$$P(O=o \mid C=c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in Vocab} \exp(u_w^T v_c)} = \hat{y}_o \in \mathbb{R}^{|x|} \qquad u_o \in \mathbb{R}^{n\times 1}, \; v_c \in \mathbb{R}^{n\times 1}$$

$$U^T = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_{Vocab} \end{bmatrix}^T \in \mathbb{R}^{n\times Vocab}$$

**(b)-(i)**
$$J_{naive-softmax}(v_c, o, U) = -\log P(O=o \mid C=c) = -\sum_{w \in Vocab} y_w \cdot \log \hat{y}_w = -\log(\hat{y}_o)$$

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

$$U \in \mathbb{R}^{Vocab \times n}$$

$$\frac{\partial J}{\partial v_c} = -u_o^T + \frac{\sum_w u_w^T \exp(u_w^T v_c)}{\sum_w \exp(u_w^T v_c)} \quad = \exp(U) \, v_c$$

$$\in \mathbb{R}^{1\times n} \qquad Vocab\times n \quad n\times 1$$

$$\hat{y}_w$$

$$= -u_o^T + \sum_w P(O=w \mid C=c) \cdot u_w^T$$
$$\qquad\qquad\qquad 1\times 1 \qquad\qquad 1\times n$$

$$= -u_o^T + \sum_w \hat{y}_w \cdot u_w^T$$
$$\qquad\qquad 1\times 1 \quad 1\times n$$

$$= \sum_w (\hat{y}_w - y_w) \cdot u_w^T = (U(\hat{y}-y))^T \in \mathbb{R}^{1\times n}$$
$$\qquad 1\times 1 \qquad 1\times n$$

$$v_c \in \mathbb{R}^{1\times n}, \; u_o \in \mathbb{R}^{1\times n} \quad \text{of } c \neq 1$$

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

$$\frac{\partial J}{\partial v_c} = (\hat{y}-y)\,U \quad = -u_o + \frac{\sum_w u_w \exp(v_c u_w^T)}{\sum_w \exp(v_c u_w^T)}$$
$$\qquad 1\times Vocab \quad Vocab\times n$$

**(b)-(ii)**
$$\frac{\partial J}{\partial v_c} = (U(\hat{y}-y))^T = \vec{0}, \; \text{The gradient is zero when the predicted distribution equals the true distribution } (\hat{y}=y)$$

**(b)-(iii)**
$$v_c \leftarrow v_c - \eta \cdot \frac{\partial J}{\partial v_c} \iff v_c \leftarrow v_c - \eta \cdot (U\hat{y} - Uy)^T$$

**(d)**
i) $w = o$
$$\frac{\partial J}{\partial w} = -v_c^T + \frac{v_c^T \cdot \exp(u_o^T v_c)}{\sum_z \exp(u_z^T v_c)} = \left(-1 + \frac{\exp(u_o^T v_c)}{\sum_z \exp(u_z^T v_c)}\right) v_c^T = (\hat{y}_o - y_o)\, v_c^T$$
$$1\times n \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad 1\times 1 \quad 1\times n$$

ii) $w \neq o$
$$\frac{\partial J}{\partial w} = \frac{v_c^T \exp(u_w^T v_c)}{\sum_z \exp(u_z^T v_c)} = \hat{y}_w \cdot v_c^T \qquad \hat{y} \in \mathbb{R}^{Vocab\times 1}$$
$$1\times n \qquad\qquad\qquad\qquad\qquad 1\times 1 \quad 1\times n$$

(e)

$$\frac{\partial J}{\partial U} = \begin{bmatrix} \frac{\partial J}{\partial u_1} \\ \vdots \\ \frac{\partial J}{\partial u_{Vocab}} \end{bmatrix} = \begin{bmatrix} \frac{\partial J}{\partial u_1} & \cdots & \frac{\partial J}{\partial u_{Vocab}} \end{bmatrix}^T$$

$\underset{Vocab \times n}{\underbrace{}}$   $\underset{1 \times n}{\underbrace{}}$

<span style="color:red">$f \in \mathbb{R}^{1 \times Vocab}$</span>

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

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

(h)-(i) Let $x = u_o^T v_c$, $y = -u_{ws}^T v_c$, $f = \log \Rightarrow J = -f(\sigma(x)) - \sum_{s=1}^{K} f(\sigma(y))$

<span style="color:red">$u_o, v_c \in \mathbb{R}^{1 \times d}$</span>

① $\frac{\partial J}{\partial v_c} = -\frac{\partial f}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial x} \cdot \frac{\partial x}{\partial v_c} - \sum_{s=1}^{K} \frac{\partial f}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial y} \cdot \frac{\partial y}{\partial v_c}$

$\underset{\in \mathbb{R}^{1 \times d}}{\underbrace{}}$

$\quad = -\frac{1}{\sigma(x)} \cdot \sigma(x) \cdot (1 - \sigma(x)) \cdot u_o^T + \sum_{s=1}^{K} \frac{1}{\sigma(y)} \cdot \sigma(y) \cdot (1 - \sigma(y)) \cdot (+u_{ws}^T)$

$\quad = (\sigma(u_o^T v_c) - 1) u_o^T + \sum_{s=1}^{K} (1 - \sigma(-u_{ws}^T v_c)) \cdot u_{ws}^T$

<span style="color:red">$(\sigma(u_o v_c^T) - 1) u_o + \sum_{s=1}^{K} (1 - \sigma(-u_{w_a} v_c^T)) u_{ws}$</span>

② $\frac{\partial J}{\partial u_o} = -\frac{\partial f}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial x} \cdot \frac{\partial x}{\partial u_o} - \sum_{s=1}^{K} \frac{\partial f}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial y} \cdot \boxed{\frac{\partial y}{\partial u_o}}$

$\quad = -\frac{1}{\sigma(x)} \cdot \sigma(x) \cdot (1 - \sigma(x)) \cdot v_c^T - \sum_{s=1}^{K} \frac{1}{\sigma(y)} \cdot \sigma(y) (1 - \sigma(y)) \cdot \textcircled{0}$

$\quad = (\sigma(u_o^T v_c) - 1) v_c^T$

③ $\dfrac{\partial J}{\partial U_{w_s}} = -\dfrac{\partial f}{\partial \sigma} \cdot \dfrac{\partial \sigma}{\partial x} \cdot \dfrac{\partial x}{\partial U_{w_s}} - \sum\limits_{S=1}^{K} \dfrac{\partial f}{\partial \sigma} \cdot \dfrac{\partial \sigma}{\partial y} \cdot \dfrac{\partial y}{\partial U_{w_s}}$

$$U_{w_s} \leftarrow U_{w_s} - \eta \cdot \dfrac{\partial J}{\partial U_{w_s}}$$

$$= -\dfrac{1}{\sigma(x)} \cdot \sigma(x) \cdot (1-\sigma(x)) \cdot \boxed{0} + \dfrac{1}{\sigma(y)} \cdot \sigma(y) \cdot (1-\sigma(y)) \cdot (+ V_c^{T})$$

$$= \left(1 - \sigma(-U_{w_s}^{T} V_c)\right) \cdot V_c^{T}$$

(h)—(ii)  $\quad U_{0, \{w_1, \cdots, w_k\}} = \left[ U_0, -w_1, \cdots, -w_k \right] \qquad \therefore \ \sigma\left( U_{0,\{w_1,\cdots, w_k\}}^{T} V_c \right)$

(i)  $\quad J_{neg\text{-}sample} = -\log\left( \sigma(U_0^{T} V_c) \right) - \sum\limits_{S=1}^{K} \log\left( \sigma(-U_{w_s}^{T} V_c) \right)$

$$= -\log\left( \sigma(U_0^{T} V_c) \right) - \left\{ \sum\limits_{i: w_i = w_s} \log\left( \sigma(-U_{w_i}^{T} V_c) \right) + \sum\limits_{j: w_j \neq w_s} \log\left( \sigma(-U_{w_j}^{T} V_c) \right) \right\}$$

$\dfrac{\partial J}{\partial U_{w_s}} = - \sum\limits_{i: w_i = w_s} \dfrac{\partial}{\partial U_{w_s}} \left\{ \log\left( \sigma(-U_{w_i}^{T} V_c) \right) \right\} = - \sum\limits_{i: w_i = w_s} \sigma(-U_{w_s}^{T} V_c) V_c$

$J_{neg\text{-}sample} = -\log\left( \sigma(U_0^{T} V_c) \right) - \sum\limits_{S=1}^{K} \log\left( \sigma(-U_{w_s}^{T} V_c) \right) \qquad U_0 \in \mathbb{R}^{1 \times d}, \ V_c \in \mathbb{R}^{1 \times d}$

$$= -\log\left( \sigma(U_0^{T} V_c) \right) - \left\{ \sum\limits_{i: w_i = w_s} \log\left( \sigma(-U_{w_i}^{T} V_c) \right) + \sum\limits_{j: w_j \neq w_s} \log\left( \sigma(-U_{w_j}^{T} V_c) \right) \right\}$$

(i)—(i)  $\quad \dfrac{\partial J_{skip\text{-}gram}(V_c, w_{t-m}, \cdots, w_{t+m}, U)}{\partial U} = \sum\limits_{-m \leq j \leq m, \ j \neq 0} \dfrac{J(V_c, w_{t+j}, U)}{\partial U}$

(j)—(ii)  $\quad \dfrac{\partial J_{skip\text{-}gram}(V_c, w_{t-m}, \cdots, w_{t+m}, U)}{\partial V_c} = \sum\limits_{-m \leq j \leq m, \ j \neq 0} \dfrac{J(\ \cdots_{t+j}, U)}{\partial V_c}$

(j)—(iii)  $\quad \dfrac{\partial J_{skip\text{-}gram}(V_c, w_{t-m}, \cdots, w_{t+m}, U)}{\partial V_w} = \sum\limits_{\substack{-m \leq j \leq m, \ j \neq 0 \\ w_{t+j} = w}} \dfrac{J(V_c, w_{t+j}, U)}{\partial V_w}$