

Date: .....

تاریخ: .....

Subject: .....

موضوع: .....

a) As mentioned before,  $y$  is a one-hot vector with a 1 for the true outside word  $o$ , so the answer of  $\Sigma$  is just  $-\log(\hat{y}_o)$ .

$$\begin{aligned}
 \text{b) i) } \frac{\partial J_{\text{naive-softmax}}(v_c, o, u)}{\partial v_c} &= \frac{\partial}{\partial v_c} \left( -\log \left( \frac{\exp(u_o^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} \right) \right) \\
 &= \frac{\partial}{\partial v_c} \left( -\log(\exp(u_o^T v_c)) + \log \left( \sum_{w \in \text{vocab}} \exp(u_w^T v_c) \right) \right) \\
 &= \frac{\partial}{\partial v_c} \left( -\exp(u_o^T v_c) + \log \left( \sum_{w \in \text{vocab}} \exp(u_w^T v_c) \right) \right) \\
 &= -u_o + \frac{\sum_{w \in \text{vocab}} u_w \exp(u_w^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} \\
 &\quad \underbrace{\hspace{10em}}_{\sum_{w \in \text{vocab}} \hat{y}_w u_w} \\
 &= u \hat{y} - u_o
 \end{aligned}$$

Date:

تاریخ:

Subject:

موضوع:

b) ii) the softmax prediction is perfect and equal to true distribution.

iii) in this gradient each time we update the center word and the loss function get minimized which means the probability get closer.

iv) Sometimes word vectors are really similar to each other according to similarity of that words in meaning. when our data is too big and we need to make it small considering that words equal is good for us but sometimes it takes useful information from us and may effect badly in our downstream task for example word vector of "very bad" and "not bad" maybe very similar to each other and by doing this normalization they will consider the same but we know that they aren't.

تاریخ:

Date:

موضوع:

Subject:

$$\begin{aligned}
 c) \quad \frac{\partial J_{\text{naive-softmax}}}{\partial U} &= \frac{\partial}{\partial U} \left( -\log \left( \frac{\exp(u_o^T \vartheta_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T \vartheta_c)} \right) \right) \\
 &= -\frac{\partial u_o^T \vartheta_c}{\partial U} + \frac{1}{\sum_{w \in \text{vocab}} \exp(u_w^T \vartheta_c)} \left( \sum_{w \in \text{vocab}} \exp(u_w^T \vartheta_c) \times \frac{\partial u_w^T \vartheta_c}{\partial U} \right) \\
 &= -\vartheta_c y^T + \sum_{w \in \text{vocab}} \underbrace{\hat{y}_w \vartheta_c \hat{y}^T}_{\vartheta_c \hat{y}^T} \times \frac{\partial u_w^T \vartheta_c}{\partial U} = \vartheta_c (\hat{y}^T - y^T) \\
 &= \vartheta_c (\hat{y} - y)^T
 \end{aligned}$$

$$\rightarrow \begin{cases} \text{if } w = 0 \rightarrow \vartheta_c \hat{y}^T - \vartheta_c \\ \text{if } w \neq 0 \rightarrow \vartheta_c \hat{y}^T \end{cases}$$

$$d) \quad \frac{\partial J_{\text{naive-softmax}}}{\partial U} = \left[ \frac{\partial J(\vartheta_c, 0, U)}{\partial u_1} \quad \frac{\partial J(\vartheta_c, 0, U)}{\partial u_2} \quad \frac{\partial J(\vartheta_c, 0, U)}{\partial u_{|\text{vocab}|}} \right]$$

$$e) \quad \begin{cases} \frac{\partial f(n)}{\partial n} = \alpha & n < 0 \\ \frac{\partial f(n)}{\partial n} = 1 & n > 0 \end{cases}$$



Date:

تاریخ:

Subject:

موضوع:

$$\begin{aligned}
 f) \quad \frac{\partial \sigma(n)}{\partial n} &= \frac{e^n(e^{n+1}) - e^n(e^n)}{(e^{n+1})^2} = \frac{e^n}{(e^{n+1})^2} \\
 &= \frac{e^n}{e^{2n+2}} \times \left( \frac{e^{n+1} - e^n}{e^{n+1}} \right) = \sigma(n) \times (1 - \sigma(n)) \\
 &= \sigma(n) - \sigma^2(n)
 \end{aligned}$$

$$\begin{aligned}
 g) \quad i) \quad \frac{\partial J_{\text{neg-sample}}}{\partial \vartheta_c} &= \frac{\partial}{\partial \vartheta_c} \left( -\log(\sigma(u_0^T \vartheta_c)) - \sum_{s=1}^K \log(\sigma(-u_{ws}^T \vartheta_c)) \right) \\
 &= \frac{-1}{\sigma(u_0^T \vartheta_c)} \times (u_0 \times (\sigma(u_0^T \vartheta_c) - \sigma(u_0^T \vartheta_c)^2)) \\
 &\quad - \sum_{s=1}^K \left( \frac{1}{\sigma(-u_{ws}^T \vartheta_c)} \times (-u_{ws} \times (\sigma(-u_{ws}^T \vartheta_c) - \sigma(-u_{ws}^T \vartheta_c)^2)) \right) \\
 &= -u_0 (1 - \sigma(u_0^T \vartheta_c)) + \sum_{s=1}^K u_{ws} (1 - \sigma(-u_{ws}^T \vartheta_c)) \\
 \frac{\partial J_{\text{neg-sample}}}{\partial u_0} &= \frac{\partial}{\partial u_0} \left( -\log(\sigma(u_0^T \vartheta_c)) - \sum_{s=1}^K \log(\sigma(-u_{ws}^T \vartheta_c)) \right) \\
 &= \frac{-1}{\sigma(u_0^T \vartheta_c)} \times (\vartheta_c \times (\sigma(u_0^T \vartheta_c) - \sigma(u_0^T \vartheta_c)^2)) - 0 \\
 &= -\vartheta_c (1 - \sigma(u_0^T \vartheta_c))
 \end{aligned}$$

Date: تاريخ:

Subject: موضوع:

$$g) \quad i) \quad \frac{\partial \mathcal{J}_{\text{neg-sample}}}{\partial u_{ws}} = \frac{\partial}{\partial u_{ws}} \left( -\log(\sigma(u_0^T \theta_c)) - \sum_{s=1}^K \log(\sigma(-u_{ws}^T \theta_c)) \right)$$

$$= 0 - \frac{1}{\sigma(-u_{ws}^T \theta_c)} \times (-\theta_c) \times (\sigma(-u_{ws}^T \theta_c) - \sigma(-u_{ws}^T \theta_c)^2)$$

$$= \theta_c (1 - \sigma(-u_{ws}^T \theta_c))$$

ii) the first one because it has both terms and we can use it in second and third one.

$$\frac{\partial \mathcal{J}_{\text{neg-sample}}}{\partial \theta_c} = -u_0 (1 - \sigma(u_0^T \theta_c))$$

iii) because in negative sampling we compute on just  $K$  words but in softmax the computation is on the whole vocabulary words, so using negative sampling is more efficient.

Date:

تاریخ:

Subject:

موضوع:

h) because this time our words are not distinct so the answer would be the  $\sum$  of the words which are equal to  $w_s$  so:

$$\frac{\partial J_{\text{neg-sampling}}}{\partial u_{w_s}} = \sum_{w=w_s} u_c (1 - \sigma(-u_{w_s}^T u_c))$$

$$i) \sum_{\substack{m=1 \\ j \neq 0}}^m \frac{\partial J(u_c, w_{t+j}, u)}{\partial u}$$

$$ii) \sum_{\substack{m=1 \\ j \neq 0}}^m \frac{\partial J(u_c, w_{t+j}, u)}{\partial u_c}$$

$$iii) = 0$$