



(a) Show that the naive-softmax loss given in Equation (2) is the same as the cross-entropy loss between y and \hat{y} .

show that $-\sum_{w \in V_{\text{vocab}}} y_w \log(\hat{y}_w) = -\log(\hat{y}_0)$

↳ Cross entropy loss between the true probability distribution p and another distribution q is $-\sum_i p_i \log(q_i)$

With given center word c , y is true empirical distribution and \hat{y} is the predicted distribution. (one hot vector with a 1 for the true center word 0, and 0 everywhere)

$$-\sum_{w \in V} y_w \log(\hat{y}_w) = -\sum_{w \neq 0, w \in V} 0 * \log(\hat{y}_w) - 1 * \log(\hat{y}_0) = -\log(\hat{y}_0)$$

(b) Compute the partial derivative of $J_{\text{naive-softmax}}(V_c, 0, U)$ with respect to V_c . Please write your answer in terms of y , \hat{y} and U .

$$\frac{\partial J(V_c, 0, U)}{\partial V_c} = - \frac{\partial \log(P(0=0|C=c))}{\partial V_c}$$

$$= \frac{-\cancel{\partial \log(\exp(U_0^T V_c))}}{\partial V_c} + \frac{\partial \log(\sum_{w=1}^V \exp(U_w^T V_c))}{\partial V_c}$$

$$= -U_0 + \sum_{w=1}^V \frac{\exp(U_w^T V_c)}{\sum_{w=1}^V \exp(U_w^T V_c)} \cdot U_w$$

$$= -U_0 + \sum_{w=1}^V \underbrace{P(0=w|C=c)}_{\hat{y}_w} U_w$$

$$= U^T (\hat{y} - y)$$



(C) Compute the partial derivatives of $J_{\text{naive-softmax}}(V_c, o, U)$ with respect to each of the 'outside' word vectors, u 's. There will be two cases: when $w=0$, the true 'outside' word vector, and $w \neq 0$ for all other words. Please write your answer in terms of y , \hat{y} , and V_c .

$$\frac{\partial J_{\text{naive-softmax}}(V_c, o, U)}{\partial u_w} = - \frac{\partial \log(\exp(u_0^T V_c))}{\partial u_0} + \frac{\partial \log\left(\sum_{w=1}^V \exp(u_w^T V_c)\right)}{\partial u_w}$$

$$\text{Case 1: } w=0, \quad \frac{\partial \log(\exp(u_0^T V_c))}{\partial u_0} + \frac{\partial \log\left(\sum_{w=1}^V \exp(u_w^T V_c)\right)}{\partial u_0}$$

$$= -V_c + \frac{1}{\sum_{w=1}^V \exp(u_w^T V_c)} \cdot \frac{\partial \exp(u_0^T V_c)}{\partial u_0}$$

$$= -V_c + \frac{\exp(u_0^T V_c)}{\sum_{w=1}^V \exp(u_w^T V_c)} V_c = \hat{y}_0 V_c - V_c$$

$$= (p(o=0|c=c) - 1) V_c$$

Case 2: $w \neq 0$

$$0 + \frac{\partial}{\partial u_w} \log \sum_{w \in V} \exp(u_w^T V_c)$$

$$= \frac{\exp(u_w^T V_c)}{\sum_{w=1}^V \exp(u_w^T V_c)} \cdot V_c = \hat{y}_w V_c = p(o=w|c=c) V_c$$

$$\frac{\partial J_{\text{naive-softmax}}(V_c, o, U)}{\partial u} = (\hat{y} - y)^T V_c$$



(d) The sigmoid function is given by Equation 4:

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

Please compute the derivative of $\sigma(x)$ with respect to x , where x is a vector.

$$\left(\frac{1}{1+e^{-x}}\right)' = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{e^x \cdot 1}{(1+e^x)(1+e^{-x})} = \sigma(x)(1-\sigma(x))$$

(e) consider Negative Sampling loss, which is an alternative to the naive Softmax loss. Assume that K negative samples are drawn from the vocabulary. refers to them as w_1, w_2, \dots, w_K , outside vectors as u_1, \dots, u_K , $o \in \{w_1, \dots, w_K\}$

$$J_{\text{neg-sample}}(v_c, o, u) = -\log(\sigma(u_o^T v_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T v_c))$$

Repeat parts (b) and (c), computing partial derivative $J_{\text{neg-sample}}$ with respect to v_c , with respect to u_o , with respect to a negative sample u_k .

$$\frac{\partial J_{\text{neg-sample}}(v_c, o, u)}{\partial v_c} = \frac{\partial (-\log(\sigma(u_o^T v_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T v_c)))}{\partial v_c}$$

$$= -\frac{\sigma(u_o^T v_c)(1-\sigma(u_o^T v_c))}{\sigma(u_o^T v_c)} \cdot \frac{\partial u_o^T v_c}{\partial v_c} - \sum_{k=1}^K \frac{\partial \log(\sigma(-u_k^T v_c))}{\partial v_c}$$

$$= -(1-\sigma(u_o^T v_c))u_o + \left(-\sum_{k=1}^K \frac{1}{\sigma(-u_k^T v_c)} \cdot \sigma(-u_k^T v_c)(1-\sigma(-u_k^T v_c))\right) \cdot -u_k$$

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



$$\begin{aligned}
 \text{e. } \frac{\partial J_{\text{reg-sample}}(V_c, 0, U)}{\partial u_0} &= \frac{-\log(\sigma(u_0^T V_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T V_c))}{\partial u_0} \\
 &= \frac{\cancel{\sigma(u_0^T V_c)} (1 - \cancel{u_0^T V_c})}{\cancel{\sigma(u_0^T V_c)}} \cdot \frac{(u_0^T V_c)}{\partial u_0} = -(1 - \sigma(u_0^T V_c)) V_c
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial J_{\text{reg-sample}}(V_c, 0, U)}{\partial u_k} &= \frac{-\log(\sigma(u_0^T V_c)) - \cancel{\sum_{k=1}^K \log(\sigma(-u_k^T V_c))}}{\partial u_k} \\
 &= \frac{\sigma(-u_k^T V_c) (1 - \sigma(-u_k^T V_c))}{\sigma(-u_k^T V_c)} \cdot \frac{-u_k^T V_c}{\partial u_k} \\
 &= (1 - \sigma(-u_k^T V_c)) V_c
 \end{aligned}$$

f.

$$\text{i) } \frac{\partial}{\partial u} J_{\text{skip-gram}}(V_c, u_{t-m}, \dots, u_{t+m}, U) = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(V_c, u_{t+j}, U)}{\partial u}$$

$$\text{ii) } \frac{\partial}{\partial V_c} J_{\text{skip-gram}}(V_c, u_{t-m}, \dots, u_{t+m}, U) = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(V_c, u_{t+j}, U)}{\partial V_c}$$

$$\text{iii) } \frac{\partial}{\partial V_w} J_{\text{skip-gram}}(V_c, u_{t-m}, \dots, u_{t+m}, U) = 0$$