

# Artificial Intelligence and Machine Learning

## Derivation of the Gradient for Softmax Regression

We compute the partial derivative of the cross-entropy cost function:

$$\begin{aligned}\frac{\partial}{\partial \theta_{ij}} \ell^{\text{CE}}(\mathbf{h}_{\Theta}(\mathbf{x}), \mathbf{y}) &= \frac{\partial}{\partial \theta_{ij}} \left( - \sum_{k=1}^K y_k \log(\zeta_k(\mathbf{z})) \right) \\ &= - \sum_{k=1}^K y_k \cdot \frac{\partial}{\partial \theta_{ij}} \log(\zeta_k(\mathbf{z}))\end{aligned}$$

[Apply chain rule]

$$= - \sum_{k=1}^K y_k \cdot \frac{\partial \log(\zeta_k(\mathbf{z}))}{\partial \zeta_k(\mathbf{z})} \cdot \frac{\partial \zeta_k(\mathbf{z})}{\partial z_j} \cdot \frac{\partial z_j}{\partial \theta_{ij}}$$

[Derivative of log (first factor produced by chain rule)]

$$= - \sum_{k=1}^K y_k \cdot \frac{1}{\zeta_k(\mathbf{z})} \cdot \frac{\partial \zeta_k(\mathbf{z})}{\partial z_j} \cdot \frac{\partial z_j}{\partial \theta_{ij}}$$

[Separate cases  $k = j$  and  $k \neq j$ ]

$$\begin{aligned}&= \overbrace{-y_j \cdot \frac{1}{\zeta_j(\mathbf{z})} \cdot \frac{\partial \zeta_j(\mathbf{z})}{\partial z_j} \cdot \frac{\partial z_j}{\partial \theta_{ij}}}^{k=j} - \sum_{\substack{k=1 \\ k \neq j}}^K \overbrace{y_k \cdot \frac{1}{\zeta_k(\mathbf{z})} \cdot \frac{\partial \zeta_k(\mathbf{z})}{\partial z_j} \cdot \frac{\partial z_j}{\partial \theta_{ij}}}^{k \neq j} \\ &= \left( -y_j \cdot \frac{1}{\zeta_j(\mathbf{z})} \cdot \frac{\partial \zeta_j(\mathbf{z})}{\partial z_j} - \sum_{\substack{k=1 \\ k \neq j}}^K y_k \cdot \frac{1}{\zeta_k(\mathbf{z})} \cdot \frac{\partial \zeta_k(\mathbf{z})}{\partial z_j} \right) \cdot \frac{\partial z_j}{\partial \theta_{ij}}\end{aligned}$$

[Derivative of the softmax function (see exercise sheet)]

$$= \left( -y_j \cdot \frac{1}{\zeta_j(\mathbf{z})} \cdot \zeta_j(\mathbf{z}) \cdot (1 - \zeta_j(\mathbf{z})) + \sum_{\substack{k=1 \\ k \neq j}}^K y_k \cdot \frac{1}{\zeta_k(\mathbf{z})} \cdot \zeta_k(\mathbf{z}) \cdot \zeta_j(\mathbf{z}) \right) \cdot \frac{\partial z_j}{\partial \theta_{ij}}$$

[Cancel terms]

$$= \left( -y_j + y_j \cdot \zeta_j(\mathbf{z}) + \sum_{\substack{k=1 \\ k \neq j}}^K y_k \cdot \zeta_j(\mathbf{z}) \right) \cdot \frac{\partial z_j}{\partial \theta_{ij}}$$

[Put the two cases  $k = j$  and  $k \neq j$  back together]

$$= \left( -y_j + \sum_{k=1}^K y_k \cdot \zeta_j(\mathbf{z}) \right) \cdot x_i$$

$[\zeta_j(\mathbf{z})$  does not depend on index  $k$ . Therefore, we can pull it out of the sum]

$$= \left( -y_j + \zeta_j(\mathbf{z}) \cdot \sum_{k=1}^K y_k \right) \cdot x_i$$

$[\mathbf{y}$  is a one-hot vector, therefore the sum of its components is equal to 1]

$$\begin{aligned} &= (-y_j + \zeta_j(\mathbf{z})) \cdot x_i \\ &= \boxed{(\zeta_j(\mathbf{z}) - y_j) \cdot x_i} \end{aligned}$$

□