

# Natural Language Processing: Homework 1

Taewook Kang

September 24, 2024

## Problem (a)

Prove that the loss  $J_{\text{naive-softmax}}$  is the same as the cross-entropy loss between  $\mathbf{y}$  and  $\hat{\mathbf{y}}$ .

*Solution.* The elements of one-hot vector  $\mathbf{y}_w$  (scalar) can be defined as:

$$\mathbf{y}_w = \begin{cases} 1 & \text{if } w = o \\ 0 & \text{else.} \end{cases} \quad (1)$$

Thus, the cross-entropy loss can be reformulated by utilizing the elements of  $\mathbf{y}_w$  as follows:

$$- \sum_{w \in \text{Vocab}} \mathbf{y}_w \log(\hat{\mathbf{y}}_w) = -\mathbf{y}_o \log(\hat{\mathbf{y}}_o) - \sum_{\substack{w \in \text{Vocab} \\ w \neq o}} \mathbf{y}_w \log(\hat{\mathbf{y}}_w) \quad (2)$$

$$= -1 \times \log(\hat{\mathbf{y}}_o) - \sum_{\substack{w \in \text{Vocab} \\ w \neq o}} 0 \times \log(\hat{\mathbf{y}}_w) \quad (3)$$

$$= -\log(\hat{\mathbf{y}}_o). \quad (4)$$

$$\therefore J_{\text{naive-softmax}} = -\log(\hat{\mathbf{y}}_o) = - \sum_{w \in \text{Vocab}} \mathbf{y}_w \log(\hat{\mathbf{y}}_w) \quad (5)$$

## Problem (b)

Compute the partial derivative of  $J_{\text{naive-softmax}}(\mathbf{v}_c, o, U)$  with respect to  $\mathbf{v}_c$ .

*Solution.*

$$J_{\text{naive-softmax}} = -\log \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (6)$$

$$= -\log \exp(\mathbf{u}_o^\top \mathbf{v}_c) + \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c) \quad (7)$$

$$= -\mathbf{u}_o^\top \mathbf{v}_c + \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c). \quad (8)$$

Using equation (8), we can express the partial derivative  $\frac{\partial J_{\text{naive-softmax}}}{\partial \mathbf{v}_c}$  as follows:

$$\frac{\partial J_{\text{naive-softmax}}}{\partial \mathbf{v}_c} = -\frac{\partial}{\partial \mathbf{v}_c} \mathbf{u}_o^\top \mathbf{v}_c + \frac{\partial}{\partial \mathbf{v}_c} \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c). \quad (9)$$

The first term of (9),  $\frac{\partial}{\partial \mathbf{v}_c} \mathbf{u}_o^\top \mathbf{v}_c$ , can be computed as follows:

$$-\frac{\partial}{\partial \mathbf{v}_c} \mathbf{u}_o^\top \mathbf{v}_c = -\mathbf{u}_o. \quad (10)$$

The second term of (9) can be computed as follows:

$$\frac{\partial}{\partial \mathbf{v}_c} \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c) = \frac{1}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \left[ \frac{\partial}{\partial \mathbf{v}_c} \sum_{x \in \text{Vocab}} \exp \mathbf{u}_x^\top \mathbf{v}_c \right] \quad (11)$$

$$= \frac{1}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \left[ \sum_{x \in \text{Vocab}} \mathbf{u}_x \exp \mathbf{u}_x^\top \mathbf{v}_c \right] \quad (12)$$

$$= \sum_{x \in \text{Vocab}} \frac{\mathbf{u}_x \exp \mathbf{u}_x^\top \mathbf{v}_c}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (13)$$

$$= \sum_{x \in \text{Vocab}} \mathbf{u}_x p(x|c). \quad (14)$$

Using (10) and (14), (9) can be expressed as the following equation:

$$\frac{\partial J_{\text{naive-softmax}}(\mathbf{v}_c, o, U)}{\partial \mathbf{v}_c} = -\mathbf{u}_o + \sum_{x \in \text{Vocab}} \mathbf{u}_x p(x|c). \quad (15)$$

Since  $\mathbf{y}$  is a one-hot vector with a 1 at the position corresponding to the true outside word  $o$  and 0 everywhere else, the following relationship holds::

$$\mathbf{u}_o = U \mathbf{y}. \quad (16)$$

Also, given that  $p(x|c) = \hat{\mathbf{y}}_x$ , the relationship can be expressed as:

$$\sum_{x \in \text{Vocab}} \mathbf{u}_x p(x|c) = U \hat{\mathbf{y}}. \quad (17)$$

Therefore, using (16) and (17), equation (15) can be reformulated in terms of  $\mathbf{y}$ ,  $\hat{\mathbf{y}}$ , and  $U$ .

$$\frac{\partial J_{\text{naive-softmax}}(\mathbf{v}_c, o, U)}{\partial \mathbf{v}_c} = U(\hat{\mathbf{y}} - \mathbf{y}). \quad (18)$$

## Problem (c)

Compute the partial derivative of  $J_{\text{naive-softmax}}(\mathbf{v}_c, o, U)$  with respect to each of the 'outside' word vectors  $\mathbf{u}_w$ 's.

*Solution.* By equation (8),  $\frac{\partial J_{\text{naive-softmax}}(\mathbf{v}_c, o, U)}{\partial \mathbf{u}_w}$  can be calculated as follows:

$$\frac{\partial J_{\text{naive-softmax}}(\mathbf{v}_c, o, U)}{\partial \mathbf{u}_w} = \frac{\partial}{\partial \mathbf{u}_w} \left\{ -\mathbf{u}_o^\top \mathbf{v}_c + \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c) \right\}. \quad (19)$$

There are two cases for  $w$ : either  $w = o$  or  $w \neq o$ .

i)  $w = o$

The first term of (19) can be calculated as follows, since  $w = o$ :

$$\frac{\partial}{\partial \mathbf{u}_w} -\mathbf{u}_o^\top \mathbf{v}_c = \frac{\partial}{\partial \mathbf{u}_o} -\mathbf{u}_o^\top \mathbf{v}_c = -\mathbf{v}_c. \quad (20)$$

The second term of (19) can be calculated as follows:

$$\frac{\partial}{\partial \mathbf{u}_w} \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c) = \frac{1}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \left\{ \frac{\partial}{\partial \mathbf{u}_w} \sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^\top \mathbf{v}_c) \right\} \quad (21)$$

$$= \frac{1}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \left\{ \frac{\partial}{\partial \mathbf{u}_o} \sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^\top \mathbf{v}_c) \right\} \quad (22)$$

$$= \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c) \mathbf{v}_c}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (23)$$

$$= p(o|c) \mathbf{v}_c \quad (24)$$

$$= \hat{\mathbf{y}}_o \mathbf{v}_c. \quad (25)$$

Therefore, (19) can be expressed as the following equation by combining the results of (20) and (25):

$$\frac{\partial J_{\text{naive-softmax}}(\mathbf{v}_c, o, U)}{\partial \mathbf{u}_w} = (\hat{\mathbf{y}}_w - 1) \mathbf{v}_c = (\hat{\mathbf{y}}_o - 1) \mathbf{v}_c. \quad (26)$$

ii)  $w \neq o$

The first term of (19) can be calculated as follows, since  $w \neq o$ :

$$\frac{\partial}{\partial \mathbf{u}_w} -\mathbf{u}_o^\top \mathbf{v}_c = 0 \quad (27)$$

The second term of (19) can be calculated as follows:

$$\frac{\partial}{\partial \mathbf{u}_w} \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c) = \frac{1}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \left\{ \frac{\partial}{\partial \mathbf{u}_w} \sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^\top \mathbf{v}_c) \right\} \quad (28)$$

$$= \frac{1}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \frac{\partial \{\dots + \exp(\mathbf{u}_w^\top \mathbf{v}_c) + \dots\}}{\partial \mathbf{u}_w} \quad (29)$$

$$= \frac{\exp(\mathbf{u}_w^\top \mathbf{v}_c) \mathbf{v}_c}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (30)$$

$$= p(w|c) \mathbf{v}_c \quad (31)$$

$$= \hat{\mathbf{y}}_w \mathbf{v}_c. \quad (32)$$

By combining the results from equations (27) and (32), we can express equation (19) as follows:

$$\frac{\partial J_{\text{naive-softmax}}(\mathbf{v}_c, o, U)}{\partial \mathbf{u}_w} = \hat{\mathbf{y}}_w \mathbf{v}_c. \quad (33)$$

## Problem (d)

Write down the partial derivative of  $J_{\text{naive-softmax}}(\mathbf{v}_c, o, U)$  with respect to  $U$ .

*Solution.* By equation 8,  $\frac{\partial J_{\text{naive-softmax}}(\mathbf{v}_c, o, U)}{\partial U}$  can be calculated as follows:

$$\frac{\partial J_{\text{naive-softmax}}(\mathbf{v}_c, o, U)}{\partial U} = \frac{\partial}{\partial U} \left[ -\mathbf{u}_o^\top \mathbf{v}_c + \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c) \right] \quad (34)$$

$$= -\frac{\partial \mathbf{u}_o^\top \mathbf{v}_c}{\partial U} + \frac{\partial}{\partial U} \left[ \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c) \right]. \quad (35)$$

The first term of (35) can be calculated as follows, since  $\mathbf{y}$  is an one-hot vector:

$$-\frac{\partial \mathbf{u}_o^\top \mathbf{v}_c}{\partial U} = -\mathbf{v}_c \mathbf{y}^\top. \quad (36)$$

The second term of (35) can be expressed as follows:

$$\frac{\partial}{\partial U} \left[ \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c) \right] = \frac{1}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \left\{ \frac{\partial}{\partial U} \sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^\top \mathbf{v}_c) \right\} \quad (37)$$

$$= \sum_{x \in \text{Vocab}} \frac{\exp(\mathbf{u}_x^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \frac{\partial \mathbf{u}_x^\top \mathbf{v}_c}{\partial U} \quad (38)$$

$$= \sum_{x \in \text{Vocab}} \hat{\mathbf{y}}_x \frac{\partial \mathbf{u}_x^\top \mathbf{v}_c}{\partial U} \quad (39)$$

$$= \mathbf{v}_c \hat{\mathbf{y}}^\top. \quad (40)$$

Therefore, using (36) and (40), equation (35) can be reformulated as a following equation:

$$\frac{\partial J_{\text{naive-softmax}}(\mathbf{v}_c, o, U)}{\partial U} = -\mathbf{v}_c \mathbf{y}^\top + \mathbf{v}_c \hat{\mathbf{y}}^\top \quad (41)$$

$$= \mathbf{v}_c (\hat{\mathbf{y}} - \mathbf{y})^\top \quad (42)$$

$$= [\hat{y}_1 \mathbf{v}_c \quad \hat{y}_2 \mathbf{v}_c \quad \cdots \quad (\hat{y}_o - 1) \mathbf{v}_c \quad \cdots \quad \hat{y}_{|\text{Vocab}|} \mathbf{v}_c]. \quad (43)$$

## Problem (e.i)

Please repeat parts (b) and (c), computing the partial derivatives of  $J_{\text{neg-sample}}$  with respect to  $\mathbf{v}_c$ , with respect to  $\mathbf{u}_o$ , and with respect to the  $s$ th negative sample  $\mathbf{u}_{w_s}$ .

*Solution.* The derivative of a sigmoide function is computed as follows:

$$\frac{\partial \sigma(\mathbf{x})}{\partial \mathbf{x}} = \sigma(\mathbf{x})(1 - \sigma(\mathbf{x})). \quad (44)$$

Based on (44), all derivatives of  $J_{\text{neg-sample}}$  can be expressed as follows:

$$\frac{\partial J_{\text{neg-sample}}(\mathbf{v}_c, o, U)}{\partial \mathbf{v}_c} = \frac{\partial}{\partial \mathbf{v}_c} \left[ -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \right] \quad (45)$$

$$= -\frac{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c))}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} \frac{\partial \sigma(\mathbf{u}_o^\top \mathbf{v}_c)}{\partial \mathbf{v}_c} \quad (46)$$

$$- \sum_{s=1}^K \frac{\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)(1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c))}{\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)} \frac{\partial \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)}{\partial \mathbf{v}_c} \quad (47)$$

$$= -(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)) \frac{\partial \sigma(\mathbf{u}_o^\top \mathbf{v}_c)}{\partial \mathbf{v}_c} - \sum_{s=1}^K (1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \frac{\partial \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)}{\partial \mathbf{v}_c} \quad (48)$$

$$= -(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)) \mathbf{u}_o + \sum_{s=1}^K (1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \mathbf{u}_{w_s}. \quad (49)$$

$$\frac{\partial J_{\text{neg-sample}}(\mathbf{v}_c, o, U)}{\partial \mathbf{u}_o} = \frac{\partial}{\partial \mathbf{u}_o} \left[ -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \right] \quad (50)$$

$$= -\frac{\partial \log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c))}{\partial \mathbf{u}_o} - \frac{\partial}{\partial \mathbf{u}_o} \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \quad (51)$$

$$= -\frac{\partial \log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c))}{\partial \mathbf{u}_o} - 0 \quad (52)$$

$$= -\frac{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c))}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} \frac{\partial \sigma(\mathbf{u}_o^\top \mathbf{v}_c)}{\partial \mathbf{u}_o} \quad (53)$$

$$= -(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)) \frac{\partial \sigma(\mathbf{u}_o^\top \mathbf{v}_c)}{\partial \mathbf{u}_o} \quad (54)$$

$$= -(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)) \mathbf{v}_c. \quad (55)$$

$$\frac{\partial J_{\text{neg-sample}}(\mathbf{v}_c, o, U)}{\partial \mathbf{u}_{w_s}} = \frac{\partial}{\partial \mathbf{u}_{w_s}} \left[ -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \right] \quad (56)$$

$$= -\frac{\partial \log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c))}{\partial \mathbf{u}_{w_s}} - \frac{\partial}{\partial \mathbf{u}_{w_s}} \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \quad (57)$$

$$= -0 - \frac{\partial}{\partial \mathbf{u}_{w_s}} \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \quad (58)$$

$$= -\frac{\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)(1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c))}{\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)} \frac{\partial \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)}{\partial \mathbf{u}_{w_s}} \quad (59)$$

$$= -(1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \frac{\partial \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)}{\partial \mathbf{u}_{w_s}} \quad (60)$$

$$= (1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \mathbf{v}_c. \quad (61)$$

## Problem (e.ii)

Describe with one sentence why this loss function is much more efficient to compute than the naive-softmax loss.

*Solution.* Negative sampling is more efficient than naive-softmax because it reduces computational complexity by sampling  $K$  negative words instead of traversing the entire vocabulary, while still approximating the naive-softmax result effectively through Monte-Carlo sampling ( $K \ll |\text{Vocab}|$ ).

**Problem (f.1)**

Write down partial derivative:  $\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial U}$ .

*Solution.*

$$\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial U} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(\mathbf{v}_c, w_{t+j}, U)}{\partial U}. \quad (62)$$

**Problem (f.ii)**

Write down partial derivative:  $\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial \mathbf{v}_c}$ .

*Solution.*

$$\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial \mathbf{v}_c} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(\mathbf{v}_c, w_{t+j}, U)}{\partial \mathbf{v}_c}. \quad (63)$$

**Problem (f.iii)**

Write down partial derivative:  $\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial \mathbf{v}_w}$  when  $w \neq c$ .

*Solution.*

$$\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial \mathbf{v}_w} = 0. \quad (64)$$