

Bakhtash Ansari

a) Based on explanations above y is one-hot vector which has a 1 for the true outside word o , and 0 everywhere else, so all terms will be zero except for $w = o$

so we have: $-\sum_{w \in V} y_w \log(\hat{y}_w) = -\log(\hat{y}_o)$

b) Dimensions of u and v is ~~vector~~ ~~the~~ number of words for shape(0) and one for shape(1) (num, 1) so they are same in dimension and we have:

$$\begin{aligned}\frac{\partial J}{\partial v_c} &= \frac{-\partial \log P(o|c)}{\partial v_c} = -\frac{\partial \log \exp(u_o^T v_c)}{\partial v_c} \\ &= -\frac{\partial (\log e^{u_o^T v_c})}{\partial v_c} + \frac{\partial (\log \sum_{w \in V} e^{u_w^T v_c})}{\partial v_c} \\ &= -\frac{\partial (u_o^T v_c)}{\partial v_c} + \sum_{w \in V} e^{u_w^T v_c} \times \frac{\partial (\sum_{w \in V} e^{u_w^T v_c})}{\partial v_c}\end{aligned}$$

$$\begin{aligned}
 &= -u_0^T + \sum_{w \in V} \frac{e^{u_w^T v_c}}{\sum_{w \in V} e^{u_w^T v_c}} \times u_w^T \\
 &= -u_0^T + [P(O=w | C=c) \times u_w^T] = -u_0^T + \hat{y} \times U^T \\
 &= U^T(\hat{y} - y)
 \end{aligned}$$

ii) when is the gradient you computed equal to zero?

based on answer in last part $U^T(\hat{y} - y) = 0$

so we have $\hat{y} = y \Rightarrow$ we can understand that gradient will be equal to zero whenever our estimated word equal to real word.

iii) When we subtract gradient from v_c we are creating Gradient Descent formula :

$$\theta^{\text{new}} = \theta^{\text{old}} - \alpha \nabla_{\theta} J(\theta)$$

that we use it when we want to minimize the value of loss function which lead to predict true word with higher probability.

i) Regarding to hint, in our data if we have some words that similar word vectors but not identical (like synonym words) by L₂ normalization these vectors consider equal and can make noise in our calculations. or for example if sentiment of a phrase is important for us and we have two phrases which have different emotions but similar word vectors like (very good, kind of good), L₂ may remove this informations. on the other hand in tasks that

high similarity is not so important, we can use it.

C) when $w=0$: $\frac{\partial(\bar{J}(\bar{v}_c, 0, \bar{U}))}{\partial \bar{v}} = \frac{\partial(-\log P(O=w|C=c))}{\partial \bar{v}}$

$$\frac{\partial}{\partial \bar{v}} \left(-\log \left(\frac{e^{u_0^T \bar{v}_c}}{\sum_{w \in \bar{V}} e^{u_w^T \bar{v}_c}} \right) \right) = -\frac{\partial \log e^{u_0^T \bar{v}_c}}{\partial \bar{v}} + \frac{\partial \log \sum_{w \in \bar{V}} e^{u_w^T \bar{v}_c}}{\partial \bar{v}}$$

$$= -\frac{\partial u_0^T \bar{v}_c}{\partial \bar{v}} + \sum_{w \in \bar{V}} \frac{e^{u_w^T \bar{v}_c}}{\sum_{w \in \bar{V}} e^{u_w^T \bar{v}_c}} \times \bar{v}_c^T$$

$$= -\frac{\partial u_0^T \bar{v}_c}{\partial \bar{v}} + \sum_{w \in \bar{V}} P(O=w|C=c) \times \bar{v}_c^T$$

$$= -\frac{\partial \bar{v}_c^T}{\partial \bar{v}} + \hat{y}^T \bar{v}_c^T = -\bar{y}^T \bar{v}_c + \hat{y}^T \bar{v}_c = \bar{v}_c^T (\hat{y} - \bar{y})$$

if $w=0 \Rightarrow y=1 \Rightarrow \text{Result} = \nabla_C^T (\hat{y}-1) = \nabla_C^T \hat{y} - \nabla_C^T$
 $0 \cdot w \Rightarrow y=0 \Rightarrow \text{Result} = \nabla_C^T \hat{y}'$

$$D) \frac{\partial J}{\partial u} = \left[\frac{\partial J(v_c, o, u)}{\partial u_1}, \frac{\partial J(v_c, o, u)}{\partial u_2}, \dots, \frac{\partial J(v_c, o, u)}{\partial u_{|\text{vocab}|}} \right]$$

e) Based on Leaky ReLU:

$$f(n) = \max(\alpha n, n) \quad (\alpha < 1) = \begin{cases} \alpha n & n < 0 \\ n & n \geq 0 \end{cases}$$

$$\text{So we have } \frac{\partial f(n)}{\partial n} (n < 0) = \alpha$$

$$\frac{\partial f(n)}{\partial n} (n \geq 0) = 1$$

$$f) \frac{\partial \alpha(n)}{\partial n} = \frac{\partial \left(\frac{e^n}{e^n + 1} \right)}{\partial n} = \frac{e^n(e^n + 1) - e^{2n}}{(e^n + 1)^2} =$$

$$\frac{e^n(e^n + 1)}{(e^n + 1)^2} - \frac{e^n e^n}{(e^n + 1)^2} = \frac{e^n}{e^n + 1} - \frac{e^n e^n}{(e^n + 1)^2} = \boxed{\alpha(n) - \alpha^2(n)}$$

$$r \leftarrow \Delta \leq v \leftarrow Q$$

d) i) $\frac{\partial (\bar{J}(r_c, 0, u))}{\partial r_c} = \frac{\partial (-\log(\alpha(u_0^T r_c)))}{\partial r_c}$

$$\sum_{s=1}^K \frac{\partial (\log(\alpha(-u_{ws}^T r_c)))}{\partial r_c} = -\frac{1}{\alpha(u_0^T r_c)} \times \alpha(u_0^T r_c)(1-\alpha(u_0^T r_c))$$

$$x u_0^T + \sum_{s=1}^K \frac{1}{\alpha(-u_{ws}^T r_c)} \times \alpha(-u_{ws}^T r_c)(1-\alpha(-u_{ws}^T r_c)) \times u_{ws}^T$$

$$= (\alpha(u_0^T r_c) - 1) u_0^T + \sum_{s=1}^K (1-\alpha(-u_{ws}^T r_c)) u_{ws}^T$$

$$= u_0^T \alpha(u_0^T r_c) - u_0^T + \sum_{s=1}^K (u_{ws}^T - u_{ws}^T \alpha(-u_{ws}^T r_c))$$

$$\frac{\partial (\bar{J}(r_c, 0, u))}{\partial u_0} = \frac{\partial (-\log(\alpha(u_0^T r_c)))}{\partial u_0} - \sum_{s=1}^K \frac{\partial (\log(\alpha(-u_{ws}^T r_c)))}{\partial u_0}$$

$$= -\frac{1}{\alpha(u_0^T r_c)} \times \alpha(u_0^T r_c)(1-\alpha(u_0^T r_c)) \times r_c^T$$

$$= (\alpha(u_0^T r_c) - 1) r_c^T$$

$$\frac{\partial (\bar{J}(r_c, 0, u))}{\partial u_{ws}} = \frac{\partial (-\log(\alpha(u_0^T r_c)))}{\partial u_{ws}}$$

$$\sum_{s=1}^K \frac{\partial (\log(\alpha(-u_{ws}^T r_c)))}{\partial u_{ws}} = 0 + \sum_{s=1}^K \frac{1}{\alpha(-u_{ws}^T r_c)} \times$$

$$\alpha(-u_{ws}^T r_c)(1-\alpha(-u_{ws}^T r_c)) \times r_c^T = (1-\alpha(-u_{ws}^T r_c)) r_c^T$$

ii) Based on calculations of last part we use computations of first derivative in others so we reuse the $\frac{\partial J_{\text{neg-sample}}}{\partial v_c}$, for representation we have

$$\frac{\partial J_{\text{neg-sample}}}{\partial v_c} = -U_0 \{w_1, \dots, w_K\} \left(\sigma(U_0^T \{w_1, \dots, w_K\} v_c - 1) \right)$$

iii) Because it only compute the gradients for small number of samples but naive-Sofmax loss needs to compute gradients for all possible samples.

h) Based on previous computations we have:

$$\frac{\partial J}{\partial w_s} = v_c^T (1 - \sigma(-u_{ws}^T v_c))$$

but now we know that for word in vocab, some w_s are equal to w_s so we should compute the derivative multiple times, so we have

$$\frac{\partial J}{\partial w_s} = \sum_{\substack{w \in V \\ w = w_s}}^K v_c^T (1 - \sigma(-u_{ws}^T v_c))$$

i)

$$\text{i) } \frac{\partial J_{\text{skip-gram}}}{\partial U} = \sum_{\substack{-m < j < m \\ j \neq 0}} \frac{\partial J(v_c, w_{t+j}, U)}{\partial U}$$

$$\text{ii) } \frac{\partial J_{\text{skip-gram}}}{\partial v_c} = \sum_{\substack{-m < j < m \\ j \neq 0}} \frac{\partial J(v_c, w_{t+j}, U)}{\partial v_c}$$

$$\text{iii) } \frac{\partial J_{\text{skip-gram}}}{\partial w} = 0$$