

CS 224n Ad Written Part.

1/a)

Prove negative-softmax loss is the same as cross-entropy loss between y and \hat{y} , i.e. show that

$$-\sum_{w \in \text{vocab}} y_w \log(\hat{y}_w) = -\log(y_0)$$

The probability for an outside word w given the center word is defined as

$$y_0 = p(u_0 | v_c) = \frac{\exp(u_0^T v_c)}{\sum_{i \in \text{vocab}} \exp(u_i^T v_c)}$$

Therefore the following equation is

$$-\sum_{w \in \text{vocab}} y_w \log(\hat{y}_w) = -y_1 \log(\hat{y}_1) - y_2 \log(\hat{y}_2) - \dots - y_n \log(\hat{y}_n)$$

Since for true label y , it's a one-hot encoded vector which only true outside context with value 1 and the rest as 0. So the previous equation could be simplified as

$$\begin{aligned} -\sum_{w \in \text{vocab}} y_w \log(\hat{y}_w) &= -1 \cdot \log(y_0) - \sum_{i \in \text{vocab}, i \neq 0} 0 \cdot \log(\hat{y}_i) \\ &= -\log(y_0) \end{aligned}$$

1/b) partial derivative of $J_{\text{naive-softmax}}$ with respect to v_c .

$$J_{\text{naive-softmax}} = -\log P(i=0 | C=c).$$

$$= -\log \frac{\exp(u_0^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)}.$$

$$= \underbrace{-u_0^T v_c}_A + \underbrace{\log \left[\sum_{w \in \text{vocab}} \exp(u_w^T v_c) \right]}_B.$$

$$\frac{\partial J}{\partial v_c} = \frac{\partial A}{\partial v_c} + \frac{\partial B}{\partial v_c}.$$

$$\frac{\partial A}{\partial v_c} = \frac{-\partial u_0^T v_c}{\partial v_c} = -u_0^T.$$

$$\frac{\partial B}{\partial v_c} = \frac{\sum u_w^T \exp(u_w^T v_c)}{\sum \exp(u_w^T v_c)}.$$

$$= \sum_{x \in \text{vocab}} \frac{\exp(u_x^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} \cdot u_x^T.$$

$$= \sum_{x \in \text{vocab}} P(u_x | v_c) \cdot u_x^T.$$

$$= \sum_{x \in \text{vocab}} \hat{y}_x u_x.$$

$$\therefore \frac{\partial J}{\partial v_c} = -u_0 + \sum_{x \in \text{vocab}} \hat{y}_x u_x.$$

The slope of loss function w.r.t center word is the difference between of the true context words and expected context words in the model.

Vc) Part I : if $u_w = u_o$

$$\frac{\partial J}{\partial u_o} = \frac{\partial A}{\partial u_o} + \frac{\partial B}{\partial u_o}$$

$$\cdot \frac{\partial A}{\partial u_o} = -v_c$$

$$\cdot \frac{\partial B}{\partial u_o} = \frac{v_c \exp(u_o^T v_c)}{\sum \exp(u_w^T v_c)}$$

$$\frac{\partial J}{\partial u_o} = -v_c + v_c \cdot \hat{y}_c = (\hat{y}_c - 1) \cdot v_c$$

Part d if $u_w \neq u_o$

$$\cdot \frac{\partial A}{\partial u_i} = 0$$

$$\cdot \frac{\partial B}{\partial u_i} = \frac{v_c \exp(u_i^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} = v_c \cdot \hat{y}_i$$

$$\frac{\partial J}{\partial u_i} = 0 + v_c \cdot \hat{y}_i$$

$$1/A) \quad b(x) = \frac{1}{e^x + 1} = \frac{e^{-x}}{1 + e^{-x}}$$

$$\frac{db}{dx} = -\frac{1}{(e^x + 1)^2} \cdot e^{-x} \cdot -1$$

$$= \frac{e^{-x}}{(e^{-x} + 1)^2}$$

$$= \frac{1}{e^{-x} + 1} \cdot \frac{e^{-x} + 1 - 1}{e^{-x} + 1}$$

$$= \frac{1}{e^{-x} + 1} \cdot \left(1 - \frac{1}{e^{-x} + 1}\right)$$

$$= b(x)(1 - b(x))$$

1/e). $f(x) = \max(0, x)$ is the ReLU

$$\text{if } x > 0, \quad \frac{df(x)}{dx} = 1$$

$$x < 0, \quad \frac{df(x)}{dx} = 0$$

1/2). $\frac{\partial J}{\partial u}$ is just the combination of partial derivatives.

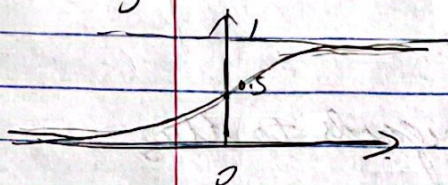
of all u_n which $n \in \text{local}$.

$$\frac{\partial J}{\partial u} = \begin{cases} (\hat{y}_i - 1) \cdot v_i & \text{if } i = 0 \\ \hat{y}_i \cdot v_i & \text{if } i \neq 0 \end{cases}$$

1/3). Negative Sample.

$$J_{\text{neg-sample}}(U_c, 0, V) = -(\log(b(\mu_0^T V_c)) - \sum (\log(b(-\mu_w^T V_c)))$$

As for Sigmoid function



$$\Rightarrow \frac{\partial J}{\partial V_c} = \frac{\partial J}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial V_c}$$

$$= - \frac{\mu_0 b(\mu_0^T V_c) [1 - b(\mu_0^T V_c)]}{b(\mu_0^T V_c)}$$

$$- \sum \frac{-\mu_w b(-\mu_w^T V_c) [1 - b(-\mu_w^T V_c)]}{b(-\mu_w^T V_c)}$$

$$= -\mu_0 [1 - b(\mu_0^T V_c)] + \sum \mu_w [1 - b(-\mu_w^T V_c)]$$

$$\text{ii)} \frac{\partial J}{\partial \mu_0} = -V_c [1 - b(\mu_0^T V_c)]$$

$$\frac{\partial J}{\partial \mu_w} = V_c [1 - b(-\mu_w^T V_c)] \quad \text{it is}$$

Negative Sample is more efficient because it takes only a part of the original sample while loss function considers all the sample.

1/h) Now without the assumption that K samples are distinct. which means we still have $\{w_1, w_2, \dots, w_K\}$ but now $i \neq j$ no longer holds.

$\frac{\partial J_{\text{neg-sample}}}{\partial v_c}$ is still the same

while for partial derivative with regards to w_s

$$\frac{\partial J_{\text{neg-sample}}}{\partial w_s} = n_{w_s} v_c [1 - \sigma(w_s^T v_c)]$$

n_{w_s} is the number of negative samples equals to w_s .

1/i).
$$\frac{\partial J_{\text{skip-gram}}}{\partial u} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(v_c, w_j, u)}{\partial u}$$

$J(v_c, w_j, u)$ could be the J_{softmax} or J_{naive} based on the setting