

outside

center

No. Date

$$P(O=o|C=c) = \frac{e^{u_o^T v_c}}{\sum_{w \in \text{Vocab}} e^{u_w^T v_c}}$$

$$U = \begin{bmatrix} u_{w_1} \\ \vdots \\ u_{w_n} \end{bmatrix} \quad V = \begin{bmatrix} v_{w_1} \\ \vdots \\ v_{w_n} \end{bmatrix} \quad w \in \text{Vocab}$$

$$J_{\text{naive-softmax}} = -\log P(O=o|C=c)$$

$$J_{\text{cross-ent}} = -\sum y_w \log(\hat{y}_w) = -\log(\hat{y}_o) ; y = \text{one hot vector,}$$

$$\hat{y} = P(O|C=c) = \frac{e^{u_o^T v_c}}{\sum_w e^{u_w^T v_c}}$$

$$J_{\text{naive-softmax}} = -\log \left[ \frac{e^{u_o^T v_c}}{\sum_{w \in \text{Vocab}} e^{u_w^T v_c}} \right] = -u_o^T v_c + \log \sum_{w \in V} e^{u_w^T v_c}$$

$$J_{\text{cross-ent}} = -\log \hat{y}_o = -\log \left[ \frac{e^{u_o^T v_c}}{\sum_{w \in V} e^{u_w^T v_c}} \right]$$

$$a) J_{\text{naive-softmax}} = -\log P(O=o|C=c) = -\log(\hat{y}_o) = -\sum_{w \in V} y_w \log(\hat{y}_w) = J_{\text{cross-entropy}}$$

Since that  $y_w$  is hot vector and will 1 if  $O=o$  we get  $-\log(\hat{y}_o)$ , and otherwise 0. While  $-\log(\hat{y}_o)$  is eq. rewrite format of  $-\log(P(O=o|C=c))$ .

$$b) \frac{\partial J}{\partial v_c} = \frac{\partial}{\partial v_c} (-\log(\hat{y}_o)) = \frac{\partial}{\partial v_c} (-u_o^T v_c + \log \sum_{w \in V} e^{u_w^T v_c}) = -u_o + \frac{1}{\sum_{w \in V} e^{u_w^T v_c}} \cdot \sum_{w \in V} e^{u_w^T v_c} \cdot u_w$$

$$= -u_o + \sum_{w \in V} \frac{e^{u_w^T v_c}}{\sum_{w \in V} e^{u_w^T v_c}} u_w = -u_o + \sum_{w \in V} P(O=w|C=c) \cdot u_w = -\sum_{w \in V} y_w u_w + \sum_{w \in V} \hat{y}_w u_w$$

$$=$$

$$\frac{\partial J}{\partial v_c} = U(-y + \hat{y}) = U(\hat{y} - y) \quad 1) \frac{\partial J}{\partial v_c} = 0, \text{ if } \hat{y} = y \text{ when predicted outside word is true outside word.}$$

2) because its outstret  $v_c$  which decrease loss.

$$c) \frac{\partial J}{\partial u_w} = \frac{\partial}{\partial u_w} \left( -u_0^T \vartheta_c + \log \sum_{w \in V} e^{u_w^T \vartheta_c} \right) = \frac{\partial}{\partial u_w} (-u_0^T \vartheta_c) + \sum_{w \in V} \frac{e^{u_w^T \vartheta_c}}{\sum_{w \in V} e^{u_w^T \vartheta_c}} \cdot \vartheta_c$$

$$\text{if } w=0: \frac{\partial J}{\partial u_w} = -\vartheta_c + P(O=0|C=c) \cdot \vartheta_c = -\vartheta_c + \hat{y} \vartheta_c \quad \Rightarrow \frac{\partial J}{\partial u_w} = (\hat{y}_w - y_w) \vartheta_c \text{ for } w=1 \dots |V|$$

$$\text{if } w \neq 0 \quad \frac{\partial J}{\partial u_w} = 0 + P(O=1|C=c) \vartheta_c = \hat{y} \vartheta_c$$

$$d) \frac{\partial J}{\partial \mathbf{u}} = \left[ \frac{\partial J}{\partial u_1}, \frac{\partial J}{\partial u_2}, \dots, \frac{\partial J}{\partial u_M} \right] = [\hat{y}_1 \vartheta_c, \hat{y}_2 \vartheta_c, \dots, -\vartheta_c + \hat{y} \vartheta_c, \dots, \hat{y}_M \vartheta_c] = (\hat{\mathbf{y}} - \mathbf{y}) \vartheta_c^T$$

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

$$f) \delta(x) = \frac{e^x}{e^x + 1} \quad \left| \frac{\partial(\delta(x))}{\partial x} = \frac{(e^x)'(e^x + 1) - e^x(e^x + 1)'}{(e^x + 1)^2} = \frac{e^{2x} + e^x - e^{2x}}{(e^x + 1)^2} = \frac{e^x}{(e^x + 1)^2} = \frac{e^x}{(e^x + 1)} \cdot \frac{1}{e^x + 1} \right.$$

$$\downarrow$$

$$e^x \cdot \delta(x) + \delta(x) = e^x \quad \left| \quad = \delta(x) \cdot \frac{1}{\frac{\delta(x)}{1-\delta(x)} + 1} = \delta(x)(1-\delta(x)) \right.$$

$$e^x = \frac{\delta(x)}{1-\delta(x)}$$

$$g) J_{\text{neg-sample}}(\vartheta_c, 0, \bar{U}) = -\log(\delta(u_0^T \vartheta_c)) - \sum_{s=1}^K \log(\delta(-u_{ws}^T \vartheta_c))$$

$$g_b) \frac{\partial J_{\text{neg-sample}}}{\partial \vartheta_c} = \frac{\partial [-\log \delta(u_0^T \vartheta_c)]}{\partial \vartheta_c} - \frac{\partial \left[ \sum_{s=1}^K \log(\delta(-u_{ws}^T \vartheta_c)) \right]}{\partial \vartheta_c}$$



(i)

$$o = u_{ws}^T v_c$$

$$g_b) \frac{\partial J_{n-s}}{\partial v_c} = -\frac{1}{\delta(u_o^T v_c)} \cdot \delta(u_o^T v_c) (1 - \delta(u_o^T v_c)) \cdot u_o - \sum_{s=1}^K \frac{1}{\delta(o)} \cdot \delta(o) (1 - \delta(o)) (-u_{ws})$$

$$\frac{\partial J_{n-s}}{\partial v_c} = -(1 - \delta(u_o^T v_c)) \cdot u_o + \sum_{s=1}^K (1 - \delta(u_{ws}^T v_c)) \cdot u_{ws} \cdot (-1)$$

$$g_c) \frac{\partial J_{n-s}}{\partial u_o} = -\frac{\partial [\log \delta(u_o^T v_c)]}{\partial u_o} - \underbrace{\frac{\partial}{\partial u_o} \left[ \sum_{s=1}^K \log(\delta(-u_{ws}^T v_c)) \right]}_0$$

$$\frac{\partial J_{n-s}}{\partial u_o} = -\frac{1}{\delta(u_o^T v_c)} \cdot \cancel{\delta(u_o^T v_c)} \cdot (1 - \delta(u_o^T v_c)) \cdot v_c = \underline{-(1 - \delta(u_o^T v_c)) v_c}$$

$$g_3) \frac{\partial J_{n-s}}{\partial u_{ws}} = -\underbrace{\frac{\partial [\log \delta(u_o^T v_c)]}{\partial u_{ws}}}_0 - \frac{\partial}{\partial u_{ws}} \left[ \sum_{s=1}^K \log[\delta(-u_{ws}^T v_c)] \right]$$

$$\frac{\partial J_{n-s}}{\partial u_{ws}} = -\frac{1}{\delta(-u_{ws}^T v_c)} \cdot \cancel{\delta(-u_{ws}^T v_c)} \cdot (1 - \delta(-u_{ws}^T v_c)) \cdot (-v_c)$$

$$\frac{\partial J_{n-s}}{\partial u_{ws}} = \underline{+(1 - \delta(-u_{ws}^T v_c)) \cdot v_c}$$

(ii) to avoid repetition we can use  $[-(1 - \delta(u_o^T v_c)) v_c]$  this expression since that it is used in the two formulas. If write in terms of  $v$  it is  $u_o$

(iii) It is efficiently use memory with fixed size  $K$ . While naive softmax use whole outside words.

$$h) \frac{\partial J_{neg-sa}}{\partial u_{ws}} = -\frac{\partial}{\partial u_{ws}} \sum_{s=1}^K \log(\delta(-u_{ws}^T \vartheta_s)) + 0 = -\frac{\partial}{\partial u_{ws}} \sum_{\substack{u_{wx} = u_{ws} \\ x=1, \dots, K}}^K \log(\delta(-u_{wx}^T \vartheta_c))$$

$$= -\frac{\partial}{\partial u_{ws}} \sum_{\substack{u_{wx} = u_{ws} \\ x=1, \dots, K}}^K \log(\delta(-u_{wx}^T \vartheta_c)) =$$

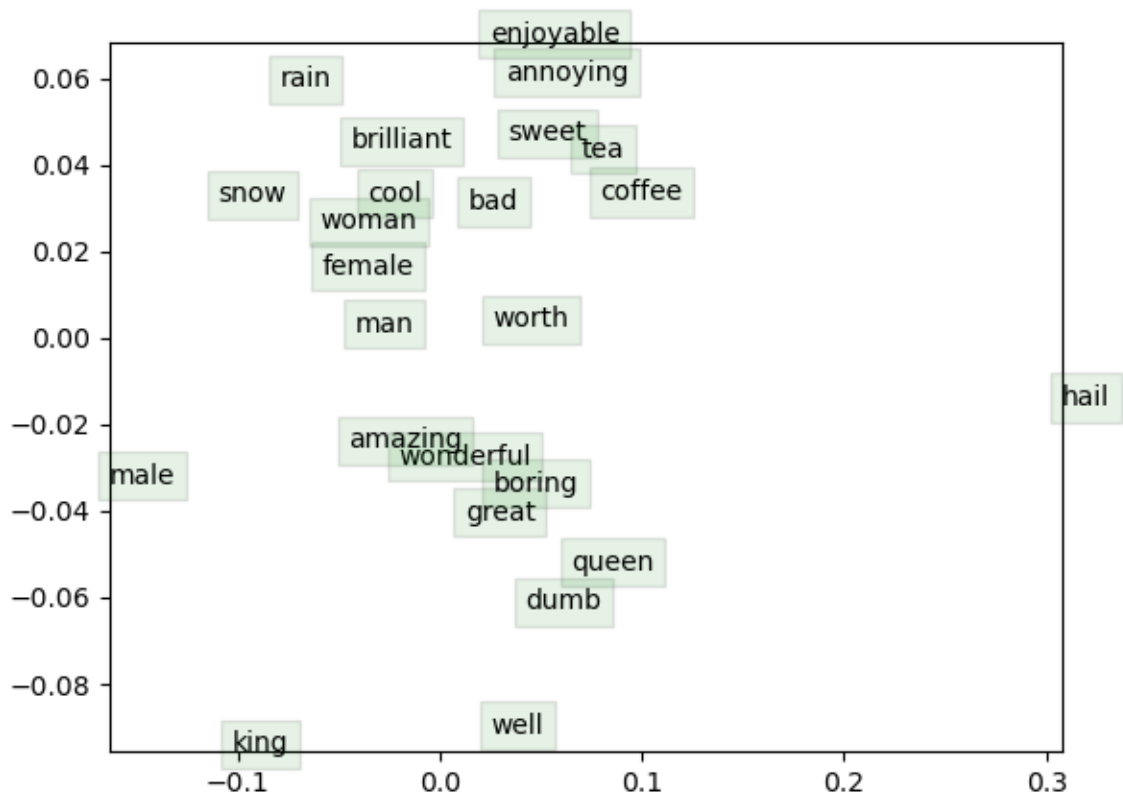
$$\frac{\partial J_{neg}}{\partial u_{ws}} = -\sum_{\substack{u_{wx} = u_{ws} \\ x=1, \dots, K}} \frac{1}{\delta(u_{wx}^T \vartheta_c)} (1 - \delta(u_{wx}^T \vartheta_c)) \cdot \delta(-u_{wx}^T \vartheta_c) \cdot (-\vartheta_c)$$

$$\frac{\partial J_{neg}}{\partial u_{ws}} = \sum_{\substack{u_{wx} = u_{ws} \\ x=1, \dots, K}} (1 - \delta(u_{wx}^T \vartheta_c)) \vartheta_c$$

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

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

$$\frac{\partial J_{skip-gram}}{\partial u_w} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(\vartheta_c, w_{t+j}, U)}{\partial u_w} = 0 \quad (w \neq c)$$



The illustration demonstrates that certain words are grouped together, but they may not always be located in the same location.

For example, the words "amazing," "wonderful," "great," and "boring" are grouped together to represent emotions, as are "cool" and "brilliant," and "enjoyable" and "annoying."

However, some of the groupings are more successful than others. For example, the words "female" and "woman," as well as "tea" and "coffee," are more closely clustered together.