

18210980068 王 1/2

Assignment 4

I word2vec

(a) derive the gradients with respect to V_c .

$$\log P(o|c) = \log \frac{\exp(u_o^T V_c)}{\sum_{w=1}^V \exp(u_w^T V_c)} = \log \exp(u_o^T V_c) - \log \sum_{w=1}^V \exp(u_w^T V_c)$$

$$\frac{\partial}{\partial V_c} \log \exp(u_o^T V_c) = \frac{\partial}{\partial V_c} u_o^T V_c = u_o$$

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

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

$$\frac{\partial \log P(o|c)}{\partial V_c} = u_o - \frac{\sum_{w=1}^V \exp(u_w^T V_c) \cdot u_w}{\sum_{w=1}^V \exp(u_w^T V_c)}$$

$$= u_o - \sum_{t=1}^V \frac{\exp(u_t^T V_c)}{\sum_{w=1}^V \exp(u_w^T V_c)} \cdot u_t$$

$$= u_o - \sum_{t=1}^V P(t|c) u_t$$

$$\frac{\partial J_{\text{softmax-CE}}(o, V_c, V)}{\partial V_c} = - \frac{\partial \log P(o|c)}{\partial V_c}$$

$$= -u_o + \sum_{t=1}^V P(t|c) u_t$$

(b) $\forall w \neq o$ 时,

$$\frac{\partial u_o^T V_c}{\partial u_w} = 0$$

$$\frac{\partial \log \sum_{w=1}^V \exp(u_w^T V_c)}{\partial u_w} = \frac{1}{\sum_{t=1}^V \exp(u_t^T V_c)} \cdot \frac{\partial \sum_{w=1}^V \exp(u_w^T V_c)}{\partial u_w}$$

(1)

$$= \frac{1}{\sum_{t=1}^V \exp(u_t^T V_c)} \cdot \frac{\partial \exp(u_w^T V_c)}{\partial u_w}$$

$$= \frac{1}{\sum_{t=1}^V \exp(u_t^T V_c)} \cdot \exp(u_w^T V_c) \cdot V_c$$

$$\frac{\partial}{\partial u_w} \log(P(o|c)) = 0 - \frac{\exp(u_w^T V_c)}{\sum_{t=1}^V \exp(u_t^T V_c)} \cdot V_c$$

$$\begin{aligned} \frac{\partial}{\partial u_w} J_{\text{softmax-CE}} &= -\frac{\partial}{\partial u_w} \log(P(o|c)) \\ &= \frac{\exp(u_w^T V_c)}{\sum_{t=1}^V \exp(u_t^T V_c)} V_c \quad (w \neq 0) \end{aligned}$$

for $w=0$,

$$\frac{\partial u_0^T V_c}{\partial u_0} = V_c$$

$$\frac{\partial}{\partial u_0} J_{\text{softmax-CE}} = -V_c + \frac{\exp(u_0^T V_c)}{\sum_{t=1}^V \exp(u_t^T V_c)} \cdot V_c \quad (w=0).$$

(c)

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

$$\frac{\partial \log(\sigma(u_0^T V_c))}{\partial V_c} = \frac{1}{\sigma(u_0^T V_c)} \cdot \frac{\partial \sigma(u_0^T V_c)}{\partial V_c} = \frac{1}{\sigma(u_0^T V_c)} \cdot \sigma(u_0^T V_c) \cdot (1 - \sigma(u_0^T V_c)) \cdot u_0$$

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

$$\frac{\partial \sum_{k=1}^K \log(\sigma(-u_k^T V_c))}{\partial V_c} = \sum_{k=1}^K \frac{\partial \log(\sigma(-u_k^T V_c))}{\partial V_c} = \sum_{k=1}^K \frac{1}{\sigma(-u_k^T V_c)} \cdot \sigma(-u_k^T V_c) \cdot (1 - \sigma(-u_k^T V_c)) \cdot (-u_k)$$

$$= \sum_{k=1}^K (1 - \sigma(-u_k^T V_c)) \cdot (-u_k)$$

(2)

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

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

(C.12)

$$\frac{\partial \log(\sigma(u_0^T V_c))}{\partial u_k} = 0$$

$k \neq 0$ 时,

$$\frac{\partial \sum_{k=1}^K \log(\sigma(-u_k^T V_c))}{\partial u_k} = \sum_{k=1}^K \frac{\partial \log(\sigma(-u_k^T V_c))}{\partial u_k} = \frac{\partial \log(\sigma(-u_k^T V_c))}{\partial u_k}$$

$$= \frac{1}{\sigma(-u_k^T V_c)} \cdot \sigma(-u_k^T V_c) (1 - \sigma(-u_k^T V_c)) (-V_c)$$

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

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

$$k=0 \text{ 时, } \frac{\partial \log(\sigma(u_0^T V_c))}{\partial u_0} = \frac{1}{\sigma(u_0^T V_c)} \cdot \sigma(u_0^T V_c) (1 - \sigma(u_0^T V_c)) \cdot V_c$$

$$= (1 - \sigma(u_0^T V_c)) \cdot V_c$$

$$\frac{\partial J_{\text{neg-sample}}}{\partial u_0} = (\sigma(u_0^T V_c) - 1) V_c + (1 - \sigma(-u_k^T V_c)) \cdot V_c$$

(C.2)

In the softmax-CE loss, the skip-gram neural network updates times as many as the size of word vocabulary. The cost of computing $P(u_0 | V_c)$ is proportional with the vocabulary size V .

In the negative sampling loss, we only update certain number weights which is a small percentage of all the possible words.

$$\frac{\partial J_{\text{skip-gram}}(\text{Word } c-m, \dots, c+m)}{\partial V_c}$$

$$= -\sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(W_{c+j}, V_c)}{\partial V_c}$$

$$\frac{\partial J_{\text{skip-gram}}(\text{Word } c-m, \dots, c+m)}{\partial U_{W_{c+j}}}$$

$$= \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(W_{c+j}, V_c)}{\partial U_{W_{c+j}}}$$

$\frac{\partial}{\partial V_c} F(W_{c+j}, V_c)$ 和 $\frac{\partial}{\partial U_{W_{c+j}}} F(W_{c+j}, V_c)$ 在 (a) ~ (c) 中已经表示出.