

Q (a) Prove -:

$$\text{softmax}(x) = \text{softmax}(x+c)$$

$$\text{softmax}(x)_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$\begin{aligned}\text{softmax}(x+c)_i &= \frac{e^{(x+c)_i}}{\sum_j e^{(x+c)_j}} = \frac{e^c e^{x_i}}{e^c \sum_j e^{x_j}} \\ &= \text{softmax}(x)_i\end{aligned}$$

2a) Derive $\frac{d\sigma(x)}{dx}$ where $\sigma(x) = \frac{1}{1+e^{-x}}$

$$\begin{aligned}\frac{d\sigma}{dx} &= -\frac{1}{(1+e^{-x})^2} \times -e^{-x} \\ &= \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1}{(1+e^{-x})} \cdot \frac{e^{-x}}{1+e^{-x}} \\ &= \sigma \left(1 - \frac{1}{1+e^{-x}}\right) \\ &= \sigma(1-\sigma)\end{aligned}$$

b) Derive $\frac{\partial CE}{\partial \theta}$ where $CE(y, \hat{y}) = -\sum_i y_i \log \hat{y}_i$
and $\hat{y}_i = \text{softmax}(\theta_i)$

Firstly let's derive

$$\frac{\partial \hat{y}}{\partial \theta}$$

$$\hat{y}_k = \frac{e^{\theta_k}}{\sum_j e^{\theta_j}}$$

$$\frac{\partial \hat{y}_k}{\partial \theta_i} \text{ when } i=k = \frac{e^{\theta_k}}{\sum_j e^{\theta_j}} + e^{\theta_k} \left(\frac{-1}{\sum_j e^{\theta_j}} \right)^2 \cdot e^{\theta_k}$$

$$= \hat{y}_k - \frac{e^{\theta_k}}{\sum_j e^{\theta_j}} \cdot \frac{e^{\theta_k}}{\sum_j e^{\theta_j}}$$

$$= \hat{y}_k (1 - \hat{y}_k)$$

when $i \neq k$

$$\frac{\partial \hat{y}_k}{\partial \theta_i} = \frac{e^{\theta_k} \cdot -1 \cdot e^{\theta_i}}{(\sum_j e^{\theta_j})^2}$$

$$= -\frac{e^{\theta_k}}{\sum_j e^{\theta_j}} \cdot \frac{e^{\theta_i}}{\sum_j e^{\theta_j}}$$

$$= -\hat{y}_k \cdot \hat{y}_i$$

$$\begin{aligned} \therefore \frac{\partial \hat{y}_k}{\partial \theta_i} &= \hat{y}_k (1 - \hat{y}_k) & \text{when } i=k \\ &= -\hat{y}_k \cdot \hat{y}_i & \text{otherwise} \end{aligned}$$

Now,

$$\frac{\partial CE}{\partial \theta_i} = -2 \left(\sum_j y_j \log \hat{y}_j \right)$$

Since, y is a one-hot vector with all but k th dimension, 0.

$$\frac{\partial CE}{\partial \theta_i} = -2 \left(y_k \log \hat{y}_k \right)$$

$$= -\frac{y_k}{\hat{y}_k} \frac{\partial (\log \hat{y}_k)}{\partial \theta_i}$$

$$= -\frac{1}{\hat{y}_k} \frac{\partial (\log \hat{y}_k)}{\partial \theta_i} \quad (\because y_k = 1)$$

when $i = k$.

$$\frac{\partial CE}{\partial \theta_i} = -\frac{1}{\hat{y}_k} \hat{y}_k (1 - \hat{y}_k) = \hat{y}_k - 1$$

when $i \neq k$

$$\frac{\partial CE}{\partial \theta_i} = -\frac{1}{\hat{y}_k} \times \hat{y}_k \hat{y}_i = -\hat{y}_i$$

$$\frac{\partial CE}{\partial \theta_i} = \hat{y}_k - 1 \quad \text{when } i = k$$

$$\text{or}$$

$$\hat{y}_k - \hat{y}_i$$

$$= \hat{y}_i$$

otherwise

Simplifying above we get

$$\frac{\partial CE}{\partial \theta_i} = \hat{y}_i - y_i$$

(c) ~~the~~ given $h = \text{sigmoid}(xW_1 + b_1)$
 $\hat{y} = \text{softmax}(hW_2 + b_2)$

let $z_1 = xW_1 + b_1$, so that
 $h = \text{sigmoid}(z_1)$

& $z_2 = hW_2 + b_2$
 so that $\hat{y} = \text{softmax}(z_2)$

$$\frac{\partial J}{\partial x} = \left(\left(\frac{\partial J}{\partial z_2} * \frac{\partial z_2}{\partial h} \right) \circ \frac{\partial h}{\partial z_1} \right) * \frac{\partial z_1}{\partial x}$$

- \circ element-wise multiplication
- $*$ dot-product

From 2(b) $\frac{\partial J}{\partial z_2} = (\hat{y} - y)$

$$\frac{\partial z_2}{\partial h} = \frac{\partial (hW_2 + b_2)}{\partial h} = W_2^T$$

$$\frac{\partial h}{\partial z_1} = \sigma'(z_1) = \sigma(z_1)(1 - \sigma(z_1))$$

$$\frac{\partial z_1}{\partial x} = W_1^T$$

$$\therefore \frac{\partial J}{\partial x} = ((\hat{y} - y)W_2^T) \circ \sigma' * W_1^T$$

(d) No. of parameters

$$x: D_x \times 1$$

$$y: D_y \times 1$$

~~W_1 must be $D_x \times H$ dimensional~~

~~W_2 must be $H \times D_y$ dimensional~~

Adding in the bias unit, hidden
No. of inputs to layer receives
inputs from $(D_x + 1)$ unit, hence
Dimensions of $W_1 \Rightarrow (D_x + 1) \times H$

Input Output layer receives inputs
from $(H + 1)$ units, hence
dimensions of $W_2 = (H + 1) \times D_y$

Total no. of parameters

$$= (D_x + 1) \times H + (H + 1) \times D_y$$

word2vec
3) (a) Given predicted word vector v_c corresponding to center word c ,

$$\hat{y} = P(o|c) = \frac{e^{u_o^T v_c}}{\sum_w e^{u_w^T v_c}}$$

u_w are "output" word vectors
& u_o is the expected word

$$\text{To find } \frac{\partial J}{\partial v_c} = \frac{\partial}{\partial v_c} \left(-\sum_y y \log \hat{y} \right)$$

assuming one-hot encoding for y

$$\frac{\partial J}{\partial v_c} = \frac{\partial}{\partial v_c} \left(-\log P(o|c) \right)$$

$$= \frac{\partial}{\partial v_c} \left(-\log e^{u_o^T v_c} + \log \left(\sum e^{u_w^T v_c} \right) \right)$$

$$= \frac{\partial}{\partial v_c} \left(u_o^T v_c \right) + \left(\frac{1}{\sum e^{u_w^T v_c}} \right) \frac{\partial}{\partial v_c} \left(\sum e^{u_w^T v_c} \right)$$

$$= u_o + \frac{1}{\sum e^{u_w^T v_c}} \left(\sum e^{u_z^T v_c} \cdot u_z \right)$$

$$= u_o + \sum_z \left(\frac{e^{u_z^T v_c}}{\sum e^{u_w^T v_c}} \right) \cdot u_z$$

$$= u_o + \sum_z P(z|c) \cdot u_z$$

(b) To find $\frac{\partial J}{\partial u_w}$

when $u_w = u_0$

$$\frac{\partial J}{\partial u_0} = -\frac{\partial}{\partial u_0} \left(\log \left(\frac{e^{u_0^T v_c}}{\sum e^{u_w^T v_c}} \right) \right)$$

$$= -\frac{\partial (u_0^T v_c)}{\partial u_0} + \frac{1}{\sum e^{u_w^T v_c}} \frac{\partial}{\partial u_0} \left(\sum e^{u_w^T v_c} \right)$$

$$= -v_c + \frac{1}{\sum e^{u_w^T v_c}} e^{u_0^T v_c} \cdot v_c$$

$$= -v_c + P(o|c) \cdot v_c$$

when $u_w \neq u_0$

$$\frac{\partial J}{\partial u_w} = -\frac{\partial}{\partial u_w} \left(\log \left(\frac{e^{u_0^T v_c}}{\sum e^{u_w^T v_c}} \right) \right)$$

$$= -\frac{\partial (u_0^T v_c)}{\partial u_w} + \frac{1}{\sum e^{u_w^T v_c}} \frac{\partial}{\partial u_w} \left(\sum e^{u_w^T v_c} \right)$$

$$= 0 + \frac{e^{u_w^T v_c}}{\sum e^{u_w^T v_c}} \cdot v_c$$

$$= P(w|c) \cdot v_c$$

hence

$$\frac{\partial J}{\partial u_w} = v_c (P(o|c) - 1)$$

$$= v_c P(w|c) \quad w \neq o$$

(d) Find $\frac{\partial J}{\partial V_c}$ & $\frac{\partial J}{\partial U_w}$ when negative sampling is applied

$$J = -\log(\sigma(U_0^T V_c)) - \sum_{k=1}^K \log(\sigma(-U_k^T V_c))$$

$$\frac{\partial J}{\partial V_c} = -\frac{1}{\sigma(U_0^T V_c)} \times \sigma \frac{\partial (\sigma(U_0^T V_c))}{\partial V_c} - \sum_{k=1}^K \frac{1}{\sigma(-U_k^T V_c)} \frac{\partial (\sigma(-U_k^T V_c))}{\partial V_c}$$

$$= -\frac{1}{\sigma_0} \times \sigma_0 (1 - \sigma_0) \frac{\partial (U_0^T V_c)}{\partial V_c} - \sum_{k=1}^K \frac{1}{\sigma_k} \times \sigma_k (1 - \sigma_k) \frac{\partial (-U_k^T V_c)}{\partial V_c}$$

$$= (\sigma_0 - 1) U_0 + \sum_{k=1}^K (1 - \sigma_k) U_k$$

$$= (\sigma(U_0^T V_c) - 1) U_0 + \sum_{k=1}^K (1 - \sigma(-U_k^T V_c)) U_k$$

For finding $\frac{\partial J}{\partial U_w}$, consider 3 cases :-

1) when $w = 0$ 2) $w = 1 \rightarrow k$ 3) $w \neq 0$ & $w \neq 1 \rightarrow k$

Case I) when $w = 0$

$$\frac{\partial J}{\partial U_0} = -\frac{1}{\sigma_0} \times \sigma_0 (1 - \sigma_0) V_c \rightarrow 0$$

$$= V_c (\sigma_0 - 1) = V_c (\sigma(U_0^T V_c) - 1)$$

Case II) when $w = 1 \rightarrow k$

$$\frac{\partial J}{\partial U_k} = -\frac{1}{\sigma_k} \times \sigma_k (1 - \sigma_k) \cdot 0 - \sum_{k=1}^K \frac{1}{\sigma_k} \times \sigma_k (1 - \sigma_k) \frac{\partial (-U_k^T V_c)}{\partial U_k}$$

$$= 0 + (1 - \sigma(-U_k^T V_c)) V_c$$

Case III) when $w \neq 0$ & $w \neq 1 \rightarrow k$

$$\frac{\partial J}{\partial U_w} = 0$$