

分类模型

第二部分 非线性分类模型

BP神经网络

张朝晖

2018-2019学年 20181009



1. 分类模块

产生式分类模型

A. 贝叶斯分类模型

判别式分类模型

线性分类模型

- B. Fisher 判别分类
- C. 感知器分类模型
- D. 大间隔分类模型 (线性 SVM)

非线性分类模型

- E. 核 SVM (非线性 SVM)
- F. 核 Fisher 判别分类
- G. 神经网络

其它分类模型

- H. KNN 分类模型
- I. 决策树分类模型
- J. Logistic 回归
- K. Softmax 回归

2. 聚类模块

- L. K-均值聚类
- M. 高斯混合聚类
- N. DBSCAN 聚类
- O. 层次聚类

3. 回归模块

- P. KNN 回归
- Q. 回归树
- R. 最小二乘线性回归
- S. 岭回归
- T. LASSO 回归

4. 集成学习

- U. Bagging
- V. 随机森林
- W. Boosting

5. 特征工程

- X. 主成分分析 (PCA)
- ...

6. 评价模块

混淆矩阵 (及其相关指标)、ROC 曲线、交叉验证

人工神经网络

--前馈神经网络与BP算法

关键词：

人工神经网络

前馈神经网络、人工神经元、激活函数

BP神经网络：梯度下降法、导数链式法则、误差反向传播

自编码器

监督式学习、非监督式学习

分类、回归、表示学习

多层BP神经网络学习

第一阶段：基于逐层表示学习的参数初始化；

第二阶段：基于BP算法的参数精调

主要内容

1. 人工神经网络基本知识、神经元与感知器

生物神经网络、生物神经元

人工神经网络、人工神经网络的基本模型、人工神经元

2. 前馈神经网络、多层感知器、及非线性分类

单层感知器、多层感知器、基于多层前馈神经网络的函数逼近

3. BP神经网络(可用于分类、回归、监督式特征学习)

监督式学习+BP算法(梯度下降法+导数链式法则)

3.1 BP网络的参数学习

3.2 BP网络的应用

4. 基于前馈神经网络的自编码器(Autoencoders)

逐层学习、非监督式学习+BP算法

已知结构的前馈神经网络的学习

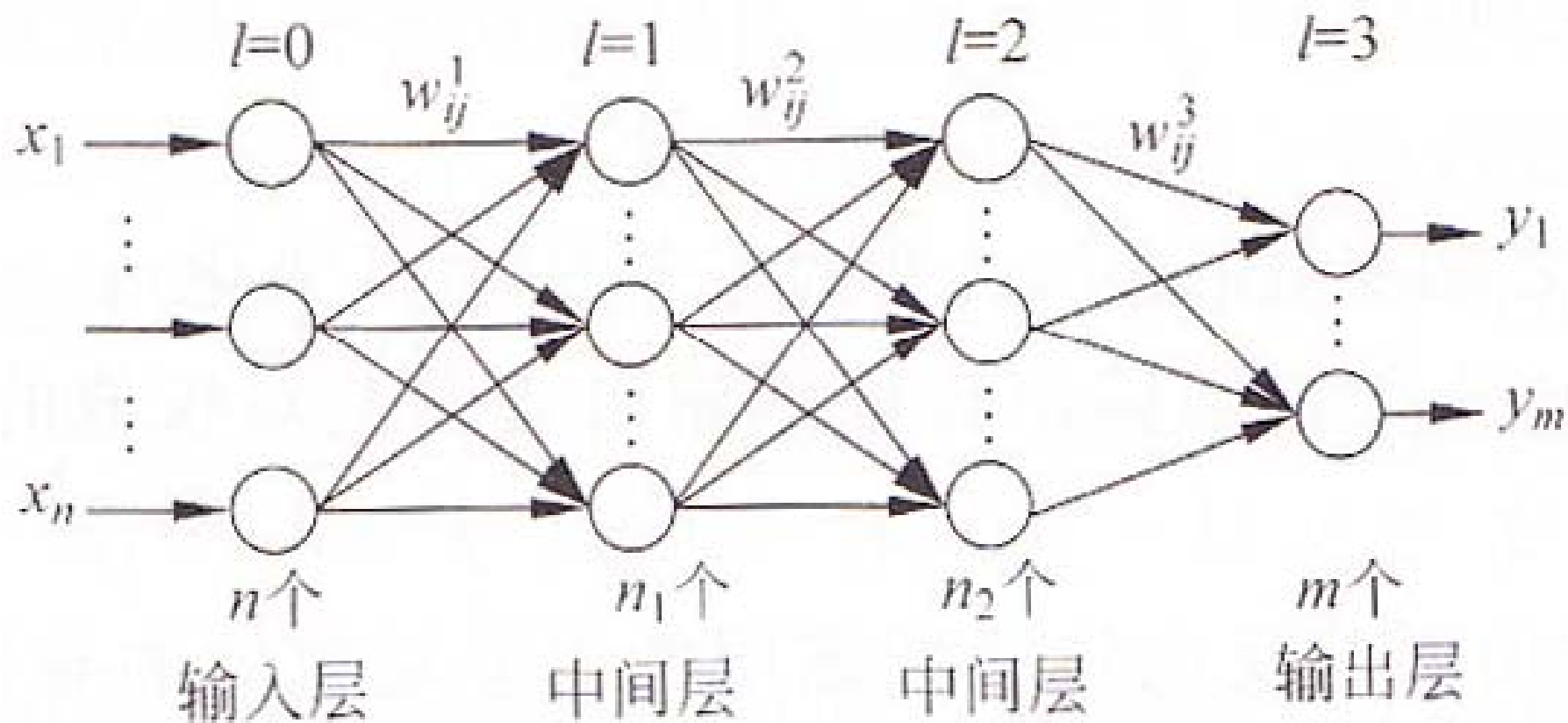
利用给定训练样本集，估计相邻层节点之间连接参数、各神经元节点的偏置量。

基于阈值神经元的前馈神经网络的不足

- 隐含层不直接与外界连接，误差无法直接估计
- 神经元节点的激活函数为阈值函数(或阶跃函数)时，无法采用梯度下降法学习神经网络的有关参数。

BP网络：基于BP(Backward Propagation)算法 进行参数学习的前馈神经网络

- 神经元节点传递函数：以连续(可微)函数代替阈值函数
- BP网络的学习：监督式学习
- 目标函数的构造：基于训练样本集的预测损失函数
- 目标函数的参数寻优：
BP算法=梯度下降法+导数链式法则
- BP网络的用途：实值函数的回归、分类、(特征提取)



多层感知器的例子和符号约定

神经网络模型训练目的在于：

基于给定的样本集 $X = \{(x_i, D_i), i = 1, \dots, N\}$

$$\text{学习参数} \left\{ \begin{array}{ll} \text{网络连接权} & \mathbf{W}_{n_{l-1} \times n_l}^l = [\omega_{ij}^l]_{n_{l-1} \times n_l} \\ \text{神经元节点偏置量} & \mathbf{b}^l = [b_j^l]_{1 \times n_l} \\ l = 1, \dots, L-1 & i = 1, \dots, n_{l-1} \quad j = 1, \dots, n_l \end{array} \right.$$

$$\Rightarrow \left\{ \begin{array}{l} \mathbf{W} = \left\{ \mathbf{W}_{n_{l-1} \times n_l}^l = [\omega_{ij}^l]_{n_{l-1} \times n_l}, l = 2, \dots, L-1 \right\} \\ \mathbf{b} = \left\{ \mathbf{b}^l = [b_j^l]_{1 \times n_l}, l = 2, \dots, L-1 \right\} \end{array} \right.$$

目标函数的构造：

$$E(W, b) = L(W, b) + \lambda \Omega(W)$$

$$\begin{cases} L(W, b) & \text{--关于训练样本集的学习能力} \\ \Omega(W) & \text{--正则项} \end{cases}$$

例: $L(\mathbf{W}, \mathbf{b})$ 的形式

T

(1) 预测残差的平均平方和(MSE)——回归, 分类:

$$L(\mathbf{W}, \mathbf{b}) = \frac{1}{2N} \sum_{i=1}^N \|Y(\mathbf{x}_i; \mathbf{W}, \mathbf{b}) - \mathbf{D}_i\|^2$$

(2) 交叉熵(CrossEntropy)——两类别分类:

$$L(\mathbf{W}, \mathbf{b}) = -\frac{1}{N} \sum_{i=1}^N [d_i \ln y_i + (1 - d_i) \ln (1 - y_i)]$$

输出层只含一个输出节点, 目标输出值 $d_i \in \{0, 1\}$, 预测输出 $y_i \in (0, 1)$

(3) 交叉熵(CrossEntropy)——C类别分类: $L(\mathbf{W}, \mathbf{b}) = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C (d_{ij} \ln y_{ij})$

$\mathbf{Y}_i = [y_{i1} \ \cdots \ y_{ic}]^T$ 为输出层关于样本 \mathbf{x}_i 的C类预测输出(概率).

$\mathbf{D}_i = [d_{i1} \ \cdots \ d_{ic}]^T$ 为输出层关于样本 \mathbf{x}_i 的C类目标输出(概率).

$d_{ij} \in \{0, 1\}$, $y_{ij} \in (0, 1)$

若训练样本 \mathbf{x}_i 来自第1类, 则 $\mathbf{D}_i = [1 \ 0 \ \cdots \ 0]$

$\Omega(\mathbf{W})$ 的形式: $\Omega(\mathbf{W}) = \|\mathbf{W}\|_F^2 = \sum_{l,i,j} (\omega_{ij}^l)^2$

$$E(W, b) = L(W, b) + \lambda \Omega(W)$$

神经网络的学习： $[\hat{W}, \hat{b}] = \arg \min_{W, b} E(W, b | X)$

梯度下降法的参数迭代估计 ($i = 1, \dots, n_{l-1} \quad j = 1, \dots, n_l$)

$$\begin{cases} \omega_{ij}^l(t+1) = \omega_{ij}^l(t) + \Delta \omega_{ij}^l(t) & \Delta \omega_{ij}^l(t) = -\eta \frac{\partial E}{\partial \omega_{ij}^l(t)} \\ b_j^l(t+1) = b_j^l(t) + \Delta b_j^l(t) & \Delta b_j^l(t) = -\eta \frac{\partial E}{\partial b_j^l(t)} \end{cases}$$

梯度下降法

{ 随机梯度下降法 (*Stochastic Gradient Descent*)

-- (也称: 单样本法、在线梯度下降法)

{ 小批量样本法 (*Mini-Batch Gradient Descent*)

{ 批量样本法 (*Batch Gradient Descent*)

{ 固定学习步长的梯度下降

{ 变学习步长的梯度下降

{ 普通梯度下降法 (不带冲量)

{ 带**冲量** (惯性项, *Momentum*)

模型参数初始化 { 参数取值应具有差异性 (高斯分布、均匀分布)
取值不能太大

给定训练样本集 $\{(x_i, D_i), i = 1, \dots, N\}$, BP
算法进行参数学习的两个阶段:

(1) 输入信号的正向传递过程

训练样本的输入部分从输入层经隐层逐层计算、正向传递, 直至到达输出层, 得到预测输出

(2) 预测误差的反向传播过程

预测输出与目标输出比较, 得到预测误差; 预测误差逐层、反向传播, 可间接计算隐层各单元的误差, 并用此误差修正网络有关参数。

BP算法训练过程描述

约定：

(1) n 维输入向量

$$\mathbf{x} = [x_1, \dots, x_n]^T$$

(2) L 层神经网络

层号 $l = 0$ 输入层

层号 $l = 1, \dots, L-2$ 隐含层

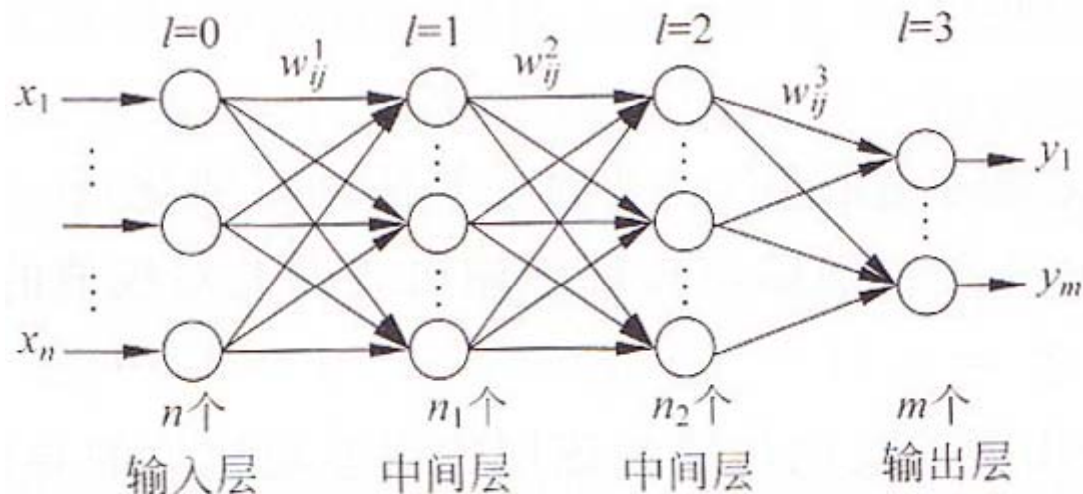
层号 $l = L-1$ 输出层

(3) 各层节点数目 n_l , $l = 0, 1, \dots, L-1$

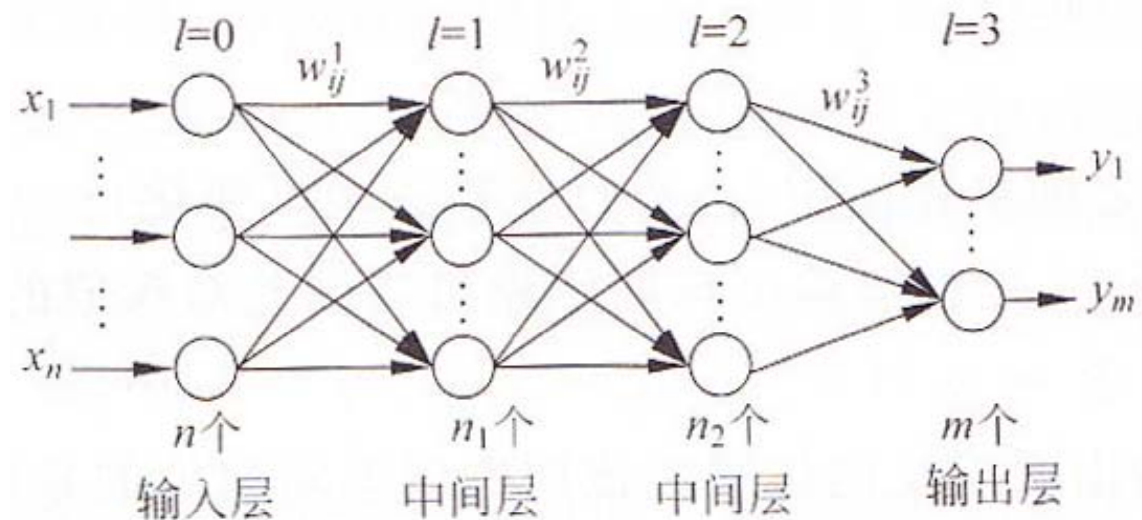
输入节点 $n_0 = n$

输出节点 $n_{L-1} = m$

(例: m 个实值函数的回归, m 个类别的分类)



多层感知器的例子和符号约定



多层感知器的例子和符号约定

(4) 相邻层连接权值

ω_{ij}^l —— 第 $l-1$ 层的节点 i 与第 l 层节点 j 的连接权值

$$\mathbf{W}_{n_{l-1} \times n_l}^l = \left[\omega_{ij}^l \right]_{n_{l-1} \times n_l} \quad \begin{cases} i = 1, \dots, n_{l-1} \\ j = 1, \dots, n_l \end{cases}$$

(5) 各计算节点处的偏置量

b_j^l —— 第 l 层节点 j 处的偏置量

$$\mathbf{b}^l = \left[b_j^l \right]_{1 \times n_l} \quad j = 1, \dots, n_l$$

BP算法训练过程描述

(6) 假定：第 l 层为**当前处理层**；

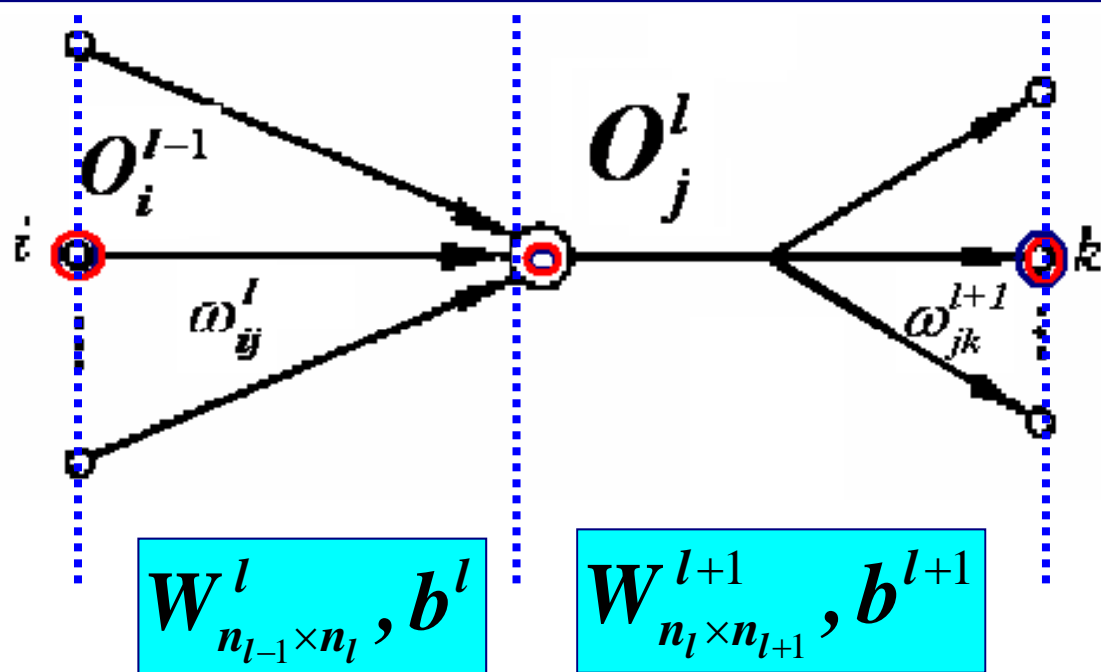
其前一层 $l-1$ 、当前层 l 、后一层 $l+1$ 的计算单元序号为 i, j, k ；

位于当前层第 j 个计算单元的预测输出为 O_j^l , $j = 1, \dots, n_l$

前层第 i 个单元到本层第 j 个单元的连接权值为 ω_{ij}^l , $i = 1, \dots, n_{l-1}$

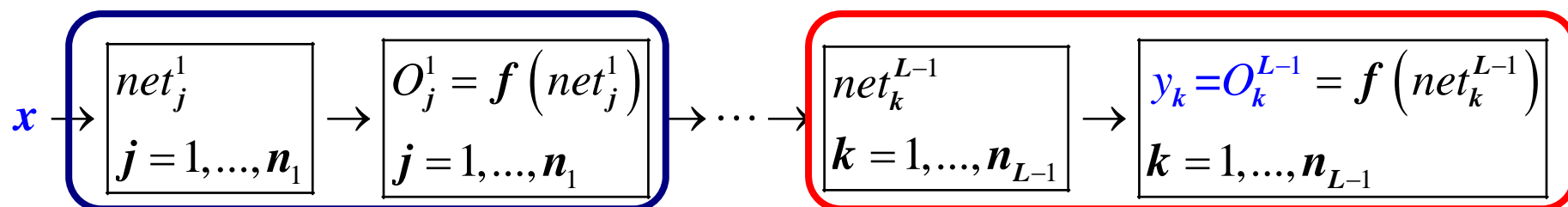
本层第 j 个单元到后层第 k 个单元的连接权值为 ω_{jk}^{l+1} , $k = 1, \dots, n_{l+1}$

注：采用**梯度法**修正权值，**计算单元输出函数**应连续可微，如：*sigmoid*函数。



第一阶段：输入信号的正向传递过程

固定步长、单样本的梯度下降法



从样本集内取出一个样本 (\mathbf{x}, \mathbf{D}) , 将 \mathbf{x} 各分量从输入层输入至网络, 由前向后逐层传递, **当前层** l 的第 j 个计算节点 ($l = 1, \dots, L-1$ $j = 1, \dots, n_l$)

$$\left\{ \begin{array}{l} \text{净输入} \quad net_j^l = \sum_{i=1}^{n_{l-1}} \omega_{ij}^l O_i^{l-1} + b_j^l \\ \text{实际输出} \quad O_j^l = f(net_j^l) = \frac{1}{1 + e^{-net_j^l}} \end{array} \right. \Rightarrow \left\{ \begin{array}{l} \text{矩阵形式:} \\ net^l = (W^l)^T O^{l-1} + b^l \\ O^l = f(net^l) \end{array} \right.$$

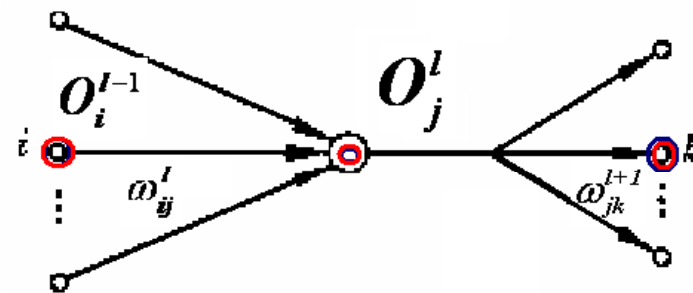
若当前层为输出层 ($l = L-1$), 则

$$\text{计算单元 } j \left\{ \begin{array}{l} \text{实际输出} \quad y_j = O_j^l \\ \text{理想输出} \quad d_j \end{array} \right. \quad \text{输出层} \left\{ \begin{array}{l} \mathbf{y} = (y_1, \dots, y_m)^T \\ \mathbf{D} = (d_1, \dots, d_m)^T \end{array} \right.$$

第二阶段：误差反向传播过程

准则函数 -- 最小误差平方和。

(I) 输出误差

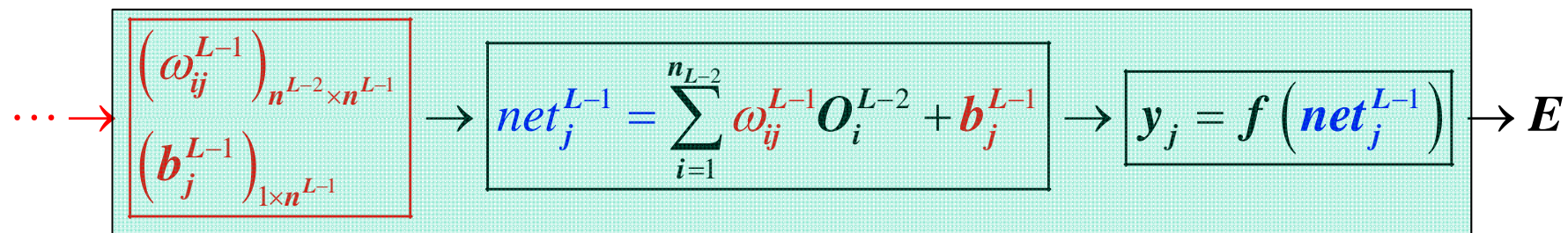


样本 $\mathbf{x} = (x_1, \dots, x_n)^T$ 在网络**输出层** (各节点) 因**输出** \mathbf{y} 与**目标****输出** \mathbf{D} 不一致而导致**输出误差**:

$$\begin{aligned} E &= \frac{1}{2} \|\mathbf{D} - \mathbf{y}\|^2 = \frac{1}{2} (\mathbf{D} - \mathbf{y})^T (\mathbf{D} - \mathbf{y}) = \frac{1}{2} (\mathbf{D}^T \mathbf{D} - 2\mathbf{y}^T \mathbf{D} + \mathbf{y}^T \mathbf{y}) \\ &= \frac{1}{2} \left\{ \mathbf{D}^T \mathbf{D} - 2 \left[f(\mathbf{net}^{L-1}) \right]^T \mathbf{D} + \left[f(\mathbf{net}^{L-1}) \right]^T f(\mathbf{net}^{L-1}) \right\} \\ &= \frac{1}{2} \sum_{j=1}^m (d_j - y_j)^2 = \frac{1}{2} \sum_{j=1}^m (d_j - O_j^{L-1})^2 = \frac{1}{2} \sum_{j=1}^m (d_j - f(\mathbf{net}_j^{L-1}))^2 \end{aligned}$$

批量样本集关于当前网络的**总输出误差** $E_{total} = \sum E$

(II) 输出层各个权值 ω_{ij}^{L-1} 调整 $\begin{cases} i = 1, \dots, n_{L-2} \\ j = 1, \dots, n_{L-1} (m) \end{cases}$



输出层的**实际输出向量** y :

$$O^{L-1} = y = f(net^{L-1}) = f\left((W^{L-1})^T O^{L-2} + b^{L-1}\right)$$

其中

$$\begin{cases} W^{L-1} = (\omega_{ij}^{L-1})_{n^{L-2} \times n^{L-1}} & b^{L-1} = [b_1^{L-1} \quad \dots \quad b_{n_{L-1}}^{L-1}]^T \\ net^{L-1} = [net_1^{L-1} \quad \dots \quad net_{n_{L-1}}^{L-1}]^T = (W^{L-1})^T O^{L-2} + b^{L-1} \\ net_j^{L-1} = \sum_{i=1}^{n_{L-2}} \omega_{ij}^{L-1} O_i^{L-2} + b_j^{L-1} \\ O^{L-1} = [O_1^{L-1} \quad \dots \quad O_{n_{L-1}}^{L-1}]^T & m = L-1 \\ y = [y_1 \quad \dots \quad y_m]^T & \text{输出端 } y_j = O_j^{L-1} = f(net_j^{L-1}) \end{cases}$$

$$f(x) = \frac{1}{1 + e^{-x}}$$

$$f'(x) = \frac{e^{-x}}{(1 + e^{-x})^2} = f(x)[1 - f(x)]$$

(II) 输出层各个权值 ω_{ij}^{L-1} 调整 $\begin{cases} i = 1, \dots, n_{L-2} \\ j = 1, \dots, n_{L-1}(m) \end{cases}$ --续1

$$\begin{aligned} \text{输出误差} \quad E &= \frac{1}{2} \|D - y\|^2 = \frac{1}{2} (D^T D - 2y^T D + y^T y) \\ &= \frac{1}{2} \left\{ D^T D - 2 \left[f(\text{net}^{L-1}) \right]^T D + \left[f(\text{net}^{L-1}) \right]^T f(\text{net}^{L-1}) \right\} \\ &= \frac{1}{2} \sum_{j=1}^m (d_j - y_j)^2 = \frac{1}{2} \sum_{j=1}^m (d_j - O_j^{L-1})^2 = \frac{1}{2} \sum_{j=1}^m (d_j - f(\text{net}_j^{L-1}))^2 \end{aligned}$$

$$\text{其中: } y = f(\text{net}^{L-1}) = \frac{1}{1 + \exp[-\text{net}^{L-1}]} \quad \text{net}^{L-1} = (\mathbf{W}^{L-1})^T \mathbf{O}^{L-2} + \mathbf{b}^{L-1}$$

$$\text{输出层误差的局部梯度} \quad \frac{\partial E}{\partial \mathbf{W}^{L-1}} = \frac{\partial E}{\partial y} \cdot \frac{\partial y}{\partial \text{net}^{L-1}} \cdot \frac{\partial \text{net}^{L-1}}{\partial \mathbf{W}^{L-1}} = (y - D) [y \cdot (1 - y)] \mathbf{O}^{L-2}$$

$$\text{其中} \begin{cases} \frac{\partial E}{\partial y} = y - D & \frac{\partial \text{net}^{L-1}}{\partial \mathbf{W}^{L-1}} = \mathbf{O}^{L-2} \\ \frac{\partial y}{\partial \text{net}^{L-1}} = \frac{e^{-\text{net}^{L-1}}}{(1 + e^{-\text{net}^{L-1}})^2} = [y \cdot (1 - y)] \end{cases}$$

$$\begin{aligned} f(x) &= \frac{1}{1 + e^{-x}} \\ f'(x) &= \frac{e^{-x}}{(1 + e^{-x})^2} = f(x) [1 - f(x)] \end{aligned}$$

$$\cdots \rightarrow \begin{bmatrix} \left(\omega_{ij}^{L-1} \right)_{n^{L-2} \times n^{L-1}} \\ \left(b_j^{L-1} \right)_{1 \times n^{L-1}} \end{bmatrix} \rightarrow \boxed{net_j^{L-1} = \sum_{i=1}^{n_{L-2}} \omega_{ij}^{L-1} o_i^{L-2} + b_j^{L-1}} \rightarrow \boxed{y_j = f \left(net_j^{L-1} \right)} \rightarrow E$$

输出层误差关于 ω_{ij}^{L-1} 的偏导数：

$$\frac{\partial E}{\partial \omega_{ij}^{L-1}} = \frac{\partial E}{\partial net_j^{L-1}} \cdot \frac{\partial net_j^{L-1}}{\partial \omega_{ij}^{L-1}} = \delta_j^{L-1} \cdot \frac{\partial net_j^{L-1}}{\partial \omega_{ij}^{L-1}},$$

$$\delta_j^{L-1} = \frac{\partial E}{\partial net_j^{L-1}} = \frac{\partial E}{\partial y_j} \cdot \frac{\partial y_j}{\partial net_j^{L-1}}$$

$$\left\{ \begin{array}{l} E = \frac{1}{2} \sum_{j=1}^m (d_j - y_j)^2 \\ y_j = f(net_j^{L-1}) = \frac{1}{1 + e^{-net_j^{L-1}}} \\ net_j^{L-1} = \sum_{i=1}^{n_{L-2}} \omega_{ij}^{L-1} o_i^{L-2} + b_j^{L-1} \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} \frac{\partial E}{\partial y_j} = y_j - d_j \\ \frac{\partial y_j}{\partial net_j^{L-1}} = f'(net_j^{L-1}) = y_j (1 - y_j) \\ \frac{\partial net_j^{L-1}}{\partial \omega_{ij}^{L-1}} = o_i^{L-2} \end{array} \right.$$

$$\Rightarrow \frac{\partial E}{\partial \omega_{ij}^{L-1}} = \delta_j^{L-1} \cdot \frac{\partial net_j^{L-1}}{\partial \omega_{ij}^{L-1}} = o_i^{L-2} y_j (1 - y_j) (y_j - d_j)$$

$$f(x) = \frac{1}{1 + e^{-x}}$$

$$f'(x) = f(x) [1 - f(x)]$$

(II) 输出层各个权值 ω_{ij}^{L-1} 调整 $\begin{cases} i = 1, \dots, n_{L-2} \\ j = 1, \dots, n_{L-1}(m) \end{cases}$ --续3

基于梯度下降法的连接权值迭代估计

ω_{ij}^{L-1} 修正方式: $\omega_{ij}^{L-1}(t+1) = \omega_{ij}^{L-1}(t) + \Delta\omega_{ij}^{L-1}(t)$

ω_{ij}^{L-1} 修正量 $\Delta\omega_{ij}^{L-1}(t) = -\eta \frac{\partial E}{\partial \omega_{ij}^{L-1}(t)}$

其中 $\frac{\partial E}{\partial \omega_{ij}^{L-1}} = \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial \text{net}_j^{L-1}} \cdot \frac{\partial \text{net}_j^{L-1}}{\partial \omega_{ij}^{L-1}} = y_j(1-y_j)(y_j - d_j)O_i^{L-2}$

矩阵形式：

W^{L-1} 修正方式: $W^{L-1}(t+1) = W^{L-1}(t) + \Delta W^{L-1}(t)$

W^{L-1} 修正量 $\Delta W^{L-1}(t) = -\eta \frac{\partial E}{\partial W^{L-1}(t)}$

其中: $\frac{\partial E}{\partial W^{L-1}} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial \text{net}^{L-1}} \cdot \frac{\partial \text{net}^{L-1}}{\partial W^{L-1}} = (y - D)[y \cdot (1 - y)]O^{L-2}$

(III) 输出层的偏置量 b_j^{L-1} 调整 $j = 1, \dots, n_{L-1}(m)$

输出端 j 的净输入 $\text{net}_j^{L-1} = \sum_{i=1}^{n_{L-2}} \omega_{ij}^{L-1} \mathbf{O}_i^{L-2} + \mathbf{b}_j^{L-1}$

输出端 j 的实际输出

$$\mathbf{O}_j^{L-1} = \mathbf{y}_j = f(\text{net}_j^{L-1}) = f\left(\sum_{i=1}^{n_{L-2}} \omega_{ij}^{L-1} \mathbf{O}_i^{L-2} + \mathbf{b}_j^{L-1}\right)$$

输出误差 $E = \frac{1}{2} \sum_{j=1}^m (d_j - \mathbf{y}_j)^2 = \frac{1}{2} \sum_{j=1}^m (d_j - f(\text{net}_j^{L-1}))^2$

偏置量 \mathbf{b}_j^{L-1} 对误差影响

$$\frac{\partial E}{\partial \mathbf{b}_j^{L-1}} = \frac{\partial E}{\partial \mathbf{y}_j} \cdot \frac{\partial \mathbf{y}_j}{\partial \text{net}_j^{L-1}} \cdot \frac{\partial \text{net}_j^{L-1}}{\partial \mathbf{b}_j^{L-1}} = (\mathbf{y}_j - d_j) \mathbf{y}_j \cdot (1 - \mathbf{y}_j) \cdot 1$$

(III) 输出层的偏置量 b_j^{L-1} 调整 (续1)

偏置量 b_j^{L-1} 对误差的影响

$$\frac{\partial E}{\partial b_j^{L-1}} = y_j (y_j - d_j) \cdot (1 - y_j)$$

输出层的偏置量 b_j^{L-1} 调整

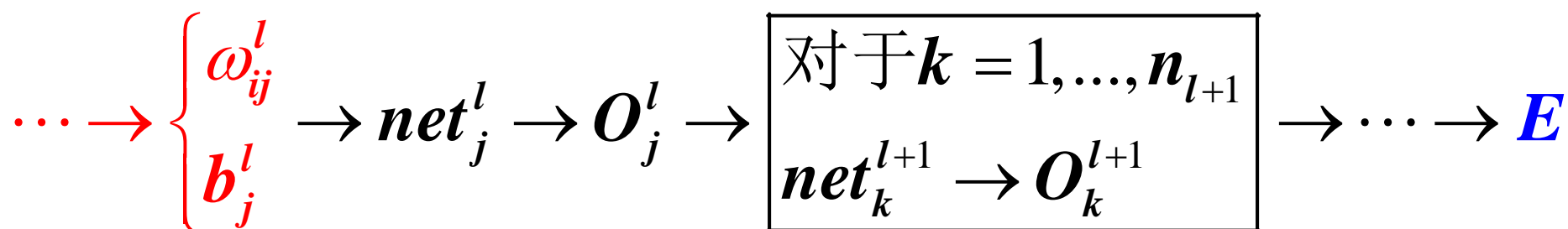
$$b_j^{L-1} \text{修正方式} \quad b_j^{L-1}(t+1) = b_j^{L-1}(t) + \Delta b_j^{L-1}(t)$$

$$b_j^{L-1} \text{修正量} \quad \Delta b_j^{L-1} = -\eta \frac{\partial E}{\partial b_j^{L-1}} = \eta y_j (d_j - y_j) \cdot (1 - y_j)$$

(IV) 隐含层(当前层 l)的权值 ω_{ij}^l 调整 $\begin{cases} i = 1, \dots, n_{l-1} \\ j = 1, \dots, n_l \end{cases}$

当前层 l 的节点 j 对后层 $(l+1)$ 全部节点 k 均有影响

后层 $(l+1)$ 的所有节点均对网络输出误差产生影响



$$\frac{\partial E}{\partial \omega_{ij}^l} = \frac{\partial E}{\partial net_j^l} \cdot \frac{\partial net_j^l}{\partial \omega_{ij}^l} = \delta_j^l \frac{\partial net_j^l}{\partial \omega_{ij}^l}$$

$$\delta_j^l = \frac{\partial E}{\partial O_j^l} \cdot \frac{\partial O_j^l}{\partial net_j^l} = \left[\sum_{k=1}^{n_{l+1}} \left(\frac{\partial E}{\partial O_k^{l+1}} \frac{\partial O_k^{l+1}}{\partial net_k^{l+1}} \frac{\partial net_k^{l+1}}{\partial O_j^l} \right) \right] \frac{\partial O_j^l}{\partial net_j^l}$$

由前面知：

$$\frac{\partial E}{\partial \omega_{ij}^l} = \delta_j^l \frac{\partial net_j^l}{\partial \omega_{ij}^l}$$

所以

$$\begin{aligned} \delta_j^l &= \frac{\partial E}{\partial net_j^l} = \frac{\partial E}{\partial O_j^l} \frac{\partial O_j^l}{\partial net_j^l} = \sum_{k=1}^{n_{l+1}} \left(\frac{\partial E}{\partial O_k^{l+1}} \frac{\partial O_k^{l+1}}{\partial net_k^{l+1}} \frac{\partial net_k^{l+1}}{\partial O_j^l} \right) \frac{\partial O_j^l}{\partial net_j^l} \\ &= \left[\sum_{k=1}^{n_{l+1}} \delta_k^{l+1} \omega_{jk}^{l+1} \right] f'(net_j^l) = \left(\sum_{k=1}^{n_{l+1}} \delta_k^{l+1} \omega_{jk}^{l+1} \right) O_j^l (1 - O_j^l) \end{aligned}$$

又

$$net_j^l = \sum_{i=1}^{n_{l-1}} \omega_{ij}^l O_i^{l-1} + b_j^l$$

所以

$$\frac{\partial net_j^l}{\partial \omega_{ij}^l} = O_i^{l-1}$$

$$\Rightarrow \frac{\partial E}{\partial \omega_{ij}^l} = \delta_j^l \frac{\partial \text{net}_j^l}{\partial \omega_{ij}^l} = \mathbf{O}_i^{l-1} \mathbf{O}_j^l (1 - \mathbf{O}_j^l) \sum_{k=1}^{n_{l+1}} \delta_k^{l+1} \omega_{jk}^{l+1}$$

隐含层节点与其前层节点连接权值 ω_{ij}^l 迭代估计

($l = L - 2, L - 3, \dots$)

ω_{ij}^l 修正方式: $\omega_{ij}^l(t+1) = \omega_{ij}^l(t) + \Delta \omega_{ij}^l(t)$

ω_{ij}^l 修正量 $\Delta \omega_{ij}^l(t) = -\eta \frac{\partial E}{\partial \omega_{ij}^l(t)}$

其中 $\frac{\partial E}{\partial \omega_{ij}^l} = \frac{\partial E}{\partial \text{net}_j^l} \cdot \frac{\partial \text{net}_j^l}{\partial \omega_{ij}^l} = \mathbf{O}_i^{l-1} \mathbf{O}_j^l (1 - \mathbf{O}_j^l) \sum_{k=1}^{n_{l+1}} \delta_k^{l+1} \omega_{jk}^{l+1}$

(V) 隐含层(当前层 l)各节点的偏置量 b_j^l 调整

$$j = 1, \dots, n_l$$

由前面分析知: $\delta_j^l = \frac{\partial E}{\partial net_j^l} = O_j^l (1 - O_j^l) \sum_{k=1}^{n_{l+1}} \delta_k^{l+1} \omega_{jk}^{l+1}$

$$net_j^l = \sum_{i=1}^{n_{l-1}} \omega_{ij}^l O_i^{l-1} + b_j^l \Rightarrow \frac{\partial net_j^l}{\partial b_j^l} = 1$$

所以: $\frac{\partial E}{\partial b_j^l} = \frac{\partial E}{\partial net_j^l} \cdot \frac{\partial net_j^l}{\partial b_j^l} = \delta_j^l \cdot \frac{\partial net_j^l}{\partial b_j^l} = \delta_j^l$

偏置量 b_j^l 修正

$$b_j^l \text{修正方式} \quad b_j^l(t+1) = b_j^l(t) + \Delta b_j^l(t)$$

$$b_j^l \text{修正量} \quad \Delta b_j^l = -\eta \frac{\partial E}{\partial b_j^l} = -\eta \delta_j^l = -\eta O_j^l (1 - O_j^l) \sum_{k=1}^{n_{l+1}} \delta_k^{l+1} \omega_{jk}^{l+1}$$

(2) 误差反向传播过程(续)

引入惯性项：为加快收敛速度，往往在权值及偏置量修正过程中加上前次的权值修正量，称为**惯性项**。

权值修正

$$\begin{aligned}\omega_{ij}(t+1) &= \omega_{ij}(t) + \Delta\omega_{ij}(t) \\ \Delta\omega_{ij}(t) &= -\eta\delta_j O_i + \alpha\Delta\omega_{ij}(t-1)\end{aligned}$$

偏置量修正

$$\begin{aligned}b_j(t+1) &= b_j(t) + \Delta b_j(t) \\ \Delta b_j(t) &= -\eta\delta_j + \alpha\Delta b_j(t-1)\end{aligned}$$

BP算法步骤

首先明确 {

- 训练样本** (输入向量, 期望输出), 如 (x_p, D_p)
- 节点参数初始化** “小随机数”
例: 输入层与第一隐含层间: $\omega_{ij}^1 \in \left(-\frac{1}{\sqrt{d}}, \frac{1}{\sqrt{d}}\right)$
- 学习步长 η** 通常固定0.1 ~ 0.3之间; 或动态调整
- 惯性(冲量)项系数 α** 通常0.9 ~ 1之间

STEP1 {

- (1) **确定神经网络结构**, 包括
输入层节点数; 隐含层数目; 各隐含层节点数目;
输出层节点数; 各神经元节点的传递函数
- (2) **样本集的标准化处理以及输出端编码**
- (3) **设定终止条件**: 最大可允许训练轮数(硬条件);
训练精度(软条件)
- (4) **以小随机数初始化网络权值**;
记训练时间 $t = 0$

BP算法步骤(续)

STEP2: 重复如下过程直至满足算法终止条件

(1) 按随机或任意顺序从训练集中抽取1个训练样本 (\mathbf{x}, \mathbf{D})

样本输入 $\mathbf{x} = [x_1, \dots, x_n]^T \in R^n$

期望输出 $\mathbf{D} = [d_1, \dots, d_m]^T \in R^m$

(2) 计算输入 \mathbf{x} 时当前网络的**实际输出** $\mathbf{y} = [y_1, \dots, y_m]^T$

$$y_r = f \left(\sum_{s=1}^{n_{L-2}} \omega_{sr}^{l=L-1} \cdots f \left(\sum_{j=1}^{n_1} \omega_{jk}^{l=2} f \left(\sum_{i=1}^n \omega_{ij}^{l=1} x_i + b_j^{l=1} \right) + b_k^{l=2} \right) \cdots + b_r^{l=L-1} \right)$$

$r = 1, \dots, m$

$$\text{其中} \begin{cases} \text{激活函数 } f(\alpha) = \frac{1}{1 + e^{-\alpha}} \\ f'(\alpha) = -\frac{-e^{-\alpha}}{(1 + e^{-\alpha})^2} = f(\alpha)(1 - f(\alpha)) \end{cases}$$

BP算法步骤

(3) 从输出层开始**调整权值**及节点**偏置量**，具体为：

$$\text{对于第 } l \text{ 层, 修正} \begin{cases} \text{权值} & \omega_{ij}^l(t+1) = \omega_{ij}^l(t) + \Delta\omega_{ij}^l(t) \\ \text{偏置量} & b_j^l(t+1) = b_j^l(t) + \Delta b_j^l(t) \end{cases}$$

$$i = 1, \dots, n_{l-1} \quad j = 1, \dots, n_l$$

$$\text{其中} \quad \Delta\omega_{ij}^l(t) = -\eta \delta_j^l x_i^{l-1} \quad \Delta b_j^l(t) = -\eta \delta_j^l$$

$$\begin{cases} \text{输出层 } (l = L-1) & \delta_j^l = y_j(1-y_j)(d_j - y_j) \\ & j = 1, \dots, m \\ \text{中间层 } (1 \leq l < L-1) & \delta_j^l = x_j^l(1-x_j^l) \sum_{k=1}^{n_{l+1}} \delta_k^{l+1} \omega_{jk}^{l+1}(t) \\ & j = 1, \dots, n_l \end{cases}$$

BP算法步骤

(4) **更新全部权值及偏置量**，对所有训练样本
重新计算输出；估计更新后网络输出误差；

检查算法**终止条件**

$\left\{ \begin{array}{l} \text{若不满足条件, 则 } t \leftarrow t + 1, \text{ 转向(1)} \\ \text{若条件满足, 则终止, 转向STEP3} \end{array} \right.$

STEP3: 算法结束. **保存**各层连接权值及偏置量。

终止条件可以是如下之一：

$\left\{ \begin{array}{l} \text{(1) 网络实际输出与期望输出总误差} < \text{阈值1} \\ \text{(2) 最近1轮训练中所有权值及偏置量变化最大值} < \text{阈值2} \\ \text{(3) 算法达到最大允许的总训练次数} = \text{阈值3} \end{array} \right.$