

CS 224N Assignment #1.

1. Softmax

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

$$\text{softmax}(x_i + c) = \frac{e^{x_i + c}}{\sum_j e^{x_j + c}} = \frac{e^c \cdot e^{x_i}}{e^c \cdot \sum_j e^{x_j}} = \frac{e^{x_i}}{\sum_j e^{x_j}} = \text{softmax}(x_i).$$

(b). CODE.

2. Neural Network Basics.

$$(a) \quad \sigma(x) = \frac{1}{1 + e^{-x}} \Rightarrow \frac{d}{dx} \sigma(x) = \frac{-(-e^{-x})}{(1 + e^{-x})^2} = (1 - \sigma(x)) \sigma(x).$$

(b)

$$\hat{y} = \text{softmax}(\theta)$$

$$\frac{\partial \mathcal{L}(y, \hat{y})}{\partial \theta} = \frac{\partial \mathcal{L}(y, \hat{y})}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \theta}.$$

The first step is to calculate the derivative of softmax() function.

$$\begin{aligned} \hat{y}_i &= \frac{e^{\theta_i}}{\sum_k e^{\theta_k}} \\ \frac{\partial \hat{y}_i}{\partial \theta_i} &= \frac{\partial}{\partial \theta_i} \frac{1}{1 + A e^{-\theta_i}} \quad (\text{where } A = \sum_k e^{\theta_k} - e^{\theta_i}) \\ &= \frac{1}{(1 + A e^{-\theta_i})^2} \left(1 - \frac{1}{1 + A e^{-\theta_i}} \right) = \hat{y}_i (1 - \hat{y}_i). \\ \frac{\partial \hat{y}_i}{\partial \theta_j} &= \frac{\partial}{\partial \theta_j} \frac{e^{\theta_i}}{A + e^{\theta_i}} \quad (\text{where } A = \sum_k e^{\theta_k} - e^{\theta_j}) \\ &= \frac{e^{\theta_i}}{A + e^{\theta_i}} \cdot \frac{-e^{\theta_j}}{A + e^{\theta_j}} = -\hat{y}_i \hat{y}_j \end{aligned}$$

$$\begin{aligned} \frac{\partial \mathcal{L}(y, \hat{y})}{\partial \theta_i} &= - \sum_k y_k \cdot \frac{\partial (\log \hat{y}_k)}{\partial \theta_i} = - \sum_k y_k \cdot \frac{1}{\hat{y}_k} \frac{\partial \hat{y}_k}{\partial \theta_i} \\ &= -(y_i \cdot \hat{y}_i \cdot \frac{1}{\hat{y}_i} (1 - \hat{y}_i) + \sum_{k \neq i} y_k \cdot \frac{1}{\hat{y}_k} (-\hat{y}_i \hat{y}_k)) \Rightarrow \frac{\partial \mathcal{L}(y, \hat{y})}{\partial \theta} = \hat{y} - y \\ &= -y_i + \sum_k \hat{y}_i y_k = (\hat{y}_i - y_i) \star \end{aligned}$$

(c).

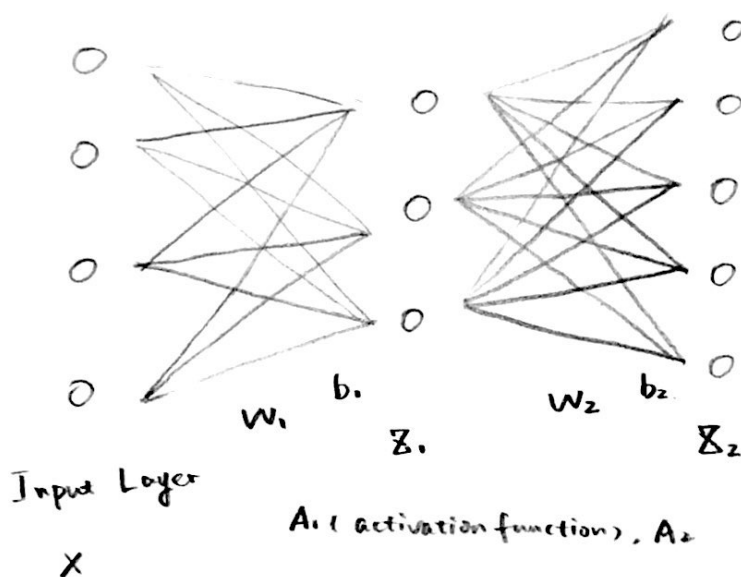
$$\frac{\partial J}{\partial x} = \frac{\partial CE(y, \hat{y})}{\partial (xw_1 + b_1)} \cdot \frac{\partial (xw_1 + b_1)}{\partial x} = (\hat{y} - y) w_1^T$$

(d). $w_1 \sim H \times D_x$ $b_1 \sim H$ $w_2 \sim D_y \times H$ $b_2 \sim D_y$

total = $H \times D_x + H + D_y \times H + D_y$

(e) ~ (g) Coding Part.

a few more details on the derivatives of Neural Networks



$$z_1 = w_1 x + b_1 \quad A_1 = A_1(z_1)$$

$$z_2 = w_2 A_1 + b_2 \quad A_2 = A_2(z_2) = \hat{Y} \text{ (predicted result).}$$

Derivatives !!

$$dz_2 = \hat{Y} - Y \quad \text{✗ (as proved in (b))}$$

$$dw_2 = dz_2 \cdot A_1^T \quad \text{or} \quad A_1^T \cdot dz_2 \quad \text{(depends on how you write it).}$$

$$db_2 = dz_2 \cdot \text{sum}() \leftarrow$$

$$dz_1 = w_2^T \cdot dz_2 / dz_2 w_2^T \cdot A_1'(z_1)$$

$$dw_1 = dz_1 \cdot x^T / x^T dz_1$$

$$db_1 = dz_1 \cdot \text{sum}() \leftarrow \text{careful for the axis (=0/1)}$$

3. word2vec

(a) skip gram - obtain center word from surroundings

first, derivative of the softmax form is:

$$\hat{y}_0 = \frac{e^{u_0^T v_c}}{\sum_w e^{u_w^T v_c}}$$

$$\frac{\partial \hat{y}_0}{\partial v_c} = \frac{u_0 e^{u_0^T v_c} (\sum_w e^{u_w^T v_c}) - e^{u_0^T v_c} (\sum_w u_w e^{u_w^T v_c})}{(\sum_w e^{u_w^T v_c})^2}$$

$$= \hat{y}_0 \frac{\sum_w (u_0 - u_w) \cdot e^{u_w^T v_c}}{\sum_w e^{u_w^T v_c}}$$

$$\frac{\partial CE(y, \hat{y})}{\partial v_c} = - \sum_k y_k \frac{\partial \log \hat{y}_k}{\partial v_c}$$

$$= - \sum_k y_k \cdot \frac{1}{\hat{y}_k} \cdot \frac{\partial \hat{y}_k}{\partial v_c}$$

$$= - \sum_k y_k \cdot \frac{1}{\hat{y}_k} \cdot \hat{y}_k \cdot \sum_w (u_0^T - u_w^T)$$

from previous section, we know that: (for classical softmax functions)

$$\frac{\partial CE(y, \hat{y})}{\partial \theta} = \hat{y} - y$$

set $\theta' = U^T \theta$, so for word2vec loss functions, it will be:

$$\frac{\partial CE(y, \hat{y})}{\partial \theta'} = \hat{y} - y$$

Thus,

$$\frac{\partial CE(y, \hat{y})}{\partial \theta} = \sum_i (\hat{y}_i - y_i) u_i^T = U(\hat{y} - y)$$

(b) first to calculate its derivative of softmax form function

$$\frac{\partial \hat{y}_0}{\partial u_0} = \frac{\partial}{\partial u_0} \frac{1}{1 + \text{constant} \cdot e^{-u_0^T v_c}} \quad (\text{where constant} = \sum_{w \neq 0} e^{u_w^T v_c}) = \hat{y}_0 (1 - \hat{y}_0) \cdot v_c$$

$$\frac{\partial \hat{y}_0}{\partial u_w (w \neq 0)} = \dots = -\hat{y}_0 \hat{y}_w \cdot v_c$$

$$\frac{\partial CE(y, \hat{y})}{\partial u_w} = - \sum_k y_k \frac{1}{\hat{y}_k} \frac{\partial \hat{y}_k}{\partial u_w} = \begin{cases} (\hat{y}_w - 1) \cdot v_c & w = 0 \\ \hat{y}_w \cdot v_c & w \neq 0 \end{cases}$$

$$(c) \quad \frac{\partial J}{\partial v_c} = - \frac{1}{\sigma(u_0^T v_c)} \sigma'(u_0^T v_c) u_0 - \sum_{k=1}^K \frac{1}{\sigma(-u_k^T v_c)} \sigma'(-u_k^T v_c) \cdot (-u_k)$$

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

$$\frac{\partial J}{\partial u_0} = - \frac{1}{\sigma(u_0^T v_c)} \sigma'(u_0^T v_c) v_c = [\sigma(u_0^T v_c) - 1] v_c$$

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

< don't know the reason why it execute more quickly...

(d)

From part (b) & (c), we already know how to derive the derivation of single cost function

(i) for skip-gram

$$\frac{\partial J_{\text{skip-gram}}(u_{t-m}, \dots, u_{t+m})}{\partial u} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(w_{t,j}, v_c)}{\partial u}$$

$$\frac{\partial J_{\text{skip-gram}}(u_{t-m}, \dots, u_{t+m})}{\partial v_c} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(w_{t,j}, v_c)}{\partial v_c}$$

$$\frac{\partial J_{\text{skip-gram}}(u_{t-m}, \dots, u_{t+m})}{\partial v_{w_{t,j}}} = 0 \quad (\text{for all } j \neq c)$$

for CBOW,

just switch the " $\sum_j \partial F(w_{t,j}, v_c)$ " into " $\partial F(w_c, \hat{v})$ ", and we're done