

## Assignment: 2

a) Show that, naive-softmax loss is same as cross entropy loss between  $y$  and  $\hat{y}$ .

$$-\sum_{w \in \text{vocab}} y_w \log(\hat{y}_w) = -\log(\hat{y}_0)$$

Given, true empirical distribution,

$$\begin{cases} w=0 \\ w \neq 0 \end{cases} \quad \begin{cases} y_w = 1 \\ y_w = 0 \end{cases}$$

$$\text{Cross entropy} = -\sum_{w \in V} y_w \log(\hat{y}_w)$$

$$= - (y_1 \log \hat{y}_1 + y_2 \log \hat{y}_2 + \dots + y_0 \log \hat{y}_0 + \dots + y_{|V|} \log \hat{y}_{|V|})$$

$$= - (0 + 0 + y_0 \log \hat{y}_0 + \dots + 0) = -1 (\log \hat{y}_0)$$

$$\because (w=0, y_w=1)$$

$$= \boxed{-\log(\hat{y}_0)} \Rightarrow \text{Proved}$$

b)  $\rightarrow$  Derivative of  $J_{\text{naivesoftmax}}(v_c, 0, v) = ?$

$\rightarrow$  write terms  $y, \hat{y}, v \sim w \cdot v + v_c$ ,

$\rightarrow$  Present in vectorize form.

$$J(v_c, \theta, U) = -\log(\hat{y}_c)$$

$$\begin{aligned}
 \frac{\partial J}{\partial v_c} &= \frac{\partial}{\partial v_c} \left[ -\log \left( \frac{\exp(u_0^T \cdot v_c)}{\sum_{w \in V} \exp(u_w^T \cdot v_c)} \right) \right] \\
 &= \frac{\partial}{\partial v_c} \left[ -\log(\exp(u_0^T \cdot v_c)) + \log \sum_{w \in V} \exp(u_w^T \cdot v_c) \right] \quad // \text{applying log property} \\
 &= \frac{\partial}{\partial v_c} \left[ -u_0^T \cdot v_c + \log \sum_{w \in V} \exp(u_w^T \cdot v_c) \right] \quad // \log \& \\
 &\quad \exp \text{ are} \\
 &= -\frac{\partial}{\partial v_c} (u_0^T \cdot v_c) + \frac{\partial}{\partial v_c} \log \sum_{w \in V} \exp(u_w^T \cdot v_c) \\
 &\because \frac{\partial x^T a}{\partial x} = \frac{\partial a^T x}{\partial x} = a \\
 &= -\mu_0^T + \frac{1}{\sum_{w \in V} \exp(u_w^T \cdot v_c)} \sum_{w \in V} (\exp(u_w^T \cdot v_c) \cdot u_w)
 \end{aligned}$$

Since,  $\frac{\exp(u_w^T \cdot v_c)}{\sum_{w \in V} \exp(u_w^T \cdot v_c)}$  is

the probability distribution over words.  
 $(\hat{y}_w)$ -

$$\Rightarrow \frac{\partial J}{\partial v_c} = -\mu_0^T + \sum_{w \in V} \hat{y}_w u_w$$

Since,  $y_w = \begin{cases} 1 & \text{if } w=0 \\ 0 & \text{if } w \neq 0 \end{cases}$ , we can write.

$$u_0 = y_w u_w$$

$$\begin{aligned} \Rightarrow \frac{\partial J}{\partial v_c} &= -\left(\sum_{w \in V} y_w u_w + \hat{y}_w u_w\right) \\ &= \sum_{w \in V} u_w (-y_w + \hat{y}_w) \\ &= \sum_{w \in V} u_w (\hat{y}_w - y) \end{aligned}$$

After vectorizing  $\Rightarrow \frac{\partial J}{\partial v_c} = U(\hat{y} - y)$

(C) Compute derivative w.r.t  $u_w$ . Consider both cases when  $w=0$  (outside word) and  $w \neq 0$  (all other words).

$$\frac{\partial J}{\partial u_w} = \frac{\partial}{\partial u_w} (-\log(\hat{y}_w))$$

$$= \frac{\partial}{\partial u_w} \left( -\log \frac{\exp(u_0^T v_c)}{\sum_{w \in V} \exp(u_w^T \cdot v_c)} \right)$$

$$\frac{\partial}{\partial u_w} \left[ -\log(\exp(u_0^T v_c)) + \log \sum_{w \in V} \exp(u_w^T \cdot v_c) \right]$$

$$\Rightarrow \frac{\partial}{\partial u_w} = \frac{\partial}{\partial u_w} \left( -u_0^T \cdot v_c + \log \sum_{w \in V} \exp(u_w^T \cdot v_c) \right)$$

$$= -\frac{\partial u_0^T(v_c)}{\partial u_w} + \frac{\partial}{\partial u_w} \left[ \log \sum_{w \in V} \exp(u_w^T \cdot v_c) \right].$$

$$= -\frac{\partial u_0^T(v_c)}{\partial u_w} + \frac{1}{\sum_{w \in V} \exp(u_w^T \cdot v_c)} \frac{\partial}{\partial u_w} \left[ \sum_{w \in V} \exp(u_w^T \cdot v_c) \right]$$

$\Rightarrow$  (derivative of sum = sum of derivatives)

$$= -\frac{\partial u_0^T(v_c)}{\partial u_w} + \frac{1}{\sum_{w \in V} \exp(u_w^T \cdot v_c)} \cdot \sum_{w \in V} \exp(u_w^T \cdot v_c) \cdot v_c$$

Since,  $\frac{\partial u_0^T}{\partial u_w} = \begin{cases} 1 & \text{if } w=0 \\ 0 & \text{if } w \neq 0 \end{cases}$  and equivalent to  $y_w$

$$\Rightarrow \sum_{w \in V} -y_w v_c + \frac{1}{\sum_{w \in V} \exp(u_w^T \cdot v_c)} \sum_{w \in V} \exp(u_w^T \cdot v_c) \cdot v_c$$

Since,  $\frac{\exp(u_w^T \cdot v_c)}{\sum_{w \in V} \exp(u_w^T \cdot v_c)}$  is conditional probability, denoted as  $\hat{y}_w$ .

$$\Rightarrow \frac{\delta J}{\delta u_w} = \sum_{v \in V} (-y_w v_c + \hat{y}_w v_c)$$

$$= \boxed{\sum_{v \in V} v_c (\hat{y}_w - y_w)}$$

case 1:  $w=0$ . (context word is the outside word)

$$\frac{\delta J}{\delta u_{w=0}} = v_c (\hat{y}_{w=0} - y_{w=0})$$

Since,  $y_{w=0} = 1$

$$\Rightarrow \boxed{\frac{\delta J}{\delta u_{w=0}} = v_c (\hat{y}_{w=0} - 1)}$$

case 2:  $w \neq 0$  Since  $y_{w \neq 0} = 0$

$$\Rightarrow \boxed{\frac{\delta J}{\delta u_{w \neq 0}} = v_c \hat{y}_{w \neq 0}}$$

$$\frac{\delta J}{\delta u_w} = \begin{cases} v_c (\hat{y}_w - 1) & \text{if } w=0 \\ v_c (\hat{y}_w) & \text{otherwise} \end{cases}$$

(d) Compute partial derivative of  $J_{\text{naive-softmax}}$  with respect to  $U$ .

Derivative of scalar  $y$  by a matrix  $A$  is given as

$$\frac{\partial y}{\partial A_{m \times n}} = \begin{bmatrix} \frac{\partial y}{\partial A_{11}} & \frac{\partial y}{\partial A_{12}} & \cdots & \frac{\partial y}{\partial A_{1n}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial y}{\partial A_{m1}} & \frac{\partial y}{\partial A_{m2}} & \cdots & \frac{\partial y}{\partial A_{mn}} \end{bmatrix}$$

Given  $ww = \text{outside word } w$ ,

$\Rightarrow$  Derivative of  $y$  w.r.t  $U$  is -

$$\frac{\partial J_{\text{naive-softmax}}(v_c, o, U)}{\partial U} = \left[ \frac{\partial J}{\partial u_1}, \frac{\partial J}{\partial u_2}, \dots, \frac{\partial J}{\partial u_m} \right]$$

where,

$$\frac{\partial J}{\partial u_w} = \begin{cases} v_c(y - 1) & \text{if } w=0 \\ v_c y & \text{otherwise} \end{cases}$$

(e) Compute derivative of sigmoid

$$\text{function } \sigma(x) = \frac{e^x}{e^x + 1} = \frac{1}{1 + e^{-x}}$$

$$\sigma = \frac{1}{1+e^{-x}} \quad (\text{given}).$$

$$\begin{aligned}
\frac{\partial \sigma}{\partial x} &= \frac{\partial}{\partial x} \left[ \frac{1}{1+e^{-x}} \right] = \frac{\partial}{\partial x} (1+e^{-x})^{-1} \\
&= -1 (1+e^{-x})^{-2} [-e^{-x}] \\
&= \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}} \\
&= \frac{1}{1+e^{-x}} \cdot \left( \frac{e^{-x} + 1 - 1}{1+e^{-x}} \right) \\
&= \frac{1}{1+e^{-x}} \left[ \left( \frac{e^{-x} + 1}{1+e^{-x}} \right) - \frac{1}{1+e^{-x}} \right] \\
&= \boxed{\sigma(x) [1 - \sigma(x)]}
\end{aligned}$$

(b) Consider negative sampling loss: for a center word  $c$ , outside word  $o$ ,

$$J_{\text{neg-sampling}}(v_c, o, U) = -\log(\sigma(u_o^T v_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T v_c)) -$$

Compute derivatives of  $J_{\text{neg-sample}}$  w.r.t  $v_c$ ,  $u_o$  and  $u_k$ . Tell why this function is efficient.

① Derivative of  $J_{\text{neg-sample}}$  w.r.t  $v_c$ :

$$\frac{\partial}{\partial v_c} = \frac{\partial}{\partial v_c} \left[ -\log(\sigma(u_0^T v_c)) - \sum_{k=1}^K \log(\sigma(-w_0^k v_c)) \right]$$

$$= \frac{\partial}{\partial v_c} \left[ -\log(\sigma(u_0^T v_c)) - \sum_{k=1}^K \frac{\partial}{\partial v_c} \log(\sigma(-w_0^k v_c)) \right]$$

$$= \frac{-1}{\sigma(u_0^T v_c)} \frac{\partial}{\partial v_c} (\sigma(u_0^T v_c)) - \sum_{k=1}^K \left[ \frac{1}{\sigma(-w_0^k v_c)} \frac{\partial}{\partial v_c} \right]$$

$$\boxed{\frac{\partial}{\partial x} (\sigma(x)) = \sigma(x)(1-\sigma(x))} \quad (\text{Part e})$$

$$\Rightarrow \frac{-1}{\sigma(u_0^T v_c)} \cdot \cancel{\sigma(u_0^T v_c)(1-\sigma(u_0^T v_c))} \cdot u_0 - \sum_{k=1}^K$$

$$\frac{1}{\sigma(-w_0^k v_c)} \cdot \cancel{\sigma(-w_0^k v_c)(1-\sigma(-w_0^k v_c))} \cdot (-w_k)$$

$$\Rightarrow -u_0(1-\sigma(u_0^T v_c)) + \sum_{k=1}^K u_k(1-\sigma(-w_k^T v_c))$$

$$\Rightarrow \boxed{u_0(1-\sigma(u_0^T v_c)) - \sum_{k=1}^K u_k(1-\sigma(-w_k^T v_c))}$$

3) Derivative of  $J_{\text{neg-sampij}}$  w.r.t  $u_0$ :

$$\frac{\partial J}{\partial u_0} = \frac{\partial}{\partial u_0} \left[ -\log(\sigma(u_0^T v_c)) - \sum_{k=1}^n \log(\sigma(-u_k^T v_c)) \right]$$

$\Rightarrow$  Since,  $a \notin \{w_1, \dots, w_K\}$

$$\Rightarrow \frac{\partial}{\partial u_0} \left[ \sum_{k=1}^n \log(\sigma(-u_k^T v_c)) \right] = 0$$

$$\Rightarrow \frac{\partial}{\partial u_0} (-\log(\sigma(u_0^T v_c)) - 0)$$

$$= -\frac{1}{\sigma(u_0^T v_c)} \frac{\partial}{\partial u_0} (\sigma(u_0^T v_c)) = -\frac{1}{\sigma(u_0^T v_c)} \cdot \frac{\sigma(u_0^T v_c)(1-\sigma(u_0^T v_c))}{v_c}$$

$$\Rightarrow -v_c(1-\sigma(u_0^T v_c))$$

$$\Rightarrow \boxed{v_c(\sigma(u_0^T v_c) - 1)}$$

3) Compute Derivative of  $J_{\text{neg-sampij}}$  w.r.t  $u_k$ :

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

Since,  $0 \notin \{w_1, w_2, \dots, w_n\}$ .

$$\Rightarrow \frac{\partial J}{\partial u_n} (-\log(\sigma(u_0^T v_c))) = 0$$

$$\Rightarrow \frac{\partial J}{\partial u_k} = 0 - \sum_{k=1}^K \frac{\partial}{\partial u_k} (\log(\sigma(-u_k^T v_c))).$$

$$\Rightarrow \sum_{k=1}^K \frac{1}{\sigma(-u_k^T v_c)} \sigma(-u_k^T v_c) (1 - \sigma(-u_k^T v_c)) \circ - v_c$$

$\Rightarrow$  all the derivative terms where  $w \neq k$ , becomes 0, while  $w=k$ , only the term  $u_w$  remains

$$\Rightarrow \frac{\partial J}{\partial u_k} = -v_c (1 - \sigma(-u_k^T v_c))$$

$$= \boxed{v_c (\sigma(-u_k^T v_c) - 1)}$$

Negative Sampling loss is much more memory and compute efficient since it requires only  $K$  sampled words.

$O(K)$  - However, Naive-Softmax uses all the words inside vocabulary  $O(V)$  and normalize probabilities

(8) compute derivative of  $J_{\text{neg-sampling}}$  w.r.t  $u_k$ . where  $K$  samples are drawn from vocabulary

Since,  $K$  words are sampled from vocabulary, which can't be distinct. Let's break sum into 2 sums

↳ ① Sum of all sampled words  $w_i = w_n$

② Sum of all sampled words  $w_i \neq w_n$ .

$\Rightarrow$  Consider Eq ① from last part

$$\begin{aligned} \frac{\partial J}{\partial u_k} &= -\frac{\partial}{\partial u_n} \left[ \sum_{i=1}^K \left( \sigma(-u_i^T v_c) \right) \right] \\ &= -\frac{\partial}{\partial u_n} \left[ \sum_{\substack{i \in \{1, \dots, K\} : \\ w_i = w_n}} \log(\sigma(-u_i^T v_c)) + \right. \\ &\quad \left. \sum_{\substack{i \in \{1, \dots, K\} : \\ w_i \neq w_n}} \log(\sigma(-u_i^T v_c)) \right] \end{aligned}$$

$$\text{Now, } \frac{\partial}{\partial u_k} \log(\sigma(-u_i^T v_c; w_i \neq w_n, v_c)) = 0.$$

$$\Rightarrow -\frac{\partial}{\partial u_n} \sum_{\substack{i \in \{1, \dots, K\} : \\ w_i = w_n}} \log(\sigma(-u_i^T v_c))$$

$$\begin{aligned} \Rightarrow -\frac{\partial}{\partial u_n} \sum_{\substack{i \in \{1, \dots, K\} : \\ w_i = w_n}} \log(\sigma(-u_i^T v_c)) &\Rightarrow -\sum_{\substack{i \in \{1, \dots, K\} : \\ w_i = w_n}} \frac{\partial J}{\partial u_n} (\log(\sigma(-u_i^T v_c))) \\ \Rightarrow -\sum_{\substack{w_i = w_n \\ i \in \{1, \dots, K\}}} \frac{1}{\log(\sigma(-u_i^T v_c))} &\cdot \cancel{\log(\sigma(-u_i^T v_c))} \cdot -v_c \cdot (1 - \sigma(-u_n^T v_c)) \end{aligned}$$

$$\Rightarrow \sum_{k=1; w_i = w_n}^K v_c (1 - \sigma(-w_R^T v_c))$$

$$\Rightarrow \left[ \sum_{i \in \{1, \dots, K\}; w_i = w_n} -v_c (\sigma(-u_k^T v_c) - 1) \right] \frac{\partial \sigma}{\partial u_k}$$

when samples  
are drawn  
from vocabulary.

(h)

①  $\frac{\partial J_{\text{skip gram}}(v_c, w_{t-m}, \dots, w_{t+n})}{\partial v} =$

$$\sum_{\substack{j \neq 0 \\ -m < j \leq m}} \frac{\partial J(v_c, w_{t+j}, v)}{\partial v}$$

②  $\frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+n})}{\partial v_c} =$

$$\sum_{\substack{-m < j \leq m \\ j \neq 0}} \frac{\partial J(v_c, w_{t+j}, v)}{\partial v_c}$$

③  $\frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+n})}{\partial v_w; w \neq c} = 0$