

② Softmax Classifier Gradient

$$(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}); \quad x^{(j)} \in \mathbb{R}^n, \quad y^{(j)} \in \{1, \dots, c\}, \quad j=1, \dots, n$$

$$\theta = \{w_i, b_i\}_{i=1, \dots, c}$$

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

$$\text{softmax}_i(x) = \frac{e^{w_i^T x + b_i}}{\sum_{k=1}^c e^{w_k^T x + b_k}}$$

$$\tilde{x} = \begin{bmatrix} x \\ 1 \end{bmatrix} \Rightarrow a_i(x) = \tilde{w}_i^T \tilde{x}$$
$$\tilde{w}_i = \begin{bmatrix} w_i \\ b_i \end{bmatrix}$$

$$p(x^{(1)}, \dots, x^{(n)}, y^{(1)}, \dots, y^{(n)} | \theta) = \prod_{i=1}^n p(x^{(i)}, y^{(i)} | \theta)$$
$$= \prod_{i=1}^n p(y^{(i)} | x^{(i)}, \theta) p(x^{(i)} | \theta)$$

$$\arg \max_{\theta} \prod_{i=1}^n p(y^{(i)} | x^{(i)}, \theta) p(x^{(i)} | \theta) = \arg \max_{\theta} \prod_{i=1}^n p(y^{(i)} | x^{(i)}, \theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^n \log \text{softmax}_{y^{(i)}}(x^{(i)}) = \arg \max_{\theta} \sum_{i=1}^n \log \left[\frac{e^{a_{y^{(i)}}(x^{(i)})}}{\sum_{j=1}^c e^{a_j(x^{(i)})}} \right]$$

$$= \arg \max_{\theta} \frac{1}{n} \sum_{i=1}^n \left[a_{y^{(i)}}(x^{(i)}) - \log \sum_{j=1}^c e^{a_j(x^{(i)})} \right]$$

$$= \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \left[\log \sum_{j=1}^c e^{a_j(x^{(i)})} - a_{y^{(i)}}(x^{(i)}) \right]$$

$$\mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \left[\log \sum_{j=1}^c e^{a_j(x^{(i)})} - a_{y^{(i)}}(x^{(i)}) \right]$$

$$\nabla_{\tilde{w}_i} \left(\log \sum_{j=1}^c e^{a_j(x)} \right) = \nabla_{\tilde{w}_i} \left(\log \left[e^{\tilde{w}_1^T \tilde{x}} + e^{\tilde{w}_2^T \tilde{x}} + \dots + e^{\tilde{w}_i^T \tilde{x}} + \dots + e^{\tilde{w}_c^T \tilde{x}} \right] \right)$$

$$\frac{\partial}{\partial \tilde{w}_{i,1}} \left(\log \left[e^{\tilde{w}_1^T \tilde{x}} + \dots + e^{\tilde{w}_{i,1}x_1 + \tilde{w}_{i,2}x_2 + \dots + \tilde{w}_{i,n}x_n + \tilde{w}_{i,0}} + \dots + e^{\tilde{w}_c^T \tilde{x}} \right] \right)$$

$$= \frac{1}{\sum_j e^{a_j(x)}} \frac{\partial}{\partial \tilde{w}_{i,1}} \left[\underbrace{e^{\tilde{w}_1^T \tilde{x}} + \dots + e^{\tilde{w}_{i,1}x_1 + \tilde{w}_{i,2}x_2 + \dots + \tilde{w}_{i,n}x_n + \tilde{w}_{i,0}}}_{0} + \underbrace{e^{\tilde{w}_c^T \tilde{x}}}_{0} \right]$$

$$= \frac{1}{\sum_j e^{a_j(x)}} x_1 e^{\tilde{w}_i^T \tilde{x}} = \frac{e^{\tilde{w}_i^T \tilde{x}}}{\sum_j e^{a_j(x)}} x_1$$

$$\frac{\partial}{\partial \tilde{w}_{i,2}} \left(\log \sum_{j=1}^c e^{a_j(x)} \right) = \frac{e^{\tilde{w}_i^T \tilde{x}}}{\sum_j e^{a_j(x)}} x_2$$

$$\left\{ \nabla_{\tilde{w}_i} \left(\log \sum_{j=1}^c e^{a_j(x)} \right) = \frac{e^{a_i(x)}}{\sum_{j=1}^c e^{a_j(x)}} \tilde{x} \right.$$

$$\nabla_{\tilde{w}_i} a_{y(k)}(x) = \begin{cases} 0 & \text{if } y(k) \neq i \\ \tilde{x} & \text{if } y(k) = i \end{cases}$$

$$\hookrightarrow y(k) \neq i \Rightarrow 0$$

$$y(k) = i \Rightarrow \nabla_{\tilde{w}_i} \tilde{w}_i^T \tilde{x} = \tilde{x}$$

$$\Rightarrow \nabla_{\tilde{w}_i} \mathcal{L} = \frac{1}{n} \sum_{k=1}^n \left[\frac{e^{a_i(x^{(k)})}}{\sum_j e^{a_j(x^{(k)})}} \tilde{x}^{(k)} - \delta_{y^{(k)}, i} \tilde{x}^{(k)} \right]$$

$$\left[\nabla_{w_i} \mathcal{L} = \frac{1}{n} \sum_{k=1}^n \left[\frac{e^{a_i(x^{(k)})}}{\sum_j e^{a_j(x^{(k)})}} - \delta_{y^{(k)}, i} \right] x^{(k)} \right]$$

$$\left[\nabla_{b_i} \mathcal{L} = \frac{1}{n} \sum_{k=1}^n \left[\frac{e^{a_i(x^{(k)})}}{\sum_j e^{a_j(x^{(k)})}} - \delta_{y^{(k)}, i} \right] \right]$$