

Extra Credit Challenge I

Singleton

$$\begin{aligned} J_{\text{naive-softmax}}(v_c, o, u) &= -\log(\hat{y}_o) = -\log\left[\frac{\exp(u_o^T v_c)}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)}\right] \\ &= -\left[\log(\exp(u_o^T v_c)) - \log\left(\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)\right)\right] \\ &= -u_o^T v_c + \log\left[\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)\right] \end{aligned}$$

$$\Rightarrow \frac{\partial J}{\partial v_c} = \frac{\partial}{\partial v_c}(-u_o^T v_c) + \frac{\partial}{\partial v_c} \log\left[\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)\right]$$

Now $\frac{\partial}{\partial v_c}(-u_o^T v_c) = -u_o$, and by the Chain Rule;

$$\begin{aligned} \frac{\partial}{\partial v_c} \log\left[\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)\right] &= \left[\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)\right]^{-1} \frac{\partial}{\partial v_c} \left[\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)\right] \\ &= \left[\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)\right]^{-1} \sum_{w \in \text{Vocab}} \exp(u_w^T v_c) u_w \end{aligned}$$

$$\Rightarrow \frac{\partial J}{\partial v_c} = -u_o + \frac{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c) u_w}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)} = \left[\sum_{w \in \text{Vocab}} \underbrace{\frac{\exp(u_w^T v_c)}{\sum_{x \in \text{Vocab}} \exp(u_x^T v_c)}}_{\hat{y}_w} \cdot u_w \right] - u_o$$

$$= \left(\sum_{w \in \text{Vocab}} \hat{y}_w u_w \right) - u_o$$

$u_o = u_y$ since y is a one-hot vector that has 1. only for word o .

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

$$= u(\hat{y} - y)$$

Extra Credit Challenge II

Singletan

$$J_{\text{naive-softmax}}(V_c, o, U) = -u_o^T V_c + \log \left[\sum_{w' \in \text{Vocab}} \exp(u_{w'}^T V_c) \right]$$

Case 1: Suppose $W=0$.

$$\frac{\partial J}{\partial u_w} = \frac{\partial J}{\partial u_o} = \frac{\partial}{\partial u_o} (-u_o^T V_c) + \frac{\partial}{\partial u_o} \log \left[\sum_{w' \in \text{Vocab}} \exp(u_{w'}^T V_c) \right]$$

$$= -V_c + \left[\sum_{x \in \text{Vocab}} \exp(u_x^T V_c) \right]^{-1} \frac{\partial}{\partial u_o} \left[\sum_{w' \in \text{Vocab}} \exp(u_{w'}^T V_c) \right]$$

But $\frac{\partial}{\partial u_o} \left[\sum_{w' \in \text{Vocab}} \exp(u_{w'}^T V_c) \right] = \exp(u_o^T V_c) V_c$ since u_o

appears in only one term of the sum. The terms where u_o does not appear vanish when $\frac{\partial}{\partial u_o}$ is differentiated w.r.t. u_o . Therefore,

$$\frac{\partial J}{\partial u_o} = -V_c + \frac{\exp(u_o^T V_c)}{\underbrace{\sum_{x \in \text{Vocab}} \exp(u_x^T V_c)}_{\hat{y}_o}} V_c$$

$$= \hat{y}_o V_c - V_c$$

$$= (\hat{y}_o - 1) V_c.$$

Case 2: Suppose $w \neq 0$.

$$\begin{aligned}\frac{\partial J}{\partial u_w} &= \frac{2}{2u_w} [-u_0^T V_c] + \frac{2}{2u_w} \log \left[\sum_{w' \in \text{vocab}} \exp(u_{w'}^T V_c) \right] \\ &= 0 + \left[\sum_{x \in \text{vocab}} \exp(u_x^T V_c) \right]^{-1} \frac{2}{2u_w} \left[\sum_{w' \in \text{vocab}} \exp(u_{w'}^T V_c) \right]\end{aligned}$$

But under the sum, $\frac{2}{2u_w} \exp(u_{w'}^T V_c) = 0$ for all values of w' except $w' = w$, in which case it equals $\exp(u_w^T V_c) V_c$. So we have

$$\begin{aligned}\frac{\partial J}{\partial u_w} &= \frac{\exp(u_w^T V_c) V_c}{\sum_{x \in \text{vocab}} \exp(u_x^T V_c)} = \frac{\exp(u_w^T V_c)}{\underbrace{\sum_{x \in \text{vocab}} \exp(u_x^T V_c)}_{\hat{y}_w}} V_c \\ &= \hat{y}_w V_c.\end{aligned}$$

Thus we have
$$\frac{\partial J}{\partial u_w} = \begin{cases} (\hat{y}_w - 1) V_c & \text{if } w=0 \\ \hat{y}_w V_c & \text{otherwise} \end{cases}$$

or equivalently
$$\frac{\partial J}{\partial u} = V_c (\hat{y} - y)^T$$