

Natural Language Processing with Deep Learning: HW 2

Author: *LIUQiaoan 520030910220*

Course: *CS224n: Natural Language Processing with Deep Learning, Stanford, Winter 2021*

Date: *October 28, 2022*

(a)

Since o is the only true outside word, then

$$\forall w \in \text{Vocab}, y_w = \begin{cases} 1, & w = o \\ 0, & \text{o.w.} \end{cases}$$

So for all $o \in \text{Vocab}$, we have

$$\begin{aligned} - \sum_{w \in \text{Vocab}} y_w \log(\hat{y}_w) &= - \sum_{w \in \text{Vocab}} 1\{y_w = 1\} \log(\hat{y}_w) \\ &= - \log(\hat{y}_o) \\ &= - \log P(O = o | C = c) \\ &= \mathbf{J}_{\text{naive_softmax}}(\mathbf{v}_c, o, \mathbf{U}) \end{aligned}$$

(b)

$$\begin{aligned} \frac{\partial \mathbf{J}_{\text{naive_softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} &= \frac{\partial \{-\mathbf{u}_o^T \mathbf{v}_c + \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)\}}{\partial \mathbf{v}_c} \\ &= -\mathbf{u}_o + \frac{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c) \mathbf{u}_w}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \\ &= -\mathbf{u}_o + \sum_{w \in \text{Vocab}} P(O = w | C = c) \mathbf{u}_w \\ &= -\mathbf{u}_o + \sum_{w \in \text{Vocab}} \hat{y}_w \mathbf{u}_w \\ &= \mathbf{U}(\hat{\mathbf{y}} - \mathbf{y}) \quad (\text{This is such a clean form!}) \end{aligned}$$

(c)

When $w = o$,

$$\begin{aligned}\frac{\partial \mathbf{J}_{\text{naive_softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} &= \frac{\partial \{-\mathbf{u}_o^T \mathbf{v}_c + \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)\}}{\partial \mathbf{u}_o} \\ &= -\mathbf{v}_c + \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c) \mathbf{v}_c}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \\ &= -\mathbf{v}_c + \hat{y}_o \mathbf{v}_c \\ &= (\hat{y}_o - y_o) \mathbf{v}_c\end{aligned}$$

When $w \neq o$,

$$\begin{aligned}\frac{\partial \mathbf{J}_{\text{naive_softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_w} &= \frac{\partial \{-\mathbf{u}_o^T \mathbf{v}_c + \log \sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c)\}}{\partial \mathbf{u}_w} \\ &= \frac{\exp(\mathbf{u}_w^T \mathbf{v}_c) \mathbf{v}_c}{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c)} \\ &= \hat{y}_w \mathbf{v}_c\end{aligned}$$

(d)

$$\frac{\partial \mathbf{J}_{\text{naive_softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{U}} = \left[\frac{\partial \mathbf{J}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_1}, \quad \frac{\partial \mathbf{J}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_2}, \quad \dots, \quad \frac{\partial \mathbf{J}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{|\text{Vocab}|}} \right]$$

(e)

$$\begin{aligned}\sigma'(x) &= \frac{e^x(e^x + 1) - e^x * e^x}{(e^x + 1)^2} \\ &= \frac{e^x}{(e^x + 1)^2} \\ &= \sigma(x)(1 - \sigma(x))\end{aligned}$$

(f)

$$\frac{\partial \mathbf{J}_{\text{neg_sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = -(1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \mathbf{u}_o + \sum_{k=1}^K (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \mathbf{u}_k$$

$$\frac{\partial \mathbf{J}_{\text{neg_sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} = -(1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \mathbf{v}_c$$

$$\frac{\partial \mathbf{J}_{\text{neg_sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_k} = (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \mathbf{v}_c$$

Compare (b) and (c) with (f), we can see that (b) requires matrix multiplication, while (f)

only needs inner product.

(g)

$$\frac{\partial \mathbf{J}_{\text{neg_sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_k} = \sum_{w=1}^K 1\{\mathbf{u}_w = \mathbf{u}_k\} (1 - \sigma(-\mathbf{u}_w^T \mathbf{v}_c)) \mathbf{v}_c$$

(h)

i.

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

ii.

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

iii.

$$\frac{\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{v}_w} = 0 \quad \text{when } w \neq c$$