

SDS 555 HW7 sh2432 Gradient

1a) $h_1 = \text{ReLU}(w_1 X + b_1)$, $h_2 = \text{ReLU}(w_2 h_1 + b_2)$, $p = \text{Softmax}(w_3 h_2 + b_3)$
 $p(Y=y_i | X) = p(Y=k | X) = \frac{\exp(f(y_i))}{\sum_{j=1}^K \exp(f(y_j))}$ for $k=1, 2, \dots, K$

$$L = \frac{1}{n} \sum_{i=1}^n -\log p(Y=y_i | X_i) = \frac{1}{n} \sum_{i=1}^n -\log \frac{\exp(f(y_i))}{\sum_{j=1}^K \exp(f(y_j))} = \frac{1}{n} \sum_{i=1}^n -f(y_i) + \log \left[\sum_{j=1}^K \exp(f(y_j)) \right]$$

therefore, $\frac{\partial L}{\partial f(y_i)} = \frac{\exp(f(y_i))}{\sum_{j=1}^K \exp(f(y_j))} - \frac{1}{n} \mathbb{1}(Y=y_i) = p(Y=y_i | X_i) - \frac{1}{n} \mathbb{1}(Y=y_i)$

Apply the chain rule:

$$\bullet \frac{\partial L}{\partial f} = \begin{pmatrix} \frac{\partial L}{\partial f_1} \\ \vdots \\ \frac{\partial L}{\partial f_K} \end{pmatrix} \quad K \times 1$$

$$\bullet \frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial f} \cdot h_2^T \in \mathbb{R}^{K \times H_2}$$

$$\bullet \frac{\partial L}{\partial b_3} = \frac{\partial L}{\partial f} \in \mathbb{R}^K$$

$$\bullet \frac{\partial L}{\partial w_2} = \mathbb{1}(h_2 > 0) \cdot \frac{\partial L}{\partial h_2} \cdot h_3^T = \mathbb{1}(h_2 > 0) \cdot w_3^T \cdot \frac{\partial L}{\partial f} \cdot h_1^T$$

$$\bullet \frac{\partial L}{\partial b_2} = \mathbb{1}(h_2 > 0) \frac{\partial L}{\partial h_2} = \mathbb{1}(h_2 > 0) \cdot w_3^T \cdot \frac{\partial L}{\partial f}$$

$$\bullet \frac{\partial L}{\partial w_1} = \mathbb{1}(h_1 > 0) \frac{\partial L}{\partial h_3} \cdot X^T = \mathbb{1}(h_1 > 0) \mathbb{1}(h_2 > 0) w_3^T \cdot \frac{\partial L}{\partial h_2} \cdot X^T$$

$$\bullet = \mathbb{1}(h_1 > 0) \mathbb{1}(h_2 > 0) w_3^T w_3^T \frac{\partial L}{\partial f} X^T$$

$$\bullet \frac{\partial L}{\partial b_1} = \mathbb{1}(h_1 > 0) \frac{\partial L}{\partial h_1} = \mathbb{1}(h_1 > 0) \mathbb{1}(h_2 > 0) w_3^T \cdot w_3^T \cdot \frac{\partial L}{\partial f}$$

$$2 a) \quad h = \text{ReLU}(W_1 x + b_1), \quad \hat{x} = \text{ReLU}(W_2 h + b_2)$$

$$L = \frac{1}{n} \sum_{i=1}^n \|x_i - \text{ReLU}(W_2 \text{ReLU}(W_1 x_i + b_1) + b_2)\|^2 = \frac{1}{n} \sum_{i=1}^n \|x_i - \hat{x}\|^2$$

$$\frac{\partial L}{\partial \hat{x}} = \frac{1}{n} \sum_{i=1}^n 2(\hat{x}_i - x_i)$$

Apply the chain Rule:

$$\bullet \quad \frac{\partial L}{\partial W_2} = \mathbb{1}(\hat{x} > 0) \cdot \frac{\partial L}{\partial \hat{x}} \cdot h^T \quad \bullet \quad \frac{\partial L}{\partial b_2} = \mathbb{1}(\hat{x} > 0) \cdot \frac{\partial L}{\partial \hat{x}}$$

$$\bullet \quad \frac{\partial L}{\partial W_1} = \mathbb{1}(h > 0) \frac{\partial L}{\partial h} \cdot X^T = \mathbb{1}(h > 0) \mathbb{1}(\hat{x} > 0) W_2^T \frac{\partial L}{\partial \hat{x}} \cdot X^T$$

$$\bullet \quad \frac{\partial L}{\partial b_1} = \mathbb{1}(h > 0) \frac{\partial L}{\partial h} = \mathbb{1}(h > 0) \mathbb{1}(\hat{x} > 0) W_2^T \frac{\partial L}{\partial \hat{x}}$$