

Subject :

written assignments

Date

a)

which is $y = [0, 0, \dots, 1, \dots, 0, 0]$

true label

which is $\hat{y} = [p_0, p_1, \dots, p_0, \dots, 0, 0]$

predicted label

$\sum_{w \in \text{vocab}} y_w \log(\hat{y}_w) = y$ because all elements of y is 0 except element at position 0 so we have just $= 1 \times p_0$

$$p_0 = \hat{y}_0 \Rightarrow - \sum_{w \in \text{vocab}} y_w \log(\hat{y}_w) = p_0 = \hat{y}_0$$

b)

$$i) \frac{\partial \bar{J}_{\text{naive-softmax}}(v_c, u, V)}{\partial v_c} = \frac{\partial}{\partial v_c} \left(- \log \left(\frac{\exp(u_0^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} \right) \right)$$

split log
division to minus $\frac{\partial}{\partial v_c} \left(- \left(\log(\exp(u_0^T v_c)) - \log \left(\sum_{w \in \text{vocab}} \exp(u_w^T v_c) \right) \right) \right)$

$$\xrightarrow{\log e^u = u} \frac{\partial}{\partial v_c} \left(- \left(\exp(u_0^T v_c) - \log \left(\sum_{w \in \text{vocab}} \exp(u_w^T v_c) \right) \right) \right)$$

apply
derivation $\frac{\sum \exp(u_w^T v_c) u_w}{\sum \exp(u_w^T v_c)} - u_0 \xrightarrow[\text{to } \hat{y}]{\text{softmax}} \sum_{w \in \text{vocab}} \hat{y}_w u_w = u_0$

$$\rightarrow U \hat{y} - u_0 \xrightarrow{u_0 = Uy} U \hat{y} - Uy \rightarrow U(\hat{y} - y)$$

because y just 0's

MICRO index is 1

ii) it's zero when $\begin{cases} U = 0 & (1) \\ \text{or} \\ y = \hat{y} & (2) \end{cases}$

① $U = 0$ \rightarrow it's not possible because vectors should vary from each other and they can't be equal to zero.

② $y = \hat{y}$ \rightarrow when \hat{y} which we predicted is completely equal to the true y . it means that all of other words probability is equal to zero and only our target U_0 has probability 1, which obviously in this case we don't need to change our vectors because we're in the place that we want to be, but in real cases is not possible this happens.

iii)

because by subtracting this gradient from our v_c we move our v_c closer to local or global minima.

how this happens relies on the concept of loss function which we want to minimize and gradient helps us in this way.

in fact gradient shows how should probability of our prediction change. by subtracting $\hat{y} - y$ we want to be closer to y by changing our \hat{y} in each iteration.

iv) when downstream applications only pay attention to cosine similarity, normalizing would be a good idea. but if they want to capture features beyond similarities, such as significance, consistency, normalizing will ignore these. so in this case it's not applicable.

● MICRO: in fact normalizing will lose length of vectors. in some cases both length and direction carries important information.

$$c) \text{ when } w = 0 \rightarrow -\frac{\delta}{\delta u_w} \log \frac{\exp(u_0^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} = -\frac{\delta}{\delta u_w} (u_w^T v_c)$$

$$+ \frac{\delta}{\delta u_w} \log \sum_{w \in \text{vocab}} \exp(u_w^T v_c)$$

$$= -v_c + \frac{v_c \exp(u_w^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} = -v_c + v_c \hat{y}$$

$$\hat{y} = \text{softmax} \leftarrow \sum_{w \in \text{vocab}} \exp(u_w^T v_c)$$

$$* \text{ when } w \neq 0 \Rightarrow -\frac{\delta}{\delta u_w} (u_0^T v_c) + \frac{\delta}{\delta u_w} \log \sum_{w \in \text{vocab}} \exp(u_w^T v_c)$$

this term is zero

so the only term remaining is the second part which we calculated above

$$v_c \cdot \text{softmax}(w) = v_c \hat{y}$$

$$\frac{\delta J_{\text{softmax}}}{\delta u_w} \bigg|_{w=0} \rightarrow -v_c + v_c \hat{y}$$

$$\frac{\delta J_{\text{softmax}}}{\delta u_w} \bigg|_{w \neq 0} \rightarrow v_c \hat{y}$$

$$d) \frac{\delta J_{\text{softmax}}(v_c, 0, U)}{\delta U} =$$

$$\left[\frac{\delta J(v_c, 0, U)}{\delta u_1} \quad \frac{\delta J(v_c, 0, U)}{\delta u} \quad \dots \quad \frac{\delta J(v_c, 0, U)}{\delta u_{|\text{vocab}|}} \right]$$

$$\left[\hat{y}_1 v_c, \quad \hat{y}_2 v_c, \quad \dots \quad \hat{y}_{|\text{vocab}|} v_c \right]$$

$$e) f(x) = \max(ax, x)$$



$$\frac{df(x)}{dx} = \begin{cases} 1 & x > 0 \\ a & x < 0 \end{cases}$$

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

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

$$g) J_{\text{neg-sample}} = -\log(\sigma(u_0^T v_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T v_c))$$

$$i) \frac{\partial J_{\text{neg sampling}}}{\partial v_c} = \frac{\sigma(u_0^T v_c)(1 - \sigma(u_0^T v_c)) \cdot \partial u_0^T v_c}{\sigma(u_0^T v_c)} + \sum_{k=1}^K \frac{\partial \log(\sigma(-u_k^T v_c))}{\partial v_c} = \sigma(u_0^T v_c - 1)u_0 + \sum_{k=1}^K (1 - \sigma(u_k^T v_c))u_k$$

$$\frac{\partial J_{\text{neg sampling}}}{\partial u_0} = -\frac{\partial}{\partial u_0} \left(\log \sigma(u_0^T v_c) - \sum_{k=1}^K \log \sigma(-u_k^T v_c) \right)$$

$$= -\frac{\partial}{\partial u_0} \left(-\log(\sigma(u_0^T v_c)) - (\sigma(u_0^T v_c) - 1)v_c \right) \quad \text{This term is 0}$$

$$\frac{\partial J_{\text{neg sampling}}}{\partial u_{ws}} = \frac{\partial (-\log \sigma(-u_{ws}^T v_c))}{\partial u_{ws}} = (1 - \sigma(-u_{ws}^T v_c))v_c$$

ii) terms which has sigmoid and inner calculations are removed
this term is repeated $\rightarrow (\sigma(u_0^T v_c) - 1)$

So we can $U = \{w_1, \dots, w_k\}$ $v_c = \begin{bmatrix} u_0^T v_c \\ -u_{w_1}^T v_c \\ \vdots \\ u_{w_k}^T v_c \end{bmatrix} \rightarrow 1 - \sigma(Uv_c) - 1 = \begin{bmatrix} \sigma(u_0^T v_c) - 1 \\ \sigma(-u_{w_1}^T v_c) - 1 \\ \vdots \\ \sigma(-u_{w_k}^T v_c) - 1 \end{bmatrix}$

iii) in naive softmax loss we obviously iterate through whole vocabulary each time which cost a lot more then here, because we are only sampling k item in negative sampling.

$$h) J_{\text{neg sampling}} = -\log(\sigma(u_o^T v_c)) + \sum_{\substack{1 \leq j \leq k \\ w_j = w_s}} \log(\sigma(-u_{w_j}^T v_c)) + \sum_{\substack{1 \leq j \leq k \\ w_j \neq w_s}} \log(\sigma(-u_{w_j}^T v_c))$$

because we want to calculate derivative from w_s so the term

which doesn't have w_s will be removed. so first and third term

are removed we just need to calculate derivative of second term.

$$\frac{\partial J_{\text{neg sampling}}}{\partial w_s} = - \sum_{\substack{1 \leq j \leq k \\ w_j = w_s}} (\sigma(-u_{w_j}^T v_c) - 1) v_c$$

i

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

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

$$iii) \frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial U_w} = 0$$

