment #2 Written Part

(a)

$$\begin{aligned} - \sum_{w \in Vocab} y_w \log(\hat{y}_w) &= -[y_1 \log(\hat{y}_1) + \dots + y_o \log(\hat{y}_o) + \dots + y_w \log(\hat{y}_w)] \\ &= -[0 \cdot \log(\hat{y}_1) + \dots + y_o \log(\hat{y}_o) + \dots + 0 \cdot \log(\hat{y}_w)] \\ &= -1 \cdot \log(\hat{y}_o) = -\log(\hat{y}_o) \end{aligned}$$

(b)

$$\begin{aligned} \frac{\partial J_{naive-softmax}(\mathbf{v}_c, o, U)}{\partial \mathbf{v}_c} &= -\frac{\partial \log(\hat{y}_o)}{\partial \mathbf{v}_c} = -\frac{\partial}{\partial \mathbf{v}_c} \log\left(\frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c)}\right) \\ &= \frac{\partial}{\partial \mathbf{v}_c} \log\left[\sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c)\right] - \frac{\partial \mathbf{u}_o^T \mathbf{v}_c}{\partial \mathbf{v}_c} \\ &= \frac{\frac{\partial}{\partial \mathbf{v}_c} \sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c)} - \mathbf{u}_o \\ &= \left(\sum_{w \in Vocab} \hat{y}_w \mathbf{u}_w\right) - \mathbf{u}_o \\ &= (\hat{y}_1 \mathbf{u}_1 - 0 \cdot \mathbf{u}_1) + \dots + (\hat{y}_o \mathbf{u}_o - 1 \cdot \mathbf{u}_o) + \dots + (\hat{y}_w \mathbf{u}_w - 0 \cdot \mathbf{u}_w) \\ &= U(\hat{\mathbf{y}} - \mathbf{y}) \end{aligned}$$

(c)

$$\begin{aligned} \frac{\partial J_{naive-softmax}(\mathbf{v}_c, o, U)}{\partial \mathbf{u}_w} &= -\frac{\partial \log(\hat{y}_o)}{\partial \mathbf{u}_w} = -\frac{\partial}{\partial \mathbf{u}_w} \log\left(\frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c)}\right) \\ &= \frac{\partial}{\partial \mathbf{u}_w} \log\left[\sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c)\right] - \frac{\partial \mathbf{u}_o^T \mathbf{v}_c}{\partial \mathbf{u}_w} \end{aligned}$$

When $w = o$,

$$\begin{aligned} &= \frac{\frac{\partial}{\partial \mathbf{u}_o} \sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c)} - \frac{\partial \mathbf{u}_o^T \mathbf{v}_c}{\partial \mathbf{u}_o} \\ &= \hat{y}_o \mathbf{v}_c - \mathbf{v}_c \end{aligned}$$

When $w \neq o$,

$$\begin{aligned}
&= \frac{\frac{\partial}{\partial \mathbf{u}_w} \sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c)} - \frac{\partial \mathbf{u}_o^T \mathbf{v}_c}{\partial \mathbf{u}_w} \\
&= \left(\sum_{w \in Vocab, w \neq o} \hat{y}_w \mathbf{v}_c \right) - 0 = \sum_{w \in Vocab, w \neq o} \hat{y}_w \mathbf{v}_c
\end{aligned}$$

Therefore,

$$\begin{aligned}
\frac{\partial J_{naive-softmax}(\mathbf{v}_c, o, U)}{\partial \mathbf{u}_w} &= \left(\sum_{w \in Vocab} \hat{y}_w \mathbf{v}_c \right) - \mathbf{v}_c \\
&= (\hat{y}_1 \mathbf{v}_c - 0 \cdot \mathbf{v}_c) + \dots + (\hat{y}_o \mathbf{v}_c - 1 \cdot \mathbf{v}_c) + \dots + (\hat{y}_w \mathbf{v}_c - 0 \cdot \mathbf{v}_c) \\
&= \mathbf{v}_c (\hat{\mathbf{y}} - \mathbf{y})^T
\end{aligned}$$

(d)

$$\frac{d\sigma(x)}{dx} = \frac{d}{dx} \frac{e^x}{e^x + 1} = 0 - \frac{d}{dx} \frac{1}{e^x + 1} = \frac{e^x}{(e^x + 1)^2}$$

(e)

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

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

$$\begin{aligned}
\frac{\partial J_{negative-sample}(\mathbf{v}_c, o, U)}{\partial \mathbf{u}_o} &= -\frac{\partial}{\partial \mathbf{u}_o} \log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \frac{\partial}{\partial \mathbf{u}_o} \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \\
&= -\frac{\frac{\partial}{\partial \mathbf{u}_o} \sigma(\mathbf{u}_o^T \mathbf{v}_c)}{\sigma(\mathbf{u}_o^T \mathbf{v}_c)} - 0 = -\sigma(\mathbf{u}_o^T \mathbf{v}_c) \cdot \mathbf{v}_c
\end{aligned}$$

$$\begin{aligned}
\frac{\partial J_{negative-sample}(\mathbf{v}_c, o, U)}{\partial \mathbf{u}_k} &= -\frac{\partial}{\partial \mathbf{u}_k} \log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \frac{\partial}{\partial \mathbf{u}_k} \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \\
&= -0 - 0 - \dots - \frac{\frac{\partial}{\partial \mathbf{u}_k} \sigma(-\mathbf{u}_k^T \mathbf{v}_c)}{\sigma(-\mathbf{u}_k^T \mathbf{v}_c)} - 0 - \dots - 0 = \sigma(\mathbf{u}_k^T \mathbf{v}_c) \cdot \mathbf{v}_c
\end{aligned}$$

It is much more efficient because with the aid of some tricks working out the derivative of sigmoid function is much easier than handling naive-softmax function.

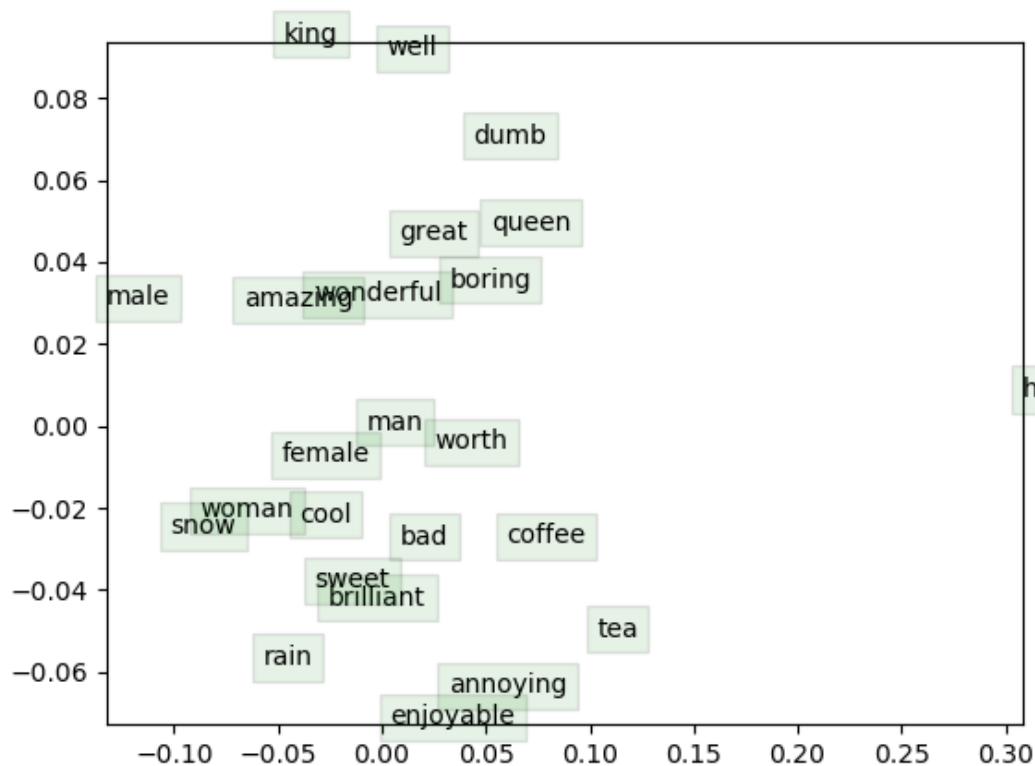
(f)

$$\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial U} = \sum_{-m < j < m, j \neq 0} \frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t+j}, U)}{\partial U}$$

$$\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial \mathbf{v}_c} = \sum_{-m < j < m, j \neq 0} \frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t+j}, U)}{\partial \mathbf{v}_c}$$

$$\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial \mathbf{v}_w} = 0$$

[My Word Vector]



Some words describing feelings or emotions like "enjoyable" and "annoying", "amazing" and "wonderful" cluster together. If we have the word vector of "male" minus "female" and the result will approximately equal the word vector of "king" minus "queen".