



Gradient Descent

① network parameters $\theta = \{w_1, \dots, b_1, b_2, \dots\}$

$$\nabla L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial w_1} \\ \vdots \\ \frac{\partial L(\theta)}{\partial b} \end{bmatrix} \quad \begin{array}{l} \text{compute } \nabla L(\theta^0) \quad \theta^1 = \theta^0 - \eta \nabla L(\theta^0) \\ \text{compute } \nabla L(\theta^1) \quad \theta^2 = \theta^1 - \eta \nabla L(\theta^1) \\ \vdots \end{array}$$

Starting parameters $\theta^0 \longrightarrow \theta^1 \longrightarrow \theta^2 \longrightarrow \dots$

But we have millions of parameters — Then we should use backpropagation

② recall that Chain Rule

Chain Rule

Case 1 $y = g(x) \quad z = h(y)$

$$\Delta x \rightarrow \Delta y \rightarrow \Delta z \quad \frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

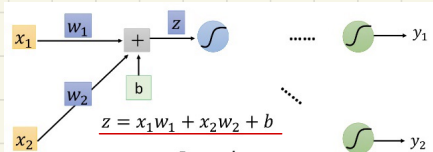
Case 2

$x = g(s) \quad y = h(s) \quad z = k(x, y)$

$$\begin{array}{c} \Delta x \\ \swarrow \quad \searrow \\ \Delta s \quad \Delta y \\ \nwarrow \quad \nearrow \\ \Delta z \end{array} \quad \frac{dz}{ds} = \frac{\partial z}{\partial x} \frac{dx}{ds} + \frac{\partial z}{\partial y} \frac{dy}{ds}$$

对于 $x^n \longrightarrow \begin{bmatrix} N \times N \\ \theta \end{bmatrix} \longrightarrow y^n \xleftrightarrow{C^n} \hat{y}^n$
 $\hookrightarrow y$ 与 \hat{y} 之间的差距

$$L(\theta) = \sum_{n=1}^N L^n(\theta) \longrightarrow \frac{\partial L(\theta)}{\partial w} = \sum_{n=1}^N \frac{\partial L^n(\theta)}{\partial w} \quad \text{逐条计算}$$



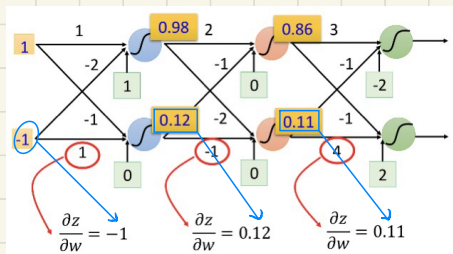
Forward pass: compute $\frac{\partial z}{\partial w}$ for all parameters

Backward pass: compute $\frac{\partial L}{\partial z}$ for all activation function inputs z .

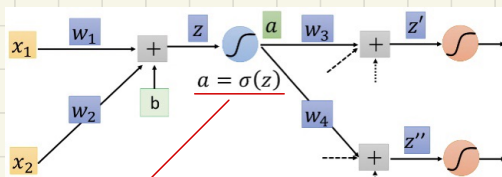
$$\frac{\partial L}{\partial w} = \frac{\partial z}{\partial w} \frac{\partial L}{\partial z}$$

① compute $\frac{\partial z}{\partial w}$ for all parameters forward pass

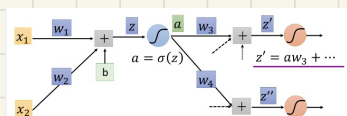
$\left. \begin{array}{l} \frac{\partial z}{\partial w_1} = x_1 \\ \frac{\partial z}{\partial w_2} = x_2 \end{array} \right\}$ 规律: The value of the input connected by the weights



② compute $\frac{\partial C}{\partial z}$ for all parameters *Backward pass*



$$\frac{\partial C}{\partial z} = \frac{\partial a}{\partial z} \cdot \frac{\partial C}{\partial a}$$



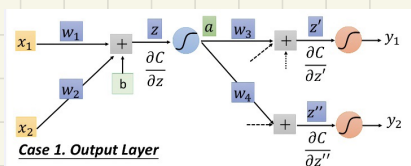
(假设只有两个 σ , 所以只有两项)

$$\begin{aligned} \frac{\partial C}{\partial a} &= \frac{\partial z'}{\partial a} \frac{\partial C}{\partial z'} + \frac{\partial z''}{\partial a} \frac{\partial C}{\partial z''} \\ &= w_3 \frac{\partial C}{\partial z'} + w_4 \frac{\partial C}{\partial z''} \end{aligned}$$

$$\frac{\partial C}{\partial z} = g'(z) [w_3 \frac{\partial C}{\partial z'} + w_4 \frac{\partial C}{\partial z''}]$$

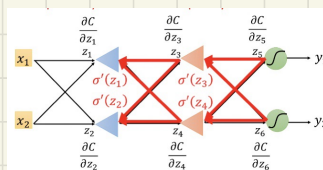


Case 1: Output layer



$$\frac{\partial C}{\partial z'} = \frac{\partial y_1}{\partial z'} \frac{\partial C}{\partial y_1} \quad \frac{\partial C}{\partial z''} = \frac{\partial y_2}{\partial z''} \frac{\partial C}{\partial y_2}$$

Case 2: Not output layer



compute $\frac{\partial C}{\partial z}$ recursively

Until we reach output layer...

反向的递推工作

Summary:

