

$$2) \quad P_{\theta}(y^{(j)} = i | x^{(j)}, \theta) = \text{softmax}_i(x^{(j)})$$

$$\text{softmax}_i(x) = \frac{e^{\tilde{w}_i^T x}}{\sum_{k=1}^C e^{\tilde{w}_k^T x}}$$

Assume: $\tilde{x} = \begin{bmatrix} x \\ 1 \end{bmatrix} \quad \tilde{w} = \begin{bmatrix} w \\ b \end{bmatrix}$

LIKELIHOOD:

$$L_H = P(\tilde{x}^{(1)}, \dots, \tilde{x}^{(n)}, y^{(1)} \dots y^{(n)} | \theta) = \prod_{j=1}^n P(x^{(j)}, y^{(j)} | \theta)$$

$$= \prod_{j=1}^n P(\tilde{x}^{(j)} | \theta) P(y^{(j)} | \tilde{x}^{(j)}, \theta)$$

$$= \prod_{j=1}^n P(y^{(j)} | \tilde{x}^{(j)}, \theta)$$

$$= \prod_{j=1}^n \text{softmax}_{y^{(j)}}(\tilde{x}^{(j)})$$

$$= \prod_{j=1}^n \left[\frac{e^{a_{y^{(j)}}(\tilde{x}^{(j)})}}{\sum_{k=1}^C e^{a_k(\tilde{x}^{(j)})}} \right]$$

where $a_k(\tilde{x}^{(j)}) = \tilde{w}_k^T \tilde{x}^{(j)}$

$$\log L_H = L = \sum_{j=1}^n \left[(a_{y^{(j)}}(\tilde{x}^{(j)})) - \log \sum_{k=1}^C e^{a_k(\tilde{x}^{(j)})} \right]$$

WE WANT TO MAXIMIZE THE LOG LIKELIHOOD:

$$\frac{\partial L}{\partial \theta} = 0$$

$$\nabla_{\tilde{w}_i} L = \frac{\partial L}{\partial \tilde{w}_i} = \frac{\partial L}{\partial a_i} \times \frac{\partial a_i}{\partial \tilde{w}_i}$$

$$= \sum_{j=1}^n \frac{-1}{\sum_{k=1}^c e^{a_k(\tilde{x}(j))}} \times e^{a_i(\tilde{x}(j))} \times \tilde{x}(j)$$

$$+ \sum_{\substack{j=1 \\ y(j)=i}}^n \tilde{x}(j) //$$

$$\text{Take } I_{ji} = \begin{cases} 1, & y(j)=i \\ 0, & y(j) \neq i \end{cases}$$

$$\nabla_{\tilde{w}_i} L = \sum_{j=1}^n \left[I_{ji} - \frac{e^{a_i(\tilde{x}(j))}}{\sum_{k=1}^c e^{a_k(\tilde{x}(j))}} \right] \cdot \tilde{x}(j)$$

$$\nabla_{w_i} L = \sum_{j=1}^m \left[I_{ji} - \frac{e^{a_k(x^{(j)})}}{\sum_{k=1}^c e^{a_k(x^{(j)})}} \right] \cdot x^{(j)}$$

$$\nabla_{b_i} L = \sum_{j=1}^m \left[I_{ji} - \frac{e^{a_k(x^{(j)})}}{\sum_{k=1}^c e^{a_k(x^{(j)})}} \right]$$

WE CAN ALSO DERIVE FOR A SINGLE SAMPLE:

$$L_j(\theta) = a_{y^{(j)}}(\tilde{x}^{(j)}) - \log \sum_{k=1}^c e^{a_k(\tilde{x}^{(j)})}$$

$$= \log \left(\frac{e^{a_{y^{(j)}}(\tilde{x}^{(j)})}}{\sum_{k=1}^c e^{a_k(\tilde{x}^{(j)})}} \right)$$

$$a_k(\tilde{x}^{(j)}) = \tilde{w}_k^T \tilde{x}^{(j)}$$

$$\text{LET } \sigma_{y^{(j)}}(\tilde{x}^{(j)}) = \frac{e^{a_{y^{(j)}}(\tilde{x}^{(j)})}}{\sum_{k=1}^c e^{a_k(\tilde{x}^{(j)})}}$$

$$\frac{\partial \sigma_{y(i)}(\tilde{x}^{(i)})}{\partial a_k(\tilde{x}^{(i)})}$$

→ if $y(i) = k$

$$\begin{aligned} \frac{\partial \sigma_k(\tilde{x}^{(i)})}{\partial a_k(\tilde{x}^{(i)})} &= \frac{\partial}{\partial a_k(\tilde{x}^{(i)})} \left(\frac{e^{a_k(\tilde{x}^{(i)})}}{\sum_{k=1}^C e^{a_k(\tilde{x}^{(i)})}} \right) \\ &= \frac{\sum_{k=1}^C e^{a_k(\tilde{x}^{(i)})} \cdot e^{a_k(\tilde{x}^{(i)})} - e^{a_k(\tilde{x}^{(i)})} \cdot e^{a_k(\tilde{x}^{(i)})}}{\left[\sum_{k=1}^C e^{a_k(\tilde{x}^{(i)})} \right]^2} \\ &= \sigma_{y(i)}(\tilde{x}^{(i)}) [1 - \sigma_{y(i)}(\tilde{x}^{(i)})] \end{aligned}$$

→ if $y(i) \neq k$

$$\frac{\partial \sigma_{y(i)}(\tilde{x}^{(i)})}{\partial a_k(\tilde{x}^{(i)})} = \frac{\partial}{\partial a_k(\tilde{x}^{(i)})} \left(\frac{e^{a_{y(i)}(\tilde{x}^{(i)})}}{\sum_{k=1}^C e^{a_k(\tilde{x}^{(i)})}} \right)$$

$$= \frac{\sum_{k=1}^C e^{a_k(\tilde{x}^{(j)})} \cdot 0 - e^{a_{y^{(j)}}(\tilde{x}^{(j)})} \cdot e^{a_k(\tilde{x}^{(j)})}}{\left[\sum_{k=1}^C e^{a_k(\tilde{x}^{(j)})} \right]^2}$$

$$= -\sigma_{y^{(j)}}(\tilde{x}^{(j)}) \sigma_k(\tilde{x}^{(j)})$$

$$\frac{\partial L_j(\theta)}{\partial a_k(\tilde{x}^{(j)})} = \frac{\partial L_j(\theta)}{\partial \sigma_{y^{(j)}}(\tilde{x}^{(j)})} \cdot \frac{\partial \sigma_{y^{(j)}}(\tilde{x}^{(j)})}{\partial a_k(\tilde{x}^{(j)})}$$

$$L_j(\theta) = \log \left(\frac{e^{a_{y^{(j)}}(\tilde{x}^{(j)})}}{\sum_{k=1}^C e^{a_k(\tilde{x}^{(j)})}} \right)$$

$$= \log(\sigma_{y^{(j)}}(\tilde{x}^{(j)}))$$

$$\frac{\partial L_j(\theta)}{\partial a_k(\tilde{x}^{(j)})} = \begin{cases} (1 - \sigma_{y^{(j)}}(\tilde{x}^{(j)})) & , y^{(j)} = k \\ -\sigma_k(\tilde{x}^{(j)}) & , y^{(j)} \neq k \end{cases}$$

$$\frac{\partial L_j(\theta)}{\partial \tilde{w}_k} = \begin{cases} (1 - \sigma_{y^{(j)}}(\tilde{x}^{(j)}))\tilde{x}^{(j)}, & y^{(j)} = k \\ -\sigma_k(\tilde{x}^{(j)})\tilde{x}^{(j)}, & y^{(j)} \neq k \end{cases} \quad \hookrightarrow \text{here } \tilde{x}^{(j)} = \begin{bmatrix} x^{(j)} \\ 1 \end{bmatrix}$$

$$\frac{\partial L_j(\theta)}{\partial w_k} = \begin{cases} (1 - \sigma_{y^{(j)}}(x^{(j)}))x^{(j)}, & y^{(j)} = k \\ -\sigma_k(x^{(j)})x^{(j)}, & y^{(j)} \neq k \end{cases} \quad \hookrightarrow \text{here } \tilde{x}^{(j)} = \begin{bmatrix} x^{(j)} \\ 1 \end{bmatrix}$$

$$\frac{\partial L_j(\theta)}{\partial b_k} = \begin{cases} (1 - \sigma_{y^{(j)}}(x^{(j)})), & y^{(j)} = k \\ -\sigma_k(x^{(j)}), & y^{(j)} \neq k \end{cases} \quad \hookrightarrow \text{here } \tilde{x}^{(j)} = \begin{bmatrix} x^{(j)} \\ 1 \end{bmatrix}$$