

The solution is in Dr. Chris Manning hint/derivation:

$$J_{\text{softmax}}(V_c, 0, U) = -u_0^T V_c + \log \left(\sum_{w \in V_{\text{vocab}}} \exp(u_w^T V_c) \right)$$

then: $\frac{\partial J}{\partial V_c} = \frac{\partial}{\partial V_c} \left[-u_0^T V_c + \log \left(\sum_{w \in V_{\text{vocab}}} \exp(u_w^T V_c) \right) \right]$

$$= -u_0 + \frac{\partial}{\partial V_c} \log \left(\sum_{w \in V_{\text{vocab}}} \exp(u_w^T V_c) \right)$$

chain-rule

$$= -u_0 + \sum \frac{1}{\sum \exp(u_w^T V_c)} \left(\exp(u_w^T V_c) \cdot u_w \right)$$

$$= -u_0 + \sum \frac{\exp(u_w^T V_c)}{\sum \exp(u_w^T V_c)} \cdot u_w$$

$$= -u_0 + \sum \hat{y}_w \cdot u_w \quad (\text{4) from Note}$$

$$= U(\hat{y} - y) \quad (\text{5) from Note}$$

~~100~~

3 Extra Credit Challenge II (2.5 Points)

The partial derivatives of $J_{\text{naive-softmax}}(\mathbf{v}_c, \mathbf{o}, \mathbf{U})$ with respect to each of the 'outside' word vectors, \mathbf{u}_w 's is given below:

$$\frac{\partial J}{\partial \mathbf{U}} = \mathbf{v}_c(\hat{\mathbf{y}} - \mathbf{y})^\top \quad (6)$$

or equivalently:

$$\frac{\partial J}{\partial \mathbf{u}_w} = \begin{cases} (\hat{y}_w - 1)\mathbf{v}_c & \text{if } w = o \\ \hat{y}_w \mathbf{v}_c & \text{otherwise} \end{cases} \quad (7)$$

Write the steps required to arrive at the partial derivative of $J_{\text{naive-softmax}}(\mathbf{v}_c, \mathbf{o}, \mathbf{U})$ with respect to each of the 'outside' word vectors, \mathbf{u}_w 's. There are two cases you need to consider: when $w = o$, the true 'outside' word vector, and $w \neq o$, for all other words. Please write your answer in terms of \mathbf{y} , $\hat{\mathbf{y}}$, and \mathbf{v}_c . The proof may take 4 or 5 steps. The loss function $J_{\text{naive-softmax}}(\mathbf{v}_c, \mathbf{o}, \mathbf{U})$ is:

$$J_{\text{naive-softmax}}(\mathbf{v}_c, \mathbf{o}, \mathbf{U}) = -\mathbf{u}_o^\top \mathbf{v}_c + \log \left(\sum_{w' \in \text{Vocab}} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_c) \right)$$

similar as previous derivation + Manning's Note.

$$\frac{\partial J}{\partial \mathbf{u}_w} = \frac{\partial}{\partial \mathbf{u}_w} \left[-\mathbf{u}_o^\top \mathbf{v}_c + \log \left(\sum_{w' \in \text{Vocab}} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_c) \right) \right]$$

$$= \frac{\partial}{\partial \mathbf{u}_w} (-\mathbf{u}_o^\top \mathbf{v}_c) + \frac{\partial}{\partial \mathbf{u}_w} \log \left(\sum_{w' \in \text{Vocab}} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_c) \right)$$

when $w = o$

$$\frac{-\mathbf{u}_o^\top \mathbf{v}_c}{\partial \mathbf{u}_o} = -\mathbf{v}_c$$

when $w \neq o$.

$$\frac{-\mathbf{u}_w^\top \mathbf{v}_c}{\partial \mathbf{u}_w} = 0.$$

(Chain Rule)

$$= \frac{1}{\sum \exp(\mathbf{u}_{w'}^\top \mathbf{v}_c)} \frac{\partial}{\partial \mathbf{u}_w} \left(\sum \exp(\mathbf{u}_{w'}^\top \mathbf{v}_c) \right)$$

$$= \frac{\exp(\mathbf{u}_w^\top \mathbf{v}_c) \cdot \mathbf{v}_c}{\sum \exp(\mathbf{u}_{w'}^\top \mathbf{v}_c)}$$

$$= \hat{y}_w \cdot \mathbf{v}_c$$

We got (7):
$$\begin{cases} (\hat{y}_w - 1)\mathbf{v}_c & \text{if } w = o \\ \hat{y}_w \mathbf{v}_c & \text{if } w \neq o. \end{cases}$$