

CS 224n Assignment #2: word2vec

1 Written: Understanding Word2vec

$$(a). - \sum_{w \in \text{Vocab}} y_w \log(\hat{y}_w) = -\log(\hat{y}_0)$$

y_w is a scalar, y is a vector

y_w is y 中第 w 个词对应的标量
 \hat{y}_w 是 y 中第 w 个词计算出的 $P(w|c)$.

$$\text{即: } - \sum_{w \in V} y_w \log(\hat{y}_w) = -y_0 \log(\hat{y}_0) = -\log(\hat{y}_0)$$

$$(b). \frac{\partial}{\partial v_c} J_{\text{naive-softmax}}(v_c, o, U)$$

answers in terms of y, \hat{y}, U

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

$$= -\log \frac{\exp(u_0^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)} = -[\log \exp(u_0^T v_c) - \log \sum_{w \in V} \exp(u_w^T v_c)]$$

$$= -u_0 + \frac{\partial}{\partial v_c} \log \sum_{w \in V} \exp(u_w^T v_c)$$

rows
columns of U : outside
vectors u_w

$$= -u_0 + \sum_{x \in V} \frac{\exp(u_x^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)} \cdot u_x$$

$$y = U^T v_c$$

$$= -u_0 + \sum_{x \in V} P(x|c) u_x$$

$$u_0 = U^T y$$

$$= -y_0^T \cdot U^T + \sum \hat{y}_x^T \cdot U^T$$

$$P(x|c) = \hat{y}_x$$

$$= \sum_{x \in V} (\hat{y}_x - y_0)^T \cdot U^T$$

$$\sum_{x \in V} P(x|c) \cdot u_x = U^T \hat{y}$$

$$= -U^T y + U^T \hat{y}$$

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

(y_1, y_2, \dots, v)

$$(c). \frac{\partial}{\partial u_w} J_{\text{naive-softmax}}(v_c, o, U)$$

answers: in terms of y, \hat{y}, v_c

$$O. w=0$$

$$J(v_c, o, U) = -u_0^T v_c + \log \sum_{w \in V} \exp(u_w^T v_c)$$

$$\frac{\partial}{\partial u_0} J(v_c, o, U) = -v_c + \frac{\partial}{\partial u_0} \log \sum_{w \in V} \exp(u_w^T v_c)$$

$$= -v_c + \frac{1}{\sum_{w \in V} \exp(u_w^T v_c)} \cdot \sum_{w \neq 0} \exp(u_w^T v_c) \cdot \frac{\partial \exp(u_w^T v_c)}{\partial u_0}$$

$$= -V_c + \frac{1}{\sum_{w \in V} \exp(u_w^T V_c)} \cdot \sum_{\substack{w \in V \\ w \neq 0}} 0$$

$$= -V_c + \sum_{\substack{w \in V \\ w \neq 0}} \frac{\exp(u_w^T V_c)}{\sum_{w \in V} \exp(u_w^T V_c)}$$

$$= -V_c + \sum_{w \in V} P(w|c) \cdot 1$$

$$= -V_c + \sum_{w \in \text{Vocab}} \hat{y}_w V_c$$

$P(c|c)$

$$(d). \frac{\partial J(V_c, 0, U)}{\partial U} \quad \text{answer}$$

$$\textcircled{2}. w \neq 0. \quad \hat{x} = u_x = u_w. \quad x \neq 0.$$

$$\frac{\partial}{\partial u_x} J_{\text{naive-softmax}}(V_c, 0, U)$$

$$J_{\text{naive-softmax}} = -u_0^T V_c + \log \sum_{w \in V} \exp(u_w^T V_c)$$

$$= -u_0^T V_c + \log [\exp(u_x^T V_c) + \sum_{\substack{w \in V \\ w \neq x}} \exp(u_w^T V_c)]$$

$$\frac{\partial J}{\partial u_x} = \frac{\partial \log [\exp(u_x^T V_c) + \sum_{\substack{w \in V \\ w \neq x}} \exp(u_w^T V_c)]}{\partial u_x}$$

$$= \frac{1}{\log \sum_{w \in V} \exp(u_w^T V_c)} \cdot \exp(u_x^T V_c) \cdot V_c$$

$$= P(x|c) \cdot V_c$$

$$= \hat{y}_x \cdot V_c$$

$$\hat{y}_p \hat{y}_w \cdot V_c.$$

$$= \frac{\partial}{\partial U} (-u_0^T V_c + \log \exp \sum_{w \in V} u_w^T V_c)$$

$$= \left(\frac{\partial J}{\partial u_1}, \frac{\partial J}{\partial u_2}, \dots, \frac{\partial J}{\partial u_{|\text{Vocab}|}} \right)^T$$

$$= (\hat{y}_1 \cdot V_c, \hat{y}_2 \cdot V_c, \dots, \hat{y}_{|\text{Vocab}|} \cdot V_c)^T = \begin{pmatrix} \hat{y}_1 \cdot V_c \\ \hat{y}_2 \cdot V_c \\ \vdots \\ \hat{y}_{|\text{Vocab}|} \cdot V_c \end{pmatrix}$$

$$(e). \quad \sigma(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x+1} \quad x \text{ is a scalar. answer: } \sigma(x) \quad \text{in terms of}$$

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

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

$$(f). \quad J_{\text{neg-sample}}(V_c, 0, U) = -\log(\sigma(u_0^T V_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T V_c))$$

answer in terms of u_0, V_c, u_k

$$\frac{\partial J_{\text{neg-sample}}(V_c, 0, U)}{\partial V_c} = -[1 - \sigma(u_0^T V_c)] u_0 + \sum_{k=1}^K [1 - \sigma(-u_k^T V_c)] \cdot u_k$$

$$= - \frac{u_0}{\exp(u_0^T V_c) + 1} + \sum_{k=1}^K \frac{u_k}{\exp(-u_k^T V_c) + 1}$$

$$\frac{\partial J_{\text{neg-sample}}(V_c, 0, U)}{\partial u_0} = -[1 - \sigma(u_0^T V_c)] V_c + \sum_{k=1}^K [1 - \sigma(-u_k^T V_c)] \cdot V_c$$

neg-sample's loss function is more efficient than naive-softmax

cuz no $\hat{y}(\Sigma)$

(g). $\frac{\partial J_{\text{neg-sample}}}{\partial u_k}$ answers in terms of $V_c, u_k, k \in [1, K]$

$$J_{\text{neg-sample}}(V_c, 0, U) = -\log(\sigma(c^T V_c)) - \sum_{\substack{x=1 \\ u_x \neq u_k}}^K \log(\sigma(c - u_k^T V_c)) - \sum_{x=1}^K \log(\sigma(-u_x^T V_c))$$

$$\frac{\partial J_{\text{neg-sample}}}{\partial u_k} = [\sigma(c - u_k^T V_c) - 1] \cdot V_c$$

skip-gram of word2vec

(h). suppose $c = w_t$. context window size = m .

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

$$J(V_c, w_{t+j}, U) : J_{\text{naive-softmax}}(V_c, w_{t+j}, U) \text{ or } J_{\text{neg-sample}}(V_c, w_{t+j}, U)$$

answers in terms of $\frac{\partial J(V_c, w_{t+j}, U)}{\partial U}$, $\frac{\partial J(V_c, w_{t+j}, U)}{\partial V_c}$

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

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

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