

$$1. \quad x \in \mathbb{R}^d \quad f(x) = \|x\|_2^2 = \left(\left(\sum_{i=1}^d x_i^2 \right)^{\frac{1}{2}} \right)^2 = \sum_{i=1}^d x_i^2$$

$$\nabla f(x) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{pmatrix} = \begin{pmatrix} 2x_1 \\ \vdots \\ 2x_d \end{pmatrix} = 2\vec{x}$$

$$2. \quad f(x) = A^T x \in \mathbb{R}^n, \quad A \in \mathbb{R}^{d \times m} \quad x \in \mathbb{R}^d$$

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_n}{\partial x_1} \\ \vdots & \vdots & & \vdots \\ \frac{\partial f_1}{\partial x_d} & \frac{\partial f_1}{\partial x_d} & \dots & \frac{\partial f_n}{\partial x_d} \end{bmatrix} = A \quad f(x) = [f_1(x), \dots, f_n(x)]^T$$

$$3. \quad g(x) = \|y\|_2^2 = \sum_{i=1}^k y_i^2$$

$$\nabla_x f(y(x)) = \underbrace{\frac{\partial g(u)}{\partial x}}_{= 2} \cdot \frac{\partial f}{\partial g(u)} = 2 \cdot \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_1} & \dots & \frac{\partial f}{\partial x_1} \\ \vdots & \vdots & & \vdots \\ \frac{\partial f}{\partial x_d} & \frac{\partial f}{\partial x_d} & \dots & \frac{\partial f}{\partial x_d} \end{bmatrix} \cdot A^T x = 2 A A^T x$$

$$\frac{\partial g(u)}{\partial u} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_n}{\partial x_1} \\ \vdots & \vdots & & \vdots \\ \frac{\partial f_1}{\partial x_d} & \frac{\partial f_1}{\partial x_d} & \dots & \frac{\partial f_n}{\partial x_d} \end{bmatrix}$$

$$\frac{\partial f}{\partial g(u)} = 2g(u)$$

$$4. \quad g(A) = A^T X \in \mathbb{R}^n \quad f(y) = \|y\|_2^2 \quad \frac{\partial f}{\partial y} = 2y \Rightarrow \boxed{\frac{\partial f}{\partial y} = 2A^T X}$$

$$\nabla_A f(g(A)) = \frac{\partial g(A)}{\partial A} \cdot \frac{\partial f}{\partial g} =$$

↓

vector by matrix derivative

$$f(g(A)) = g(A)^T g(A) =$$

$$= (A^T X)^T A^T X = X^T A A^T X = y \text{ scalar}$$

derivative of scalar by matrix:

$$\nabla_A g(A) = X e_1^T$$

combining the result

$$\nabla_A X^T A A^T X = \begin{bmatrix} \frac{\partial y}{\partial A_{11}} & \dots & \frac{\partial y}{\partial A_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial A_{11}} & \dots & \frac{\partial y}{\partial A_{nn}} \end{bmatrix} =$$

more explicitly I use the trace for

the scalar: Since $f(g(A))$ is a scalar only y then:

$$y = X^T A A^T X$$

$$y = \text{trace}(y) = \text{trace}(X^T A A^T X) = \text{trace}(A^T X X^T A)$$

$$X = X^T - \text{Symmetric}$$

$$\frac{\partial}{\partial A} \text{tr}(A^T X X^T A) = \frac{\partial}{\partial A} \text{tr}(A^T X A) =$$

$\underbrace{\hspace{1cm}}$
and $\underbrace{\hspace{1cm}}$

$n \times n$

$$= 2 X A = \boxed{2 X X^T A}$$

B.

$$J = -y \log g - (1-y) \log (1-g)$$

1. $\hat{y} = G\left(x + \underbrace{w_1 x + b_1}_{b} + b_2\right)$

$b = w_2 \cdot \max(0, d)$ $\text{Relu}(w_2 \cdot x + b_2)$

$d = w_1 \cdot x + b_1$

$\hat{y} = \frac{1}{1 + e^{-(w_2 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2) + x}}$

2.

$$\frac{\partial J}{\partial w_2} = \frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2} = \left[\frac{-y}{\hat{y}} + \frac{(1-y)}{1-\hat{y}} \right] \cdot \frac{e^{-(w_2 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2) + x}}{(1 + e^{-(w_2 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2) + x})^2} \cdot \text{Relu}'(w_1 \cdot x + b_1)$$

(1): $\frac{-y}{\hat{y}} + \frac{(1-y)}{1-\hat{y}}$

(2): $\frac{w_2}{\hat{y}} \cdot \frac{1}{1 + e^{-(w_2 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2) + x}} = G'(w_2 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2 + x) \cdot \text{Relu}'(w_1 \cdot x + b_1) =$

$G'(x) = G(x) \cdot (1 - G(x))$

$$= \frac{1}{1 + e^{-(w_2 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2) + x}} \cdot \frac{e^{-(w_2 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2) + x}}{1 + e^{-(w_2 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2) + x}} \cdot \text{Relu}'(w_1 \cdot x + b_1) =$$

$$= \frac{e^{-(w_2 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2) + x}}{(1 + e^{-(w_2 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2) + x})^2} \cdot \text{Relu}'(w_1 \cdot x + b_1)$$

$$\frac{\partial J}{\partial b_2} = \frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b_2} = \left[\frac{-y}{\hat{y}} + \frac{(1-y)}{1-\hat{y}} \right] \cdot \frac{e^{(w_1 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2) \cdot x}}{(1 + e^{(w_1 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2) \cdot x})^2}$$

$$\textcircled{1} : \frac{-y}{\hat{y}} + \frac{(1-y)}{1-\hat{y}}$$

$$\textcircled{2} \quad g'(w_1 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2 \cdot x) = \frac{e^{(w_1 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2 \cdot x)}}{(1 + e^{(w_1 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2 \cdot x)})^2}$$

$$3. \quad \frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_1} = \left[\frac{-y}{\hat{y}} + \frac{(1-y)}{1-\hat{y}} \right] \cdot \frac{e^{(w_1 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2 \cdot x)}}{(1 + e^{(w_1 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2 \cdot x)})^2} - w_2 \cdot x, \quad w_1 \cdot x + b_1 > 0 \\ \textcircled{1}, \quad \textcircled{2}, \quad \text{otherwise}$$

$$\textcircled{1} : \frac{-y}{\hat{y}} + \frac{(1-y)}{1-\hat{y}}$$

$$\textcircled{2} : \frac{\partial}{\partial w_1} \frac{1}{1 + e^{(w_1 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2 \cdot x)}} = \frac{e^{(w_1 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2 \cdot x)}}{(1 + e^{(w_1 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2 \cdot x)})^2} \cdot \frac{\partial}{\partial w_1} [w_2 \cdot \text{Relu}(w_1 \cdot x + b_1)]$$

$$\frac{\partial}{\partial w_1} [w_2 \cdot \text{Relu}(w_1 \cdot x + b_1)] = w_2 \cdot x \cdot \underbrace{\text{Relu}'[w_1 \cdot x + b_1]}$$

$$\begin{cases} 0, & w_1 \cdot x + b_1 \leq 0 \\ 1, & w_1 \cdot x + b_1 > 0 \end{cases}$$

$$\textcircled{2} : \frac{e^{(w_1 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2 \cdot x)}}{(1 + e^{(w_1 \cdot \text{Relu}(w_1 \cdot x + b_1) + b_2 \cdot x)})^2} - w_2 \cdot x, \quad w_1 \cdot x + b_1 > 0 \\ 0, \quad \text{otherwise}$$

$$3.b. \frac{\partial J}{\partial b_1} = \frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b_1} = \left[\frac{-y}{\hat{y}} + \frac{(1-y)}{1-\hat{y}} \right] \cdot \frac{e^{(w_1 \text{ReLu}(w_1 x + b_1) + b_2 - x)}}{\left(1 + e^{(w_1 \text{ReLu}(w_1 x + b_1) + b_2 - x)} \right)^2} \cdot w_2, \quad w_1 \neq b_1 \geq 0$$

0 other wise

$$(1) \quad \frac{\partial J}{\partial b_1} = \frac{\partial}{\partial b_1} \frac{1}{1 + e^{(w_1 \text{ReLu}(w_1 x + b_1) + b_2 - x)}} = \frac{e^{(w_1 \text{ReLu}(w_1 x + b_1) + b_2 - x)}}{\left(1 + e^{(w_1 \text{ReLu}(w_1 x + b_1) + b_2 - x)} \right)^2} \cdot \frac{\partial}{\partial b_1} [w_2 \cdot \text{ReLu}(w_1 x + b_1)]$$

$$\frac{\partial}{\partial b_1} [w_2 \cdot \text{ReLu}(w_1 x + b_1)] = w_2 \cdot 1 \cdot \underbrace{\text{ReLu}'(w_1 x + b_1)}_{\begin{cases} 0, & w_1 x + b_1 \leq 0 \\ 1, & w_1 x + b_1 > 0 \end{cases}}$$

3.c. $\frac{\partial J}{\partial x} =$ going to repeat same steps as before with
the difference of the derivative of ReLu

$$\frac{\partial}{\partial x} [w_2 \cdot \text{ReLu}(w_1 x + b_1) + b_2 - x] = w_2 \cdot w_1 \cdot \underbrace{\text{ReLu}'(w_1 x + b_1)}_{\begin{cases} 0, & w_1 x + b_1 \leq 0 \\ 1, & w_1 x + b_1 > 0 \end{cases}} + 1$$

$$\frac{\partial J}{\partial x} = \left[\frac{-y}{\hat{y}} + \frac{(1-y)}{1-\hat{y}} \right] \cdot \frac{e^{(w_1 \text{ReLu}(w_1 x + b_1) + b_2 - x)}}{\left(1 + e^{(w_1 \text{ReLu}(w_1 x + b_1) + b_2 - x)} \right)^2} \cdot (w_2 \cdot w_1 + 1), \quad w_1 \neq b_1 \geq 0$$

$$\left[\frac{-y}{\hat{y}} + \frac{(1-y)}{1-\hat{y}} \right] \cdot \frac{e^{(w_1 \text{ReLu}(w_1 x + b_1) + b_2 - x)}}{\left(1 + e^{(w_1 \text{ReLu}(w_1 x + b_1) + b_2 - x)} \right)^2}, \quad \text{other wise}$$

4. We need to cache: x , $z_1 = w_1x + b_1$, $z_2 = w_2 \cdot \text{ReLU}(z_1) + b_2 + x$

$$G(x) \text{ and } G'(x) = G(x) \cdot (1 - G(x)), \quad \frac{\partial J}{\partial g} = \frac{-y}{g} + \frac{1-y}{1-g},$$

$G'(w_2 \text{ReLU}(w_1x+b_1)+b_2+x) = \frac{e^{(w_1x+b_1)(w_1x+b_1+y)}}{(1+e^{(w_1x+b_1)(w_1x+b_1+y)})^2}$

$$\text{Rely}'(x) = \begin{cases} 1, & w_1x + b_1 > 0 \\ 0, & \text{otherwise} \end{cases}$$

$$\bullet \frac{\partial g}{\partial k}$$

(2): where g is a function of k

Variable by which we want to differentiate