

△ 反向传播.

KEYWORDS: Forward Pass . Backward Pass .

train 个神经网络实际是在优化个 loss function .

e.g. 数字识别. input 为 256-dim vector, output 为 10-dim vector.

用交叉熵 $C(y, \hat{y}) = -\sum_{i=1}^n \hat{y}_i \ln y_i$ 定义 loss function:

相当于 $L(\theta) = \sum_{n=1}^N C^n(\theta)$. 目标

"1": $y = \begin{bmatrix} y_1 \\ \vdots \\ y_{10} \end{bmatrix}_{(10 \times 1)}$ $\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_{10} \end{bmatrix}_{(10 \times 1)}$ $\theta = \begin{bmatrix} w_1 \\ \vdots \\ b_1 \end{bmatrix}$ $\nabla L(\theta) = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \vdots \\ \frac{\partial L}{\partial b_1} \end{bmatrix}$

$\theta^* = \arg \min_{\theta} L(\theta)$.

gradient descent: $\theta^i \leftarrow \theta^{i-1} - \eta \nabla L(\theta^{i-1})$
 \uparrow learning rate.

Backpropagation: 是 $\frac{\partial L}{\partial w_1} \dots \frac{\partial L}{\partial b_1}$ 的向量形式 (因为参数可能有自己的直接导数).

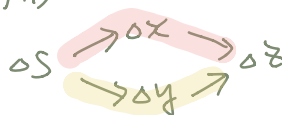
Chain Rule:

① $y = g(x)$ $z = h(y)$ $\Delta x \rightarrow \Delta y \rightarrow \Delta z$

$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$ (按路径从后往前算)

② $x = g(s)$ $y = h(s)$ $z = f(x, y)$

$\frac{dz}{ds} = \frac{\partial z}{\partial x} \frac{dx}{ds} + \frac{\partial z}{\partial y} \frac{dy}{ds}$

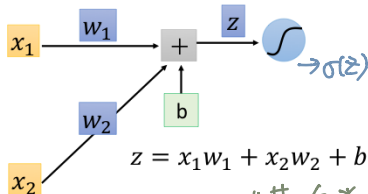


这里是偏导, 因为 z 有 2 个参数.

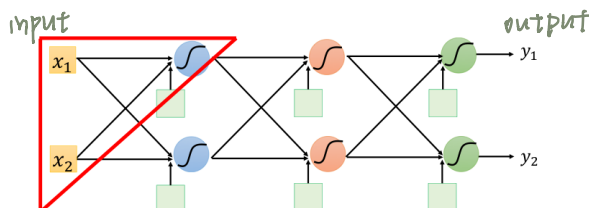
$L(\theta) = \sum_{n=1}^N C^n(\theta)$ $\frac{\partial L(\theta)}{\partial w} = \sum_{n=1}^N \frac{\partial C^n(\theta)}{\partial w}$

$C(\theta)$ 相当于 $C(y_1, \dots, y_p)$
 我们的目标: 算 C 对每个参数的偏微分.

1 neuron:



w 为某一参数:



$z = z(w)$
 $C = C(z)$

$\Delta w \rightarrow \Delta z \rightarrow \Delta C$
 $\therefore \frac{\partial C}{\partial w} = \frac{\partial C}{\partial z} \frac{\partial z}{\partial w}$

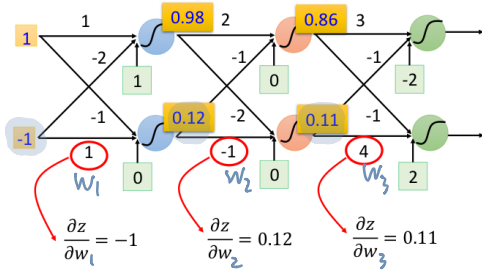
难算 \rightarrow 好算 (有表达式)

Forward Pass: 算 $\frac{\partial z}{\partial w}$

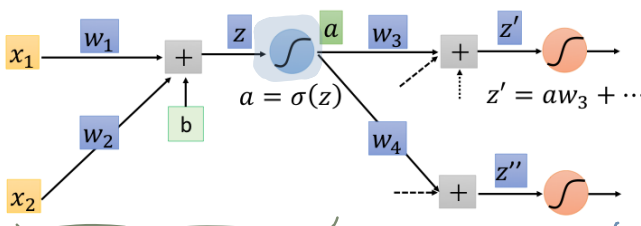
$\frac{\partial z}{\partial w_1} = x_1$, $\frac{\partial z}{\partial w_2} = x_2$. 即看参数 w 前面的 input 是什么, 结果就是什么.

正向计算出所有 neuron 的 output (z) 后

∴ 给定了一组 input (x) 和参数 w , 可以直接读出 $\frac{\partial z}{\partial w}$.
(z 是相对每个 neuron 来算的).



Backward Pass: 算 $\frac{\partial C}{\partial z}$



用 chain rule 表示时, 不一定要把中间每层函数都拆出来, 因为想拆多少都是拆出来的. 重要的是能拆出一部分可求导的.
是常数, 因为 z 已被计算出值.

$$a = \sigma(z) \Rightarrow \frac{\partial C}{\partial z} = \frac{\partial C}{\partial a} \frac{\partial a}{\partial z} = \sigma'(z)$$

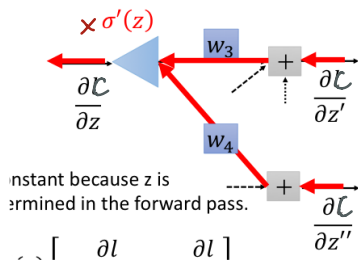
$$z' = z'(a) \quad C = C(z', z'') \quad \Delta a \rightarrow \Delta z' \rightarrow \Delta C$$

$$\frac{\partial C}{\partial a} = \frac{\partial C}{\partial z'} \frac{\partial z'}{\partial a} + \frac{\partial C}{\partial z''} \frac{\partial z''}{\partial a} = w_3 + w_4$$

$$\therefore \frac{\partial C}{\partial a} = w_3 \frac{\partial C}{\partial z'} + w_4 \frac{\partial C}{\partial z''}$$

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

用一个倒过来的 neuron 表示: ① 当 z', z'' 不在 output layer:

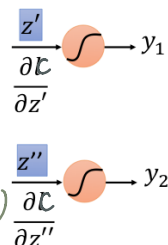


求 $\frac{\partial C}{\partial z'}$, $\frac{\partial C}{\partial z''}$ 和求 $\frac{\partial C}{\partial z}$ 相同.

② 当 z', z'' 在 output layer:

$$y_1 = \sigma(z') \quad C = C(y_1) \quad \frac{\partial C}{\partial z'} = \frac{\partial C}{\partial y_1} \frac{\partial y_1}{\partial z'} = \sigma'(z') \frac{\partial C}{\partial y_1}$$

$$\text{同理, } \frac{\partial C}{\partial z''} = \frac{\partial C}{\partial y_2} \frac{\partial y_2}{\partial z''} = \sigma'(z'') \frac{\partial C}{\partial y_2}$$



$\frac{\partial C}{\partial y}$ 用 C 的定义直接求.

△ Framework

$$\therefore L(\theta) = \sum_{n=1}^N C^n(\theta) \quad \therefore \frac{\partial L(\theta)}{\partial W} = \sum_{n=1}^N \frac{\partial C^n(\theta)}{\partial W}$$

$$\frac{\partial C(\theta)}{\partial W} = \underbrace{\frac{\partial C}{\partial z}}_{\text{Forward Pass, 直接读 input.}} \underbrace{\frac{\partial z}{\partial W}}_{\text{Forward Pass, 直接读 input.}}$$

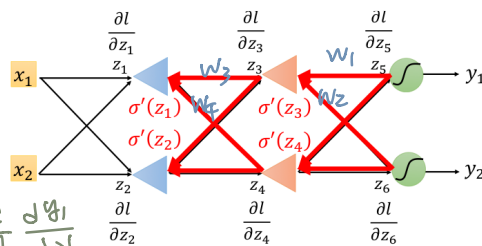
Backward Pass

$$\frac{\partial C}{\partial z} = \underbrace{\frac{\partial C}{\partial a}}_{\sigma'(z)} \underbrace{\frac{\partial a}{\partial z}}_{=1} = \sigma'(z) \quad (z, z' \text{ 为 } w \text{ 所在 neuron 的线性组合值})$$

$$\frac{\partial C}{\partial z'} = \underbrace{\frac{\partial C}{\partial a}}_{\sigma'(z')} \underbrace{\frac{\partial a}{\partial z'}}_{=w_3} + \underbrace{\frac{\partial C}{\partial z''}}_{\sigma'(z'')} \underbrace{\frac{\partial a}{\partial z'}}_{=w_4}$$

$$\begin{cases} z, z' \text{ 在 output layer: } \frac{\partial C}{\partial z'} = \frac{\partial C}{\partial y_1} \frac{\partial y_1}{\partial z'} \\ \frac{\partial C}{\partial z''} = \frac{\partial C}{\partial y_2} \frac{\partial y_2}{\partial z''} \end{cases} \text{ 均需求.}$$

z, z'' 不在 output layer: 连续递归



可推出:

$$\frac{\partial L}{\partial z_5} = \frac{\partial L}{\partial y_1} \frac{\partial y_1}{\partial z_5}$$

$$\frac{\partial L}{\partial z_3} = \sigma'(z_3) \left(w_1 \frac{\partial L}{\partial z_5} + w_2 \frac{\partial L}{\partial z_6} \right)$$

$$\frac{\partial L}{\partial z_1} = \sigma'(z_1) \left(w_3 \frac{\partial L}{\partial z_3} + w_4 \frac{\partial L}{\partial z_4} \right)$$