

①

## 1. Understanding word2vec

(a) Since word  $o$  is the expected word, the true distribution  $y$  is a one-hot vector with the value 1 for word  $o$  and 0 otherwise.  
 Therefore  $-\sum_{w \in V} y_w \log(\hat{y}_w) = -y_o \log(\hat{y}_o) = -\log(\hat{y}_o)$

(b) Compute  $\partial J_{\text{naive-softmax}}(v_c, o, U) / \partial v_c$ , write answer in terms of  $y$ ,  $\hat{y}$ , and  $U$

First,  $J_{\text{naive-softmax}}(v_c, o, U)$

$$= -\log P(o|c)$$

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

$$\text{Then, } \frac{\partial J}{\partial v_c} = -\frac{\partial}{\partial v_c} \log \frac{\exp(u_o^T v_c)}{\sum_w \exp(u_w^T v_c)}$$

$$= - \left[ \frac{\partial(u_o^T v_c)}{\partial v_c} - \frac{\partial \log \sum_w \exp(u_w^T v_c)}{\partial v_c} \right]$$

$$= - \left[ u_o - \frac{\sum_w \exp(u_w^T v_c)}{\sum_w \exp(u_w^T v_c)} \frac{\partial \sum_w \exp(u_w^T v_c)}{\partial v_c} \right]$$

$$= - \left[ u_o - \frac{\sum_x \exp(u_x^T v_c) u_x}{\sum_w \exp(u_w^T v_c)} \right]$$

$$= - \left[ u_o - \sum_x \frac{\exp(u_x^T v_c)}{\sum_w \exp(u_w^T v_c)} u_x \right]$$

$$= - \left[ u_o - \sum_x P(x|c) u_x \right]$$

$$= - \left[ u_o - \sum_{x=1}^V \hat{y}_x u_x \right] = U(\hat{y} - y)$$

(2)

(C) Compute  $\frac{\partial J_{\text{naive-softmax}}(V_C, 0, U)}{\partial u_W}$ , write answer in terms of  $y$ ,  $\hat{y}$ , and  $V_C$

when  $w=0$

$$\frac{\partial J}{\partial u_0} = - \frac{\partial}{\partial u_0} \log \frac{\exp(u_0^T V_C)}{\sum_w \exp(u_w^T V_C)}$$

$$= - \left[ \frac{\partial (u_0^T V_C)}{\partial u_0} - \frac{\partial \log \sum_w \exp(u_w^T V_C)}{\partial u_0} \right]$$

$$= - \left( V_C - \frac{1}{\sum_w \exp(u_w^T V_C)} \frac{\partial \exp(u_0^T V_C)}{\partial u_0} \right)$$

$$= - \left( V_C - \frac{\exp(u_0^T V_C) \cdot V_C}{\sum_w \exp(u_w^T V_C)} \right)$$

$$= -(V_C - \hat{y}_0 \cdot V_C)$$

$$= (\hat{y}_0 - 1) V_C.$$

when  $w \neq 0$

$$\frac{\partial J}{\partial u_w} = - \frac{\partial}{\partial u_w} \log \frac{\exp(u_0^T V_C)}{\sum_x \exp(u_x^T V_C)}$$

$$= - \left( 0 - \frac{\partial \log \sum_x \exp(u_x^T V_C)}{\partial u_w} \right)$$

$$= \frac{1}{\sum_x \exp(u_x^T V_C)} \frac{\partial \sum_x \exp(u_x^T V_C)}{\partial u_w}$$

$$= \frac{\exp(u_w^T V_C) V_C}{\sum_x \exp(u_x^T V_C)} = \hat{y}_w \cdot V_C.$$

Combining above

$$\frac{\partial J}{\partial u_W} = \begin{cases} (\hat{y}_w - 1) V_C & w=0 \\ \hat{y}_w \cdot V_C & w \neq 0 \end{cases}$$

$$= V_C (\hat{y} - y)^T$$

(3)

(d) Compute  $\frac{\partial \sigma(x)}{\partial x}$  where  $\sigma(x)$  is the sigmoid function.

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

(e) Compute  $\frac{\partial J_{\text{neg-sampling}}}{\partial v_c}$ ,  $\frac{\partial J}{\partial u_0}$ ,  $\frac{\partial J}{\partial u_k}$

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

$$= -\frac{1}{\sigma(u_0^T v_c)} \frac{\partial \sigma(u_0^T v_c)}{\partial v_c} - \sum_k \frac{1}{\sigma(-u_k^T v_c)} \frac{\partial \sigma(-u_k^T v_c)}{\partial v_c}$$

$$= -\frac{\sigma(u_0^T v_c)(1-\sigma(u_0^T v_c)) u_0}{\sigma(u_0^T v_c)} - \sum_k \frac{\sigma(-u_k^T v_c)(1-\sigma(-u_k^T v_c)) (-u_k)}{\sigma(-u_k^T v_c)}$$

$$= (\sigma(u_0^T v_c) - 1) u_0 - \sum_k (\sigma(-u_k^T v_c) - 1) u_k$$

$$\text{Second } \frac{\partial J}{\partial u_0} = \frac{\partial [-\log \sigma(u_0^T v_c)]}{\partial u_0}$$

note that  $u_k \{k=1\dots K\}$  do not contain  $u_0$

$$= -\frac{\sigma(u_0^T v_c)(1-\sigma(u_0^T v_c)) v_c}{\sigma(u_0^T v_c)}$$

$$= (\sigma(u_0^T v_c) - 1) v_c$$

$$\text{Third. } \frac{\partial J}{\partial u_k} = \frac{\partial [-\sum_{x=1}^K \log(\sigma(-u_x^T v_c))]}{\partial u_k}$$

$$= -\frac{\partial \log(\sigma(-u_k^T v_c))}{\partial u_k} = -\frac{\sigma(-u_k^T v_c)(1-\sigma(-u_k^T v_c))}{\sigma(-u_k^T v_c) \cdot (-v_c)}$$

$$= -(\sigma(-u_k^T v_c) - 1) v_c \quad \text{for all } k=1\dots K$$

This is more efficient to compute than naive-softmax loss, which need to compute  $\hat{y}$ , which is normalized by  $\sum_{w \in V} \exp(u_w^T v_c)$ , a sum across the whole vocabulary.

(4)

$$(f) \text{ i) } \frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m} \dots w_{t+m}, u)}{\partial u}$$

$$= \frac{\partial \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} J(v_c, w_{t+j}, u)}{\partial u}$$

$$= \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(v_c, w_{t+j}, u)}{\partial u}$$

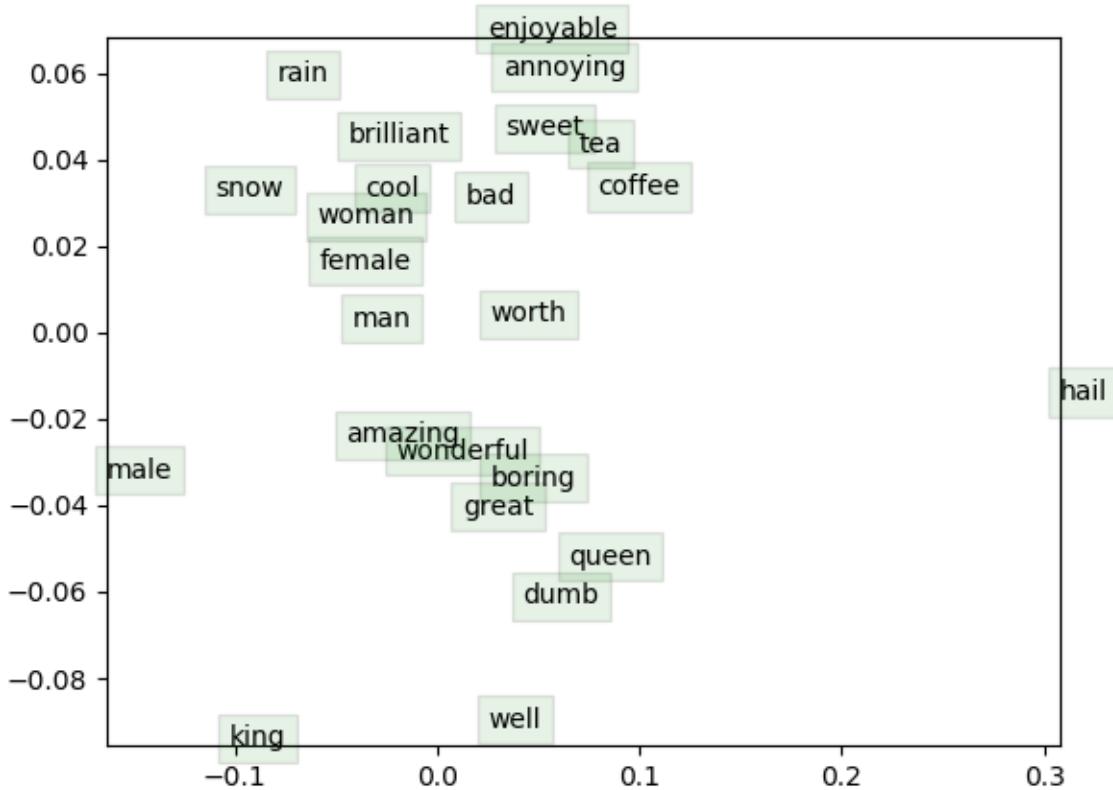
$$\text{ii) } \frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots w_{t+m}, u)}{\partial v_c}$$

$$= \frac{\partial \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} J(v_c, w_{t+j}, u)}{\partial v_c}$$

$$= \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(v_c, w_{t+j}, u)}{\partial v_c}$$

$$\text{iii) } \frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots w_{t+m}, u)}{\partial v_w} (w \neq c)$$

$$= 0$$



The above figure shows the output of training the word vectors. We can see analogies from the figure, i.e. the linearity between (queen: king) :: (female: male). We can also see synonyms that tend to have similar contexts cluster together, e.g. ("amazing," "wonderful," "great") and ("coffee," "tea"). The 2D visualization is not showing everything though—some word vectors may be close to each other in the higher dimensional space, but is not captured here.