Name: Bhavnick Minhas

CS 224 N

Assignment 2 [Written]

Understanding Word2Vec

Given:

$$P(O=o \mid C=c) = \frac{\exp(u_o^T v_c)}{\sum_{\omega \in Vocab} \exp(u_\omega^T v_c)}$$

$$J_{naive-softmax}(v_c, o, U) = -\log P(O=o \mid C=c)$$

a)   To prove:

$$-\sum_{\omega \in Vocab} y_\omega \log(\hat{y}_\omega) = -\log(\hat{y}_o)$$

Since, $\vec{y}_o = [0 \cdots \underset{o^{th}\, index}{1} \cdots 0]$

as $y_\omega$ is just one-hot pointing to the true outer word, hence.

$$y_\omega = \begin{cases} 0 & \text{for } \omega \neq o \\ 1 & \text{for } \omega = o \end{cases}$$

$$\Rightarrow -\sum y_\omega \log(\hat{y}_\omega) = 0 + 0 \cdots \\ -\log(\hat{y}_o) \\ +0 \cdots \\ = -\log(\hat{y}_o)$$

b) To get: $\dfrac{\partial J_{nawe-softmax}}{\partial v_c}$

First, lets simplify $J_{ns}$ a bit.

$$J_{ns} = -\log\left(\frac{\exp(u_o^T v_c)}{\sum_w \exp(u_w^T v_c)}\right)$$

$$= -u_o^T v_c + \log\left(\sum_w \exp(u_w^T v_c)\right)$$

Then, we have,

$$\frac{\partial J_{ns}}{\partial v_c} = -u_o + \frac{1}{\sum_w \exp(u_w^T v_c)} \cdot \sum_w \exp(u_w^T v_c) \cdot u_w$$

$$\underset{d \times m}{\phantom{x}} \qquad \underset{m \times d}{\phantom{x}} \quad \underset{d \times 1}{u_o}$$

$$\equiv -u_o + \frac{U \cdot \exp(U^T v_c)}{\underset{1 \times m}{1} \cdot \exp(U^T v_c)}$$

where $U \in \mathbb{R}^{d \times m}$
$v_c \in \mathbb{R}^{d \times 1}$
$\bar{1} = [1, 1, \ldots]$

$$\equiv -U \cdot y + U \cdot \hat{y}$$

$$= U \cdot (\hat{y} - y)$$

c) To compute: $\dfrac{\partial J_{ns}}{\partial \vec{u}_\omega}$

for $\omega \in \{1, \dots m\}$

case I: $\omega \neq 0$, $u_\omega$

$$\dfrac{\partial}{\partial u_\omega}\left(-\vec{u}_\omega^T v_c + \log\left(\sum \exp(u_\omega^T v_c)\right)\right)$$

$$= 0 + \dfrac{1}{\sum_\omega \exp(u_\omega^T v_c)} \cdot \exp(u_\omega^T v_c) \cdot v_c$$

$$= \hat{y}_\omega \cdot \vec{v}_c$$

case II: $\omega = 0$

$$\dfrac{\partial}{\partial u_0}\left(-u_0^T v_c + \log\left(\sum \exp(u_\omega^T v_c)\right)\right)$$

$$= -v_c + \dfrac{1}{\sum \exp(u_\omega^T v_c)} \cdot \exp(u_0^T v_c) \cdot v_c$$

$$= -\vec{v}_c + \hat{y}_0 \cdot \vec{v}_c$$

$$= \vec{v}_c (\hat{y}_0 - y_0).$$

d)
$$\frac{\partial J_{no}}{\partial U} = \left[ \frac{\partial J}{\partial u_o} \; \frac{\partial J}{\partial u_1} \; \cdots \; \frac{\partial J}{\partial u_m} \right]$$

$$= \left[ \hat{y}_o \bar{v}_c, \; \hat{y}_1 \bar{v}_c \cdots \hat{y}_j(\hat{y}_o - y_o) \; \cdots \; \hat{y}_m v_c \right]$$

$$= \bar{v}_c \cdot (\hat{y} - y)^T$$

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

we need: $\dfrac{\partial \sigma(x)}{\partial x} \quad = \quad \dfrac{1}{(1 + e^{-x})^2} \cdot + e^{-x}$

$$= \frac{e^{-x}}{(1 + e^{-x})^2} = \sigma(x) \left[ \frac{e^{-x}}{1 + e^{-x}} \right]$$

$$= \sigma(x) \left[ \frac{1 + e^{-x} - 1}{1 + e^{-x}} \right]$$

$$= \sigma(x) \left[ 1 - \frac{1}{1 + e^{-x}} \right]$$

$$= \sigma(x) (1 - \sigma(x))$$

f)

$$J_{neg\text{-}sampling} = -\log(\sigma(u_o^T v_c))$$

$$-\sum_{k=1}^{K} \log(\sigma(-u_k^T v_c))$$

i) To get: $\dfrac{\partial J_{neg\text{-}sampling}}{\partial v_c}$

$$\frac{\partial J_{ns}}{\partial v_c} = -\frac{1}{\sigma(u_o^T v_c)} \sigma(u_o^T v_c)(1 - \sigma(u_o^T v_c)) \cdot u_o$$

$$+ \sum_k \frac{1}{\sigma(-u_k^T v_c)} \cdot \sigma(-u_k^T v_c) \cdot (1 - \sigma(u_k^T v_c)) \cdot u_k$$

$$= -1(1 - \sigma(u_o^T v_c))\, \bar{u}_o$$
$$+ \sum \bar{u}_k (1 - \sigma(-u_k^T v_c))$$

ii) To get: $\dfrac{\partial J_{ns}}{\partial u_o}$

$$\frac{\partial J_{ns}}{\partial u_o} = -\frac{1}{\sigma(u_o^T v_c)} \cdot \sigma(u_o^T v_c) \cdot (1 - \sigma(u_o^T v_c))$$
$$\cdot v_c$$

$$= -(1 - \sigma(u_o^T v_c))\, v_c$$

iii) To get: $\dfrac{\partial J_{ns}}{\partial u_k}$

$$\dfrac{\partial J_{ns}}{\partial u_k} = \dfrac{1}{\sigma(-u_k^T v_c)} \cdot \sigma(-u_k^T v_c)(1-\sigma(u_k^T v_c))$$

$$\cdot \; +v_c$$

$$= (1-\sigma(-u_k^T v_c))\, v_c$$

g) In the case that the indices are distinct, the indices can be devided into two parts:

    i) One contaning same index k
    ii) Other not being k

$$\Rightarrow \quad \dfrac{\partial J_{ns}}{\partial u_k} = +\sum_{d=1}^{d_k} \dfrac{1}{\sigma(-u_k^T v_c)} \cdot \dfrac{\sigma(-u_k^T v_c)}{(1-\sigma(-u_k^T v_c))}$$

$$\cdot \; +v_c$$

$$= d_k (1-\sigma(-u_k^T v_c))\, v_c$$

where $d_k$ are the no. of $k^{th}$ indices in samples.

b) For skip gram, since,

$$J_{sg}(v_c, w_{t-m}, \ldots, w_{t+m}, U)$$

$$= \sum_{\substack{-m \le j \le m \\ j \ne 0}} J(v_c, w_{t+j}, U)$$

i) $$\frac{\partial J_{sg}}{\partial U} = \sum_{\substack{-m \le j \le m \\ j \ne 0}} \frac{\partial J(v_c, w_{t+j}, U)}{\partial U}$$

ii) $$\frac{\partial J_{sg}}{\partial v_c} = \sum_{\substack{-m \le j \le m \\ j \ne 0}} \frac{\partial J(v_c, w_{t+j}, U)}{\partial v_c}$$

iii) $$\frac{\partial J_{sg}}{\partial v_w} = 0 \quad \left[\text{since the term } v_w \text{ doesn't occur at all in the expression.}\right]$$