

chap-4 前馈神经网络.

4.1 神经元

净输入 z 一个神经元获得输入信号 x 的加权总和

$$z = \sum w_i x_i + b$$

$$= \bar{w}^T \bar{x} + b$$

活性值

$$a = f(z)$$

\Rightarrow 4.1.1 Sigmoid 型激活函数

- logistic

- tanh

$$\sigma'(x) = \sigma(x) (1 - \sigma(x))$$

$$\sigma'(\bar{x}) = \text{diag}(\sigma(\bar{x})) \odot (1 - \sigma(\bar{x}))$$

Logistic 函数

$$\sigma(x) = \frac{1}{1 + \exp(-x)} \quad \in (0, 1)$$

Tanh 函数

$$\tanh(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)} \quad \in (-1, 1)$$

$$\tanh(x) = 2\sigma(2x) - 1$$

The Tanh 的输出是零中心化的. logistic 不是, 总是大于 0.

非零中心化 使得其上一层神经网络产生偏置 (bias shift), 并进一步使 GD 收敛速度变慢.

⇒ ⇒ 4.1.1.1 Hard-logistic 和 Hard-Tanh 函数

用另两个函数近似

$$\text{hard-logistic}(x) = \max(\min(0.25x + 0.5, 1), 0)$$

$$\text{hard-tanh}(x) = \max(\min(x, 1), -1)$$

⇒ 4.1.1 修正线性单元

$$\text{ReLU}(x) = \max(0, x)$$

优点 1. 计算高效

2. 生物学解释: 单侧抑制, 兴奋有边界

3. 稀疏性 50% 的神经元

4. 缓解梯度消失

缺点 1. 非零中心化, 偏差偏移

2. 死 ReLU 问题

No.

Date.

⇒⇒ 带泄露的 ReLU

$$\text{Leaky ReLU}(x) = \begin{cases} x & \text{if } (x > 0) \\ rx & \text{if } (x < 0) \end{cases}$$

$$= \max(0, x) + r \min(0, x)$$

当 $r < 1$ 时, r 一般等于 0.01.

$$\text{Leaky ReLU}(x) = \max(x, rx)$$

相当于一个简单的 max out 单元

⇒⇒ ~~P.T. 2.2~~ 带参数的 ReLU

引入一个可学习的参数, 不同神经元可以有不同参数。

$$\text{PReLU}_i(x) = \begin{cases} x & \text{if } x > 0 \\ r_i x & \text{if } x < 0 \end{cases}$$

$$= \max(0, x) + r_i \min(0, x)$$

⇒⇒ ELU 指数线性单元

$$\text{ELU}(x) = \begin{cases} x & \text{if } x > 0 \\ r(\exp(x) - 1) & \text{if } x < 0 \end{cases}$$

$$= \max(0, x) + \min(0, r(\exp(x) - 1))$$

$\Rightarrow \Rightarrow$ softplus 函数

rectifier 的平滑版本, 没有系统收敛性

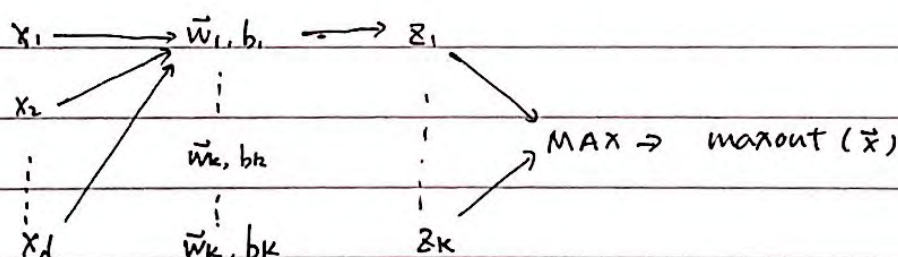
$$\text{softplus}(x) = \log(1 + \exp(x))$$

④ p.1.3 Maxout 单元

maxout 的输入是上层神经元的全部输出 $\vec{x} = [x_1; \dots; x_d]$

k 个权重向量 $\vec{w}_k \in \mathbb{R}^d$ 和偏置 b_k , k 个净输出 z_k

$$z_k = \vec{w}_k^T \vec{x} + b_k$$



$$\text{maxout}(\vec{x}) = \max(z_k)$$

整体学习输入到输出之间的非线性映射关系。

No.

Date.

4.2 网络结构

⇒ 4.2.1 前馈网络

全连接前馈网络

卷积神经网络

⇒ 4.2.2 反馈网络

循环神经网络

Hopfield网络. 玻尔兹曼机

递归增强网络

⇒ 4.2.3 图网络

4.1 前馈神经网络

L : 神经网络的层数

n^L : 第 L 层神经网络的个数

f_{act} : 激活函数

$W^{(l)}$: $(l-1)$ 到 l 层的权重矩阵

$b^{(l)}$: l 层的偏置

$\vec{z}^{(l)}$: l 层神经元的净输入

$\vec{a}^{(l)}$: l 层的净输出

$$\vec{z}^{(l)} = W^{(l)} \cdot \vec{a}^{(l-1)} + b^{(l)}$$

$$\vec{a}^{(l)} = f_{\text{act}}(\vec{z}^{(l)})$$

$$a^{(l)} \longrightarrow z^{(l)} \longrightarrow a^{(l)}$$

⇒ 4.3.1 通用近似定理

只需要一个包含足够多神经元的隐藏层，多层前馈网络能以任意精度逼近任意复杂度的函数。

⇒ 4.3.2 应用到机器学习

特征抽取 $x \rightarrow \varphi(x)$

No.

Date.

⇒ 4.3.3 参数学习

$$L(\vec{y}, \hat{\vec{y}}) = -\vec{y}^T \log \hat{\vec{y}}$$

EMF

$$R(w, \vec{b}) = \frac{1}{N} \sum L(\vec{y}^{(w)}, \hat{\vec{y}}^{(w)}) + \frac{1}{2} \lambda \|W\|_F^2$$

Frobenius 范数

$$\|W\|_F^2 = \sum_i \sum_j \sum_k (w_{ij}^{(k)})^2$$

第 L 层参数更新公式

$$w^{(L)} \leftarrow w^{(L)} - \alpha \cdot \frac{\partial R(w, \vec{b})}{\partial w^{(L)}}$$

$$\vec{b}^{(L)} \leftarrow \vec{b}^{(L)} - \alpha \cdot \frac{\partial R(w, \vec{b})}{\partial \vec{b}^{(L)}}$$

4.4 反向传播算法

因为 $\frac{\partial L(\vec{y}, \hat{\vec{y}})}{\partial w^{(l)}}$ 计算涉及矩阵微分, 较为繁琐, 故我们先计算 $\frac{\partial L(\vec{y}, \hat{\vec{y}})}{\partial w^{(l)}}$



4.4.1 $X \rightarrow \vec{y} \rightarrow z$

$$\frac{\partial z}{\partial x_{ij}} = \left(\frac{\partial z}{\partial \vec{y}} \right)^T \frac{\partial \vec{y}}{\partial x_{ij}}$$

$$\frac{\partial L(\vec{y}, \hat{\vec{y}})}{\partial w_{ij}^{(l)}} = \left(\frac{\partial \vec{z}^{(l)}}{\partial w_{ij}^{(l)}} \right)^T \cdot \frac{\partial L(\vec{y}, \hat{\vec{y}})}{\partial \vec{z}^{(l)}} \quad (4.46)$$

$$\frac{\partial L(\vec{y}, \hat{\vec{y}})}{\partial b^{(l)}} = \left(\frac{\partial \vec{z}^{(l)}}{\partial b^{(l)}} \right)^T \cdot \frac{\partial L(\vec{y}, \hat{\vec{y}})}{\partial \vec{z}^{(l)}} \quad (4.47)$$

$$\frac{\partial L(\vec{y}, \hat{\vec{y}})}{\partial \vec{z}^{(l)}}$$

目标函数关于第 l 层神经元 $z^{(l)}$ 的偏导数

称为 误差项

$$\frac{\partial \vec{z}^{(l)}}{\partial b^{(l)}}$$

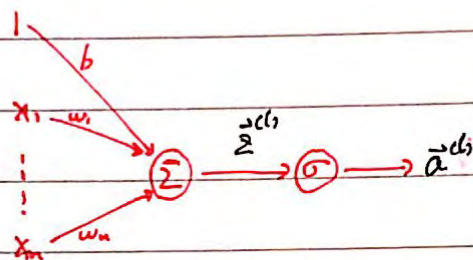
$$\frac{\partial L(\vec{y}, \hat{\vec{y}})}{\partial \vec{z}^{(l)}} \cdot \frac{\partial \vec{z}^{(l)}}{\partial w_{ij}^{(l)}}$$

No.

Date.

$\Rightarrow \Rightarrow \Rightarrow$ 计算偏导数 $\frac{\partial \vec{z}^{(l)}}{\partial w_{ij}^{(l)}}$

由 $\vec{z}^{(l)}$ 与 $w_{ij}^{(l)}$ 的函数关系为 $\vec{z}^{(l)} = w^{(l)} \vec{a}^{(l-1)} + \vec{b}^{(l)}$



$$\frac{\partial \vec{z}^{(l)}}{\partial w_{ij}^{(l)}} = \frac{\partial (w_{ij}^{(l)} \vec{a}^{(l-1)} + \vec{b}^{(l)})}{\partial w_{ij}^{(l)}}$$

$$= \begin{bmatrix} \frac{\partial (w_{1j}^{(l)} \vec{a}^{(l-1)} + \vec{b}^{(l)})}{\partial w_{ij}^{(l)}} \\ \vdots \\ \frac{\partial (w_{mj}^{(l)} \vec{a}^{(l-1)} + \vec{b}^{(l)})}{\partial w_{ij}^{(l)}} \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ a_j^{(l-1)} \\ \vdots \\ 0 \end{bmatrix} = \mathbb{I}_i(a_j^{(l-1)})$$

$\Rightarrow \Rightarrow \Rightarrow$ 计算偏导数 $\frac{\partial \vec{z}^{(l)}}{\partial \vec{b}^{(l)}}$

由 $\vec{z}^{(l)} = w^{(l)} \vec{a}^{(l-1)} + \vec{b}^{(l)}$

$$\frac{\partial \vec{z}^{(l)}}{\partial \vec{b}^{(l)}} = \mathbb{I}_{m^{(l)}} \text{ 为 } m^{(l)} \times m^{(l)} \text{ 的单位矩阵}$$

$\Rightarrow \Rightarrow \Rightarrow$ 计算误差项 $\frac{\partial L(\vec{y}, \hat{\vec{y}})}{\partial \vec{z}^{(l)}}$

第 l 层神经元的误差项

表示 l 层神经元对最终误差的响应。

$$\delta^{(l)} = \frac{\partial L(\vec{y}, \hat{\vec{y}})}{\partial \vec{z}^{(l)}}$$

根据 $\vec{z}^{(l+1)} = W^{(l+1)} \vec{a}^{(l)} + b^{(l+1)}$

$$\frac{\partial \vec{z}^{(l+1)}}{\partial \vec{a}^{(l)}} = (W^{(l+1)})^T$$

根据 $\vec{a}^{(l)} = f(\vec{z}^{(l)})$, $f(\cdot)$ 为逐点计算函数

$$\frac{\partial \vec{a}^{(l)}}{\partial \vec{z}^{(l)}} = \frac{\partial f(\vec{z}^{(l)})}{\partial \vec{z}^{(l)}}$$

$$= \text{diag}(f'_i(\vec{z}^{(l)}))$$

根据链式法则

$$\delta^{(l)} = \frac{\partial L(\vec{y}, \hat{\vec{y}})}{\partial \vec{z}^{(l)}}$$

$$= \frac{\partial \vec{a}^{(l)}}{\partial \vec{z}^{(l)}} \cdot \frac{\partial \vec{z}^{(l+1)}}{\partial \vec{a}^{(l)}} \cdot \frac{\partial L(\vec{y}, \hat{\vec{y}})}{\partial \vec{z}^{(l+1)}}$$

$$= \text{diag}(f'_i(\vec{z}^{(l)})) \cdot (W^{(l+1)})^T \cdot \delta^{(l+1)}$$

$$= f'_{l+1}(\vec{z}^{(l)}) \odot ((W^{(l+1)})^T \delta^{(l+1)})$$

No.

Date.

$$\frac{\partial L(\vec{y}, \hat{\vec{y}})}{\partial w_{ij}^{(1)}} = \Pi_i (a_j^{(1)})^T \delta^{(6)} = \delta_i^{(6)} a_j^{(6)}$$

第(6)层权重 $w^{(6)}$ 的梯度为

$$\frac{\partial L(\vec{y}, \hat{\vec{y}})}{\partial w^{(6)}} = \delta^{(6)} \cdot (a^{(6)})^T$$

同理

$$\frac{\partial L(\vec{y}, \hat{\vec{y}})}{\partial b^{(6)}} = \delta^{(6)}$$

4.5 自动微分

⇒ 4.5.1 数值微分

计算不准

计算复杂

⇒ 4.5.2 符号微分

优点：平台无关

处理数学表达式

缺点：编译时间长

查询语言

难以调试

⇒ 4.5.3 自动微分

处理一个数学表达式

$$f(x; w, b) = \frac{1}{\exp(-wx + b) + 1} \Rightarrow \text{求解或计算图}$$

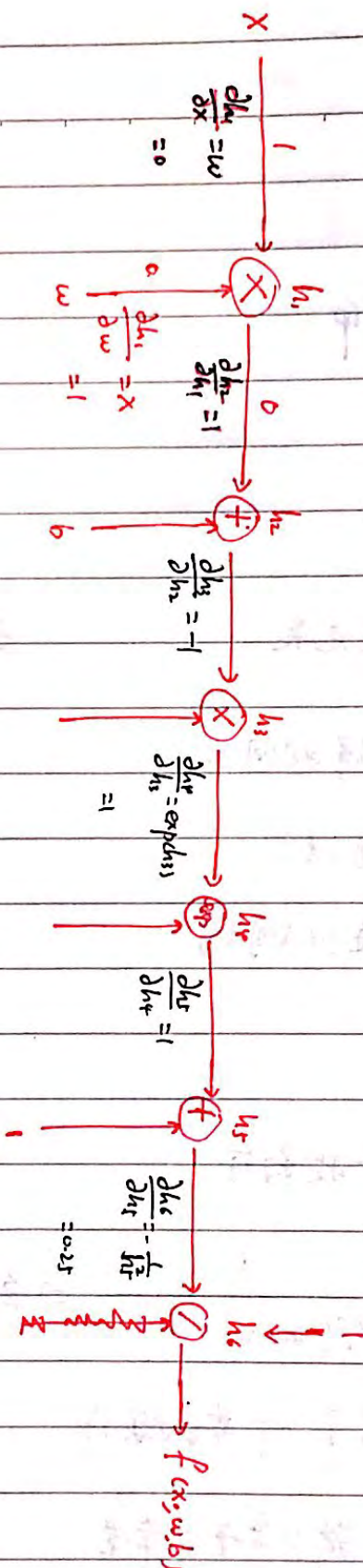
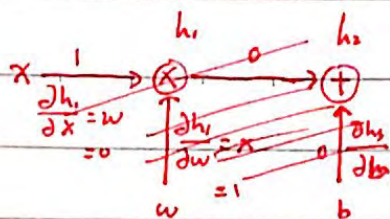
非叶子节点表示一个基本操作

叶子节点为一个输入变量或常量

复合函数 $f(x; w, b)$ 由 6 个基本函数 h_i 组成

No.

Date.



$$\frac{\partial f(x; w, b)}{\partial w} = \frac{\partial f(x; w, b)}{\partial h_6} \frac{\partial h_6}{\partial h_5} \frac{\partial h_5}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w}$$

$$\frac{\partial f(x; w, b)}{\partial b} = \frac{\partial f(x; w, b)}{\partial h_6} \frac{\partial h_6}{\partial h_5} \frac{\partial h_5}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial b}$$

如果参数步数之间有多条路径, 可以将多条路径上的参数再相加, 得到最终梯度。

前向模式 $\frac{\partial h_1}{\partial w} \rightarrow \frac{\partial f(x; w, b)}{\partial w}$

反向模式 $\frac{\partial f(x; w, b)}{\partial w} \rightarrow \frac{\partial h_1}{\partial w}$

符号微分与自动微分

静态计算图与动态计算图。

No.

Date.

4.6 优化问题

⇒ 4.6.1 非凸优化问题

⇒ 4.6.2 梯度消失问题

$$s^{(L)} = f_L(\tilde{z}^{(L)}) \odot (W^{(L+1)})^T s^{(L+1)}$$

sigmoid 饱和已导致接近 0.