

David Freilich SkipGram Derivation

$$J_{\text{softmax-CE}}(O, V_c, V) = - \sum_{j=1}^V \log(\hat{y}_j) y_j, \quad \hat{y}_0 = p(O|c) = \frac{\exp(v_0^T v_c)}{\sum \exp(v_w^T v_c)}$$

$$\frac{\partial \text{soft}}{\partial v_c} = v_w^T (y - \hat{y})$$

$$\frac{\partial \text{soft}}{\partial v_w} = \frac{\partial J}{\partial v_w} = -y \cdot \frac{\partial \log \hat{y}}{\partial v_w} \Rightarrow \text{we only need to calculate for } y=1, \text{ at } y_0$$

$$= -y_0 \cdot \frac{\partial \log \hat{y}}{\partial v_w} \Rightarrow \text{we'll split federations now, for } v_0 \neq v_k$$

$$\boxed{y_0=1}$$

$$\frac{\partial J}{\partial v_0} = - \frac{\partial \log \hat{y}}{\partial v_0} = - \left(\frac{\partial v_0^T v_c}{\partial v_0} - \frac{\partial \sum \exp(v_w^T v_c)}{\partial v_0} \right) = - \left(\frac{\partial v_0^T v_c}{\partial v_0} - \frac{1}{\sum \exp(v_w^T v_c)} \cdot \frac{\partial \sum \exp(v_w^T v_c)}{\partial v_0} \right)$$

Expanding out for each k , we'd get:

$$= - \left(\frac{\exp(v_0^T v_c)}{\sum \exp(v_w^T v_c)} \cdot \frac{\partial v_0^T v_c}{\partial v_0} \right)$$

which would be 0 for all $k \neq 0$, so we simplify as:

$$= - \left(\frac{\partial v_0^T v_c}{\partial v_0} - \frac{\exp(v_0^T v_c)}{\sum \exp(v_w^T v_c)} \cdot \frac{\partial v_0^T v_c}{\partial v_0} \right) = -(v_c^T - p(O|c) \cdot v_c^T) = \boxed{(\hat{y} - 1) v_c^T}$$

$$\frac{\partial J}{\partial v_k} = -y_0 \cdot \frac{\partial \log \hat{y}}{\partial v_k} = - \left(\frac{\partial v_0^T v_c}{\partial v_k} - \frac{\partial \log \sum \exp(v_w^T v_c)}{\partial v_k} \right)$$

$$= - \left(0 - \frac{1}{\sum \exp(v_w^T v_c)} \cdot \sum \exp(v_w^T v_c) \cdot \frac{\partial v_w^T v_c}{\partial v_k} \right)$$

The only non zero one will be the k we are looking at

$$= p(k|c) \cdot v_c^T \quad \text{for each } k \text{ where } k \neq 0$$

$$= \boxed{\hat{y}_k v_c^T}$$

David Freilich

Negative Sampling

$$J_{\text{neg-sam}} = -\log(\sigma(u_o^T v_c)) - \sum_k \log(\sigma(-u_k^T v_c))$$

$$\frac{\partial J_{\text{neg}}}{\partial v_c} = -2 \log \sigma(u_o^T v_c) - \sum_k 2 \log \sigma(-u_k^T v_c)$$

$$= -\left(\frac{1}{\sigma(u_o^T v_c)} \cdot \frac{\sigma(u_o^T v_c)(1-\sigma(u_o^T v_c)) \cdot u_o^T}{1} \right) - \sum_k \left(\frac{1}{\sigma(-u_k^T v_c)} \cdot \frac{\sigma(-u_k^T v_c)(1-\sigma(-u_k^T v_c)) \cdot (-u_k^T)}{1} \right)$$

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

Assumptions:

$$\frac{\partial \log x}{\partial x} = \frac{1}{x}$$

$$\frac{\partial \sigma x}{\partial x} = \frac{\sigma x (1 - \sigma x)}{1}$$

$$\frac{\partial J_{\text{neg}}}{\partial v_k} = -2 \log \sigma(u_o^T v_c) - \sum_k \log \sigma(-u_k^T v_c)$$

$$= -\left(\frac{1}{\sigma(u_o^T v_c)} \cdot \frac{\sigma(u_o^T v_c)(1-\sigma(u_o^T v_c)) \cdot 0}{1} \right) - \sum_k \frac{1}{\sigma(-u_k^T v_c)} \cdot \frac{\sigma(-u_k^T v_c)(1-\sigma(-u_k^T v_c)) \cdot (-v_c)}{1}$$

$$\frac{\partial J}{\partial v_k} = 0 - \sum_k (\sigma(-u_k^T v_c) - 1) \cdot v_c$$

Since we're ~~der~~ ^{der} gradient of each k alone, the summation disappears, so

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