

$$3(b) \frac{\partial(CE(\mathbf{y}, \hat{\mathbf{y}}))}{\partial \mathbf{U}}$$

$= -\frac{\partial}{\partial \mathbf{U}} \sum_{w=1}^W \left(y_w \log(\hat{y}_w) \right)$	(1)
$= \begin{bmatrix} -\frac{\partial}{\partial \mathbf{u}_l} \left((y_l \log \hat{y}_l) + \dots + (y_o \log \hat{y}_o) + \dots + (y_w \log \hat{y}_w) \right) \\ \vdots \\ -\frac{\partial}{\partial \mathbf{u}_o} \left((y_l \log \hat{y}_o) + \dots + (y_o \log \hat{y}_o) + \dots + (y_w \log \hat{y}_w) \right) \\ \vdots \\ -\frac{\partial}{\partial \mathbf{u}_w} \left((y_l \log \hat{y}_l) + \dots + (y_o \log \hat{y}_o) + \dots + (y_w \log \hat{y}_w) \right) \\ \vdots \end{bmatrix}$	(2)
$= \begin{bmatrix} -\frac{\partial}{\partial \mathbf{u}_l} \left(\left(y_l \log \left(\frac{\exp(\mathbf{u}_l^T \mathbf{v}_c)}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) + \dots + y_o \log \left(\frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) + \dots + y_w \log \left(\frac{\exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) \right) \\ \vdots \\ -\frac{\partial}{\partial \mathbf{u}_o} \left(\left(y_l \log \left(\frac{\exp(\mathbf{u}_l^T \mathbf{v}_c)}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) + \dots + y_o \log \left(\frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) + \dots + y_w \log \left(\frac{\exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) \right) \\ \vdots \\ -\frac{\partial}{\partial \mathbf{u}_w} \left(\left(y_l \log \left(\frac{\exp(\mathbf{u}_l^T \mathbf{v}_c)}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) + \dots + y_o \log \left(\frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) + \dots + y_w \log \left(\frac{\exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) \right) \\ \vdots \end{bmatrix}$	(3)

If we look at the w^{th} row of Equation (3), it is of the form below:

$= -\frac{\partial}{\partial \mathbf{u}_w} \left(\sum_{i=1}^{i=W} y_i \log \left(\frac{\exp(\mathbf{u}_i^T \mathbf{v}_c)}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) \right)$	(4)
$= -\sum_{i=1}^{i=W} \frac{\partial}{\partial \mathbf{u}_w} \left(y_i \log \left(\frac{\exp(\mathbf{u}_i^T \mathbf{v}_c)}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) \right)$	(5)
$= -\sum_{i=1}^{i=W} \frac{\partial}{\partial \mathbf{u}_w} \left(y_i \log(\exp(\mathbf{u}_i^T \mathbf{v}_c)) - \log \sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c) \right)$	(6)
$= -\sum_{i=1}^{i=W} y_i \left(\frac{\partial(\mathbf{u}_i^T \mathbf{v}_c)}{\partial \mathbf{u}_w} - \frac{\partial}{\partial \mathbf{u}_w} \left(\log \sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c) \right) \right)$	(7)
$= -\sum_{i=1}^{i=W} y_i \left(\frac{\partial(\mathbf{u}_i^T \mathbf{v}_c)}{\partial \mathbf{u}_w} - \left(\frac{1}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) \frac{\partial}{\partial \mathbf{u}_w} \left(\sum_{x=1}^{x=W} \exp(\mathbf{u}_x^T \mathbf{v}_c) \right) \right)$	(8)
$= -\sum_{i=1}^{i=W} y_i \left(\frac{\partial(\mathbf{u}_i^T \mathbf{v}_c)}{\partial \mathbf{u}_w} - \sum_{x=1}^{x=W} \frac{\frac{\partial}{\partial \mathbf{u}_w} (\exp(\mathbf{u}_x^T \mathbf{v}_c))}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right)$	(9)

When $i = w$ in Equation (9):

$= - \left(\dots + y_w \left(\frac{\partial(\mathbf{u}_w^T \mathbf{v}_c)}{\partial \mathbf{u}_w} - \frac{\frac{\partial}{\partial \mathbf{u}_w} (\exp(\mathbf{u}_1^T \mathbf{v}_c) + \dots + \exp(\mathbf{u}_w^T \mathbf{v}_c) + \dots + \exp(\mathbf{u}_W^T \mathbf{v}_c))}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) + \dots \right)$	(10)
$= - \left(\dots + y_w \left(\mathbf{v}_c - \frac{\exp(\mathbf{u}_w^T \mathbf{v}_c) \mathbf{v}_c}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) + \dots \right)$	(11)
$= - \left(\dots + y_w (1 - \hat{y}_w) \mathbf{v}_c + \dots \right)$	(12)
$= - \left(\dots + y_w \mathbf{v}_c - y_w \hat{y}_w \mathbf{v}_c + \dots \right)$	(13)

Now let's consider all cases where $i \neq w$ in Equation (9):

=	$-\left(\dots + y_{i \neq w} \left(\frac{\partial(\mathbf{u}_{i \neq w}^T \mathbf{v}_c)}{\partial \mathbf{u}_w} - \frac{\frac{\partial}{\partial \mathbf{u}_w}(\exp(\mathbf{u}_l^T \mathbf{v}_c) + \dots + \exp(\mathbf{u}_w^T \mathbf{v}_c) + \dots + \exp(\mathbf{u}_w^T \mathbf{v}_c))}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) + \dots \right)$	(14)
---	---	------

=	$-\left(\dots + y_{i \neq w} \left(0 - \frac{\exp(\mathbf{u}_w^T \mathbf{v}_c) \mathbf{v}_c}{\sum_{m=1}^{m=W} \exp(\mathbf{u}_m^T \mathbf{v}_c)} \right) + \dots \right)$	(15)
---	---	------

=	$-\left(\dots + y_{i \neq w} (0 - \hat{y}_w) \mathbf{v}_c + \dots \right)$	(16)
---	---	------

=	$-\left(\dots + -y_{i \neq w} \hat{y}_w \mathbf{v}_c + \dots \right)$	(17)
---	--	------

Combining Equations (13) and (17):

=	$-\left(\dots + y_w \mathbf{v}_c - y_w \hat{y}_w \mathbf{v}_c - y_{i \neq w} \hat{y}_w \mathbf{v}_c + \dots \right)$	(18)
---	---	------

=	$-y_w \mathbf{v}_c + y_w \hat{y}_w \mathbf{v}_c + y_{i \neq w} \hat{y}_w \mathbf{v}_c + \dots$	(19)
---	--	------

=	$-y_w \mathbf{v}_c + \hat{y}_w \mathbf{v}_c \sum_{n=1}^{n=W} y_n$	(20)
---	---	------

Now, $\sum_n y_n = 1$ because \mathbf{y} is a one-hot label vector. Substituting in Equation (20):

=	$-y_w \mathbf{v}_c + \hat{y}_w \mathbf{v}_c (1)$	(21)
---	--	------

=	$\left(\hat{y}_w - y_w \right) \mathbf{v}_c$	(22)
---	---	------

Substituting in Equation (3):

$\frac{\partial(CE(\mathbf{y}, \hat{\mathbf{y}}))}{\partial \mathbf{U}}$	$\begin{bmatrix} \left(\hat{y}_l - y_l \right) \mathbf{v}_c \\ \vdots \\ \left(\hat{y}_o - y_o \right) \mathbf{v}_c \\ \vdots \\ \left(\hat{y}_w - y_w \right) \mathbf{v}_c \end{bmatrix}$	(23)
--	---	------

Since only **o** is the expected word ($y_o = 1$):

$\frac{\partial(CE(\mathbf{y}, \hat{\mathbf{y}}))}{\partial \mathbf{U}} = \begin{cases} \left(\hat{y}_w - 1\right) \mathbf{v}_c, & w = o \\ \hat{y}_w \mathbf{v}_c, & \text{otherwise} \end{cases}$	(24)
---	------