

CS224N HW2

1. Understanding word2vec.

(a) For all $w \neq 0$, $y_w = 0$; For $w=0$ $y_w = y_0 = 1$; thus:

$$-\sum_{w \in V} y_w \log(\hat{y}_w) = 0 - 1 \cdot \log(\hat{y}_0) = -\log(\hat{y}_0)$$

$$(b) J_{\text{naive-softmax}}(v_c, 0, U) = -\log\left(\frac{\exp(u_0^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}\right)$$

$$\frac{\partial}{\partial v_c} J_{\text{naive-softmax}}(v_c, 0, U) = -\frac{\exp(u_0^T v_c) \cdot u_0}{\exp(u_0^T v_c)} + \frac{\sum_{w \in V} \exp(u_w^T v_c) u_w}{\sum_{w \in V} \exp(u_w^T v_c)}$$

$$= -u_0 + \sum_{w \in V} (\hat{y}_w \cdot u_w)$$

(1) Note that $u_0 = U \cdot y$, where $U = [u_1, u_2, \dots, u_V]$ and y is true label for word v_c and 0.

(2) Note that $\sum_{w \in V} (\hat{y}_w \cdot u_w) = U \cdot \hat{y}$, where \hat{y} is predicted values for each word

$$\therefore \frac{\partial}{\partial v_c} J_{\text{naive-softmax}}(v_c, 0, U) = -U \cdot y + U \cdot \hat{y} \\ = U(\hat{y} - y)$$

$$(c) J_{\text{naive-softmax}}(v_c, 0, U) = -\log(\exp(u_0^T v_c)) + \log\left(\sum_{w \in V} \exp(u_w^T v_c)\right)$$

• Case 1. when $u_w \neq u_0$:

$$\frac{\partial}{\partial u_w} J_{\text{naive-softmax}}(v_c, 0, U) = \frac{\exp(u_w^T v_c) \cdot v_c}{\sum_{x \in V} \exp(u_x^T v_c)} \\ = \hat{y}_w \cdot v_c \quad (d,)$$

• Case 2. When $u_w = u_0$:

$$\frac{\partial}{\partial u_0} J_{\text{naive-softmax}}(v_c, 0, U) = -\frac{\exp(u_0^T v_c) \cdot v_c}{\exp(u_0^T v_c)} + \frac{\exp(u_0^T v_c) v_c}{\sum_{w \in V} \exp(u_w^T v_c)} \\ = (\hat{y}_0 - y_0) v_c$$

$$(d) \frac{\partial}{\partial U} J_{\text{naive-softmax}}(V_c, o, U) = \left[\frac{\partial J(V_c, o, U)}{\partial u_1}, \frac{\partial J(V_c, o, U)}{\partial u_2}, \dots, \frac{\partial J(V_c, o, U)}{\partial u_{|Vocab|}} \right]$$

(e) Let $z = e^{-x}$, then:

$$\begin{aligned} \frac{d}{dx} \sigma(x) &= \frac{d}{dz} \sigma(z) \cdot \frac{dz}{dx} = -\frac{1}{(1+z)^2} \cdot (-e^{-x}) \\ &= \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1}{1+e^x} \cdot \sigma(x) = \sigma(x)(1-\sigma(x)) \end{aligned}$$

$$(f) J_{\text{neg-sample}}(V_c, o, U) = -\log(\sigma(u_o^T V_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T V_c))$$

$$\begin{aligned} ① \quad \frac{\partial}{\partial V_c} J(V_c, o, U) &= -\frac{\sigma(u_o^T V_c)(1-\sigma(u_o^T V_c))}{\sigma(u_o^T V_c)} \cdot u_o^T - \sum_{k=1}^K \frac{\sigma(-u_k^T V_c)(1-\sigma(-u_k^T V_c))}{\sigma(-u_k^T V_c)} \cdot (-u_k^T) \\ &= -u_o^T (1-\sigma(u_o^T V_c)) + \sum_{k=1}^K u_k^T (1-\sigma(-u_k^T V_c)) \end{aligned}$$

$$② \quad \frac{\partial}{\partial u_o} J(V_c, o, U) = -\frac{\sigma(u_o^T V_c)}{\sigma(u_o^T V_c)} (1-\sigma(u_o^T V_c)) \cdot V_c^T = -(1-\sigma(u_o^T V_c)) \cdot V_c^T$$

$$③ \quad \frac{\partial}{\partial u_k} J(V_c, o, U) = -\frac{\sigma(-u_k^T V_c)}{\sigma(-u_k^T V_c)} (1-\sigma(-u_k^T V_c)) \cdot (-V_c^T) = (1-\sigma(-u_k^T V_c)) \cdot V_c^T$$

$$(g) \frac{\partial}{\partial U} J(V_c, o, U) = \sum_{u_w = u_k} \frac{\sigma(-u_w^T V_c)}{\sigma(-u_w^T V_c)} (1-\sigma(-u_w^T V_c)) \cdot (-V_c^T)$$

Assume there are n samples that equals to u_k .

$$= n(1-\sigma(-u_k^T V_c)) \cdot (-V_c^T)$$

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

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

$$(iii) \frac{\partial}{\partial V_w} J_{\text{skip-gram}} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0 \\ \text{if } t+j = w}} \frac{\partial}{\partial V_w} J(V_c, W_{t+j}, U)$$

(c) Some analogies can be seen like "man" - "king" v.s "woman" - "queen". Some clusters of adj., such as "great", "wonderful" and "amazing" are clustered together. Positive adjectives and negative ones like "boring" are not divided in low dimension visualization.