

# Assignment #1

## 1. Softmax

(a) for all dimensions  $i \in [1, \dim(n)]$

$$\begin{aligned}
 \text{softmax}(n+c)_i &= \frac{e^{(n_i+c)}_i}{\sum_j e^{(n_j+c)}_j} \\
 &= \frac{e^{n_i} \cdot e^c}{e^c \sum_j e^{n_j}} \\
 &= \frac{e^{n_i}}{\sum_j e^{n_j}} \\
 &= \text{softmax}(n)_i
 \end{aligned}$$

## 2. Neural Network Basics

(a)  $\sigma(n) = \frac{1}{1+e^{-n}}$

$\sigma'(n) = \frac{-(-e^{-n})}{(1+e^{-n})^2}$

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

$$= \sigma(x) (1 - \sigma(x))$$

(b)

$$CE(y, \hat{y}) = - \sum_i y_i \log(\hat{y}_i)$$

$$= - \sum_i y_i \log \left( \frac{e^{\theta_i}}{\sum_j e^{\theta_j}} \right)$$

Let  $k^{th}$  element of  $y=1$  & all others are zero

$$\begin{aligned} \frac{\partial CE(y, \hat{y})}{\partial \theta_l} &= -y_l \frac{\partial}{\partial \theta_l} \left( \log(e^{\theta_l}) - \log \left( \sum_i e^{\theta_i} \right) \right) \\ &\quad - \frac{\partial}{\partial \theta_l} \sum_{i \neq l} y_i (\log(e^{\theta_i}) - \log \left( \sum_j e^{\theta_j} \right)) \\ &= -y_l \left( 1 - \frac{1}{\sum_i e^{\theta_i}} e^{\theta_l} \right) + \sum_{i \neq l} \frac{y_i \cdot e^{\theta_l}}{\sum_j e^{\theta_j}} \\ &= -y_l (1 - \hat{y}_l) + \hat{y}_l \sum_{i \neq l} y_i \end{aligned}$$

$$\cancel{\nabla_{\theta_l} CE(y, \hat{y}) = }$$

~~$\hat{y}_l - y_l$~~  if  ~~$l=k$~~

if  $l = k$ ,

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

if  $l \neq k$ ,

$$\frac{\partial}{\partial \theta_l} CE(y, \hat{y}) = \hat{y}_l$$

So,

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

(c)  $J = - \sum_i y_i \log(\hat{y}_i)$

Let  $z_2 = h w_2 + b_2$

$z_1 = x w_1 + b_1$

then

$$h = \text{sigmoid}(z_1) = \sigma(z_1)$$

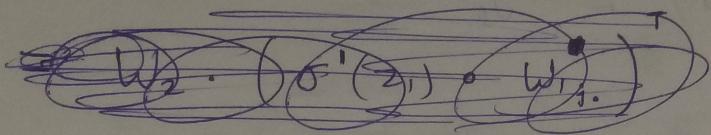
$$\hat{y} = \text{softmax}(z_2)$$

$$\frac{\partial J}{\partial x_j} = \frac{\partial J}{\partial z_2} \cdot \frac{\partial z_2}{\partial x_j}$$

$$= (\hat{y} - y) \cdot \frac{\partial z_2}{\partial x_j}$$

$$\delta_1 = \hat{y} - y$$

$$\begin{aligned}\frac{\partial z_2}{\partial x_j} &= \sum_k \frac{\partial z_2}{\partial h_k} \cdot \frac{\partial h_k}{\partial x_j} \\ &= \sum_k w_{2, k}^T \cdot \sigma'(z_{1k}) \cdot w_{1, jk}\end{aligned}$$



$$\frac{\partial J}{\partial x_j} = \sum_k \underbrace{(w_1^T \cdot w_2^T \cdot \sigma'(z_{1k}) \cdot w_{1, jk})}_{\delta_1}$$

$$= \sum_k \underbrace{\delta_1 w_{2, k}^T \sigma'(z_{1k}) w_{1, jk}}_{\delta_2}$$

$$\delta_2 = \delta_1 w_2^T$$

$$\frac{\partial J}{\partial x_j} = \sum_k \underbrace{\delta_2 \underbrace{\sigma'(z_{1k})}_{\delta_3} w_{1, jk}}_{\delta_3} = \sum_k \delta_3 w_{1, jk}$$

$$\frac{\partial J}{\partial x} = \underbrace{\delta_3 \wedge \dots \wedge \delta_3}_{\text{6x6}} \cdot \underbrace{w_1^T}_{\text{7x6}}$$

$$\frac{256}{6 \times 6} \cdot (B - \bar{B})$$

$$(d) \quad n \in \mathbb{R}^{D_n}$$

$$\hat{y} \in \mathbb{R}^{D_y}$$

H hidden units.

Parameters  $\Rightarrow w_1, w_2, b_1, b_2$

$$w_1 \in \mathbb{R}^{H \times D_n}$$

$$b_1 \in \mathbb{R}^H$$

$$w_2 \in \mathbb{R}^{H \times D_y}$$

$$b_2 \in \mathbb{R}^{D_y}$$

$$\# \text{ Parameters} = H(D_n + 1) + D_y(H + 1)$$

### 3. word2vec

$$(a) \quad J(\theta) = CE(y, \hat{y})$$

$$= - \sum_{\omega=1}^W y_\omega \log(\hat{y}_\omega)$$

$\therefore o$  is the expected word.

$$\therefore y_o = 1 \quad \& \quad y_i = 0 \quad \forall i \neq o$$

$$J(\theta) = - \log(\hat{y}_o)$$

$$= -\log \left( \frac{\exp(u_0^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)} \right)$$

$$= -u_0^T v_c + \log \sum_{w=1}^W \exp(u_w^T v_c)$$

$$\frac{\partial J}{\partial v_c} = -u_0 + \frac{\sum_{w=1}^W \exp(u_w^T v_c) u_w}{\sum_{i=1}^W \exp(u_i^T v_c)}$$

$$\boxed{\frac{\partial J}{\partial v_c} = -u_0 + \sum_{w=1}^W \hat{y}_w u_w}$$

Let

$$U = \begin{bmatrix} | & | & | \\ u_1 & u_2 & \cdots & u_W \\ | & | & | \end{bmatrix}$$

then

$$\frac{\partial J}{\partial v_c} = U \hat{y} - U y$$

$$\boxed{\frac{\partial J}{\partial v_c} = U(\hat{y} - y)}$$

(b) from (a)

$$J(\theta) = -u_0^T v_c + \log \sum_{w=1}^W \exp(u_w^T v_c)$$

$$\begin{aligned}\frac{\partial J}{\partial u_i} &= \mathbb{1}_{\{i=0\}} (-v_c) + \frac{\cancel{\exp(u_0^T v_c)} v_c}{\sum_{j=1}^W \exp(u_j^T v_c)} \\ &= \mathbb{1}_{\{i=0\}} (-v_c) + \hat{y}_i v_c\end{aligned}$$

$$\frac{\partial J}{\partial u_w} = \hat{y}_w v_c - y_w v_c$$

$$\boxed{\frac{\partial J}{\partial u_w} = (\hat{y}_w - y_w) v_c}$$

$$\frac{\partial J}{\partial v} = (\hat{y} - y) v^T$$

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

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

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

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

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

~~$\sigma(u_k^T v_c) (1 - \sigma(u_k^T v_c))$~~

~~$\log \left( \frac{1}{1 + \exp(-u_0^T v_c)} \right)$~~

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

~~$- \log (1 + \exp(-u_0^T v_c))$~~

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

$$\frac{\partial J}{\partial u_i} = \mathbb{1}_{\{i=0\}} \left[ - \frac{\sigma(u_0^T v_c) (1 - \sigma(u_0^T v_c))}{\sigma(u_0^T v_c)} v_c \right] - \sum_{k=1}^K \mathbb{1}_{\{k=i\}} \left[ \frac{\sigma(-u_k^T v_c) (1 - \sigma(-u_k^T v_c))}{\sigma(-u_k^T v_c)} (-v_k) \right]$$

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

$$\frac{\partial J}{\partial u_k} = -(\sigma(-u_k^T v_c) - 1) v_c \quad \forall k=1,2,\dots,K$$

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

$$\forall i \notin \{1, 2, \dots, k\} \cup \{0\}$$

The negative sampling loss function requires only  $K$  ( $\approx 100$ ) training examples to compute while softmax-CE loss function requires all training examples (words) which is of the order of millions. i.e.

$$\text{speed-up ratio} \approx \frac{w}{K}$$

(d) for skip-gram

$$J_{\text{skip-gram}}(\text{word}_{c-m\dots c+m}) = \sum_{\substack{-m \leq j \leq m, \\ j \neq 0}} f(w_{c+j}, v_c)$$

$$= \sum_{\substack{-m \leq j \leq m, \\ j \neq 0}} J_{\text{softmax-CE}}(o, v_c, u)$$

From previous parts, we know to compute the gradients  $\frac{\partial F(w_i, v)}{\partial v}$  and  $\frac{\partial F(w_i, v)}{\partial u}$ .

$$\frac{\partial J_{\text{skip-gram}}(\text{word}_{c-m \dots c+m})}{\partial u} = \sum_{\substack{-m \leq j \leq m, \\ j \neq 0}} \frac{\partial F(w_{cj}, u_c)}{\partial u}$$

$$\frac{\partial J_{\text{skip-gram}}(\text{word}_{c-m \dots c+m})}{\partial u_c} = \sum_{\substack{-m \leq j \leq m, \\ j \neq 0}} \frac{\partial F(w_{cj}, u_c)}{\partial u_c}$$

$$\frac{\partial J_{\text{skip-gram}}(\text{word}_{c-m \dots c+m})}{\partial u_j} = 0 \quad \forall j \neq c$$

For CBOW

$$\frac{\partial J_{\text{CBOW}}(\text{word}_{c-m \dots c+m})}{\partial u} = \frac{\partial F(w_c, \hat{v})}{\partial u}$$

$$\frac{\partial J_{\text{CBOW}}(\text{word}_{c-m \dots c+m})}{\partial u_j} = \frac{\partial F(w_c, \hat{v})}{\partial \hat{v}}, \quad \forall j \in \{c-m, \dots, c-1, c+1, \dots, c+m\}$$

$$\frac{\partial J_{\text{CBOW}}(\text{word}_{c-m \dots c+m})}{\partial u_j} = 0, \quad \forall j \notin \{c-m, \dots, c-1, c+1, \dots, c+m\}$$

## 4. Sentiment Analysis

- (b) To avoid overfitting to the training data and make the model generalize well on unseen examples.
- (c) To chose best regularization parameter, we tried ~~all~~ the regularization parameter values from  $0.0, 0.00001, 0.00003, 0.0001, 0.0003, 0.001, 0.003, 0.01$  and choose the one that maximizes the dev set accuracy.
- This strategy gives the best regularization parameter with value of  $1e-4$  with train, dev and test accuracies of 29.12, 31.43 and 27.56.