

ECE C247 HW #2

2) Softmax Classifier Gradient

Let $w_i^T x^{(j)} + b_i = 0_i$

$$\text{Likelihood} = \prod_{j=1}^m \prod_{n=1}^c \text{Pr}(y^{(j)} = n | x^{(j)}, \theta) \delta_n(y^{(j)})$$

where $\delta_n(y^{(j)}) = 1$ if $y^{(j)} = n$
0 otherwise.

We are given that $\text{Pr}(y^{(j)} = n | x^{(j)}, \theta) = \text{softmax}_n(x^{(j)}) = \frac{e^{\theta_n}}{\sum_{k=0}^c e^{\theta_k}}$

$$\text{So, likelihood} = \prod_{j=1}^m \prod_{n=1}^c \text{softmax}_n(x^{(j)}) \delta_n(y^{(j)})$$

$$\mathcal{L} = \text{Log-likelihood} = \sum_{j=1}^m \sum_{n=1}^c \delta_n(y^{(j)}) \log(\text{softmax}_n(x^{(j)}))$$

$$\frac{d\mathcal{L}}{dw_i} = \frac{d\mathcal{L}}{d\text{softmax}_n} \cdot \frac{d\text{softmax}_n}{d\theta_i} \cdot \frac{d\theta_i}{dw_i}$$

$$\frac{dL}{do_i} = \sum_{j=1}^m \delta_i(y^{(j)}) \cdot \frac{1}{\text{Softmax}_i(x^{(j)})} \cdot \frac{d \text{Softmax}_i(x^{(j)})}{do_i} + \sum_{n \neq i} \delta_n(y^{(j)}) \cdot \frac{1}{\text{Softmax}_n(x^{(j)})} \cdot \frac{d \text{Softmax}_n(x^{(j)})}{do_i}$$

$$\begin{aligned} \frac{d \text{Softmax}_i(x^{(j)})}{do_i} &= \frac{d e^{o_i}}{\sum_{k=1}^c e^{o_k}} = \frac{e^{o_i} \cdot \sum_{k=1}^c e^{o_k} - e^{o_i} \cdot e^{o_i}}{\left(\sum_{k=1}^c e^{o_k} \right)^2} \\ &= \frac{e^{o_i}}{\sum_k e^{o_k}} - \left(\frac{e^{o_i}}{\sum_k e^{o_k}} \right)^2 \end{aligned}$$

$$= \text{Softmax}_i(x^{(j)}) [1 - \text{Softmax}_i(x^{(j)})]$$

$$\frac{d \text{Softmax}_n(x^{(j)})}{do_i} \text{ when } n \neq i = \frac{d e^{o_n}}{\sum_k e^{o_k}}$$

$$= \frac{0 \cdot \sum_k e^{o_k} - e^{o_n} \cdot e^{o_i}}{\left(\sum_k e^{o_k} \right)^2}$$

$$= -\text{Softmax}_n(x^{(j)}) \cdot \text{Softmax}_i(x^{(j)})$$

Plugging in to $\frac{dL}{do_i}$ gives:

$$\frac{dL}{do_i} = \sum_{j=1}^m f_i(y^{(j)}) \cdot \frac{\text{softmax}_i(x^{(j)}) \cdot [1 - \text{softmax}_i(x^{(j)})]}{\text{softmax}_i(x^{(j)})} + \sum_{n \neq i} f_n(y^{(j)}) \cdot \frac{-\text{softmax}_n(x^{(j)}) \cdot \text{softmax}_i(x^{(j)})}{\text{softmax}_n(x^{(j)})}$$

$$= \sum_{j=1}^m f_i(y^{(j)}) \cdot [1 - \text{softmax}_i(x^{(j)})] + \sum_{n \neq i} f_n(y^{(j)}) \cdot -\text{softmax}_i(x^{(j)})$$

$$= \sum_{j=1}^m f_i(y^{(j)}) - \text{softmax}_i(x^{(j)}) \cdot \sum_{n=1}^n f_n(y^{(j)})$$

We know that $\sum_{n=1}^n f_n(y^{(j)}) = 1$ since $f_n(y^{(j)})$ is 1 when $y^{(j)} = n$ or 0 otherwise.

$$\text{So, } \frac{dL}{do_i} = \sum_{j=1}^m f_i(y^{(j)}) - \text{softmax}_i(x^{(j)})$$

$$\frac{do_i}{dw_i} = \frac{dw_i^T x^{(j)} + b_i}{dw_i} = x^{(j)}$$

$$\frac{do_i}{db_i} = \frac{dw_i^T x^{(j)} + b_i}{db_i} = 1$$

$$\text{So, } \frac{dL}{dw_i} = \frac{dL}{do_i} \cdot \frac{do_i}{dw_i} = \sum_{j=1}^m [f_i(y^{(j)}) - \text{softmax}_i(x^{(j)})] \cdot x^{(j)}$$

$$\frac{dL}{db_i} = \frac{dL}{do_i} \cdot \frac{do_i}{db_i} = \sum_{j=1}^m f_i(y^{(j)}) - \text{softmax}_i(x^{(j)})$$