

CS224n : Assignment2 , written part

[question website](#)

Variables notation

Attention: All the variables' dimensions here are consistent with the code part in Assignment 2 for easy understanding.

\mathbf{U} , matrix of shape (vocab_size,embedding_dim) ,all the 'outside' vectors .

\mathbf{V} , matrix of shape (vocab_size,embedding_dim) ,all the 'center' vectors .

\mathbf{y} , vector of shape (vocab_size,1), the true empirical distribution \mathbf{y} is a one-hot vector with a 1 for the true outside word o , and 0 everywhere else .

$\hat{\mathbf{y}}$, vector of shape (vocab_size,1), the predicted distribution $\hat{\mathbf{y}}$ is the probability distribution $P(O|C = c)$ given by our model .

question a

Given outside word o and context word c .

The distribution of \mathbf{y} is as follows:

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

$$-\sum_{w=1}^V y_w \log(\hat{y}_w) = -y_o \log(\hat{y}_o) = -\log(\hat{y}_o)$$

Here , V represents the vocab_size.

question b

$$\begin{aligned}
& \frac{\partial J_{naive-softmax}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} \\
&= - \frac{\partial \log(P(O = o | C = c))}{\partial \mathbf{v}_c} \\
&= - \frac{\partial \log(\exp(\mathbf{u}_o^T \mathbf{v}_c))}{\partial \mathbf{v}_c} + \frac{\partial \log(\sum_{w=1}^V \exp(\mathbf{u}_w^T \mathbf{v}_c))}{\partial \mathbf{v}_c} \\
&= -\mathbf{u}_o + \sum_{w=1}^V \frac{\exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{w=1}^V \exp(\mathbf{u}_w^T \mathbf{v}_c)} \mathbf{u}_w \\
&= -\mathbf{u}_o + \sum_{w=1}^V P(O = w | C = c) \mathbf{u}_w \\
&= \mathbf{U}^T (\hat{\mathbf{y}} - \mathbf{y})
\end{aligned}$$

question c

$$\begin{aligned}
& \frac{\partial J_{naive-softmax}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_w} \\
&= - \frac{\partial \log(\exp(\mathbf{u}_o^T \mathbf{v}_c))}{\partial \mathbf{u}_w} + \frac{\partial \log(\sum_{w=1}^V \exp(\mathbf{u}_w^T \mathbf{v}_c))}{\partial \mathbf{u}_w}
\end{aligned}$$

when $w = o$,

$$\begin{aligned}
& \frac{\partial J_{naive-softmax}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_w} \\
&= -\mathbf{v}_c + \frac{1}{\sum_{w=1}^V \exp(\mathbf{u}_w^T \mathbf{v}_c)} \frac{\partial \sum_{w=1}^V \exp(\mathbf{u}_w^T \mathbf{v}_c)}{\partial \mathbf{u}_o} \\
&= -\mathbf{v}_c + \frac{1}{\sum_{w=1}^V \exp(\mathbf{u}_w^T \mathbf{v}_c)} \frac{\partial \exp(\mathbf{u}_o^T \mathbf{v}_c)}{\partial \mathbf{u}_o} \\
&= -\mathbf{v}_c + \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w=1}^V \exp(\mathbf{u}_w^T \mathbf{v}_c)} \mathbf{v}_c \\
&= (P(O = o | C = c) - 1) \mathbf{v}_c
\end{aligned}$$

when $w \neq o$,

$$\begin{aligned} & \frac{\partial J_{naive-softmax}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_w} \\ &= \frac{\exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{w=1}^V \exp(\mathbf{u}_w^T \mathbf{v}_c)} \mathbf{v}_c \\ &= P(O = w | C = c) \mathbf{v}_c \end{aligned}$$

In summary,

$$\begin{aligned} & \frac{\partial J_{naive-softmax}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{U}} \\ &= (\hat{\mathbf{y}} - \mathbf{y})^T \mathbf{v}_c \end{aligned}$$

question d

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

question e

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

$$\begin{aligned}
& \frac{\partial J_{neg-sample}(v_c, o, U)}{\partial u_o} \\
&= \frac{\partial(-\log(\sigma(u_o^T v_c)))}{\partial u_o} = -(1 - \sigma(u_o^T v_c))v_c \\
& \frac{\partial J_{neg-sample}(v_c, o, U)}{\partial u_k} \\
&= \frac{\partial(-\log(\sigma(-u_k^T v_c)))}{\partial u_k} = (1 - \sigma(-u_k^T v_c))v_c
\end{aligned}$$

question f

i)

$$\begin{aligned}
& \frac{\partial J_{skip-gram}(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial U} \\
&= \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(v_c, w_{t+j}, U)}{\partial U}
\end{aligned}$$

ii)

when $w=c$,

$$\begin{aligned}
& \frac{\partial J_{skip-gram}(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial v_c} \\
&= \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(v_c, w_{t+j}, U)}{\partial v_c}
\end{aligned}$$

iii)

when $w \neq c$,

$$\begin{aligned}
& \frac{\partial J_{skip-gram}(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial v_w} \\
&= \mathbf{0}
\end{aligned}$$