

# Written part Linsen Li

## (a) Softmax

Assume  $\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}_{n \times 1}$ , then  $\text{softmax}(\vec{x}) = \begin{bmatrix} \frac{e^{x_1}}{\sum_j e^{x_j}} \\ \frac{e^{x_2}}{\sum_j e^{x_j}} \\ \vdots \\ \frac{e^{x_n}}{\sum_j e^{x_j}} \end{bmatrix}_{n \times 1}$ , that is,  $\text{softmax}(\vec{x})_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$

$$\text{So, } \text{softmax}(\vec{x} + c)_i = \frac{e^{x_i + c}}{\sum_j e^{x_j + c}} = \frac{e^{x_i} \cdot e^c}{\sum_j e^{x_j} \cdot e^c} = \frac{e^c \cdot e^{x_i}}{e^c \cdot \sum_j e^{x_j}}$$

$$\text{Since } e^c \neq 0, \text{ then } \text{softmax}(\vec{x} + c)_i = \frac{e^{x_i}}{\sum_j e^{x_j}} = \text{softmax}(\vec{x})_i$$

Then,  $\text{softmax}(\vec{x}) = \text{softmax}(\vec{x} + c)$  for any vector  $\vec{x}$  and constant  $c$ .

## (b) Sigmoid

$$\text{Set } M = 1 + e^{-x}, \text{ then } G(x) = \frac{1}{M}, \quad \frac{\partial G}{\partial M} = -\frac{1}{M^2} = -\frac{1}{(1+e^{-x})^2}$$

$$\text{then } \frac{\partial M}{\partial x} = e^{-x} \cdot (-1) = -e^{-x}$$

$$\therefore \frac{\partial G}{\partial x} = \frac{\partial M}{\partial x} \cdot \frac{\partial G}{\partial M} = -e^{-x} \times \left(-\frac{1}{(1+e^{-x})^2}\right)$$

$$= \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1}{1+e^{-x}} \times \frac{e^{-x}}{1+e^{-x}}$$

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

$$= G(x) \times (1 - G(x))$$

$$\therefore G'(x) = G(x)(1 - G(x))$$

(c) Word2vec

(1) For  $v_c$ :

Since  $y$  is a one-hot vector, that is  $y = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix}$  -  $i$ -th

then when  $o$  is the expected word,  $\sum_i y_i \log(\hat{y}_i) = 1 \times \log(\hat{y}_o) = \log(\hat{y}_o)$

$$\therefore J_{CE}(u, v_c, U) = -\log(\hat{y}_o) = -\log \frac{\exp(u_o^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)}$$

$$= -u_o^T v_c + \log \sum_{w=1}^W \exp(u_w^T v_c)$$

$$\frac{\partial u_o^T v_c}{\partial v_c} = u_o, \text{ which is trivial}$$

$$\text{Set } M = \sum_{w=1}^W \exp(u_w^T v_c), \quad Y = \log M.$$

$$\text{then } \frac{\partial M}{\partial v_c} = \sum_{w=1}^W \exp(u_w^T v_c) \cdot u_w, \quad \frac{\partial Y}{\partial M} = \frac{1}{M}$$

$$\therefore \frac{\partial Y}{\partial v_c} = \frac{\partial M}{\partial v_c} \cdot \frac{\partial Y}{\partial M} = \sum_{w=1}^W \exp(u_w^T v_c) \cdot u_w / \sum_{w=1}^W \exp(u_w^T v_c)$$

$$\therefore \frac{\partial J}{\partial v_c} = -u_o + \sum_{x=1}^W \frac{\exp(u_x^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)} \cdot u_x$$

$$= -u_o + \sum_{x=1}^W p(x|c) \cdot u_x$$

Since  $\hat{y} = \begin{bmatrix} p(1|c) \\ p(2|c) \\ \vdots \\ p(w|c) \end{bmatrix}, \quad U = [u_1, u_2, \dots, u_W], \quad u_o = [u_1, u_2, \dots, u_W] \cdot \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix}$  -  $o$ -th.

$$\therefore \sum_{x=1}^W p(x|c) \cdot u_x = U \cdot \hat{y}, \quad -u_o = U \cdot y$$

$$\therefore \frac{\partial J_{CE}(u, v_c, U)}{\partial v_c} = U (\hat{y} - y)$$

(2) For  $u_w$

$$\text{Since } J_{CE}(o, v_c, u) = -u^T v_c + \log \sum_{w=1}^W \exp(u_w^T v_c)$$

When  $w=0$ :

$$\begin{aligned}\frac{\partial J_{CE}(o, v_c, u)}{\partial u_w} &= -v_c + \frac{1}{\sum_{x=1}^W \exp(u_x^T v_c)} \times \exp(u_w^T v_c) \times v_c \\ &= -v_c + v_c \cdot \frac{\exp(u_w^T v_c)}{\sum_{x=1}^W \exp(u_x^T v_c)} \\ &= -v_c + v_c \cdot P(w|c) \\ &= v_c(\hat{y}_w - 1)\end{aligned}$$

when  $w \neq 0$ : the only difference is the first term

$$\frac{\partial J_{CE}(o, v_c, u)}{\partial u_w} = 0 + v_c \cdot P(w|c) = v_c \cdot \hat{y}_w$$

$$\therefore \frac{\partial J_{CE}(o, v_c, u)}{\partial u_w} = \begin{cases} (\hat{y}_w - 1) v_c, & w=0 \\ \hat{y}_w v_c, & w \neq 0 \end{cases}$$

(3)

For negative sampling, and  $G(x) = G(x)(1-G(x))$

$$J_{\text{neg-sample}}(o, v_c, u) = -\log(G(u_o^T v_c)) - \sum_{k=1}^K \log(G(-u_k^T v_c))$$

① For  $v_c$ :

$$\frac{\partial J_{\text{neg-sample}}}{\partial v_c} = -\frac{G(u_o^T v_c) \cdot (1 - G(u_o^T v_c))}{G(u_o^T v_c)} \cdot u_o - \sum_{k=1}^K \frac{G(-u_k^T v_c) \cdot (1 - G(-u_k^T v_c))}{G(-u_k^T v_c)} \cdot (-u_k)$$

Since  $G(\cdot) \neq 0$  for all domain.

$$= u_0 [G(u_0^T v_c) - 1] - \sum_{k=1}^K (G(-u_k^T v_c) - 1) \cdot u_k$$

(2) For  $u_w$ :

When  $w=0$ : since  $0 \notin \{1, 2, \dots, K\}$ , so  $u_w$  now is irrelevant with 2nd term.

$$\begin{aligned}\therefore \frac{\partial J_{\text{neg-sample}}}{\partial u_w} &= - \frac{G(u_0^T v_c)(1 - G(u_0^T v_c))}{G(u_0^T v_c)} \cdot v_c - 0 \\ &= (G(u_0^T v_c) - 1) \cdot v_c\end{aligned}$$

When  $w \neq 0$ : then  $u_w$  is irrelevant with 1st term

$$\begin{aligned}\therefore \frac{\partial J_{\text{neg-sample}}}{\partial u_w} &= \frac{\partial}{\partial u_w} \left( - \sum_{k=1}^K \log(G(-u_k^T v_c)) \right) \\ &= \frac{\partial}{\partial u_w} \left[ -\log(G(-u_w^T v_c)) \right] + \frac{\partial}{\partial u_w} \left[ - \sum_{k \neq w} \log(G(-u_k^T v_c)) \right] \\ &= - \frac{G(-u_w^T v_c)(1 - G(-u_w^T v_c))}{G(-u_w^T v_c)} \cdot (-v_c) + 0 \\ &= -(G(-u_w^T v_c) - 1) \cdot v_c, \text{ for all } w = 1, 2, \dots, K\end{aligned}$$

In conclusion:

$$\frac{\partial J_{\text{neg-sample}}}{\partial v_c} = [G(u_0^T v_c) - 1] \cdot u_0 - \sum_{k=1}^K (G(-u_k^T v_c) - 1) \cdot u_k$$

$$\frac{\partial J_{\text{neg-sample}}}{\partial u_w} = \begin{cases} (G(u_0^T v_c) - 1) \cdot v_c, & \text{when } w=0 \\ -[G(-u_w^T v_c) - 1] \cdot v_c, & \text{for all } w=1, 2, \dots, K \\ & \text{and } w \neq 0 \end{cases}$$

(4)

Since we have computed

$$\frac{\partial F(w_i, v_c)}{\partial U} \text{ and } \frac{\partial F(w_i, v_c)}{\partial V_c}$$

So, for skip-gram, the gradients for the cost of one context window are

$$\frac{\partial J_{\text{skip-gram}}(w_{t-m}, \dots, w_{t+m})}{\partial U} = \sum_{-m < j < m, j \neq 0} \frac{\partial F(w_{t+j}, v_c)}{\partial U}$$

$$\frac{\partial J_{\text{skip-gram}}(w_{t-m}, \dots, w_{t+m})}{\partial V_c} = \sum_{-m < j < m, j \neq 0} \frac{\partial F(w_{t+j}, v_c)}{\partial V_c}$$

$$\frac{\partial J_{\text{skip-gram}}(w_{t-m}, \dots, w_{t+m})}{\partial v_{w_{t+j}}} = 0, \text{ for all } j \neq l.$$