

CS224n Assignment 2 Written

(a) $y_w = 0$ everywhere except $w = 0 \rightarrow -\sum y_w \log(\hat{y}_w) = -y_0 \log(\hat{y}_0)$.
 Since $y_0 = 1$, $-y_0 \log(\hat{y}_0) = -\log(\hat{y}_0) \sum_{w \in \text{Vocab}}$

(b) $J = -\log\left(\frac{\exp(u_0^T v_c)}{\sum_w \exp(u_w^T v_c)}\right)$; \sum_w denotes $\sum_{w \in \text{Vocab}}$

$$= -(u_0^T v_c) + \log\left(\sum_w \exp(u_w^T v_c)\right)$$

$$\frac{\partial J}{\partial v_c} = -u_0 + \frac{1}{\sum_w \exp(u_w^T v_c)} \cdot \sum_w \exp(u_w^T v_c) \cdot u_w$$

$$u_0 = v_y \leftarrow -u_0 + \sum_w u_w \cdot \left(\frac{\exp(u_w^T v_c)}{\sum_w \exp(u_w^T v_c)} \right) \rightarrow \hat{y}_w$$

$$= -v_y + \sum_w u_w \cdot \hat{y}_w = -v_y + v_{\hat{y}} = \boxed{v(\hat{y} - y)}$$

(1) Gradient is zero when $y = \hat{y}$

(2) Since we're optimizing v_c to minimize J , we're attempting to maximize $\frac{\exp(u_0^T v_c)}{\sum_w \exp(u_w^T v_c)}$ which primarily involves increasing

the dot product between u_0 and v_c , which requires bringing v_c closer to u_0 . This is also evident in the intermediate steps of deriving $\frac{\partial J}{\partial v_c}$, where subtracting $-u_0$ implies bringing v_c closer to u_0 .

(c) $\frac{\partial J}{\partial u_i} = \begin{cases} i=0: -v_c + \left(\frac{1}{\sum_w \exp(u_w^T v_c)} \cdot \exp(u_i^T v_c) \right) \cdot v_c = (\hat{y}_i - 1)v_c \\ i \neq 0: \frac{1}{\sum_w \exp(u_w^T v_c)} \cdot \exp(u_i^T v_c) \cdot v_c = \hat{y}_i v_c \end{cases}$
 $= (\hat{y}_i - y_i) v_c$

(d) $\frac{\partial J}{\partial v} = [v_c(\hat{y}_1 - y_1), v_c(\hat{y}_2 - y_2), \dots, v_c(\hat{y}_{|\text{Vocab}|} - y_{|\text{Vocab}|})]$
 $= v_c (\hat{y} - y)^T$

$$(e) \frac{df}{dx} = \begin{cases} 0, & x < 0 \\ 1, & x > 0 \end{cases}$$

$$(f) \frac{d\sigma}{dx} = \frac{(e^x+1) \cdot e^x - e^x \cdot e^x}{(e^x+1)^2}$$

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

$$(g)(i) \frac{\partial J}{\partial v_c} = - \frac{1}{\sigma(u_0^T v_c)} \cdot \cancel{\sigma(u_0^T v_c)} \cdot (1-\sigma(u_0^T v_c)) \cdot u_0$$

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

$$= \boxed{-(1-\sigma(u_0^T v_c)) u_0 + \sum_{s=1}^K (1-\sigma(-u_{w_s}^T v_c)) u_{w_s}}$$

$$\frac{\partial J}{\partial u_0} = - \frac{1}{\cancel{\sigma(u_0^T v_c)}} \cdot \cancel{\sigma(u_0^T v_c)} \cdot (1-\sigma(u_0^T v_c)) \cdot v_c$$

$$= -(1-\sigma(u_0^T v_c)) v_c = \boxed{(\sigma(u_0^T v_c) - 1) v_c}$$

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

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

$$(ii) 1 - \sigma(v^T v_c)$$

(iii) It only needs to backpropagate to K instead of $|\text{Vocab}|$ outside vectors.

$$(h) \frac{\partial J}{\partial u_{w_i}} = \frac{\partial}{\partial u_{w_i}} (-\log(\sigma(u_0^T v_c))) - \frac{\partial}{\partial u_{w_i}} \sum_{\substack{s=1 \\ w_i \neq w_s}}^K \log(\sigma(-u_{w_s}^T v_c))$$

$$- \frac{\partial}{\partial u_{w_i}} \sum_{\substack{s=1 \\ w_i \neq w_s}}^K \log(\sigma(-u_{w_0}^T v_c))$$

$$= 0 - \sum_{\substack{s=1 \\ w_i \neq w_s}}^K \frac{1}{\cancel{\sigma(-u_{w_s}^T v_c)}} \cdot \cancel{\sigma(-u_{w_s}^T v_c)} \cdot (1-\sigma(-u_{w_s}^T v_c)) \cdot (-v_c)$$

$$= \sum_{\substack{s=1 \\ w_i \neq w_s}}^K 0$$

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

$$(i) \quad \frac{\partial J(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial U} \quad (iii) \quad \frac{\partial J(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial v_w}$$

$$= \left[\sum_{\substack{-n \leq j \leq m \\ j \neq 0}} \frac{\partial J(v_c, w_{t+j}, U)}{\partial U} \right] = \boxed{0}$$

$$(ii) \quad \frac{\partial J(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial v_c}$$

$$= \left[\sum_{\substack{-n \leq j \leq m \\ j \neq 0}} \frac{\partial J(v_c, w_{t+j}, U)}{\partial v_c} \right]$$