a). Because $y$ is a hot vectors with 0 if $w \neq 0$ and 1 if $w = 0$.

Therefore $-\sum_{w \in Vocab} y_w \log(\hat{y}_w) = -(y_1 \log(\hat{y}_1) + y_2 \log(\hat{y}_2) + \cdots + y_0 \log(\hat{y}_0) + \cdots)$

$$= -y_0 \log(\hat{y}_0)$$

$\because y_0 = 1$ when $w = 0$

$\therefore = -\log(\hat{y}_0)$

b). $J_{naive-softmax}(v_0, o, U) = -\log \dfrac{\exp(u_0^T v_c)}{\sum_{w \in Vocab} \exp(u_w^T v_c)}$

$$\frac{\partial J}{\partial v_c} = -\frac{\partial}{\partial v_c}\left( \log(\exp(u_0^T v_c)) - \log(\sum_{w \in Vocab} \exp(u_w^T v_c)) \right)$$

$$= -\frac{\partial}{\partial v_c}( u_0^T v_c) + \frac{1}{\sum_{w \in Vocab}\exp(u_w^T v_c)} \cdot \sum_{x \in Vocab} u_x \cdot \exp(u_x^T v_c)$$

$$= -u_0 + \sum_{x \in Vocab} \frac{\exp(u_x^T v_c)}{\sum_{w \in Vocab}\exp(u_w^T v_c)} \cdot u_x$$

$$= -u_0 + \sum_{x \in Vocab} P(u_x | v_c) \cdot u_x$$

$$= -u_0 + \sum_{x \in Vocab} \hat{y}_x \cdot u_x$$

c). Case 1: $w = 0$

$$\frac{\partial J}{\partial u_{w=0}} = -\frac{\partial}{\partial u_{w=0}}\left( \log(\exp(u_0^T v_c) - \log \sum_{w \in Vocab}((\exp(u_w^T v_c)) \right)$$

$$= -\frac{\partial}{\partial u_{w=c}}(u_0^T v_c) + \frac{\partial}{\partial u_{w=c}}(\log \sum_{w \in Vocab}(\exp(u_w^T v_c)))$$

$$= -v_c + \frac{\exp(u_0^T v_c)}{\sum_{w \in Vocab}((\exp(u_w^T v_c))} v_c$$

$$= -v_c + \hat{y}_0 v_c$$

$$= v_c \cdot (\hat{y}_0 - 1)$$

Case 2: $w \neq 0$

$$\frac{\partial J}{\partial u_{w \neq 0}} = \frac{-\partial}{\partial u_{w \neq 0}}\left( \log(\exp(u_0^T v_c)) - \log \sum_{w \in Vocab}((\exp(u_w^T v_c)) \right)$$

$$= 0 - \frac{\exp(u_{w\neq0}^T v_c)}{\sum_{w\in vocab}(\exp(u_w^T v_c))} \cdot v_c$$

$$= y_{w\neq0} \, v_c$$

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

$\qquad = -\dfrac{e^{-x}}{(1+e^{-x})^2}$

$\qquad = -\sigma(x) \cdot \dfrac{e^{-x}}{(1+e^{-x})}$

$\qquad = \sigma(x) \cdot (1 - \sigma(x))$

e). 1. $\dfrac{\partial J}{\partial v_c} = -\dfrac{1}{\sigma(u_0^T v_c)} \cdot (\sigma(u_0^T v_c) \cdot (1-\sigma(u_0^T v_c)) \cdot u_0$

$\qquad\quad + \sum_{k=1}^{K} \dfrac{1}{\sigma(u_k^T v_c)} \cdot (\sigma(u_k^T v_c) \cdot (1-\sigma(u_k^T v_c)) \cdot u_k$

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

2. $\dfrac{\partial J}{\partial u_0} = -v_c \dfrac{1}{\sigma(u_0^T v_c)} \cdot (\sigma(u_0^T v_c) \cdot (1-\sigma(u_0^T v_c)) - 0$

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

3. $\dfrac{\partial J}{\partial u_k} = \sum_{k=1}^{K} v_c \cdot (1-\sigma(-u_k^T v_c))$

Instead of calculating all the words in the Vocab, we now
only need to calculate k samples.

f). i) $\dfrac{\partial J_{skip\,c}(v_c, w_{t-m}\ldots w_{t+m}, U)}{\partial U} = \sum_{\substack{-m\leq j\leq m \\ j\neq 0}} \dfrac{\partial J(v_c, w_{t+j}, U)}{\partial U}$

ii) $= \sum_{\substack{-m\leq j\leq m \\ j\neq 0}} \dfrac{\partial J(v_c, w_{t+j}, U)}{\partial v_c}$

iii) $= 0$