

## Overview

- Backpropagation

### Differential circuits / networks

composition of arithmetic operations and elementary functions



(+, -, ×, /)

(cos, sin, exp, log, relu, sigmoid)

### Theorem (informally stated)

Suppose a differentiable circuit of size  $N$ ,

computes a real-valued function  $f: \mathbb{R}^l \rightarrow \mathbb{R}$

Then the gradient of  $\nabla f$  can also be computed in time  $O(N)$

by a circuit of size  $O(N)$

$$\cdot f \rightarrow J^{(1)}(\theta) \quad l = \# \text{ of parameters}$$

$$N \approx O(\# \text{ of parameters})$$

$$\nabla J^{(1)}(\theta) \quad \text{time } \approx O(N)$$

### Preliminary: Chain rule

Suppose  $J$  is a function of  $\theta_1, \dots, \theta_p$

$$\text{intermediate } g_j = g_j(\theta_1, \dots, \theta_p) \quad j=1, \dots, k$$

$$\text{output var } J = J(g_1, \dots, g_k)$$

$$\frac{\partial J}{\partial \theta_i} = \sum_{j=1}^k \frac{\boxed{\frac{\partial J}{\partial g_j}}}{\boxed{\frac{\partial g_j}{\partial \theta_i}}} \cdot \frac{\boxed{\frac{\partial g_j}{\partial \theta_i}}}{\boxed{\frac{\partial g_j}{\partial \theta_i}}}$$

### One-neuron neuron network example

$$\Sigma = w^T x + b$$

$$h_\theta(x) \rightarrow 0 = \text{relu}(\Sigma)$$

input:  $x$ , label:  $y$

$w \in \mathbb{R}^d$ ,  $x \in \mathbb{R}^d$ ,  $b \in \mathbb{R}$ ,  $\Sigma \in \mathbb{R}$

$$\text{Loss: } J(\theta) = \frac{1}{2} (y - \hat{y})^2$$

$\hat{y}$  as intermediate variable

$$\begin{aligned}\frac{\partial J}{\partial w_i} &= \boxed{\frac{\partial J}{\partial \hat{y}}} \frac{\partial \hat{y}}{\partial w_i} = (0 - \hat{y}) \frac{\partial \hat{y}}{\partial w_i} \\ &= (0 - \hat{y}) \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_i} \\ &= (0 - \hat{y}) \text{relu}'(z) \cdot x_i \\ &= (0 - \hat{y}) \mathbb{1}(z \geq 0) \cdot x_i\end{aligned}$$

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial b} = \frac{\partial J}{\partial \hat{y}} \cdot \boxed{\frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial b}}$$

Two layer neural network

$$w_j^{[1]} \in \mathbb{R}^d \quad z_j = w_j^{[1] T} x + b_j^{[1]}$$

$$b_j^{[1]} \in \mathbb{R} \quad a_j = \text{relu}(z_j)$$

$$a = [a_1, \dots, a_m]^T \in \mathbb{R}^m$$

$$0 = w^{[2] T} a + b^{[2]}$$

$$J = \frac{1}{2} (y - 0)^2$$

$$w_j^{[1]} = ((w_j^{[1]})_1, \dots, (w_j^{[1]})_d)^T \in \mathbb{R}^d$$

$$\frac{\partial J}{\partial (w_j^{[1]})_k} = \sum_{l=1}^m \frac{\partial J}{\partial z_k} \frac{\partial z_k}{\partial (w_j^{[1]})_k}$$

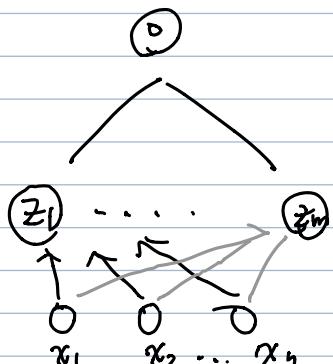
$z_k$  and  $w_j$  has No dependency  
if  $k \neq j$

$$= \frac{\partial J}{\partial z_j} \boxed{\frac{\partial z_j}{\partial (w_j^{[1]})_k}}$$

$$= \boxed{\frac{\partial J}{\partial z_j}} x_l$$

$$= \frac{\partial J}{\partial a_j} \frac{\partial a_j}{\partial z_j} \cdot x_l$$

$$= \frac{\partial J}{\partial a_j} \text{relu}'(z_j) \cdot x_l$$



$$= \frac{\partial J}{\partial o} \cdot \frac{\partial o}{\partial a_j} \text{relu}'(z_j) x_l$$

$$= (o - y) \underbrace{(w^{[2]})_j}_{\text{relu}'(z_j)} \underbrace{\text{relu}'(z_j) \cdot x_l}$$

Computing  $\frac{\partial J}{\partial (w_j^{[2]})_l}$  takes  $O(\# \text{parameters})$  time

full gradient  $\Rightarrow O(\# \text{parameters}^2)$  time.

$\forall j, l$   $(o - y)$  computation is shared

$\forall l$ , fixed  $j$   $\text{relu}(z_j)$  shared  $(\frac{\partial J}{\partial z_1}, \dots, \frac{\partial J}{\partial z_m})$

For backprop two layer neural network with vectorized notations

$$\delta^{[2]} \triangleq \frac{\partial J}{\partial o} = (o - y) \in \mathbb{R}$$

# Note  $A \odot B = \begin{bmatrix} A \cdot B_1 \\ \vdots \\ A \cdot B_n \end{bmatrix}$

$$\delta^{[1]} \triangleq \begin{bmatrix} \frac{\partial J}{\partial z_1} \\ \vdots \\ \frac{\partial J}{\partial z_m} \end{bmatrix} \in \mathbb{R}^m \quad \frac{\partial J}{\partial z} = \begin{bmatrix} \frac{\partial J}{\partial o} \cdot \frac{\partial o}{\partial a_j} \cdot \frac{\partial a_j}{\partial z_j} \\ \vdots \end{bmatrix}$$

$$= (o - y) (w^{[2]})^\top \odot \text{relu}'(z)$$

① Compute  $\frac{\partial J}{\partial o} = (o - y)$

② Compute  $\begin{bmatrix} \frac{\partial J}{\partial z_1} \\ \vdots \\ \frac{\partial J}{\partial z_m} \end{bmatrix} \quad \frac{\partial J}{\partial z_j} = \frac{\partial J}{\partial o} \cdot \frac{\partial o}{\partial a_j} \cdot \frac{\partial a_j}{\partial z_j} = (o - y) w_j^{[2]} \text{relu}'(z_j)$

③  $\frac{\partial J}{\partial (w_j^{[2]})_l} = \frac{\partial J}{\partial z_j} \cdot \frac{\partial z_j}{\partial (w_j^{[2]})_l} = \frac{\partial J}{\partial z_j} \cdot x_l$

$$\frac{\partial J}{\partial w_j^{[2]}} = \begin{bmatrix} \frac{\partial J}{\partial z_j} \cdot \frac{\partial z_j}{\partial (w_j^{[2]})_1} \\ \vdots \\ \frac{\partial J}{\partial z_j} \cdot \frac{\partial z_j}{\partial (w_j^{[2]})_m} \end{bmatrix} = \frac{\partial J}{\partial z} \cdot x^\top$$

$$\frac{\partial J}{\partial b_j^{[2]}} = \frac{\partial J}{\partial z_j} \cdot \frac{\partial z_j}{\partial b_j^{[2]}} \stackrel{(w^T x + b)}{=} 1$$

$$\frac{\partial J}{\partial (b_j^{[2]})_j} = \frac{\partial J}{\partial (w_j^{[2]})_j}$$

## Some notations

$$\frac{\partial A}{\partial B}$$

$A$  is a real value variable

$B$  is a vector / matrix

$$B \in \mathbb{R}^m \quad \frac{\partial A}{\partial B} = \begin{bmatrix} \frac{\partial A}{\partial B_1} \\ \vdots \\ \frac{\partial A}{\partial B_m} \end{bmatrix} \in \mathbb{R}^m$$

$$B \in \mathbb{R}^{m \times n} \quad \frac{\partial A}{\partial B} = \begin{bmatrix} \frac{\partial A}{\partial B_{ij}} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

## Vectorized Gradients

$$g^{[2]} = o - y \left( \frac{\partial J}{\partial o} \right)$$

$$g^{[1]} = (o - y) (w^{[2]})^T \circ \text{relu}(z) \in \mathbb{R}^m \left( \frac{\partial J}{\partial z} \right)$$

$$\frac{\partial J}{\partial w^{[1]}} = \frac{\partial J}{\partial z} \cdot \underbrace{x^T}_{\mathbb{R}^m \times \mathbb{R}^{1 \times d}} = g^{[1]} \in \mathbb{R}^{m \times d}$$

$$\frac{\partial J}{\partial b^{[1]}} = \frac{\partial J}{\partial z} \cdot \frac{\partial z}{\partial b^{[1]}} = g^{[1]} \in \mathbb{R}$$

$$\frac{\partial J}{\partial w^{[0]}} = \frac{\partial J}{\partial o} \cdot \frac{\partial o}{\partial w^{[0]}} = (o - y) \cdot a = g^{[2]} \cdot a^T \in \mathbb{R}^{1 \times m}$$

$$\frac{\partial J}{\partial b^{[0]}} = \frac{\partial J}{\partial o} \cdot \frac{\partial o}{\partial b^{[0]}} = g^{[1]}$$

## Chain rule for matrix multiplication - 1

$$z = w u + b$$

$$J = f(z)$$

$$u \in \mathbb{R}^d$$

$$w \in \mathbb{R}^{m \times d}, b \in \mathbb{R}^d$$

$$\frac{\partial J}{\partial u} = w^T \cdot \frac{\partial J}{\partial z}$$

$$\frac{\partial J}{\partial z} \in \mathbb{R}^m$$

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial z} \cdot u^T \quad w^T \in \mathbb{R}^{d \times m}, \quad \frac{\partial J}{\partial u} \in \mathbb{R}^d$$

two layer NN case

$$\frac{\partial J}{\partial w^{(2)}} = \frac{\partial J}{\partial z} \cdot x^T$$

$$\begin{aligned} W^{(2)} &\leftrightarrow w \\ J &\leftrightarrow J \\ z &\leftrightarrow z \\ x &\leftrightarrow u \\ b^{(2)} &\leftrightarrow b \end{aligned}$$

$$\frac{\partial J}{\partial b^{(2)}} = \frac{\partial J}{\partial z} \cdot \frac{\partial z}{\partial b^{(2)}} = \frac{\partial J}{\partial z}$$

Chain rule for matrix multiplication - 2

$$a = \sigma(z) \quad a, z \in \mathbb{R}^m$$

$$J = J(a) \quad \sigma \text{ is element-wise application}$$

$$\frac{\partial J}{\partial z} = \frac{\partial J}{\partial a} \odot \sigma'(z)$$

$$\frac{\partial J}{\partial z} = \frac{\partial J}{\partial a} \odot \sigma'(z)$$

$$= \frac{\partial J}{\partial a} \cdot \frac{\partial a}{\partial z} \odot \sigma'(z)$$

$$= \frac{\partial J}{\partial a} \cdot (w^{(2)})^T \odot \sigma'(z) \quad (\int^{(2)} \text{from above})$$

Multiple layer example

$$a^{(1)} = \text{relu}(w^{(1)}x + b^{(1)})$$

$$a^{(2)} = \text{relu}(w^{(2)}a^{(1)} + b^{(2)})$$

$$\vdots \\ a^{(r)} \Rightarrow h(x) = w^{(r)}a^{(r-1)} + b^{(r)}$$

$$J = \frac{1}{2}(a^{(r)} - y)^2$$

$$\text{#Note } \delta^{[k]} = \frac{\partial J}{\partial z^{[k]}} \in \mathbb{R}^{m_k} \quad \text{if } a^{[k]} \in \mathbb{R}^{m_k}$$

$$\delta^{[0]} = \frac{\partial J}{\partial z^{[0]}} = (a^{[0]} - y) = (z^{[0]} - y)$$

$$\begin{aligned} \delta^{[r-1]} &= \frac{\partial J}{\partial z^{[r-1]}} = \left( w^{[r-1]T} \delta^{[r]} \right) \circ \text{relu}'(z^{[r-1]}) \\ &\vdots \\ \frac{\partial J}{\partial z^{[r-1]}} &= \frac{\partial J}{\partial z^{[0]}} \cdot \frac{\partial z^{[r-1]}}{\partial a^{[r-1]}} \cdot \frac{\partial a^{[r-1]}}{\partial z^{[r-1]}} \end{aligned}$$

$$\delta^{[k]} = \frac{\partial J}{\partial z^{[k]}} = \left( w^{[k+1]T} \delta^{[k+1]} \right) \circ \text{relu}'(z^{[k]})$$

$$\Rightarrow \frac{\partial J}{\partial w^{[r+1]}} = \delta^{[r+1]} \cdot z^{[r]T} \quad \boxed{z^{[0]} = x}$$

⋮

$$\frac{\partial J}{\partial b^{[k+1]}} = \delta^{[k+1]}$$