# 1. Softmax

(a) we want to prove : $softmax(X+c) = softmax(X)$

we know that: $softmax(X)_i = \dfrac{e^{x_i}}{\sum_{j=1}^{Dim(X)} e^{x_j}}$

So: $softmax(X+c)_i = \dfrac{e^{x_i+c}}{\sum_{j=1}^{Dim(X)} e^{x_j+c}} = \dfrac{e^{x_i} \times e^{c}}{\sum_{j=1}^{Dim(X)} e^{x_j} \times \underbrace{e^{c}}_{constant}} = \dfrac{e^{x_i} \times \cancel{e^{c}}}{\cancel{e^{c}}\sum_{j=1}^{Dim(X)} e^{x_j}}$

$\Rightarrow softmax(X+c) = softmax(X)$

---

# 2. Neural Network Basics:

(a) we want to derive the gradients of the sigmoid function and show that
it can be rewritten as a function of the function value :

$\sigma(x) = \dfrac{1}{1+e^{-x}} = (1+e^{-x})^{-1} \xrightarrow[\text{using chain rule}]{\text{take derivative}} \sigma'(x) = -(1+e^{-x})^{-2} \times -e^{-x}$

$\Rightarrow \sigma'(x) = \dfrac{e^{-x}}{(1+e^{-x})^2} = \underbrace{\dfrac{e^{-x}}{1+e^{-x}}}_{1-\sigma(x)} \times \underbrace{\dfrac{1}{1+e^{-x}}}_{\sigma(x)} = \sigma(x)(1-\sigma(x))$

(b) we want to derive the gradient w.r.t the inputs of a softmax
function when cross entropy loss is used for evaluation :

$\begin{cases} \hat{y} = softmax(\theta) \\ CE(y,\hat{y}) = -\sum_i y_i \log(\hat{y}_i) \end{cases}$

$\dfrac{\partial CE}{\partial \theta} = -\sum_i y_i \dfrac{\partial}{\partial \theta} \log(\hat{y}_i) \Rightarrow$ we know that $y$ is a one-hot label vector. we assume that only the k-th dimension of $y$ is one, So:

$\dfrac{\partial CE}{\partial \theta} = -\overset{1}{\cancel{y_K}} \dfrac{\partial}{\partial \theta} \log(\hat{y}_K) = -\dfrac{\partial}{\partial \theta} \log\left(\dfrac{e^{\theta_K}}{\sum_j e^{\theta_j}}\right) \Rightarrow$ $\theta$ is a vector. So we expand this derivative

Gradient
w.r.t $\theta_i$ $\Rightarrow$ $\dfrac{\partial}{\partial \theta_i}\left[\log \dfrac{e^{\theta_K}}{\sum_j e^{\theta_j}}\right] = \dfrac{\partial}{\partial \theta_i}\left[\log e^{\theta_K} - \log \sum_j e^{\theta_j}\right] = \dfrac{\partial}{\partial \theta_i}\left[\theta_K - \log \sum_j e^{\theta_j}\right]$

Now there is
two options $\Rightarrow$
$\left\{\begin{array}{l}\text{if } i = K \Rightarrow \dfrac{\partial}{\partial \theta_K}\left(\theta_K - \log \sum_j e^{\theta_j}\right) = 1 - \dfrac{e^{\theta_K}}{\sum_j e^{\theta_j}} = 1 - \hat{y}_K = 1 - \hat{y}_i = y_i - \hat{y}_i \\ \qquad\qquad\qquad \boxed{\hat{y}_i \text{ when } i = K} \\[20pt] \text{if } i \neq K \Rightarrow \dfrac{\partial}{\partial \theta_i}\left[\theta_K - \log \sum_j e^{\theta_j}\right] = 0 - \dfrac{e^{\theta_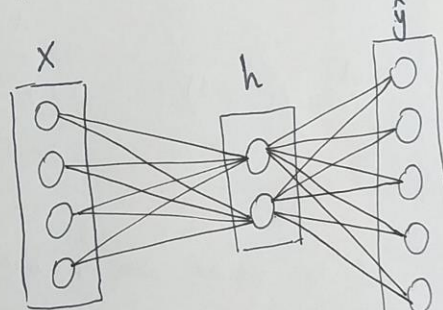i}}{\sum_j e^{\theta_j}} = 0 - \hat{y}_i = y_i - \hat{y}_i \\ \qquad\qquad\qquad \boxed{\hat{y}_i \text{ when } i \neq K}\end{array}\right.$

So anyway we
can rewritten
the gradient like: $\qquad \dfrac{\partial}{\partial \theta}\left[\log \dfrac{e^{\theta_K}}{\sum_j e^{\theta_j}}\right] = y - \hat{y}$

finally $\Rightarrow$ $\dfrac{\partial CE}{\partial \theta} = -\dfrac{\partial}{\partial \theta}\left[\log \dfrac{e^{\theta_K}}{\sum_j e^{\theta_j}}\right] = \hat{y} - y$

(C) we want to move backward through the neural network and take
all derivatives step by step.

let's assume $\left\{\begin{array}{l} z^2 = h w^2 + b^2 \\ z^1 = x w^1 + b^1 \end{array}\right.$



$h = \text{sigmoid}(x w^1 + b^1)$
$\hat{y} = \text{softmax}(h w^2 + b^2)$

from part (b) we know that:

$\dfrac{\partial CE}{\partial z^2} = (\hat{y} - y) = dz^2$

$\underbrace{\dfrac{\partial CE}{\partial w^2}}_{H \times D_y} = \dfrac{\partial z^2}{\partial w^2} \times \dfrac{\partial CE}{\partial z^2} = \underbrace{\dfrac{\partial z^2}{\partial w^2}}_{\text{it should be } H \times 1} \times \underbrace{(\hat{y} - y)}_{1 \times D_y}$ $\qquad$ let's compute
this term

According to the neural network architecture we can conclude

$$\begin{cases} h_1 = \dfrac{\partial z_1^2}{\partial w_{11}^2} = \dfrac{\partial z_2^2}{\partial w_{12}^2} = \dfrac{\partial z_3^2}{\partial w_{13}^2} = \cdots = \dfrac{\partial z_{D_y}^2}{\partial w_{1D_y}^2} \\[3mm] h_2 = \dfrac{\partial z_1^2}{\partial w_{21}^2} = \dfrac{\partial z_2^2}{\partial w_{22}^2} = \dfrac{\partial z_3^2}{\partial w_{23}^2} = \cdots = \dfrac{\partial z_{D_y}^2}{\partial w_{2D_y}^2} \\[3mm] \vdots \\[2mm] h_H = \dfrac{\partial z_1^2}{\partial w_{H1}^2} = \dfrac{\partial z_2^2}{\partial w_{H2}^2} = \dfrac{\partial z_3^2}{\partial w_{H3}^2} = \cdots = \dfrac{\partial z_{D_y}^2}{\partial w_{HD_y}^2} \end{cases}$$

$$\frac{\partial CE}{\partial w^2} = \begin{bmatrix} \dfrac{\partial CE}{\partial w_{11}^2} & \dfrac{\partial CE}{\partial w_{12}^2} & \cdots & \dfrac{\partial CE}{\partial w_{1D_y}^2} \\[3mm] \dfrac{\partial CE}{\partial w_{21}^2} & \dfrac{\partial CE}{\partial w_{22}^2} & \cdots & \dfrac{\partial CE}{\partial w_{2D_y}^2} \\[3mm] \vdots & \vdots & & \vdots \\[3mm] \dfrac{\partial CE}{\partial w_{H1}^2} & \dfrac{\partial CE}{\partial w_{H2}^2} & \cdots & \dfrac{\partial CE}{\partial w_{HD_y}^2} \end{bmatrix} = \begin{bmatrix} \dfrac{\partial z_1^2}{\partial w_{11}^2}\times\dfrac{\partial CE}{\partial z_1^2} & \dfrac{\partial z_2^2}{\partial w_{12}^2}\times\dfrac{\partial CE}{\partial z_2^2} & \cdots & \dfrac{\partial z_{D_y}^2}{\partial w_{1D_y}^2}\times\dfrac{\partial CE}{\partial z_{D_y}^2} \\[3mm] \dfrac{\partial z_1^2}{\partial w_{21}^2}\times\dfrac{\partial CE}{\partial z_1^2} & \dfrac{\partial z_2^2}{\partial w_{22}^2}\times\dfrac{\partial CE}{\partial z_2^2} & \cdots & \dfrac{\partial z_{D_y}^2}{\partial w_{2D_y}^2}\times\dfrac{\partial CE}{\partial z_{D_y}^2} \\[3mm] \vdots & & & \vdots \\[3mm] \dfrac{\partial z_1^2}{\partial w_{H1}^2}\times\dfrac{\partial CE}{\partial z_1^2} & \dfrac{\partial z_2^2}{\partial w_{H2}^2}\times\dfrac{\partial CE}{\partial z_2^2} & \cdots & \dfrac{\partial z_{D_y}^2}{\partial w_{HD_y}^2}\times\dfrac{\partial CE}{\partial z_{D_y}^2} \end{bmatrix}_{H\times D_y}$$

$$\Rightarrow \frac{\partial CE}{\partial w^2} = \begin{bmatrix} h_1\times(\hat{y}_1-y_1) & h_1(\hat{y}_2-y_2) & \cdots & h_1(\hat{y}_{D_y}-y_{D_y}) \\[2mm] h_2\times(\hat{y}_1-y_1) & h_2\times(\hat{y}_2-y_2) & \cdots & h_2\times(\hat{y}_{D_y}-y_{D_y}) \\[2mm] \vdots & \vdots & & \vdots \\[2mm] h_H\times(\hat{y}_1-y_1) & h_H\times(\hat{y}_2-y_2) & \cdots & h_H\times(\hat{y}_{D_y}-y_{D_y}) \end{bmatrix}_{H\times D_y} = \begin{bmatrix} h_1 \\ h_2 \\ h_3 \\ \vdots \\ h_H \end{bmatrix}_{H\times 1}\begin{bmatrix} \hat{y}_1-y_1 & \hat{y}_2-y_2 & \cdots & \hat{y}_{D_y}-y_{D_y} \end{bmatrix}_{1\times D_y}$$

$$= h^T = \frac{\partial z^2}{\partial w^2}$$

finally

$$\Rightarrow \frac{\partial CE}{\partial w^2} = \frac{\partial z^2}{\partial w^2}\times(\hat{y}-y) = h^T(\hat{y}-y) = h^T dz^2 = dw^2$$

Gradient w.r.t $b^2$

$$\Rightarrow \frac{\partial CE}{\partial b^2} = \begin{bmatrix} \dfrac{\partial CE}{\partial b_1^2} & \dfrac{\partial CE}{\partial b_2^2} & \cdots & \dfrac{\partial CE}{\partial b_{D_y}^2} \end{bmatrix}$$

$$\frac{\partial CE}{\partial b^2} = \begin{bmatrix} \underbrace{\dfrac{\partial z_1^2}{\partial b_1^2}}_{1}\times\dfrac{\partial CE}{\partial z_1^2} & \underbrace{\dfrac{\partial z_2^2}{\partial b_2^2}}_{1}\times\dfrac{\partial CE}{\partial z_2^2} & \cdots & \underbrace{\dfrac{\partial z_{D_y}^2}{\partial b_{D_y}^2}}_{1}\times\dfrac{\partial CE}{\partial z_{D_y}^2} \end{bmatrix}$$

$$\frac{\partial CE}{\partial b^2} = \begin{bmatrix} \dfrac{\partial CE}{\partial z_1^2} & \dfrac{\partial CE}{\partial z_2^2} & \cdots & \dfrac{\partial CE}{\partial z_{D_y}^2} \end{bmatrix} = \frac{\partial CE}{\partial z^2} = (\hat{y}-y) = db^2$$

Gradient $\Rightarrow$ $\dfrac{\partial CE}{\partial h} = \underset{\text{1x1Dy}}{\dfrac{\partial CE}{\partial Z^2}} \times \underset{\text{should be: Dy×H}}{\dfrac{\partial Z^2}{\partial h}} = dZ^2 \times \dfrac{\partial Z^2}{\partial h}$, let's compute
w.r.t $h$ $\underset{\text{Dim=1×H}}{\quad}$ this term

we know
that $\longrightarrow$

$$z_1^2 = h_1 w_{11}^2 + h_2 w_{21}^2 + \cdots + h_H w_{H1}^2 + b_1^2$$

$$z_2^2 = h_1 w_{12}^2 + h_2 w_{22}^2 + \cdots + h_H w_{H2}^2 + b_2^2 \qquad \searrow w_{ij}^2 = \dfrac{\partial z_j^2}{\partial h_i}$$

$$\vdots$$

$$z_{Dy}^2 = h_1 w_{1Dy}^2 + h_2 w_{2Dy}^2 + \cdots + h_H w_{HDy}^2 + b_{Dy}^2$$

SO: $\dfrac{\partial Z^2}{\partial h} = \begin{bmatrix} \dfrac{\partial z_1^2}{\partial h_1} & \dfrac{\partial z_1^2}{\partial h_2} & \cdots & \dfrac{\partial z_1^2}{\partial h_H} \\ \dfrac{\partial z_2^2}{\partial h_1} & \dfrac{\partial z_2^2}{\partial h_2} & \cdots & \dfrac{\partial z_2^2}{\partial h_H} \\ \vdots & \vdots & & \vdots \\ \dfrac{\partial z_{Dy}^2}{\partial h_1} & \dfrac{\partial z_{Dy}^2}{\partial h_2} & \cdots & \dfrac{\partial z_{Dy}^2}{\partial h_H} \end{bmatrix}_{Dy×H} = \begin{bmatrix} w_{11}^2 & w_{21}^2 & \cdots & w_{H1}^2 \\ w_{12}^2 & w_{22}^2 & \cdots & w_{H2}^2 \\ \vdots & & & \\ w_{1Dy}^2 & w_{2Dy}^2 & \cdots & w_{HDy}^2 \end{bmatrix}_{Dy×H}$

$\Rightarrow \dfrac{\partial Z^2}{\partial h} = w^{2T}$

finally : $\dfrac{\partial CE}{\partial h} = \left(\dfrac{\partial CE}{\partial Z^2}\right) \times \left(\dfrac{\partial Z^2}{\partial h}\right) = (\hat{y}-y) w^{2T} = dh$

Gradient $\Rightarrow$ $\dfrac{\partial CE}{\partial Z^1} = \dfrac{\partial CE}{\partial h} \times \dfrac{\partial h}{\partial Z^1} = dh \times \dfrac{\partial h}{\partial Z^1}$, let's compute
w.r.t $Z^1$ this term

we know that $\Rightarrow$ $h = \underset{\sigma(Z^1)}{\underbrace{sigmoid(Z^1)}} \Rightarrow \dfrac{\partial h}{\partial Z^1} = \sigma'(Z^1)$

$\Rightarrow$ SO: $\underset{\text{Dim=1×H}}{\dfrac{\partial CE}{\partial Z^1}} = \dfrac{\partial CE}{\partial h} \times \underset{\text{1×H}}{\dfrac{\partial h}{\partial Z^1}} = \underset{\text{1×H}}{\underbrace{(\hat{y}-y) w^{2T}}} \circ \underset{\text{1×H}}{\underbrace{\sigma'(Z^1)}} = dz^1$

Gradient
w.r.t $w^1$ → $\underbrace{\dfrac{\partial CE}{\partial w^1}}_{\text{Dim} = D_x \times H} = \underbrace{\dfrac{\partial z^1}{\partial w^1}}_{\downarrow} \times \underbrace{\dfrac{\partial CE}{\partial z^1}}_{1 \times H} = \dfrac{\partial z^1}{\partial w^1} \times dz^1$ ⟶ let's compute this term

Should be: $D_x \times 1$

According to the neural network, we can conclude

$$
\begin{cases}
x_1 = \dfrac{\partial z_1^1}{\partial w_{11}^1} = \dfrac{\partial z_2^1}{\partial w_{12}^1} = \cdots = \dfrac{\partial z_H^1}{\partial w_{1H}^1} \\[3mm]
x_2 = \dfrac{\partial z_1^1}{\partial w_{21}^1} = \dfrac{\partial z_2^1}{\partial w_{22}^1} = \cdots = \dfrac{\partial z_H^1}{\partial w_{2H}^1} \\[3mm]
\vdots \\[2mm]
x_{D_x} = \dfrac{\partial z_1^1}{\partial w_{D_x 1}^1} = \dfrac{\partial z_2^1}{\partial w_{D_x 2}^1} = \cdots = \dfrac{\partial z_H^1}{\partial w_{D_x H}^1}
\end{cases}
$$

$$
\frac{\partial CE}{\partial w^1} =
\begin{bmatrix}
\dfrac{\partial CE}{\partial w_{11}^1} & \dfrac{\partial CE}{\partial w_{12}^1} & \cdots & \dfrac{\partial CE}{\partial w_{1H}^1} \\[3mm]
\dfrac{\partial CE}{\partial w_{21}^1} & \dfrac{\partial CE}{\partial w_{22}^1} & \cdots & \dfrac{\partial CE}{\partial w_{2H}^1} \\[3mm]
\vdots & \vdots & & \vdots \\[3mm]
\dfrac{\partial CE}{\partial w_{D_x 1}^1} & \dfrac{\partial CE}{\partial w_{D_x 2}^1} & \cdots & \dfrac{\partial CE}{\partial w_{D_x H}^1}
\end{bmatrix}_{D_x H}
=
\begin{bmatrix}
\dfrac{\partial z_1^1}{\partial w_{11}^1}\times\dfrac{\partial CE}{\partial z_1^1} & \dfrac{\partial z_2^1}{\partial w_{12}^1}\times\dfrac{\partial CE}{\partial z_2^1} & \cdots & \dfrac{\partial z_H^1}{\partial w_{1H}^1}\times\dfrac{\partial CE}{\partial z_H^1} \\[3mm]
\dfrac{\partial z_1^1}{\partial w_{21}^1}\times\dfrac{\partial CE}{\partial z_1^1} & \dfrac{\partial z_2^1}{\partial w_{22}^1}\times\dfrac{\partial CE}{\partial z_2^1} & \cdots & \dfrac{\partial z_H^1}{\partial w_{2H}^1}\times\dfrac{\partial CE}{\partial z_H^1} \\[3mm]
\vdots & \vdots & & \vdots \\[3mm]
\dfrac{\partial z_1^1}{\partial w_{D_x 1}^1}\times\dfrac{\partial CE}{\partial z_1^1} & \dfrac{\partial z_2^1}{\partial w_{D_x 2}^1}\times\dfrac{\partial CE}{\partial z_2^1} & \cdots & \dfrac{\partial z_H^1}{\partial w_{D_x H}^1}\times\dfrac{\partial CE}{\partial z_H^1}
\end{bmatrix}
$$

$$
\Rightarrow \frac{\partial CE}{\partial w^1} =
\begin{bmatrix}
x_1 \dfrac{\partial CE}{\partial z_1^1} & x_1 \dfrac{\partial CE}{\partial z_2^1} & \cdots & x_1 \dfrac{\partial CE}{\partial z_H^1} \\[3mm]
x_2 \dfrac{\partial CE}{\partial z_1^1} & x_2 \dfrac{\partial CE}{\partial z_2^1} & \cdots & x_2 \dfrac{\partial CE}{\partial z_H^1} \\[3mm]
\vdots & & & \\[2mm]
x_{D_x} \dfrac{\partial CE}{\partial z_1^1} & x_{D_x} \dfrac{\partial CE}{\partial z_2^1} & \cdots & x_{D_x} \dfrac{\partial CE}{\partial z_H^1}
\end{bmatrix}
=
\underbrace{\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{D_x} \end{bmatrix}}_{\substack{D_x \times 1 \\ = X^T = \frac{\partial z^1}{\partial w^1}}}
\underbrace{\begin{bmatrix} \dfrac{\partial CE}{\partial z_1^1} & \dfrac{\partial CE}{\partial z_2^1} & \cdots & \dfrac{\partial CE}{\partial z_H^1} \end{bmatrix}}_{1 \times H}
$$

$$
\Rightarrow \frac{\partial CE}{\partial w^1} = X^T \frac{\partial CE}{\partial z^1} = X^T (\hat{y} - y) w^{2T} \circ \sigma'(z^1) = dw^1
$$

Gradient → w.r.t $b'$

$$\frac{\partial CE}{\partial b'} = \left[ \frac{\partial CE}{\partial b'_1} \quad \frac{\partial CE}{\partial b'_2} \quad \cdots \quad \frac{\partial CE}{\partial b'_H} \right]$$

$$\frac{\partial CE}{\partial b'} = \left[ \underbrace{\frac{\partial z'_1}{\partial b'_1}}_{1} \times \frac{\partial CE}{\partial z'_1} \quad \underbrace{\frac{\partial z'_2}{\partial b'_2}}_{1} \times \frac{\partial CE}{\partial z'_2} \quad \cdots \quad \underbrace{\frac{\partial z'_H}{\partial b'_H}}_{1} \times \frac{\partial CE}{\partial z'_H} \right]$$

$$\frac{\partial CE}{\partial b'} = \left[ \frac{\partial CE}{\partial z'_1} \quad \frac{\partial CE}{\partial z'_2} \quad \cdots \quad \frac{\partial CE}{\partial z'_H} \right] = \frac{\partial CE}{\partial z'} = (\hat{y}-y) w^{2T} \circ \sigma'(z') = db'$$

And finally we can compute

$$\frac{\partial CE}{\partial x} = \frac{\partial CE}{\partial z'} \times \frac{\partial z'}{\partial x} \quad \longrightarrow \quad$$ similar to computation of $\frac{\partial z^2}{\partial h}$, it is equal to: $w^{1T}$

$$\Rightarrow \quad \frac{\partial CE}{\partial x} = (\hat{y}-y) w^{2T} \circ \sigma'(z') w^{1T}$$

(d) How many parameters are there in this neural network? let's suppose:

$\begin{cases} D_x : \text{input vector dimension} \\ D_y : \text{output vector dimension} \\ H : \text{number of hidden units} \end{cases}$

\# parameters = \# parameter in $w^2, b^2, w', b' =$

$$= H \times D_y + 1 \times D_y + D_x \times H + 1 \times H$$

$$= (H+1) D_y + (D_x+1) H$$

---

## 3. Word2Vec:

(a) With following assumptions, we want to derive the gradient of cost function with respect to $V_c$:

- $V_c$: predicted word vector corresponding to the center word $c$ for skipgram

- $U_o$: output word vector corresponding to the expected word $o$

(b) Derive the gradient w.r.t all output vectors: ⑧

$$\frac{\partial J}{\partial u_i} = \frac{\partial}{\partial u_i} \left[ -u_0^T v_c + \log \left[ \sum_{w \geq 1}^{W} \exp(u_w^T v_c) \right] \right]$$

① if $i \neq 0$ $\Rightarrow$ $\frac{\partial J}{\partial u_i} = \left[ -(0) + \frac{1}{\sum_{w \geq 1}^{W} \exp(u_w^T v_c)} \times v_c \exp(u_i^T v_c) \right]$

$$\Rightarrow \frac{\partial J}{\partial u_i} = \left[ -(0) + v_c \times \underbrace{\frac{\exp(u_i^T v_c)}{\sum_{w \geq 1}^{W} \exp(u_w^T v_c)}}_{} \right] = \left[ -(0) + v_c \overset{\overset{v_c \hat{y}_i}{}}{\hat{y}_i} \right]$$

$$\Rightarrow \frac{\partial J}{\partial u_i} = v_c (\hat{y}_i - y_i)$$

② if $i = 0$ $\Rightarrow$ $\frac{\partial J}{\partial u_0} = \left[ -\overset{v_c y_0}{v_c} + v_c \hat{y}_0 \right] = v_c(\hat{y}_0 - y_0) = v_c(\hat{y}_i - y_i)$

$\Rightarrow$ SO $\Rightarrow$ $\frac{\partial J}{\partial u_w} = v_c(\hat{y}_w - \hat{y}_w)$ $\Rightarrow$ or more generally $\Rightarrow$ $\frac{\partial J}{\partial U} = v_c(\hat{y} - y)^T$

any word
in vocabulary

Dim $= N \times W = (N \times 1) \times (1 \times W)$

(c) Repeat part (a) and (b) assuming we are using the negative sampling loss function ( Note that $0 \notin \{1, 2, \dots K\}$ ):

$$J_{negative-sample}(0, v_c, U) = -\log(\sigma(u_0^T v_c)) - \sum_{k=1}^{K} \log(\sigma(-u_k^T v_c))$$

$$\frac{\partial J}{\partial v_c} = -\frac{\partial}{\partial v_c} \log(\sigma(u_0^T v_c)) - \sum_{k=1}^{K} \frac{\partial}{\partial v_c} \log(\sigma(-u_k^T v_c))$$

$$\frac{\partial J}{\partial v_c} = -\frac{\sigma(u_0^T v_c)(1 - \sigma(u_0^T v_c))}{\sigma(u_0^T v_c)} u_0 - \sum_{k=1}^{K} \frac{-\sigma(-u_k^T v_c)(1 - \sigma(-u_k^T v_c))}{\sigma(-u_k^T v_c)} u_k$$

$$\frac{\partial J}{\partial v_c} = (\sigma(u_0^T v_c) - 1) u_0 - \sum_{k=1}^{K} (\sigma(-u_k^T v_c) - 1) u_k$$

Gradient $\Rightarrow$ $\dfrac{\partial J}{\partial u_0} = -\dfrac{\partial}{\partial u_0} \log\left(\sigma(u_0^T v_c)\right) = -\dfrac{\sigma(u_0^T v_c)(1-\sigma(u_0^T v_c))}{\sigma(u_0^T v_c)} \times v_c$
w.r.t $u_0$

$$\dfrac{\partial J}{\partial u_0} = v_c\left(\sigma(u_0^T v_c) - 1\right)$$

Gradient $\Rightarrow$ $\dfrac{\partial J}{\partial u_K} = -\dfrac{\partial}{\partial u_K} \sum_{K=1}^{K} \log\left(\sigma(-u_K^T v_c)\right)$
w.r.t $u_K$

$$\dfrac{\partial J}{\partial u_K} = -\sum_{K \neq t}^{K} \dfrac{\sigma(-u_K^T v_c)(1-\sigma(-u_K^T v_c))}{\sigma(-u_K^T v_c)} \times -v_c$$

$$\dfrac{\partial J}{\partial u_K} = \left(1 - \sigma(-u_K^T v_c)\right) v_c$$

(d) for skip-gram, the cost for a context centered around $c$ is:

$$J_{skip-gram}\left(word_{c-m\ldots c+m}\right) = \sum_{-m \leq j \leq m, j \neq 0} F(w_{c+j}, v_c)$$

Then the gradients for the cost of one context window are:

$$\dfrac{\partial J_{skip-gram}\left(word_{c-m\ldots c+m}\right)}{\partial U} = \sum_{-m \leq j \leq m, j \neq 0} \dfrac{\partial F(w_{c+j}, v_c)}{\partial U}$$

$$\dfrac{\partial J_{skip-gram}\left(word_{c-m\ldots c+m}\right)}{\partial v_c} = \sum_{-m \leq j \leq m, j \neq 0} \dfrac{\partial F(w_{c+j}, v_c)}{\partial v_c}$$

$$\dfrac{\partial J_{skipgram}\left(word_{c-m\ldots c+m}\right)}{\partial v_j} = 0 \qquad \text{for all the other } j$$

for CBOW ~~B~~ instead of using $v_c$ as the predicted vector, we use $\hat{v}$ ⑩
defined below: $\qquad \hat{v}_c = \sum\limits_{-m \le j \le m, j \ne 0} v_{c+j}$

then the CBOW cost is:

$$J_{CBOW}(word_{c-m\cdots c+m}) = F(w_c, \hat{v}_c)$$

for CBOW then we have:

$$\left\{ \begin{array}{l} \dfrac{\partial J_{CBOW}(word_{c-m\cdots c+m})}{\partial U} = \dfrac{\partial F(w_c, \hat{v}_c)}{\partial U} \\[4mm] \dfrac{\partial J_{CBOW}(word_{c-m\cdots c+m})}{\partial v_j} = \dfrac{\partial F(w_c, \hat{v}_c)}{\partial \hat{v}_c} \qquad \text{for all } j \in \{c-m\cdots c+m\} \\[4mm] \dfrac{\partial J_{CBOW}(word_{c-m\cdots c+m})}{\partial v_j} = 0 \qquad \text{for all } j \notin \{c-m\cdots c+m\} \end{array} \right.$$