

$$1. a) J_{\text{naive-softmax}} = -\log P(o=0 | C=c) = -\log \hat{y}_0 = -\sum_w y_w \log \hat{y}_w$$

$$\text{Since } y_w = \begin{cases} 1 & \text{if } w=0 \\ 0 & \text{otherwise} \end{cases}$$

$$b) \frac{\partial J_{\text{naive-softmax}}}{\partial v_c} = -\frac{\partial}{\partial v_c} \log \hat{y}_0 = -\frac{\partial}{\partial v_c} \log \frac{e^{u_0^T v_c}}{\sum_w e^{u_w^T v_c}}$$

$$= -\frac{\partial}{\partial v_c} \left(u_0^T v_c - \log \sum_w e^{u_w^T v_c} \right)$$

$$= -u_0 + \frac{1}{\sum_w e^{u_w^T v_c}} \sum_w e^{u_w^T v_c} u_w$$

$$= -u_0 + \sum_w \hat{y}_w u_w = -U y + U \hat{y}$$

$$= U(\hat{y} - y) \quad \text{where columns of } U \text{ are outputs word vectors}$$

$$c) \frac{\partial J}{\partial u_w} = -\frac{\partial}{\partial u_w} \left(u_0^T v_c - \log \sum_w e^{u_w^T v_c} \right)$$

$$= -1(w=0)v_c + \frac{1}{\sum_w e^{u_w^T v_c}} e^{u_w^T v_c} v_c$$

$$= -1(w=0)v_c + \hat{y}_w v_c$$

$$= (\hat{y}_w - 1(w=0))v_c$$

$$\begin{aligned}
 d) \frac{d\sigma}{dx} &= \frac{(e^x + 1)e^x - e^x e^x}{(e^x + 1)^2} = \frac{e^x(e^x + 1 - e^x)}{(e^x + 1)^2} \\
 &= \frac{e^x}{e^x + 1} \cdot \frac{1}{e^x + 1} = \sigma(x) \left(\frac{e^x + 1}{e^x + 1} - \frac{e^x}{e^x + 1} \right) \\
 &= \sigma(x) (1 - \sigma(x))
 \end{aligned}$$

$$\begin{aligned}
 e) \frac{\partial J_{\text{neg-sample}}}{\partial v_c} &= - \frac{1}{\sigma(u_0^T v_c)} \sigma(u_0^T v_c) (1 - \sigma(u_0^T v_c)) u_0 \\
 &\quad - \sum_{k=1}^K \frac{1}{\sigma(u_k^T v_c)} \sigma(u_k^T v_c) (1 - \sigma(u_k^T v_c)) u_k \\
 &= (\sigma(u_0^T v_c) - 1) u_0 - \sum_{k=1}^K (\sigma(u_k^T v_c) - 1) u_k
 \end{aligned}$$

$$\frac{\partial J_{\text{neg-sample}}}{\partial u_0} = (\sigma(u_0^T v_c) - 1) v_c$$

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

Only need to calculate $\sigma(u_w^T v_c)$ for $k \ll |\text{Vocab}|$ word vectors,
 so $J_{\text{neg-sample}}$ is much more efficient to calculate
 than $J_{\text{naive-softmax}}$