

1. a) when  $w \neq 0$ ,  $y_w = 0$ , so  $-\sum_{w \in V \times C} y_w \log \hat{y}_w = -y_0 \log \hat{y}_0 = -\log \hat{y}_0$ .

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

$$\begin{aligned} \frac{\partial \bar{J}}{\partial v_c} &= -u_0 + \frac{1}{\sum_{w \in V} \exp(u_w^T v_c)} \cdot \left[ \sum_{w \in V} u_w \cdot \exp(u_w^T v_c) \right] \\ &= -u_0 + \frac{\sum_{w \in V} u_w \cdot \exp(u_w^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)} = -u_0 + \sum_{w \in V} u_w \cdot P(0=w | c=c) \\ &= -u_0 + \sum_{w \in V} u_w \cdot \hat{y}_w = V(\hat{y} - y) \end{aligned}$$

(1)  $\frac{\partial \bar{J}}{\partial v_c} = 0 \Rightarrow \hat{y} = y.$

(2) If zero, then the parameter updating stops.

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

$w=0$ :

$$\begin{aligned} \frac{\partial \bar{J}}{\partial u_w} &= -v_c + \frac{1}{\sum_{w \in V} \exp(u_w^T v_c)} \cdot v_c \cdot \exp(u_w^T v_c) \\ &= -v_c + v_c \cdot \frac{\exp(u_w^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)} \\ &= -v_c + v_c \cdot P(w|c) = -v_c + v_c \cdot \hat{y}_w = v_c(\hat{y}_w - 1) \end{aligned}$$

$w \neq 0$ :  $\frac{\partial \bar{J}}{\partial u_w} = v_c \cdot \hat{y}_w.$

d) 
$$\nabla \bar{J} = \begin{pmatrix} \frac{\partial \bar{J}}{\partial u_0} \\ \vdots \\ \frac{\partial \bar{J}}{\partial u_w} \end{pmatrix}$$

e) 
$$\frac{\partial f(x)}{\partial x} = \begin{cases} 1 & x > 0 \\ 0 & x < 0 \end{cases}$$

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

g)  $J = -\log(\sigma(u_s^T v_c)) - \sum_{s=1}^K \log(\sigma(-u_{w_s}^T v_c))$

(i)  $\frac{\partial J}{\partial v_c} = -\frac{1}{\sigma(u_s^T v_c)} \cdot \sigma(u_s^T v_c)(1 - \sigma(u_s^T v_c)) \cdot u_s - \sum_{s=1}^K (-u_{w_s}) \cdot \frac{1}{\sigma(-u_{w_s}^T v_c)} \cdot \sigma(-u_{w_s}^T v_c)(1 - \sigma(-u_{w_s}^T v_c))$

$= -(1 - \sigma(u_s^T v_c))u_s + \sum_{s=1}^K u_{w_s} (1 - \sigma(-u_{w_s}^T v_c))$

$\frac{\partial J}{\partial u_s} = -\frac{1}{\sigma(u_s^T v_c)} \cdot \sigma(u_s^T v_c) \cdot (1 - \sigma(u_s^T v_c)) \cdot v_c - 0 = -(1 - \sigma(u_s^T v_c)) \cdot v_c$

$\frac{\partial J}{\partial u_{w_s}} = -\frac{1}{\sigma(-u_{w_s}^T v_c)} \cdot \sigma(-u_{w_s}^T v_c) \cdot (1 - \sigma(-u_{w_s}^T v_c)) \cdot (-v_c)$

$= v_c (1 - \sigma(-u_{w_s}^T v_c))$

(ii). we can reuse  $\sigma(V^T v_c)$

(iii) Because it takes only  $K$  negative samples into account instead of the whole corpus while computing the gradient. it's faster.

(h)  $J = -\log(\sigma(u_s^T v_c)) - \sum_{s=1}^K \log(\sigma(-u_{w_s}^T v_c))$

$\frac{\partial J}{\partial u_{w_s}} = \sum_{\substack{i=1 \\ w_s = w_i}}^K v_c (1 - \sigma(-u_{w_i}^T v_c))$

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

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

$\frac{\partial J(v_c, w_{t-m}, \dots, w_{t+m}, v)}{\partial v_w} = 0, \text{ when } w \neq c.$

2. see codes.