

yqkong@dlut.edu.cn

目录(CONTENT)



- 01 多层感知机 (Multi-layer Perceptron)
- 02 前馈神经网络 (Feedforward Neural Network)
- 03 梯度反向传播 (Gradient Back Propagation)
- 04 闭形式解 (Close-form Solution)



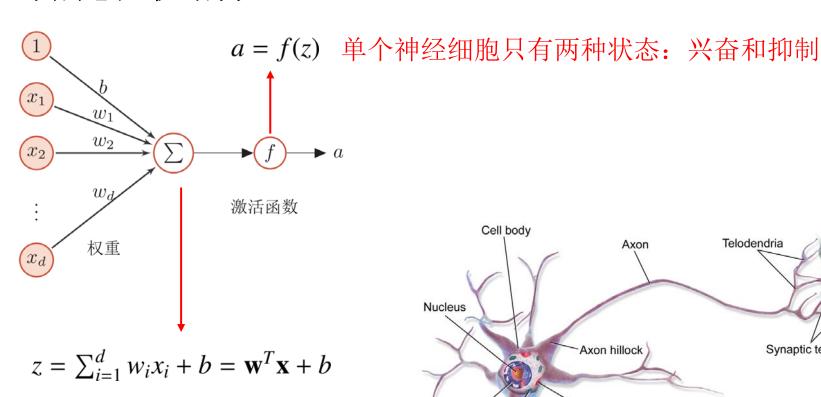


多层感知机 Multi-layer Perceptron

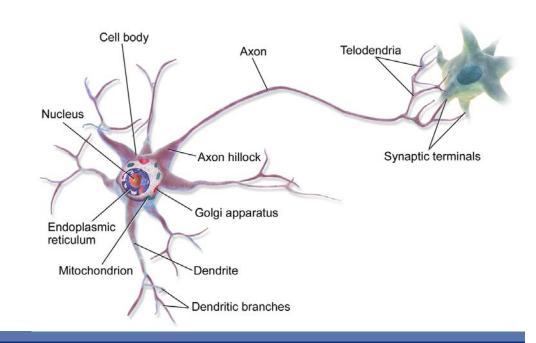
单层感知机



■ 单层感知机结构



一个简单的线性模型!



单层感知机



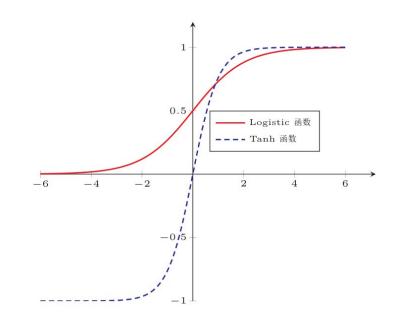
■ 常见激活函数

■ Sigmoid激活函数

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

□ Tanh激活函数

$$\tanh(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$$



□性质

- ✔ 饱和函数: 自变量绝对值较大的时候, 值接近;
- ✓ tanh函数是零中心化的,而logistic函数的输出恒大于0;
- ✓ 非零中心化的输出会使得其后一层的神经元的输入发生偏置 偏移(bias shift),并进一步使得梯度下降的收敛速度变慢。

单层感知机



■ 常见激活函数

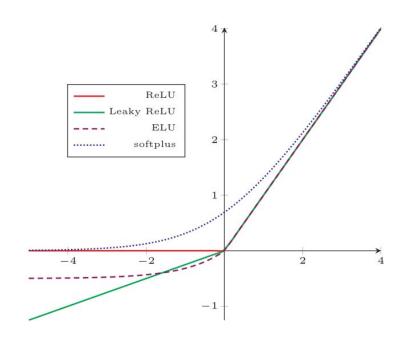
$$ReLU(x) = \begin{cases} x & x \ge 0 \\ 0 & x < 0 \end{cases}$$

LeakyReLU(x) =
$$\begin{cases} x & \text{if } x > 0 \\ \gamma x & \text{if } x \le 0 \end{cases}$$

$$PReLU_i(x) = \begin{cases} x & \text{if } x > 0\\ \gamma_i x & \text{if } x \le 0 \end{cases}$$

$$ELU(x) = \begin{cases} x & \text{if } x > 0\\ \gamma(\exp(x) - 1) & \text{if } x \le 0 \end{cases}$$
$$= \max(0, x) + \min(0, \gamma(\exp(x) - 1))$$

$$softplus(x) = log(1 + exp(x))$$



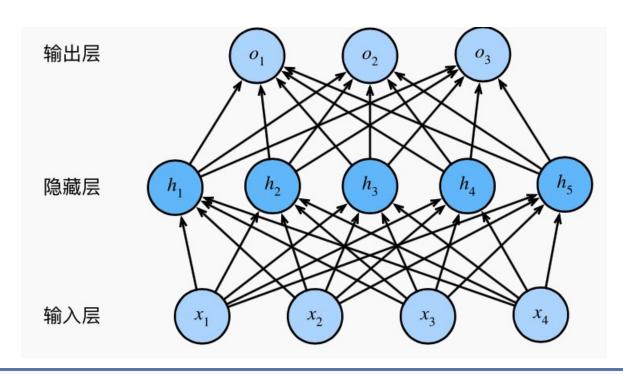
死亡ReLU问题(Dying ReLU Problem)

- □ 计算上更加高效
- □ 生物学合理性
- □ 单侧抑制、宽兴奋边界
- □ 在一定程度上缓解梯度消失问题

多层感知机



- □ 多层感知机在单层感知机的基础上引入了一到多个隐藏层(hidden layer)。隐藏层位于输入层和输出层之间。不难发现,即便再添加更多的隐藏层,以上设计依然只能与仅含输出层的单层感知机等价。
- □ 线性变换下,等价于一个单层感知机







前馈神经网络 Feedforward NNs



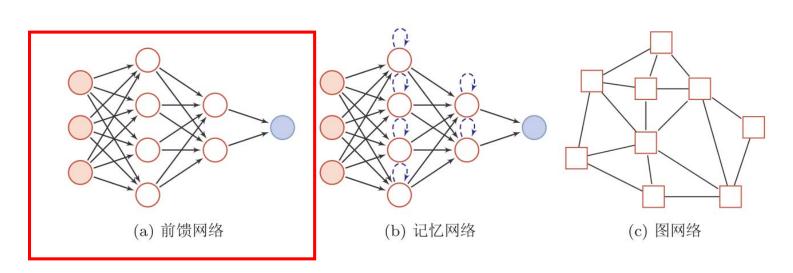
■ 人工神经网络 (Artificial Neural Network)

人工神经网络主要由大量的神经元以及它们之间的有向连接构成。 因此考虑三方面:

- □ 神经元的激活规则
 - 主要是指神经元输入到输出之间的映射关系,一般为非线性函数。
- □ 网络的拓扑结构
 - 不同神经元之间的连接关系
- □ 学习算法
 - 通过训练数据来学习神经网络的参数



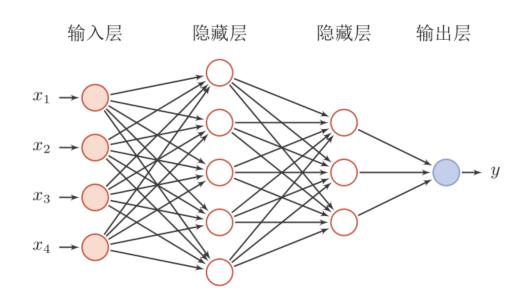
- 前馈神经网络 (Feedforward Neural Networks)
 - □ 人工神经网络由神经元模型构成,这种由许多神经元组成的信息处理网络具有并行分布结构。
 - □ **前馈神经网络**是人工神经网络的一种形式,信息前向传递,实现逐层编码处理的网络结构。



圆形节点表示一个神经元, 方形节点表示一组神经元。



- 前馈神经网络(全连接神经网络、多层感知器)
 - □ 各神经元分别属于不同的层,层内无连接。
 - □ 相邻两层之间的神经元全部两两连接。
 - 整个网络中无反馈,信号从输入层向输出层单向传播,可用一个有向无环图表示。



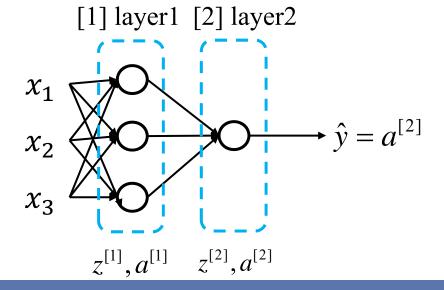
记号	含义
L	神经网络的层数
M_l	第1层神经元的个数
$f_l(\cdot)$	第1层神经元的激活函数
$\pmb{W}^{(l)} \in \mathbb{R}^{M_l imes M_{l-1}}$	第 $l-1$ 层到第 l 层的权重矩阵
$\pmb{b}^{(l)} \in \mathbb{R}^{M_l}$	第 <i>l</i> -1层到第 <i>l</i> 层的偏置
$\pmb{z}^{(l)} \in \mathbb{R}^{M_l}$	第1层神经元的净输入(净活性值)
$\pmb{a}^{(l)} \in \mathbb{R}^{M_l}$	第1层神经元的输出(活性值)



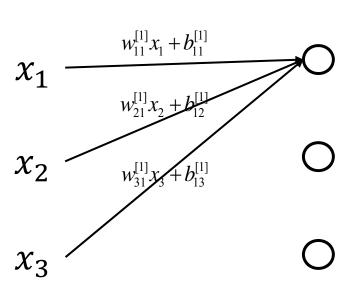
■ 逻辑回归模型

$$\begin{cases} x \\ w \\ b \end{cases} \rightarrow z = w^{T} x + b \rightarrow \hat{y} = \sigma(z) \rightarrow L(\hat{y}, y)$$

■ 神经网络模型



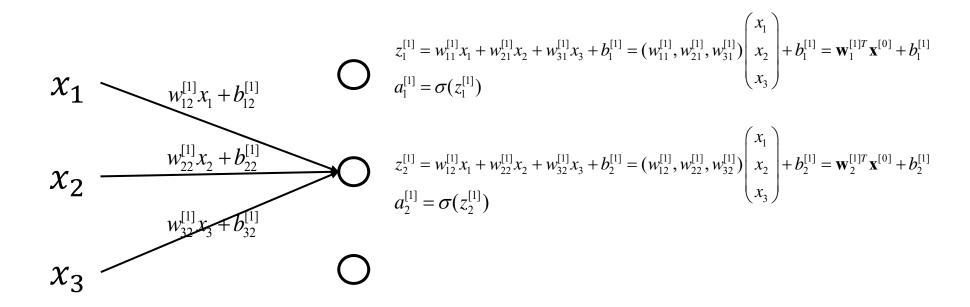




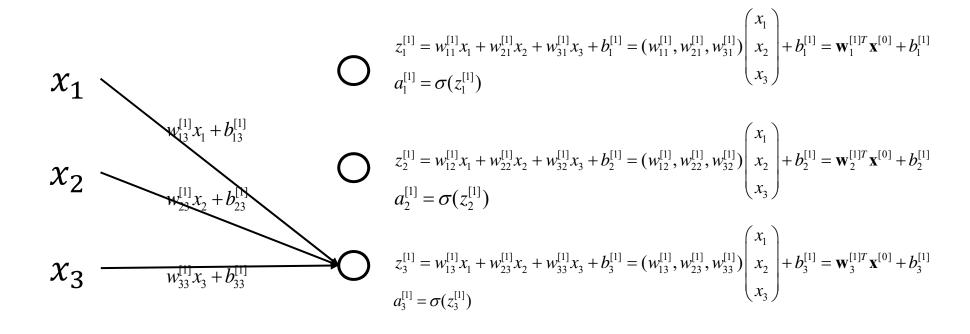
$$z_{1}^{[1]} = w_{11}^{[1]} x_{1} + w_{21}^{[1]} x_{2} + w_{31}^{[1]} x_{3} + b_{1}^{[1]} = (w_{11}^{[1]}, w_{21}^{[1]}, w_{31}^{[1]}) \begin{pmatrix} x_{1} \\ x_{2} \\ x_{3} \end{pmatrix} + b_{1}^{[1]} = \mathbf{w}_{1}^{[1]T} \mathbf{x}^{[0]} + b_{1}^{[1]}$$

$$a_{1}^{[1]} = \sigma(z_{1}^{[1]})$$





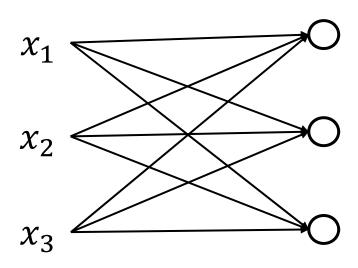






$$z_{1}^{[1]} = w_{11}^{[1]} x_{1} + w_{21}^{[1]} x_{2} + w_{31}^{[1]} x_{3} + b_{1}^{[1]} = (w_{11}^{[1]}, w_{21}^{[1]}, w_{31}^{[1]}) \begin{pmatrix} x_{1} \\ x_{2} \\ x_{3} \end{pmatrix} + b_{1}^{[1]} = \mathbf{w}_{1}^{[1]T} \mathbf{x}^{[0]} + b_{1}^{[1]}$$

$$a_{1}^{[1]} = \sigma(z_{1}^{[1]})$$



$$z_{2}^{[1]} = w_{12}^{[1]} x_{1} + w_{22}^{[1]} x_{2} + w_{32}^{[1]} x_{3} + b_{2}^{[1]} = (w_{12}^{[1]}, w_{22}^{[1]}, w_{32}^{[1]}) \begin{pmatrix} x_{1} \\ x_{2} \\ x_{3} \end{pmatrix} + b_{2}^{[1]} = \mathbf{w}_{2}^{[1]T} \mathbf{x}^{[0]} + b_{2}^{[1]}$$

$$a_{2}^{[1]} = \sigma(z_{2}^{[1]})$$

$$z_{3}^{[1]} = w_{13}^{[1]} x_{1} + w_{23}^{[1]} x_{2} + w_{33}^{[1]} x_{3} + b_{3}^{[1]} = (w_{13}^{[1]}, w_{23}^{[1]}, w_{33}^{[1]}) \begin{pmatrix} x_{1} \\ x_{2} \\ x_{3} \end{pmatrix} + b_{3}^{[1]} = \mathbf{w}_{3}^{[1]T} \mathbf{x}^{[0]} + b_{3}^{[1]}$$

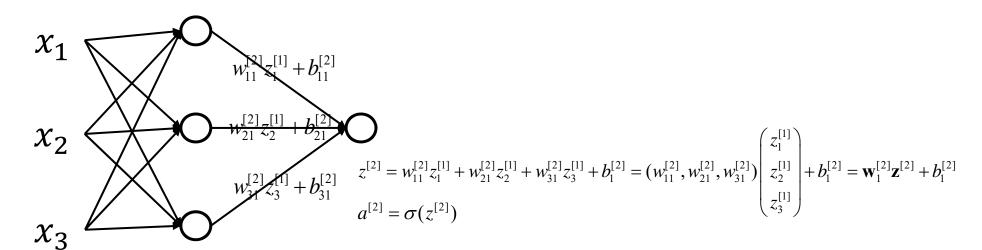
$$a_{3}^{[1]} = \sigma(z_{3}^{[1]})$$

$$\begin{pmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \end{pmatrix} = \begin{pmatrix} w_{11}^{[1]} & w_{21}^{[1]} & w_{31}^{[1]} \\ w_{12}^{[1]} & w_{22}^{[1]} & w_{32}^{[1]} \\ w_{13}^{[1]} & w_{23}^{[1]} & w_{33}^{[1]} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} + \begin{pmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \end{pmatrix}$$

$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]T} \mathbf{z}^{[0]} + \mathbf{b}^{[1]}$$
$$\mathbf{a}^{[1]} = \sigma(\mathbf{z}^{[1]})$$



■ 第二层

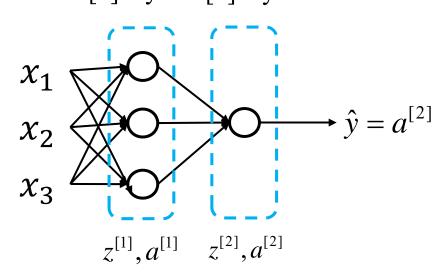


$$\mathbf{z}^{[2]} = \mathbf{W}^{[2]T} \mathbf{a}^{[1]} + \mathbf{b}^{[2]}$$



■ 第L层

[1] layer1 [2] layer2



$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]T} \mathbf{z}^{[0]} + \mathbf{b}^{[1]}$$
 第一层

$$\mathbf{z}^{[2]} = \mathbf{W}^{[2]T} \mathbf{a}^{[1]} + \mathbf{b}^{[2]}$$
 第二层

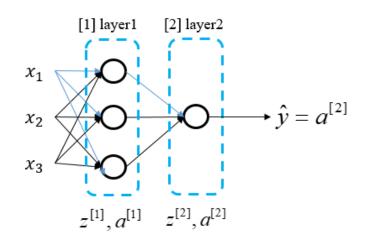
$$\mathbf{z}^{[L]} = \mathbf{W}^{[L]T} \mathbf{a}^{[L-1]} + \mathbf{b}^{[L]}$$
 第L层



■ 第一层

$$\mathbf{W}^{[1]} \rightarrow \mathbf{z}^{[1]} = \mathbf{W}^{[1]T} \mathbf{x} + \mathbf{b}^{[1]} \rightarrow \mathbf{a}^{[1]} = \sigma(\mathbf{z}^{[1]})$$

$$\mathbf{b}^{[1]}$$



■ 第二层

$$\mathbf{a}^{[1]} = \sigma(\mathbf{z}^{[1]}) \\
\mathbf{W}^{[2]} \\
\mathbf{b}^{[2]}$$

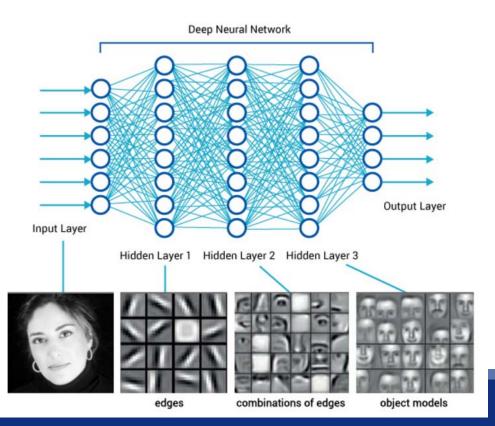
$$\rightarrow \mathbf{z}^{[2]} = \mathbf{W}^{[2]T} \mathbf{a}^{[1]} + \mathbf{b}^{[2]} \rightarrow \mathbf{a}^{[2]} = \sigma(\mathbf{z}^{[2]}) \rightarrow L(\mathbf{a}^{[2]}, y)$$



■ 前馈计算—信息向前传递

$$z^{(l)} = W^{(l)}a^{(l-1)} + b^{(l)},$$

 $a^{(l)} = f_l(z^{(l)}).$
 $x = a^{(0)} \to z^{(1)} \to a^{(1)} \to z^{(2)} \to \cdots \to a^{(L-1)} \to z^{(L)} \to a^{(L)} = \phi(x; W, b))$





■ 通用近似定理

定理 4.1 – 通用近似定理(Universal Approximation Theorem) [Cybenko, 1989, Hornik et al., 1989]: $\Diamond \varphi(\cdot)$ 是一个非常数、有界、单调递增的连续函数, \mathcal{I}_d 是一个d维的单位超立方体 $[0,1]^d$, $C(\mathcal{I}_d)$ 是定义在 \mathcal{I}_d 上的连续函数集合。对于任何一个函数 $f \in C(\mathcal{I}_d)$,存在一个整数 m,和一组实数 $v_i, b_i \in \mathbb{R}$ 以及实数向量 $\mathbf{w}_i \in \mathbb{R}^d$, $i = 1, \cdots, m$,以至于我们可以定义函数

$$F(\mathbf{x}) = \sum_{i=1}^{m} v_i \varphi(\mathbf{w}_i^{\mathrm{T}} \mathbf{x} + b_i), \qquad (4.33)$$

作为函数 f 的近似实现,即

$$|F(\mathbf{x}) - f(\mathbf{x})| < \epsilon, \forall \mathbf{x} \in \mathcal{I}_d.$$
 (4.34)

其中 $\epsilon > 0$ 是一个很小的正数。

根据通用近似定理,对于具有线性输出层和至少一个使用"挤压"性质的激活函数的隐藏层组成的前馈神经网络,只要其隐藏层神经元的数量足够,它可以以任意的精度来近似任何从一个定义在实数空间中的有界闭集函数。



■ 通用近似定理

□ 神经网络可以作为一个"万能"函数来使用,可以用来进行 复杂的特征转换,或逼近一个复杂的条件分布

$$\hat{y} = g(\varphi(\mathbf{x}), \theta)$$

神经网络



■ 通用近似定理

□ 神经网络可以作为一个"万能"函数来使用,可以用来进行 复杂的特征转换,或逼近一个复杂的条件分布

$$\hat{y} = g(\varphi(\mathbf{x}), \theta)$$

分类器

山 如果 $g(\cdot)$ 为Logistic回归,那么Logistic回归分类器可以看成神经网络的最后一层



■ 分类问题

□ 如果使用Softmax回归分类器,相当于网络最后一层设置C个神经元,其输出经过Softmax函数进行归一化后可以作为每个类的条件概率。

$$\hat{\mathbf{y}} = \operatorname{softmax}(\mathbf{z}^{(L)})$$

□ 采用交叉熵损失函数,对于样本(x,y),其损失函数为

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = -\mathbf{y}^{\mathrm{T}} \log \hat{\mathbf{y}}$$



■ 交叉熵损失函数

□ 信息量:信息是用来消除随机不确定性的东西。信息量的大小与信息 发生的概率成反比。概率越大,信息量越小。概率越小,信息量越大。 设某一事件发生的概率为P(x),其信息量表示为

$$I(x) = -\log(P(x))$$

□ 信息熵: 所有信息量的期望 $\mathbf{X} = (x_1, x_2, \dots, x_n)$

$$H(\mathbf{X}) = \sum_{i=1}^{n} P(x_i) \log(P(x_i))$$

□ 相对熵 (KL散度):

$$D_{KL}(p \parallel q) = \sum_{i=1}^{n} p(x_i) \log(\frac{p(x_i)}{q(x_i)})$$
真实分布 预测分布



■ 交叉熵损失函数

□ 交叉熵:交叉熵等于KL散度加上一个常量(信息熵),且公式相比 KL散度更加容易计算,所以在机器学习中常常使用交叉熵损失函数 来计算loss

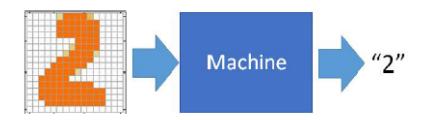
$$D_{KL}(p \| q) = \sum_{i=1}^{n} p(x_i) \log(\frac{p(x_i)}{q(x_i)})$$

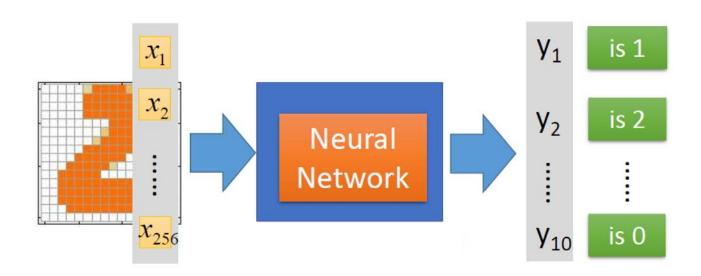
$$= \sum_{i=1}^{n} p(x_i) \log(p(x_i)) - \sum_{i=1}^{n} p(x_i) \log(q(x_i))$$

$$= -H(p(x)) + \left[-\sum_{i=1}^{n} p(x_i) \log(q(x_i))\right]$$
信息熵 交叉熵



■ 例: 手写数字识别



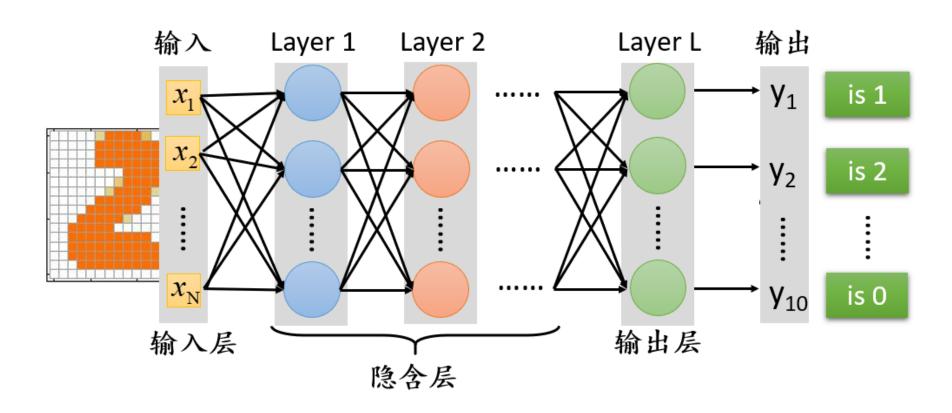


输入: 256维向量

输出: 10维向量



■ 例: 手写数字识别

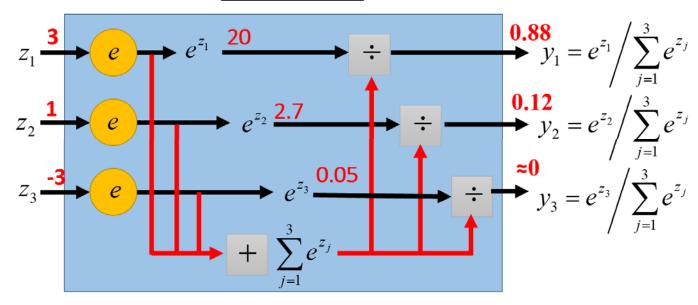


对网络结构进行合理设计,使其函数集合中能够有一个较好的函数



- 最后的输出层—Softmax Layer
 - □ 图像分类:将输入图像的信息转换成其属于各个类别的概率

Softmax Layer



Probability:

■
$$1 > y_i > 0$$

$$\blacksquare \sum_i y_i = 1$$





Gradient Back Propagation



■ 参数学习 (Parameter Learning)

□ 给定训练集为 $D = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^{N}$, 将每个样本 $\mathbf{x}^{(n)}$ 输入给前馈神经网络,得到网络输出为 $\hat{y}^{(n)}$, 骑在数据集D上的结构化风险函数为:

$$\mathcal{R}(W, \mathbf{b}) = \frac{1}{N} \sum_{n=1}^{N} \mathcal{L}(\mathbf{y}^{(n)}, \hat{\mathbf{y}}^{(n)}) + \frac{1}{2} \lambda ||W||_F^2$$

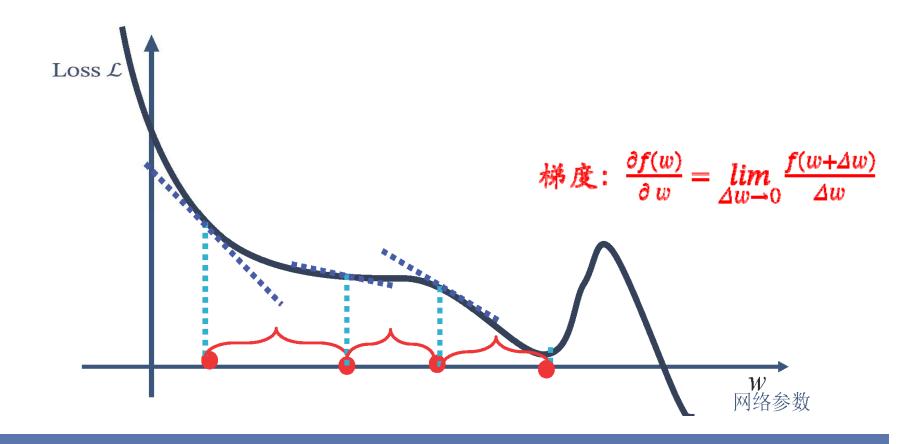
□梯度下降:

$$W^{(l)} \leftarrow W^{(l)} - \alpha \frac{\partial \mathcal{R}(W, \mathbf{b})}{\partial W^{(l)}}$$

$$\mathbf{b}^{(l)} \leftarrow \mathbf{b}^{(l)} - \alpha \frac{\partial \mathcal{R}(W, \mathbf{b})}{\partial \mathbf{b}^{(l)}}$$



- 梯度下降
 - □梯度下降类似盲人下山





■ 梯度下降

□神经网络是一个复杂的复合函数(链式法则)

$$y = f^{5}(f^{4}(f^{3}(f^{2}(f^{1}(x))))) \rightarrow \frac{\partial y}{\partial x} = \frac{\partial f^{5}}{\partial f^{4}} \frac{\partial f^{4}}{\partial f^{3}} \frac{\partial f^{3}}{\partial f^{2}} \frac{\partial f^{2}}{\partial f^{1}} \frac{\partial f^{1}}{\partial x}$$

- □ 梯度反向传播算法(Back Propagation): 根据前馈网络的特点而设计的 高效方法
- 更加通用的计算方法: 自动微分(Automatic Differentiation)



■ 链式法则

□链式法则(Chain Rule)是在微积分中求复合函数导数的一种常用方法。

(1)若
$$x \in \mathbb{R}$$
, $u = u(x) \in \mathbb{R}^s$, $g = g(u) \in \mathbb{R}^t$, 则

$$\frac{\partial \boldsymbol{g}}{\partial x} = \frac{\partial \boldsymbol{u}}{\partial x} \frac{\partial \boldsymbol{g}}{\partial \boldsymbol{u}} \in \mathbb{R}^{1 \times t}.$$

$$(2)$$
若 $oldsymbol{x} \in \mathbb{R}^p$, $oldsymbol{y} = g(oldsymbol{x}) \in \mathbb{R}^s$, $oldsymbol{z} = f(oldsymbol{y}) \in \mathbb{R}^t$,则

$$rac{\partial oldsymbol{z}}{\partial oldsymbol{x}} = rac{\partial oldsymbol{y}}{\partial oldsymbol{x}} rac{\partial oldsymbol{z}}{\partial oldsymbol{y}} \quad \in \mathbb{R}^{p imes t}.$$

(3) 若
$$X \in \mathbb{R}^{p \times q}$$
 为矩阵, $\mathbf{y} = g(X) \in \mathbb{R}^s$, $z = f(\mathbf{y}) \in \mathbb{R}$,则

$$\frac{\partial z}{\partial X_{ij}} = \frac{\partial \boldsymbol{y}}{\partial X_{ij}} \frac{\partial z}{\partial \boldsymbol{y}} \quad \in \mathbb{R}$$



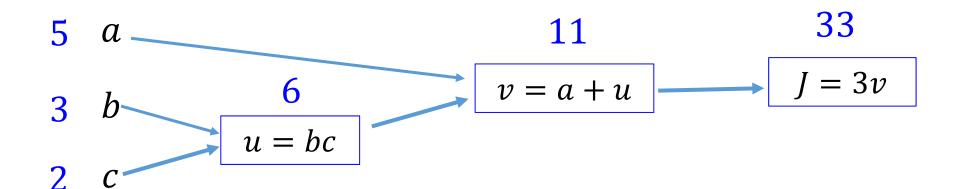
■ 梯度下降—计算图

$$J(a,b,c) = 3 (a + bc) = 3(5 + 3 \times 2) = 33$$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$





■ 梯度下降—计算图

$$y = f(x) = 3x, f'(x) = \frac{dy}{dx} = 3$$

5

$$\frac{dJ}{db} = \frac{dJ}{dv} \cdot \frac{dv}{du} \cdot \frac{du}{db} = 6$$

u = bc



$$\frac{dJ}{dc} = 9 \qquad \frac{dJ}{du} = \frac{dJ}{dv} \cdot \frac{dv}{du} = 3$$

 $b = 3 \rightarrow 3.001$

$$u = 6 \rightarrow bc = 6.002$$

$$v = 11 \rightarrow 11.002$$

$$J = 3v = 33 \rightarrow 33.006$$

$$\frac{du}{db} = 2, \frac{dv}{du} = 1, \frac{dJ}{db} = 6$$

 $u = 6 \to 6.001$

 $v = 11 \rightarrow 11.001$

 $J = 3v = 33 \rightarrow 33.003$

 $\frac{dv}{du} = 1, \qquad \frac{dJ}{dv} = 3$

11

v = a + u

33

J = 3v

$$\frac{dJ}{dv} = 3$$

$$v = 11 \rightarrow 11.001$$

$$J = 3v = 33 \rightarrow 33.003$$

$$\frac{dJ}{dv} = 3$$

$$a = 5 \to 5.001$$

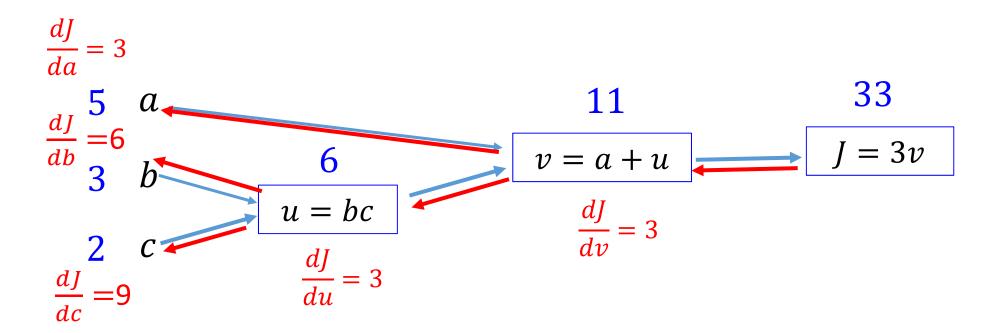
$$v = 11 \rightarrow 11.001$$

$$J = 3v = 33 \rightarrow 33.003$$

$$\frac{dv}{da} = 1, \qquad \frac{dJ}{da} = 3$$



■ 梯度下降—计算图





■ 梯度下降—计算图

- **口**逻辑回归算法: $\hat{y} = f(w^T x + b)$, 其中 $z = w^T x + b$, $f(z) = \frac{1}{1 + e^{-z}}$
- □损失函数:

$$L(\hat{y}^{(i)}, y^{(i)}) = -y^{(i)}log(\hat{y}^{(i)}) - (1 - y^{(i)})log(1 - \hat{y}^{(i)})$$

□代价函数:

$$J(w,b) = \frac{1}{m} \sum_{i=1}^{m} L(\hat{y}^{(i)}, y^{(i)})$$

□梯度下降法参数更新:

$$w \coloneqq w - \eta \frac{\partial L(w,b)}{\partial w}, b \coloneqq b - \eta \frac{\partial L(w,b)}{\partial b}$$



- 梯度下降—计算图
 - 假设单个样本,样本只有两个特征 x_1 和 x_2 ;参数 w_1 、 w_2 和b:

$$z = w_1 x_1 + w_2 x_2 + b$$
$$\hat{y} = f(z)$$

□损失函数:

$$L(\hat{y}, y) = -y\log(\hat{y}) - (1 - y)\log(1 - \hat{y})$$



■ 梯度下降—计算图

$$z = w_1 x_1 + w_2 x_2 + b$$

$$\hat{y} = f(z) = \frac{1}{1 + e^{-z}}$$

$$L(\hat{y}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

$$f'(z) = \frac{e^{-z}}{(1 + e^{-z})^2} = f(z)(1 - fz)$$



■ 梯度下降—计算图

$$x_1$$

$$w_1$$

$$dL(\hat{y},y))$$

$$w_2$$

$$dL(\hat{y},y))$$

$$dz$$

$$dL(\hat{y},y)$$

$$dz$$

$$dL(\hat{y},y))$$

$$dz$$

$$dL(\hat{y},y)$$

$$dL(\hat{$$



■ 反向传播 (Back Propagation)

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

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial w_{ij}^{(l)}} = \frac{\partial \mathbf{z}^{(l)}}{\partial w_{ij}^{(l)}} \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}}
\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{b}^{(l)}} = \frac{\partial \mathbf{z}^{(l)}}{\partial \mathbf{b}^{(l)}} \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}}$$



■ 反向传播 (Back Propagation)

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

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial w_{ij}^{(l)}} = \frac{\partial \mathbf{z}^{(l)}}{\partial w_{ij}^{(l)}} \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}}
\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{b}^{(l)}} = \frac{\partial \mathbf{z}^{(l)}}{\partial \mathbf{b}^{(l)}} \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}}$$

误差项

$$\delta^{(l)} = \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}} \in \mathbb{R}^{m^{(l)}}$$



■ 反向传播 (Back Propagation)

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

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial w_{ij}^{(l)}} = \frac{\partial \mathbf{z}^{(l)}}{\partial w_{ij}^{(l)}} \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}}
\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{b}^{(l)}} = \frac{\partial \mathbf{z}^{(l)}}{\partial \mathbf{b}^{(l)}} \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}}$$

误差项

$$\delta^{(l)} = \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}} \in \mathbb{R}^{m^{(l)}}$$

权重梯度

$$\frac{\partial \mathbf{z}^{(l)}}{\partial w_{ij}^{(l)}} = \left[\frac{\partial z_{1}^{(l)}}{\partial w_{ij}^{(l)}}, \cdots, \frac{\partial z_{i}^{(l)}}{\partial w_{ij}^{(l)}}, \cdots, \frac{\partial z_{m^{(l)}}^{(l)}}{\partial w_{ij}^{(l)}} \right] \\
= \left[0, \cdots, \frac{\partial (\mathbf{w}_{i:}^{(l)} \mathbf{a}^{(l-1)} + b_{i}^{(l)})}{\partial w_{ij}^{(l)}}, \cdots, 0 \right] \\
= \left[0, \cdots, a_{j}^{(l-1)}, \cdots, 0 \right] \\
\triangleq \mathbb{I}_{i}(a_{j}^{(l-1)}) \in \mathbb{R}^{m^{(l)}},$$



■ 反向传播 (Back Propagation)

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

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial w_{ij}^{(l)}} = \frac{\partial \mathbf{z}^{(l)}}{\partial w_{ij}^{(l)}} \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}}$$
$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{b}^{(l)}} = \frac{\partial \mathbf{z}^{(l)}}{\partial \mathbf{b}^{(l)}} \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}}$$

误差项

$$\delta^{(l)} = \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}} \in \mathbb{R}^{m^{(l)}}$$

权重梯度

$$\frac{\partial \mathbf{z}^{(l)}}{\partial w_{ij}^{(l)}} = \left[\frac{\partial z_{1}^{(l)}}{\partial w_{ij}^{(l)}}, \cdots, \frac{\partial z_{i}^{(l)}}{\partial w_{ij}^{(l)}}, \cdots, \frac{\partial z_{m^{(l)}}^{(l)}}{\partial w_{ij}^{(l)}} \right] \\
= \left[0, \cdots, \frac{\partial (\mathbf{w}_{i:}^{(l)} \mathbf{a}^{(l-1)} + b_{i}^{(l)})}{\partial w_{ij}^{(l)}}, \cdots, 0 \right] \\
= \left[0, \cdots, a_{j}^{(l-1)}, \cdots, 0 \right] \\
\triangleq \mathbb{I}_{i}(a_{j}^{(l-1)}) \in \mathbb{R}^{m^{(l)}},$$

偏置梯度

$$\frac{\partial \mathbf{z}^{(l)}}{\partial \mathbf{b}^{(l)}} = \mathbf{I}_{m^{(l)}} \quad \in \mathbb{R}^{m^{(l)} \times m^{(l)}}$$



■ 反向传播 (Back Propagation)

神经网络结构

$$\mathbf{a}^{(l)} = f_l(\mathbf{z}^{(l)})$$

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

逐层反向传播

$$\frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{z}^{(l)}} = \frac{\partial f_l(\mathbf{z}^{(l)})}{\partial \mathbf{z}^{(l)}} \\
= \operatorname{diag}(f'_l(\mathbf{z}^{(l)}))$$

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

$$\frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{z}^{(l)}} = \frac{\partial f_l(\mathbf{z}^{(l)})}{\partial \mathbf{z}^{(l)}} \\
= \operatorname{diag}(f'_l(\mathbf{z}^{(l)})) \\
= \frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{z}^{(l)}} \cdot \frac{\partial \mathbf{z}^{(l+1)}}{\partial \mathbf{z}^{(l)}} \cdot \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l+1)}} \\
= \frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{z}^{(l)}} \cdot \frac{\partial \mathbf{z}^{(l+1)}}{\partial \mathbf{a}^{(l)}} \cdot \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l+1)}} \\
= \operatorname{diag}(f'_l(\mathbf{z}^{(l)})) \cdot (W^{(l+1)})^{\mathrm{T}} \cdot \delta^{(l+1)} \\
= f'_l(\mathbf{z}^{(l)}) \odot ((W^{(l+1)})^{\mathrm{T}} \delta^{(l+1)}),$$



■ 反向传播 (Back Propagation)

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial w_{ij}^{(l)}} = \frac{\partial \mathbf{z}^{(l)}}{\partial w_{ij}^{(l)}} \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{z}^{(l)}} = \mathbb{I}_i(a_j^{(l-1)}) \delta^{(l)} = \delta_i^{(l)} a_j^{(l-1)}$$

 \square 进一步, $\mathcal{L}(\mathbf{y},\hat{\mathbf{y}})$ 关于第 l 层权重 $\mathbf{W}^{(l)}$ 的梯度为

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial W^{(l)}} = \delta^{(l)} (\mathbf{a}^{(l-1)})^{\mathrm{T}}$$

□ 同理, $\mathcal{L}(\mathbf{y},\hat{\mathbf{y}})$ 关于第 l 层偏置 $\mathbf{b}^{(l)}$ 的梯度为

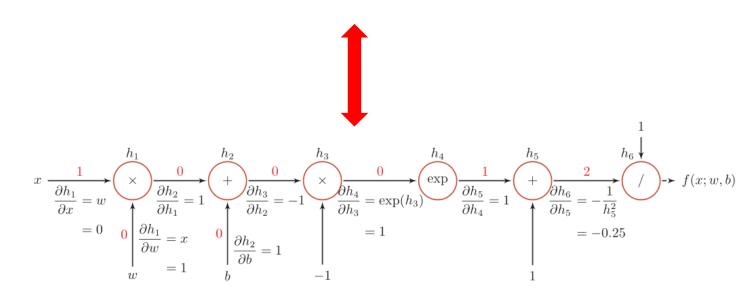
$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{b}^{(l)}} = \delta^{(l)}$$



■ 自动微分 (Automatic Differentiation)

自动微分是利用链式法则来自动计算一个复合函数的梯度。

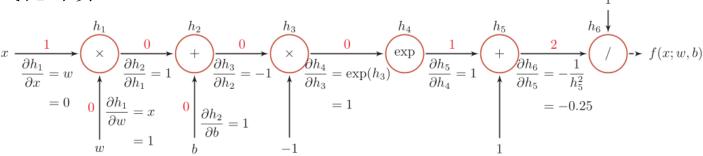
$$f(x; w, b) = \frac{1}{\exp(-(wx + b)) + 1}.$$





■ 自动微分 (Automatic Differentiation)

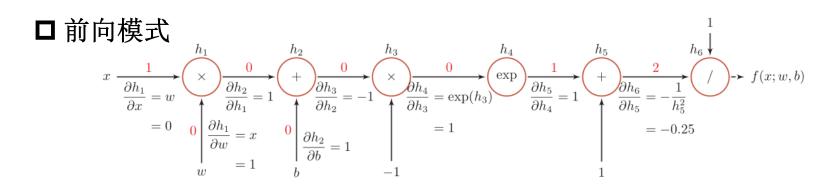
□ 形式化计算



函数	导数	
$h_1 = x \times w$	∂w	$\frac{\partial h_1}{\partial x} = w$
$h_2 = h_1 + b$	∂h_1	$\frac{\partial h_2}{\partial b} = 1$
$h_3 = h_2 \times -$	∂h_2	$2f(\dots, h)$ $2f(\dots, h)$ $2h$ $2h$ $2h$ $2h$ $2h$ $2h$ $2h$
$h_4 = \exp(h$	$\frac{\partial h_4}{\partial h_3} = \exp(h_3)$	$\frac{\partial f(x; w, b)}{\partial w} _{x=1, w=0, b=0} = \frac{\partial f(x; w, b)}{\partial h_6} \frac{\partial h_6}{\partial h_5} \frac{\partial h_5}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w}$
$h_5 = h_4 + 1$	$\frac{\partial h_5}{\partial h_4} = 1$	$= 1 \times -0.25 \times 1 \times 1 \times -1 \times 1 \times 1$
$h_6 = 1/h_5$	$rac{\partial h_6}{\partial h_5} = -rac{1}{h_5^2}$	= 0.25.



■ 自动微分 (Automatic Differentiation)



□ 反向模式(与反向传播的计算梯度的方式相同)

前向计算每一层的状态和激活值,直到最后一层

反向计算每一层的参数的偏导数

更新参数

□ 如果函数和参数之间有多条路径,可以将这多条路径上的导数 再进行相加,得到最终的梯度。



■ 自动微分 (Automatic Differentiation)

□ 静态计算图

在编译时构建计算图,计算图构建好之后在程序运行时不能改变。 代表性工具箱Theano和Tensorflow

□ 动态计算图

在程序运行时动态构建,可自动根据函数的结构调整计算图。代表性工具箱DyNet,Chainer和PyTorch

□ 优缺点

静态计算图在构建时可以进行优化,并行能力强,但灵活性比较差低。动态计算图则不容易优化,当不同输入的网络结构不一致时,难以并行计算,但是灵活性比较高。



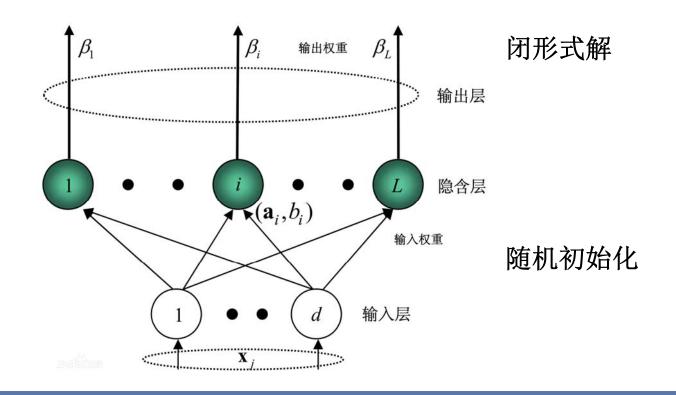


Close-form Solution



■ 极端学习机 (Extreme Learning Machine)

快速学习算法,对于单隐层神经网络(三层),ELM可以随机初始化输入权重和偏置,并利用最小二乘得到相应的输出权重。





■ 极端学习机 (Extreme Learning Machine)

□ 在一个 ELM 的单隐层结构中, 第 i 个隐藏节点的输出为:

$$h_i(\mathbf{x}) = G(\mathbf{a}_i, \mathbf{b}_i, \mathbf{x})$$

□ 整个隐藏层的输出映射为:

$$\mathbf{h}(\mathbf{x}) = [G(h_1(\mathbf{x})), G(h_2(\mathbf{x})), \dots, G(h_L(\mathbf{x}))]$$

□ 给定 N 个训练样本, ELM 的隐藏层输出矩阵 H 给出为:

$$\mathbf{H} = \begin{bmatrix} \mathbf{h}(\mathbf{x}_1) \\ \vdots \\ \mathbf{h}(\mathbf{x}_N) \end{bmatrix} = \begin{bmatrix} G(\mathbf{a}_1, \mathbf{b}_1, \mathbf{x}_1) & \cdots & G(\mathbf{a}_L, \mathbf{b}_L, \mathbf{x}_1) \\ \vdots & \vdots & \vdots \\ G(\mathbf{a}_1, \mathbf{b}_1, \mathbf{x}_N) & \cdots & G(\mathbf{a}_L, \mathbf{b}_L, \mathbf{x}_N) \end{bmatrix}$$

□ 具有 L 个隐藏节点的 ELM 结构的输出为:

$$f_L(\mathbf{x}) = \sum_{i=1}^{L} \beta_i \mathbf{h}_i(\mathbf{x})$$



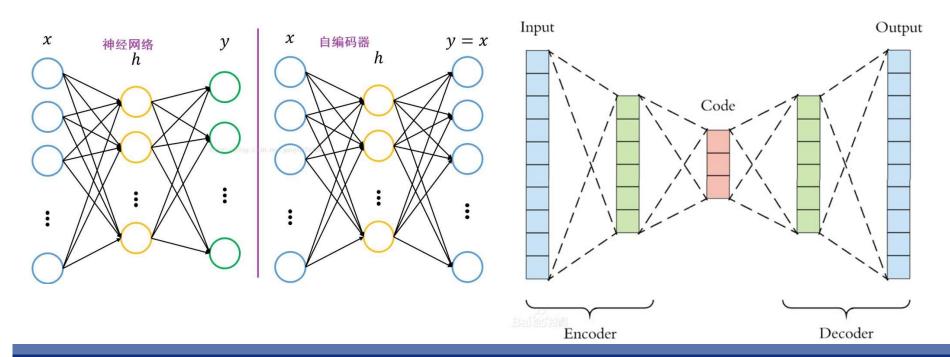
- 极端学习机 (Extreme Learning Machine)
 - □ 整个隐藏层的输出为: **H**β
 - □ 目标矩阵 T (标签矩阵)为: $\mathbf{T} = [\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_N]^T$
 - □ 目标函数为: $Minimize: ||\beta||_p + C||\mathbf{H}\beta \mathbf{T}||_q$
 - **□** 求解矩阵H的Moore-Penrose广义逆: $\hat{\beta} = \mathbf{H}^{+}\mathbf{T}$

闭形式解!



■ 自编码器 (Auto-Encoder)

- □ 自编码器是一类在半监督学习和无监督学习中使用的人工神经网络, 其功能是通过将输入信息作为学习目标,对输入信息进行表征学习
- □ 自编码器包含编码器(Encoder)和解码器(Decoder)两部分





■ 自编码器 (Auto-Encoder)

- □ 按学习范式,自编码器可以被分为收缩自编码器、正则自编码器和变分自编码器,其中前两者是判别模型、后者是生成模型
- □ 按构建类型,自编码器可以是前馈结构或递归结构的神经网络
- □ 自编码器具有一般意义上表征学习算法的功能,被应用于降维和异常 值检测
- □ 包含卷积层结构的自编码器可被应用于计算机视觉问题,包括图像降 噪、图像风格迁移等



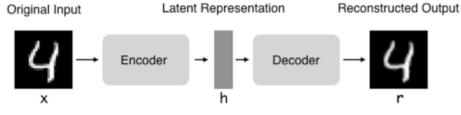
■ 自编码器 (Auto-Encoder)

□ 自编码器是一个输入和输出相同的神经网络,给定输入空间和特征空间, 自编码器求解两者的映射使输入特征的重建误差达成最小:

$$y = h(x)$$

 $r = f(y) = f(h(x))$
重建误差最小

□ 求解完成后,由编码器输出的隐含层特征,即"编码特征"可视为输入数据的表征 Original Input Latent Representation Reconstructed Output

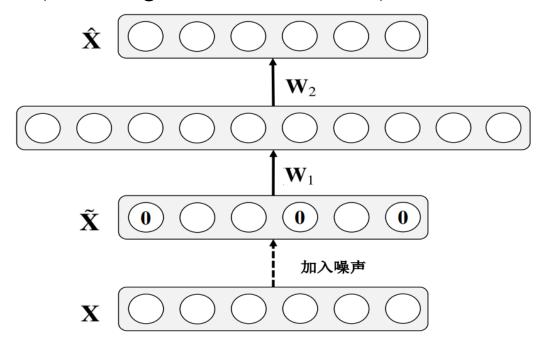


Architecture of an Autoencode

□ 当自编码器的编码、解码能力很强时,可以确保对训练样本的完美重建,但却难以发现数据的内在分布和结构规律



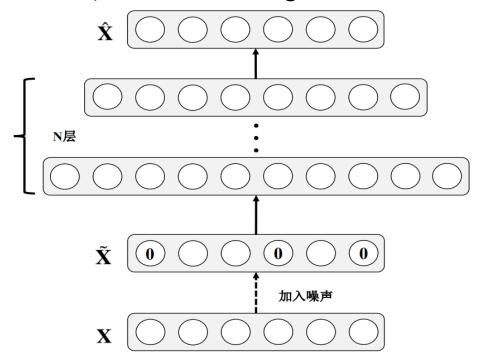
- 自编码器 (Auto-Encoder)
 - 去噪自编码器(Denoising Auto-Encoder, DAE)



通常还可约束编码部分和解码部分的权重保持某种关系,比如捆绑约束(Tied Constraint)、稀疏约束(Sparse Constraint)等。



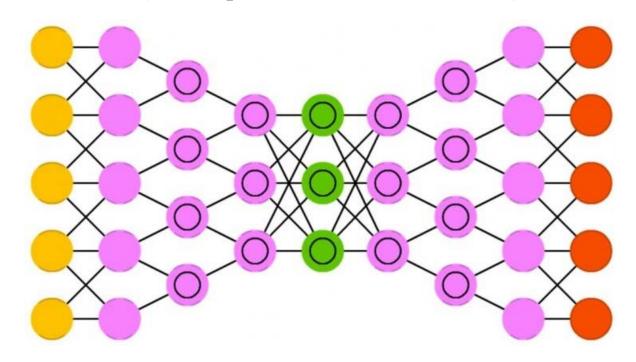
- 自编码器 (Auto-Encoder)
 - □ 堆叠去噪自编码器(Stacked Denoising Auto-Encoder, DAE)



引入了更多的隐藏层,反复堆叠,从而获得了更多输入数据的深层次和抽象化特征。



- 自编码器 (Auto-Encoder)
 - □ 欠完备自编码器(Uncomplete Auto-Encoders, UAE)



编码维度小于输入维度的自编码器。学习欠完备的表示将强制自编码器捕捉训练数据中最显著的特征。

小结



■ 多层感知机

- □ 单层感知机(连接方式,激活函数)
- □ 多层感知机(等价形态,多层扩展)

■ 前馈神经网络

- □ 人工神经网络
- □ 前馈神经网络(连接模式,数据传递)
- □ 通用近似定理
- □ 分类问题

■ 梯度反向传播

- □ 参数学习
- □ 梯度下降(链式法则,反向传播,自动微分)

■ 闭形式解

- □ 极端学习机
- □ 自编码器(降噪自编码器,堆叠降噪自编码器,欠完备降噪自编码器)



```
class FCNet(nn.Module):
    def init (self, input size, hidden size, output size):
        super(FCNet, self). init ()
        self.fc1 = nn.Linear(input size, hidden size)
        self.relu = nn.ReLU()
        self.fc2 = nn.Linear(hidden size, output size)
                                class student(object):
                                    def init (self, name, age, score):
    def forward(self, x):
                                        self.name = name
        out = self.fc1(x)
                                        self.age = age
        out = self.relu(out)
                                        self.score = score
        out = self.fc2(out)
        return out
                                    def print age(self):
                                        print("%s's name: %s" % (self.name, self.age))
                                    def print score(self):
                                        print("%s's age: %s" % (self.name, self.score))
                                stu1 = student('Amy', 13, 86)
```

stu2 = student('Tom', 15, 79)



```
class FCNet(nn.Module):
    def __init__(self, input_size, hidden_size, output_size):
        super(FCNet, self).__init__()
        self.fc1 = nn.Linear(input_size, hidden_size)
        self.relu = nn.ReLU()
        self.fc2 = nn.Linear(hidden_size, output_size)
    def forward(self, x):
        out = self.fc1(x)
        out = self.relu(out)
        out = self.fc2(out)
        return out
```



```
class MyTestData(Dataset):
   def init (self, params):
       super(MyTestData, self). init ()
       self.test_im_path = params['test_im_path']
       self.test lb path = params['test lb path']
       self.test im num = 10000
       self.test labels = open(self.test lb path, 'r').read().splitlines()
   def len (self):
       return self.test im num
   def getitem (self, index):
       # load image
       img file = os.path.join(self.test im path, str(index) + '.png')
       img = Image.open(img file)
       im = self.transform(img)
       lb = int(self.test labels[index])
       return im, 1b
   def transform(self, img):
       transform img = transforms.Compose([transforms.Resize((28, 28)),
                                           transforms.ToTensor()])
       im = transform_img(img)
       return im
```

初始化 参数设置

读入图像和标签

对图像进行变换 (转化为tensor)



基于今许控网络的王军粉宁识别



device = torch.device('cuda' if torch.cuda.is available() else 'cpu')

```
# 参数设置
params = \{\}
params['train im path'] = 'D:\\myproject\\mnist pytorch\\mnist\\train-images\\'
params['test im path'] = 'D:\\myproject\\mnist pytorch\\mnist\\t10k-images\\'
params['train lb path'] = 'D:\\myproject\\mnist pytorch\\mnist\\train-labels.txt'
params['test lb path'] = 'D:\\myproject\\mnist pytorch\\mnist\\t10k-labels.txt'
params['INPUT SIZE'] = 784
params['HIDDEN SIZE'] = 500
params['OUTPUT SIZE'] = 10
params['NUM_EPOCHES'] = 10
params['BATCH SIZE'] = 100
                                                                        参数设置
params['LEARNING RATE'] = 0.001
# load dataset
                                                            加载训练数据和测试数据
train loader = DataLoader(MyData(params),
                     shuffle=True,
                     batch_size=params['BATCH_SIZE'])
test loader = DataLoader(MyTestData(params),
                        shuffle=False,
                        batch_size=1)
                                                                       定义网络
model = FCNet(params['INPUT SIZE'], params['HIDDEN SIZE'], params['OUTPUT SIZE']).to(device)
                                                      定义损失函数(交叉熵损失函数)
criterion = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(model.parameters(), lr=params['LEARNING RATE'])
```

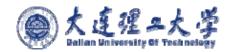
定义优化方法(ADAM)



■ 网络训练

```
for epoch in range(params['NUM_EPOCHES']):
   for i, (images, labels) in enumerate(train_loader):
       # Move tensors to the configured device
       images = images.reshape(-1, 28*28).cuda(device=0) # # -1 是指模糊控制的意思,即固定784列,不知道多少行
       labels = labels.cuda(device=0)
       # images = images.reshape(-1, 28 * 28).to(device) # # -1 是指模糊控制的意思,即固定784列,不知道多少行
       # labels = labels.to(device)
       # 前向传播
       outputs = model(images)
       loss = criterion(outputs, labels)
       # 反向传播和优化
       optimizer.zero_grad()
       loss.backward()
       optimizer.step()
       if (i+1) \% 100 == 0:
           print('Epoch [{}/{}], Step [{}/{}], Loss: {:.4f}'
                 .format(epoch+1, params['NUM_EPOCHES'], i+1, total_step, loss.item()))
           losses.append(loss.item())
```





```
测试模型
# 在测试阶段,不用计算梯度
with torch.no grad():
   correct = 0
   total = 0
   model = model.eval()
   for images, labels in test loader:
      images = images.reshape(-1, 28 * 28).to(device)
      labels = labels.to(device)
      outputs = model(images)
      , predict = torch.max(outputs.data, 1)
      ##这里返回两组数据,最大image_data和最大值索引,可以用torch.argmax()更为直观,这里去理解其
      ## 这个 _ , predicted是python的一种常用的写法,表示后面的函数其实会返回两个值
      # 但是我们对第一个值不感兴趣,就写个_在那里,把它赋值给_就好,我们只关心第二个值predicted
      # 比如 _ ,a = 1,2 这中赋值语句在python中是可以通过的,你只关心后面的等式中的第二个位置的值是
      total += labels.size(0) ##更新测试图片的数量 size(0),返回行数
      correct += (predict == labels).sum().item()
   accuracy = correct / total
   accs.append(accuracy)
   print('Accuracy of the network on the 10000 test images: {} %'.format(100 * accuracy))
```



■ 保存网络模型

保存模型 torch.save(model.state_dict(), 'model.pth')

基于全连接

```
Epoch [1/10], Step [100/600], Loss: 0.6180
Epoch [1/10], Step [200/600], Loss: 0.2788
Epoch [1/10], Step [300/600], Loss: 0.3504
Epoch [1/10], Step [400/600], Loss: 0.2611
Epoch [1/10], Step [500/600], Loss: 0.3334
Epoch [1/10], Step [600/600], Loss: 0.1110
Accuracy of the network on the 10000 test images: 93.2 %
Epoch [2/10], Step [100/600], Loss: 0.2007
Epoch [2/10], Step [200/600], Loss: 0.1345
Epoch [2/10], Step [300/600], Loss: 0.1862
Epoch [2/10], Step [400/600], Loss: 0.2383
Epoch [2/10], Step [500/600], Loss: 0.1317
Epoch [2/10], Step [600/600], Loss: 0.1385
Accuracy of the network on the 10000 test images: 95.5 %
Epoch [3/10], Step [100/600], Loss: 0.0772
Epoch [3/10], Step [200/600], Loss: 0.1712
Epoch [3/10], Step [300/600], Loss: 0.1935
Epoch [3/10], Step [400/600], Loss: 0.2046
Epoch [3/10], Step [500/600], Loss: 0.0839
Epoch [3/10], Step [600/600], Loss: 0.2274
Accuracy of the network on the 10000 test images: 95.67999999999999 %
Epoch [4/10], Step [100/600], Loss: 0.1628
Epoch [4/10], Step [200/600], Loss: 0.0509
Epoch [4/10], Step [300/600], Loss: 0.1581
Epoch [4/10], Step [400/600], Loss: 0.0813
Epoch [4/10], Step [500/600], Loss: 0.1619
Epoch [4/10], Step [600/600], Loss: 0.0967
Accuracy of the network on the 10000 test images: 96.12 %
Epoch [5/10], Step [100/600], Loss: 0.0991
Epoch [5/10], Step [200/600], Loss: 0.0900
Epoch [5/10], Step [300/600], Loss: 0.1194
Epoch [5/10], Step [400/600], Loss: 0.1203
Epoch [5/10], Step [500/600], Loss: 0.0709
Epoch [5/10], Step [600/600], Loss: 0.1261
Accuracy of the network on the 10000 test images: 96.72 %
Epoch [6/10], Step [100/600], Loss: 0.0240
Epoch [6/10], Step [200/600], Loss: 0.0739
Epoch [6/10], Step [300/600], Loss: 0.0865
Epoch [6/10], Step [400/600], Loss: 0.0295
Epoch [6/10], Step [500/600], Loss: 0.0645
```

Epoch [6/10], Step [600/600], Loss: 0.1076







```
# plot loss figure and accuracy figure
plt.figure()
plt.title('loss')
plt.plot(losses)
plt.figure()
plt.title('accuracy')
plt.plot(accs)
plt.show()
torch.save(model.state_dict(), 'model.pth')
```



