

Solution to Problem (a):

Since $y_o = 1$ and $y_w = 0$ for any $w \in \text{Vocab} \setminus \{o\}$, The equation (3) is obvious. ■

Solution to Problem (b):

According to the equations (1) and (2), note that

$$\begin{aligned} J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) &= -\log \left(\frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \right) \\ &= -\mathbf{u}_o^T \mathbf{v}_c + \log \left(\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c) \right). \end{aligned}$$

Thus, the desired partial derivative is,

$$\begin{aligned} \frac{\partial}{\partial \mathbf{v}_c} J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) &= -\mathbf{u}_o + \sum_{x \in \text{Vocab}} \left(\frac{\exp(\mathbf{u}_x^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \mathbf{u}_x \right) \\ &= -\mathbf{u}_o + \sum_{x \in \text{Vocab}} (\hat{y}_x \mathbf{u}_x) \\ &= \sum_{x \in \text{Vocab}} ((\hat{y}_x - y_x) \mathbf{u}_x) \\ &= \mathbf{U}(\hat{\mathbf{y}} - \mathbf{y}). \end{aligned}$$

Note that this solution strictly follows “the shape convention”. (Transposed answer could be possible in some case.) ■

Solution to Problem (c):

We may take advantage of the first equation in the previous solution. First, when $w = o$,

$$\begin{aligned} \frac{\partial}{\partial \mathbf{u}_o} J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) &= -\mathbf{v}_c + \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \mathbf{v}_c \\ &= -\mathbf{v}_c + \hat{y}_o \mathbf{v}_c. \end{aligned}$$

On the other hand, when $w \neq o$,

$$\begin{aligned} \frac{\partial}{\partial \mathbf{u}_w} J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) &= \frac{\exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c)} \mathbf{v}_c \\ &= \hat{y}_w \mathbf{v}_c. \end{aligned}$$

In short, for any $w \in \text{Vocab}$,

$$\frac{\partial}{\partial \mathbf{u}_w} J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = (\hat{y}_w - y_w) \mathbf{v}_c.$$

■

Solution to Problem (d):

$$\frac{\partial}{\partial \mathbf{U}} J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = \begin{bmatrix} \frac{\partial J(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_1} & \frac{\partial J(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_2} & \dots & \frac{\partial J(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{|\text{Vocab}|}} \end{bmatrix}.$$

■

Solution to Problem (e):

$$\sigma'(x) = \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{1}{1 + e^{-x}} \cdot \frac{1}{e^x + 1} = \sigma(x) \cdot (1 - \sigma(x)).$$

■

Solution to Problem (f):

Note that $\sigma(-x) = 1 - \sigma(x)$ and $\frac{d}{dx} \log(\sigma(x)) = \frac{\sigma(x)(1-\sigma(x))}{\sigma(x)} = 1 - \sigma(x)$.

$$\begin{aligned} \frac{\partial}{\partial \mathbf{v}_c} J_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) &= -(1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \mathbf{u}_o - \sum_{k=1}^K (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c)) (-\mathbf{u}_k) \\ &= (\sigma(\mathbf{u}_o^T \mathbf{v}_c) - 1) \mathbf{u}_o + \sum_{k=1}^K \sigma(\mathbf{u}_k^T \mathbf{v}_c) \mathbf{u}_k, \end{aligned}$$

$$\frac{\partial}{\partial \mathbf{u}_o} J_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = (\sigma(\mathbf{u}_o^T \mathbf{v}_c) - 1) \mathbf{v}_c,$$

$$\frac{\partial}{\partial \mathbf{u}_k} J_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = \sigma(\mathbf{u}_k^T \mathbf{v}_c) \mathbf{v}_c.$$

The negative sampling loss function is much more efficient than the naive-softmax loss, because it uses only a portion of vocabulary, so the computational burden gets significantly decreased.

■

Solution to Problem (g):

- (i) $\frac{\partial}{\partial \mathbf{U}} J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) = \sum_{-m \leq j \leq m; j \neq 0} \frac{\partial}{\partial \mathbf{U}} J(\mathbf{v}_c, w_{t+j}, \mathbf{U}),$
- (ii) $\frac{\partial}{\partial \mathbf{v}_c} J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) = \sum_{-m \leq j \leq m; j \neq 0} \frac{\partial}{\partial \mathbf{v}_c} J(\mathbf{v}_c, w_{t+j}, \mathbf{U}),$ and
- (iii) $\frac{\partial}{\partial \mathbf{v}_w} J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) = \mathbf{0},$ when $w \neq c.$

■