

数据挖掘与机器学习

潘斌

panbin@nankai.edu.cn

范孙楼227

1

实验5：基于BP的信用卡欺诈预测

- 给定信用卡数据，判断是否存在欺诈
- 用BP神经网络算法实现
- 本周四（5.20）上午1、2节，二主楼B403

大作业汇报安排

- 5.24、5.31汇报，3-5人一组，25分钟汇报+5分钟提问，说明每个人的贡献
- 自行分组，由组长将组员信息、拟选的题目、期望汇报的时间本周四前发到我的邮箱
- 汇报顺序优先按照意愿，不能满足时抽签

5 非线性分类器

- 5.1 多类问题概述
- 5.2 最小距离分类器
- 5.3 分段线性分类器概述
- 5.4 决策树
- 5.5 近邻法分类器
- 5.6 人工神经网络
- 5.7 SVM

5.1 多类问题概述

- 5.1.1 多类问题
- 5.1.2 解决方案

5.1.1 多类问题概述

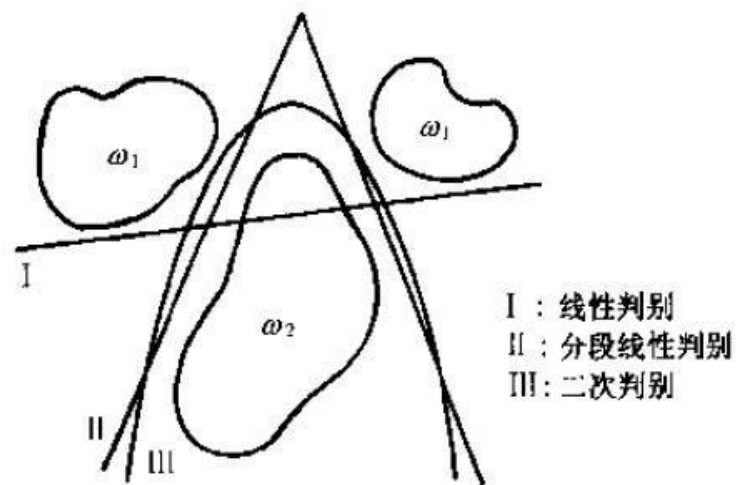
- 多类问题包括
 - 多类情况（类别数 $C > 2$ ）
 - 单峰分布
 - 多峰分布

5.1.1 多类问题

- 多类问题包括
 - 两类情况 ($C = 2$)
 - 样本集具有多峰分布

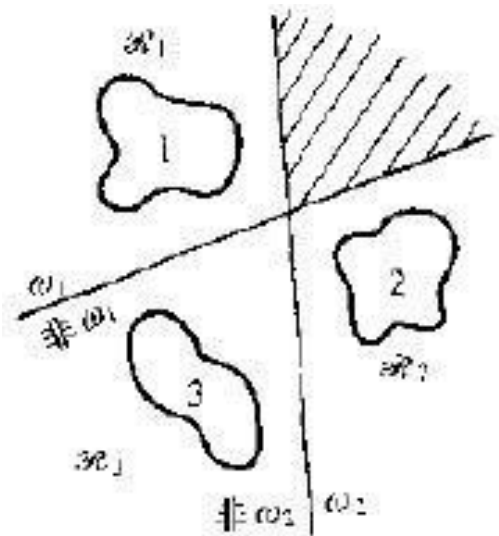
5.1.2 解决方案

- 如何解多类问题
 - Bayes分类器
 - 二次型判别函数
 - 线性分类器
 - 其它分类器



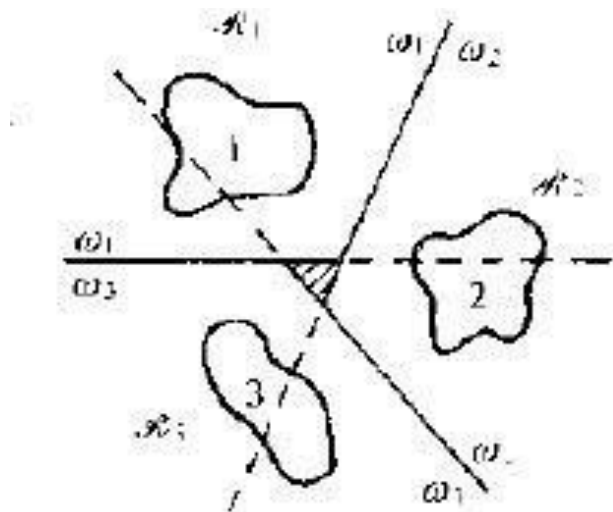
5.1.2 解决方案

- 多类问题是否可以用线性分类器解？
 - 思路一
 - 对“类与类的非”进行线性分类
 - 只需要 $C - 1$ 个线性分类器就可以



5.1.2 解决方案

- 多类问题是否可以用线性分类器解？
 - 思路二
 - “两两分类”进行线性分类
 - 需要 $C(C-1)/2$ 个线性分类器就可以

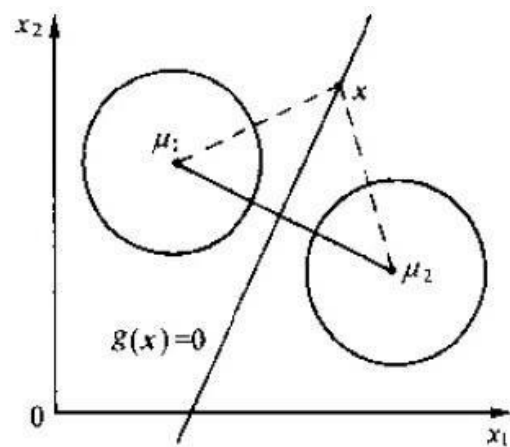


5.2 最小距离分类器

- 5.2.1 最小距离分类器原理
- 5.2.2 分段最小距离分类器
- 5.2.3 特点

5.2.1 最小距离分类器原理

- 回顾两类单峰线性分类器
 - 垂直平分 / 最小距离分类器
 - 基于两类样本均值点作垂直平分线



5.2.1 最小距离分类器原理

- 其最小距离形式

- 判别函数

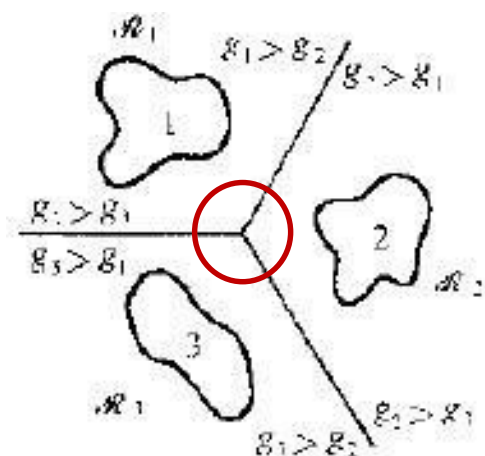
- $G_1(x) = d_1(x) = \|x - m_1\|$
 - $G_2(x) = d_2(x) = \|x - m_2\|$

- 决策规则

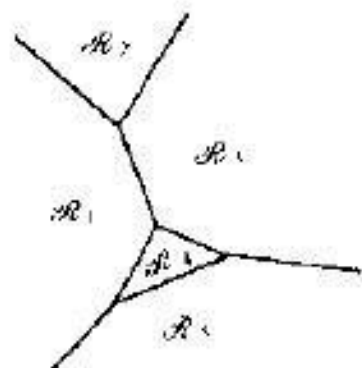
- 对于未知样本 x ，若 $d_1(x) < d_2(x)$ ，则 x 决策为 ω_1 类
 - 若 $d_1(x) > d_2(x)$ ，则 x 决策为 ω_2 类

5.2.1 最小距离分类器原理

- 直接使用可以解决多类问题
 - 解决C类单峰问题



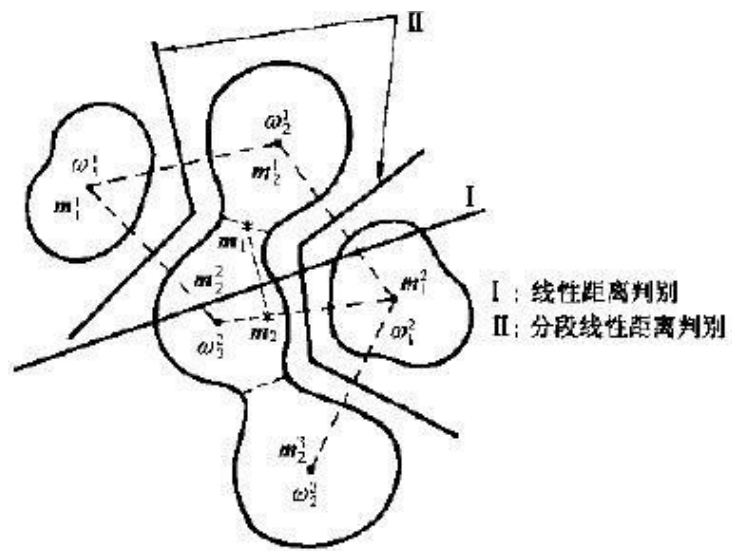
(a) 三类



(b) 五类

5.2.1 最小距离分类器原理

- 直接使用可以解决多类问题
 - 解决两类多峰问题



5.2.2 分段最小距离分类器

- 问题

- 已知各类及其子类
- 求分段最小距离分类器

5.2.2 分段最小距离分类器

- 分类器设计
 - 先求各子类均值
 - m_{ij} (ω_i 类的第j子类)
 - 定义各类判别函数
 - $G_i(x) = \min_j \|x - m_{ij}\|$
 - 决策规则
 - 对于未知样本 x ，若 $G_k(x) = \min_i G_i(x)$ ，则 x 决策为 ω_k 类

5.2.3 特点

- 分类器特点
 - 解决两类多峰或多类问题的分段线性分类器
 - 可以解决几乎所有分类问题但要已知各类子类
 - 概念直观简单，未经优化
 - 分类器设计简单容易
 - （无重叠区或空白区）

5.3 分段线性分类器概述

- 5.3.1 问题与思路
- 5.3.2 设计说明

5.3.1 问题与思路

- 思路
 - 参考分段最小距离分类器
 - 定义判别函数
 - 定义决策规则

5.3.1 问题与思路

- 针对不同已知条件
 - 1、已知各类子类个数及子类分布区域
 - 2、已知各类子类个数（分布区域不知）
 - 3、一般情况（子类个数和分布区域均不知）

5.3.2 设计说明

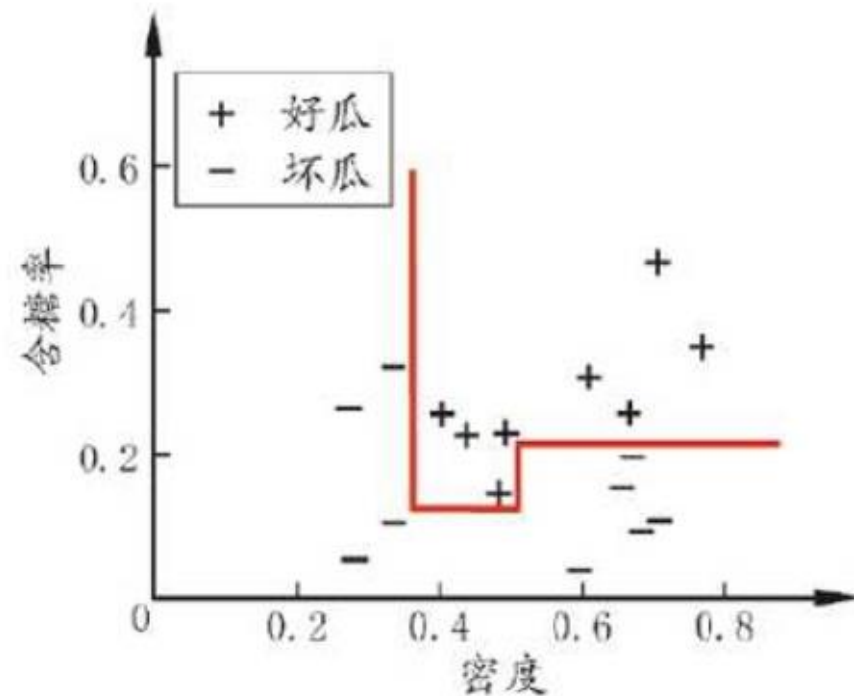
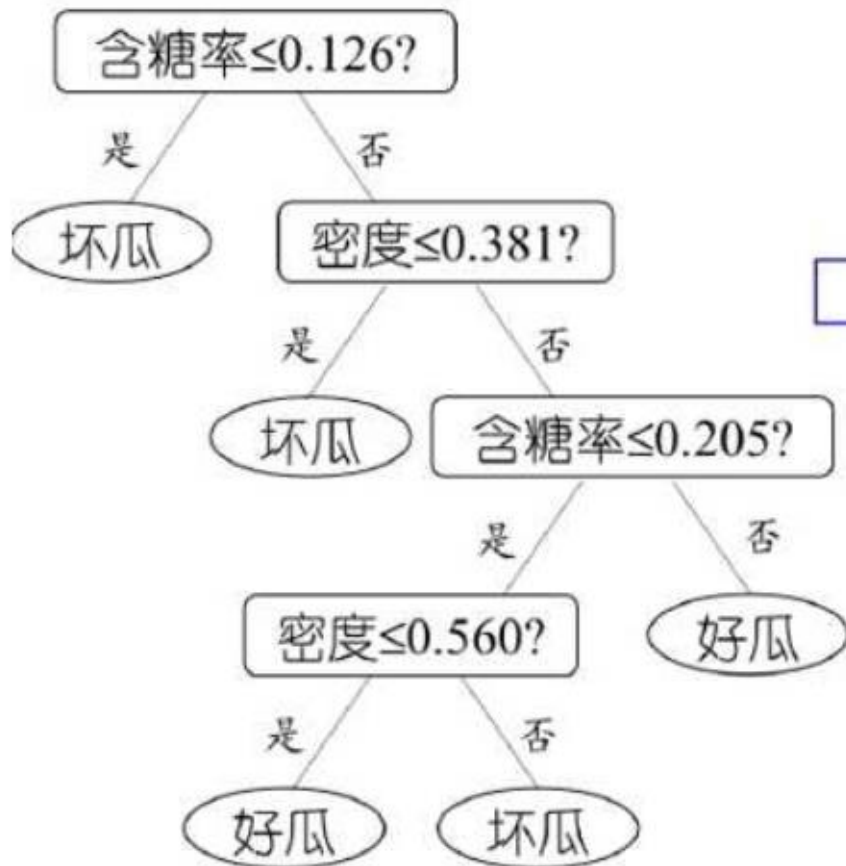
- 两种判别函数的区别
 - 小写g函数（分界面）
 - 每段设计一个g函数，容易做
 - 多个分段，如何判断正负侧，需要特殊规则
 - 大写G函数（计算值）
 - 每个子类设计一个G函数，需要知道类别分布区域
 - 直接计算Max或Min，判别规则简单

5.3.2 设计说明

- 设计关键

- 如何确定各类的子类个数
- 如何确定子类的分布区域
- 如何求解各子类的权向量和阈值权
- 若采用小写g函数，决策规则如何

5.4 分段线性分类器——决策树



5.5 分段线性分类器——近邻法分类器

- 5.5.1 近邻法原理
- 5.5.2 最近邻法
- 5.5.3 k-近邻法

5.5.1 近邻法原理

- 近邻法分类是一种简单实用的分类方法
 - 线性分类器
 - Bayes分类器
 - 分段线性分类器

5.5.1 近邻法原理

- 最小距离分类器

- 用均值点作为代表点，按最近均值点进行决策
- 有时候均值点不具有很好的代表性

- 近邻法分类器思路

- 全部训练样本都是代表点，按最近邻点进行决策

5.5.2 最近邻法

- 问题
 - 设C类问题: $\omega_1, \omega_2, \dots, \omega_C$
 - ω_i 类样本集 $Z_i = \{\dots, x_{ik}, \dots\}$
 - 求近邻法分类器

5.5.2 最近邻法

- 判别函数

- 定义 $G_i(x) = \min \|x - x_{ik}\|$ $i = 1, 2, \dots, C$

- 决策规则

- 对于未知样本 x ，若 $G_j(x) = \min G_i(x)$ ，则 $x \in \omega_j$

- 决策面

- 分段线性（画一下）

5.5.2 最近邻法

- 实例

- 甲类: $[0\ 3]^T$ 、 $[2\ 4]^T$ 、 $[1\ 3]^T$ 、 $[2\ 3]^T$ 、 $[0\ 2]^T$
- 乙类: $[4\ 1]^T$ 、 $[3\ 2]^T$ 、 $[2\ 1]^T$ 、 $[3\ 0]^T$ 、 $[3\ 1]^T$
- 待分类样本为 $\mathbf{x} = [5\ 0]^T$, 问 \mathbf{x} 应决策为哪一类?

5.5.2 最近邻法

- 近邻法特点

- 可以解决几乎所有分类问题
- 概念直观简单未经优化，但错误率并不高
- 分类器设计容易
- 运算量大，需要设计快速算法

5.5.2 最近邻法

- 快速算法
 - 剪辑法
 - 进行预分类
 - 剪辑掉错分样本
 - 剪辑法可以重复进行
 - 清理两类见的边界，去掉类别混杂的样本，使两类边界更清晰
 - 剪掉边际处的混淆样本

5.5.2 最近邻法

- 快速算法
 - 剪辑法例一

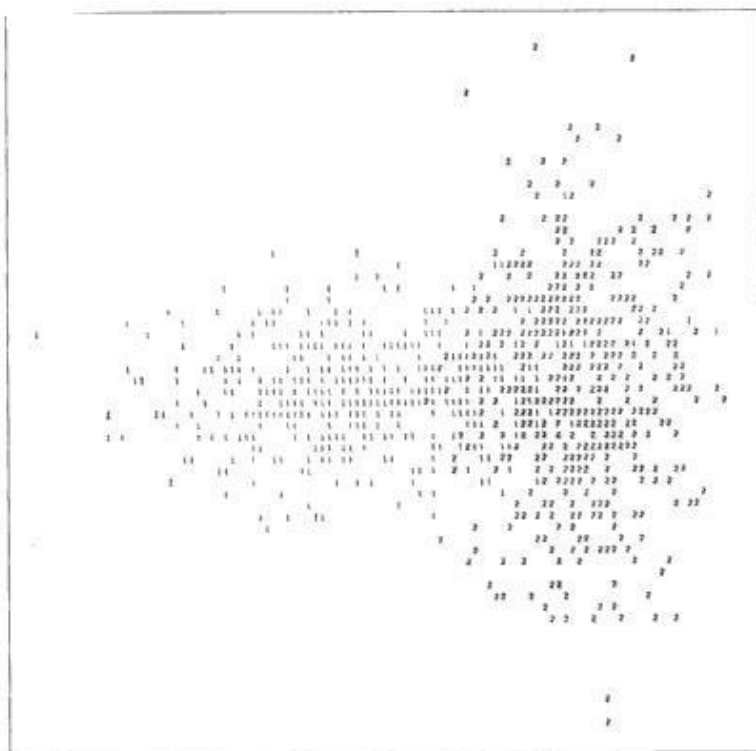


图 6.5 MULTIEDIT 算法实验:原始样本集

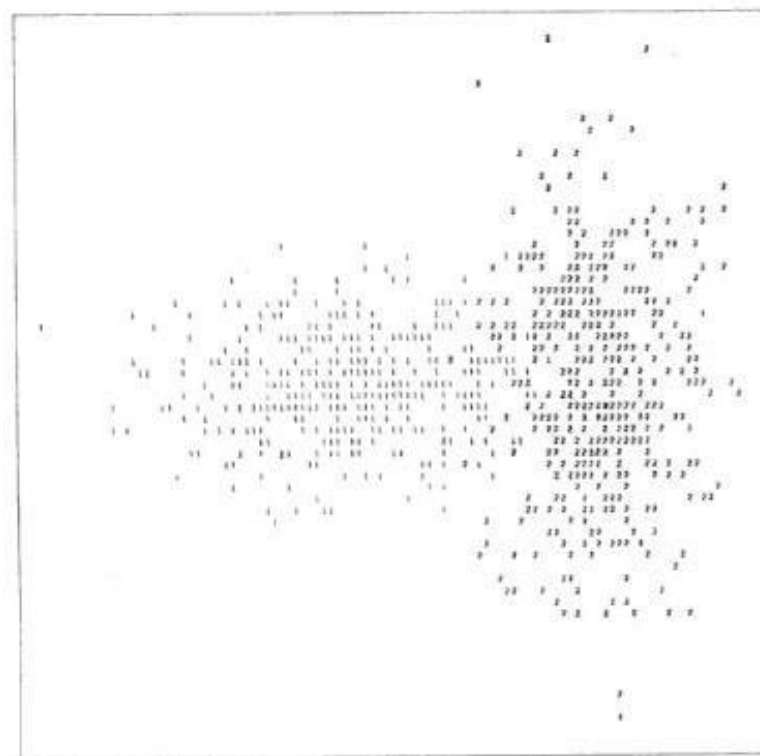


图 6.6 MULTIEDIT 算法实验:第一次迭代后留下的样本

5.5.2 最近邻法

- 快速算法
 - 剪辑法例一

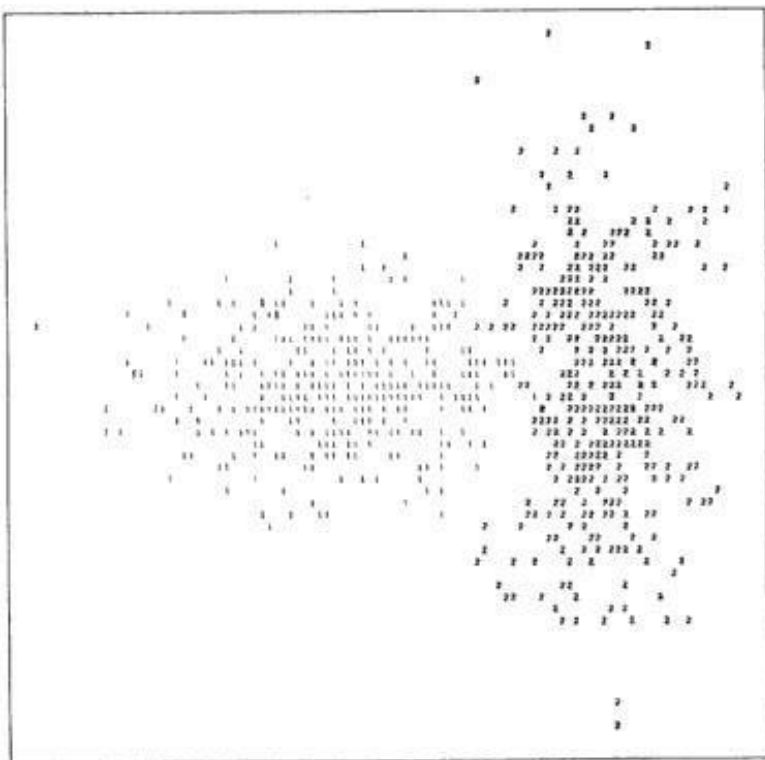


图 6.7 MULTIEDIT 算法实验:经三次迭代后留下的样本

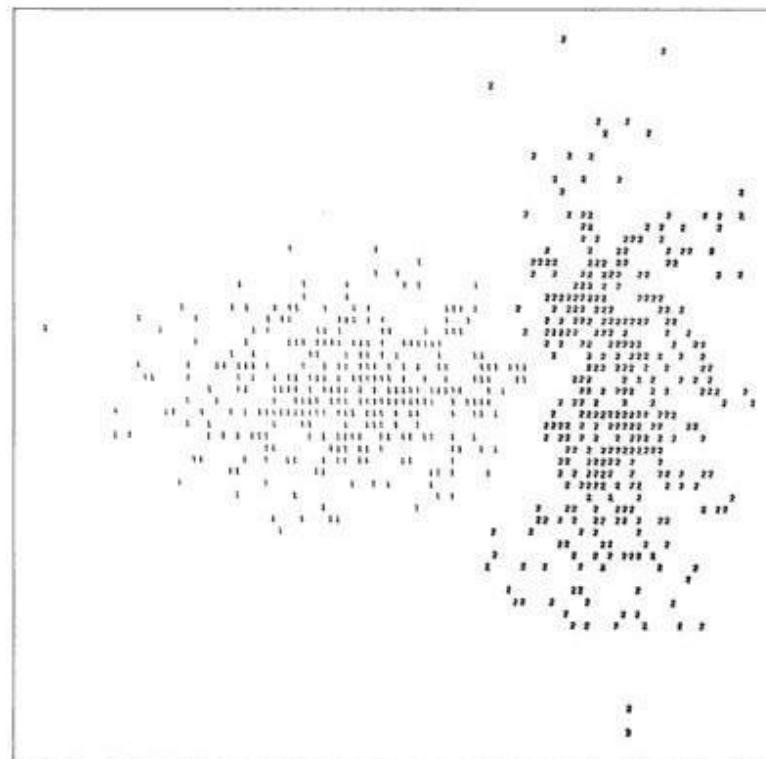


图 6.8 MULTIEDIT 算法实验:算法终止时留下的样本

5.5.2 最近邻法

- 快速算法
 - 剪辑法例二

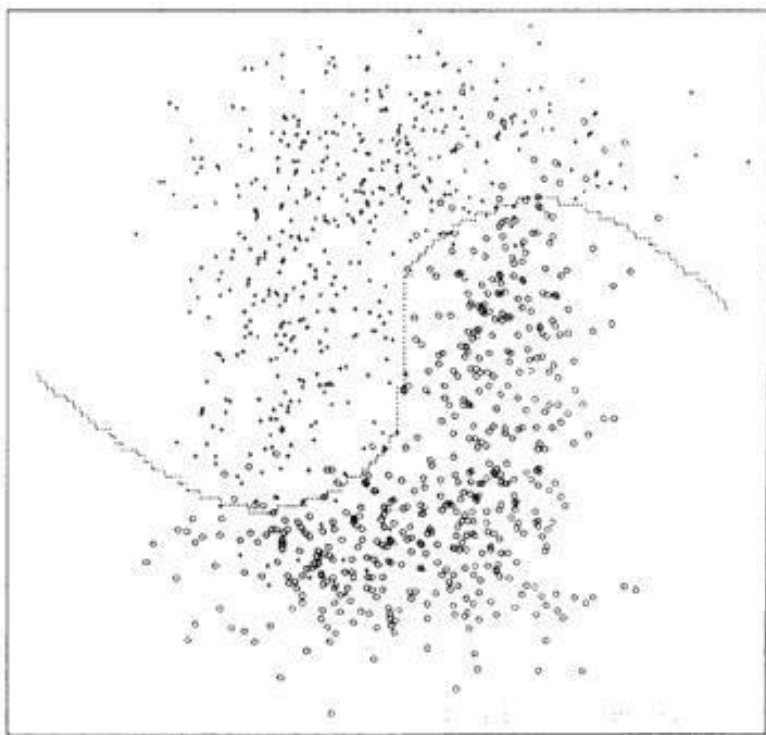


图 6.9 非正态分布下 MULTIEDIT 重复剪辑实验:初始样本集

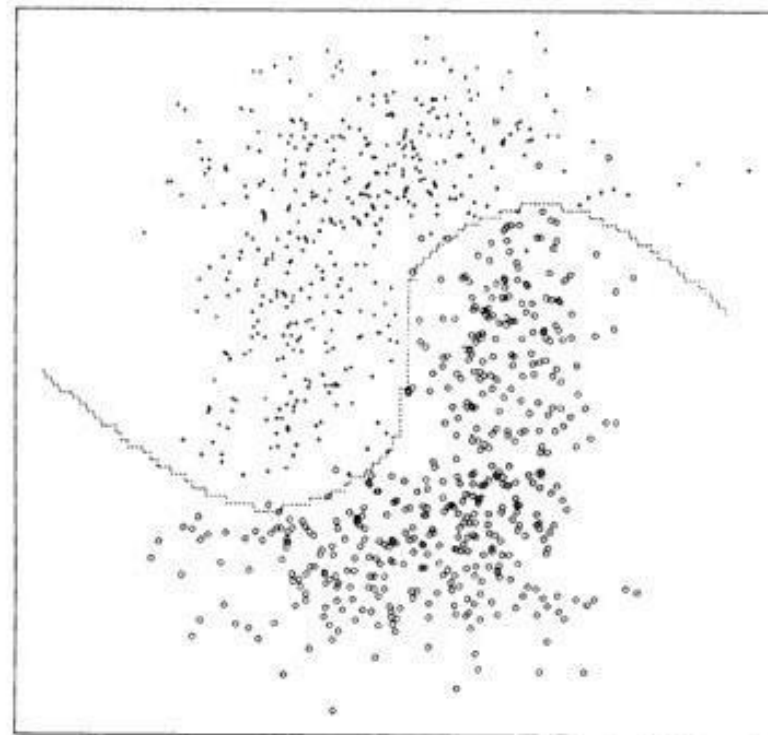


图 6.10 非正态分布下 MULTIEDIT 重复剪辑实验:第一次剪辑后的样本集

5.5.2 最近邻法

- 快速算法
 - 剪辑法例二

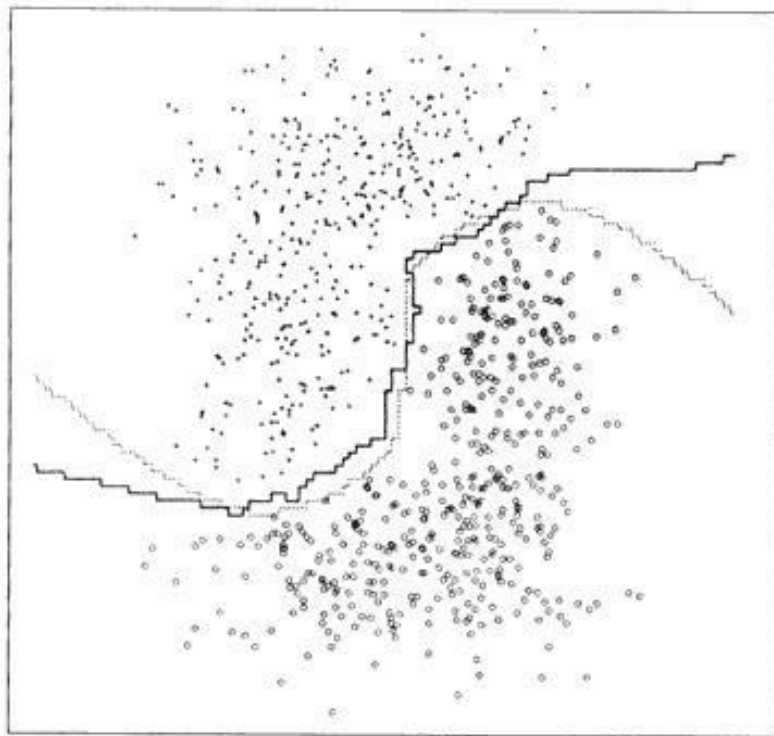


图 6.11 非正态分布下 MULTIEDIT 重复剪辑实验:最终结果

5.5.2 最近邻法

- 快速算法

- 压缩法

- 大多数样本参与计算，但却不起决定作用，是多余样本
 - 剪掉容易分的样本

- 压缩法步骤

- 每类各取一个代表样本（例如均值点近邻），组成压缩样本集 Z_s 。
 - 以当前 Z_s 对样本集做最近邻分类，然后将错分样本放入 Z_s 。
 - 重复上述步骤，直至无样本放入为止。

5.5.2 最近邻法

- 快速算法
 - 压缩法示例

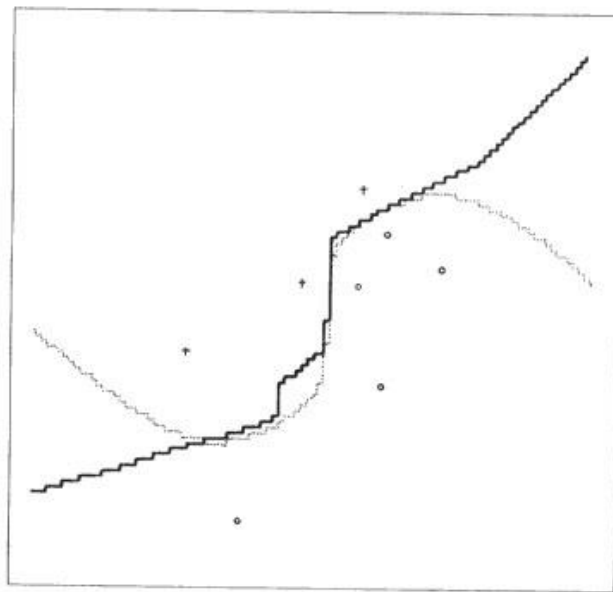


图 6.12 图 6.9 的数据经 MULTIEDIT 算法剪辑之后再使用 CONDENSING 压缩近邻算法的结果

5.5.3 k-近邻法

- 最近邻法的问题
 - 噪点干扰
- **k-近邻法**
 - 如果一个样本附近的k个最近样本的大多数属于某一个类别，则该样本也属于这个类别

5.5.3 k-近邻法

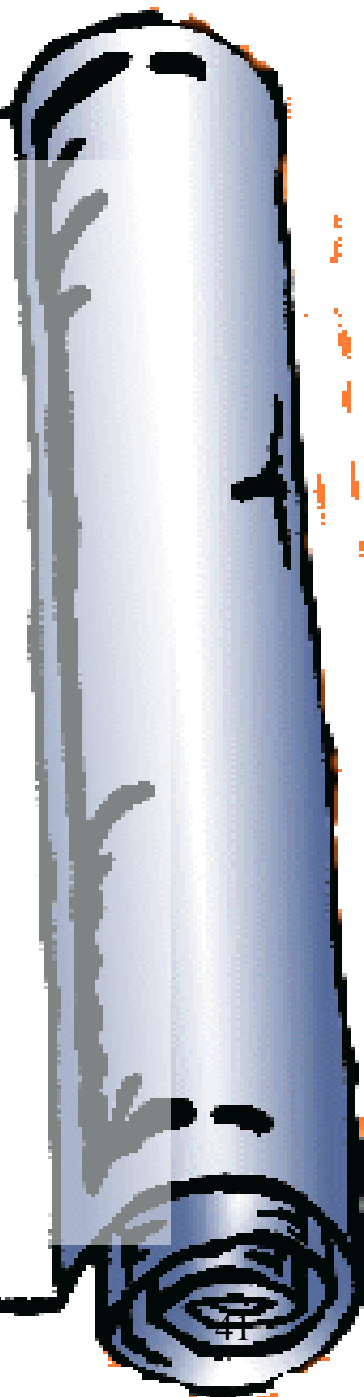
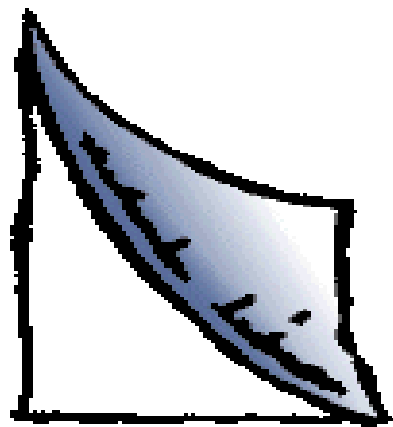
- 判别函数

- 定义 $G_i(x) = k_i$ $i = 1, 2, \dots, C$

- 决策规则

- 对于未知样本 x ，若 $G_j(x) = \max G_i(x)$ ，则 $x \in \omega_j$

人工神经网络



Artificial Neural Networks

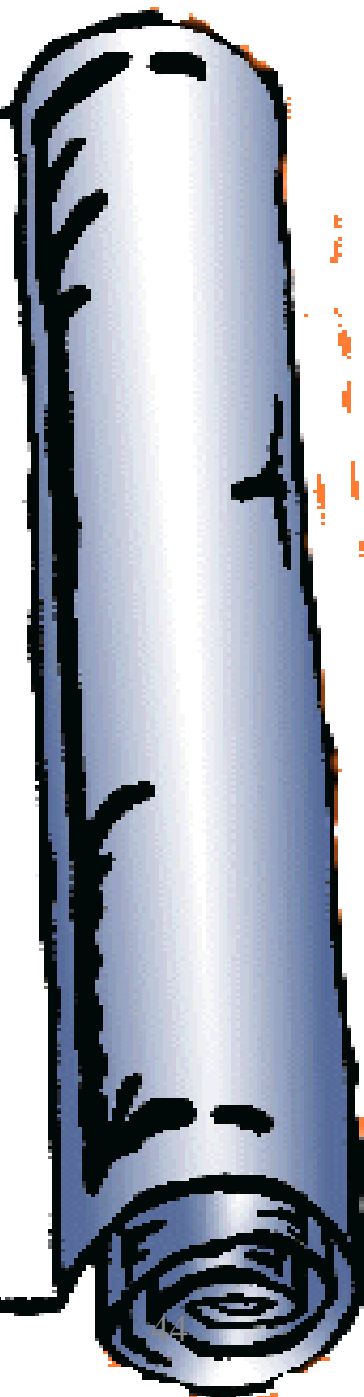
- 本章关注：
 - 感知器
 - 感知器法则
 - 梯度下降（delta）法则
 - 多层网络和BP算法
 - 多层网络
 - BP算法

Artificial Neural Networks

- 人工神经网络（Artificial Neural Network, ANN）

–神经网络是由具有适应性的简单单元组成的广泛并行互连的网络，它的组织能够模拟生物神经系统对真实世界物体所作出的交互反应。（Kohonen,1988）

感知器 (Perceptron)



感知准则函数

- 感知准则函数

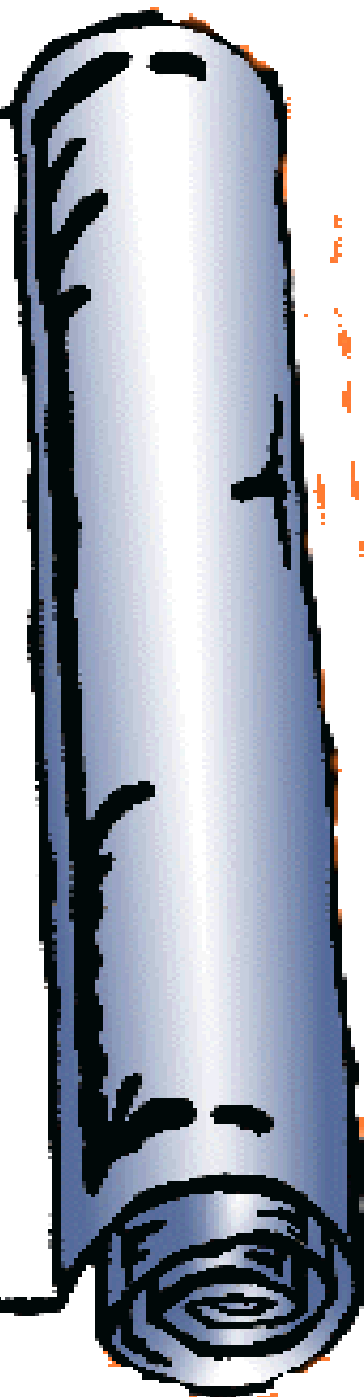
- 对于规范化的增广样本集

- $a^T y_i < 0$ ——错误分类

- 定义感知准则函数，作为优化准则函数

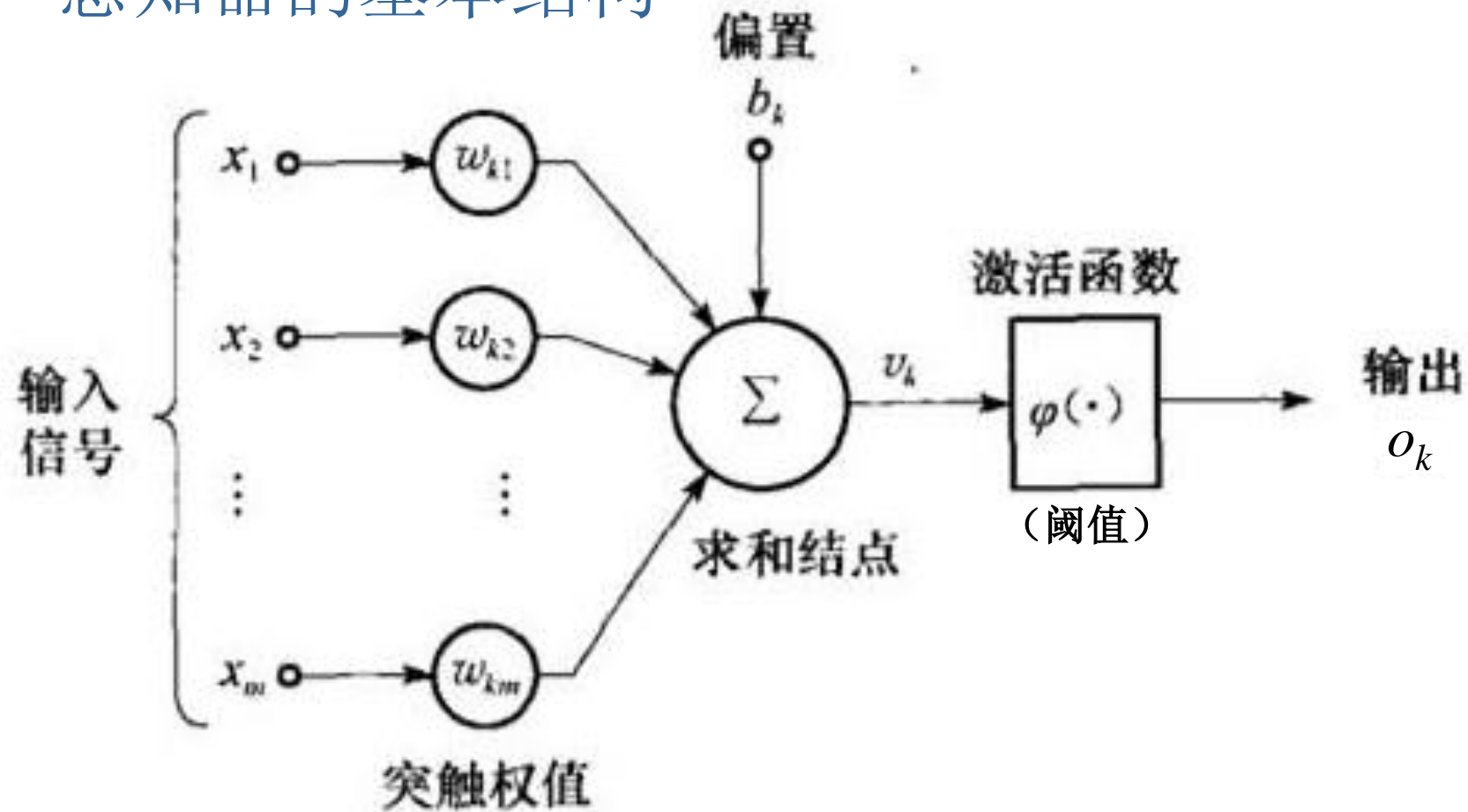
$$\min J_p(a) = \sum_{y \in Z_E} (-a^T y)$$

- 求解向量（或解区）



Artificial Neural Networks

- 感知器的基本结构



Artificial Neural Networks

- 感知器训练方法

- 单个感知器的学习任务：决定一个权向量，它可以使感知器对于给定的训练样例输出正确的1或-1。
- 感知器训练的方法有：感知器法则、delta法则，这两种方法是训练多层神经网络的基础。

Artificial Neural Networks

- 感知器法则
- 算法过程：
 - 随机选取权值
 - 反复应用这个感知器到每个训练样例，只要它误分类样例就修改感知器的权值： $w_i \leftarrow w_i + \Delta w_i$ $\Delta w_i = \eta (t - o)x_i$
 - 其中： η 为学习速率、 t 为当前样本目标输出、 o 为感知器输出。
 - 重复这个过程，直到感知器正确分类所有的训练样例。
 - 如果训练样例线性可分，并且使用充分小的 η ，感知器法则会收敛到一个能正确分类所有训练样例的权向量

Artificial Neural Networks

- 为什么这个更新法则会成功收敛到正确的权值呢？为了得到直观的感觉，考虑一些特例。
 - 假定训练样本已被感知器正确分类。这时， $(t-o)$ 是0，这使 Δw_i 为0，所以没有权值被修改。
 - 而如果当目标输出是+1时感知器输出一个-1，这种情况为使感知器输出一个+1而不是-1，权值必须被修改以增大 Δw_i 的值(x 恒正)。
 - 例如，如果 $x_i > 0$ ，那么增大 w_i 会使感知器更接近正确分类这个实例。
 - 如果 $x_i=0.8$ ， $\eta=0.1$ ， $t=1$ ，并且 $o=-1$ ，那么权更新就是 $\Delta w_i = \eta(t-o)x_i = 0.1(1-(-1))0.8 = 0.16$ 。
 - 另一方面，如果 $t=-1$ 而 $o=1$ ，那么和正的 x_i 关联的权值会被减小而不是增大。

Artificial Neural Networks

- 例：训练数据集中的正例为 $\underline{x}_1=(3,3)'$, $\underline{x}_2=(4,3)'$, 反例为 $\underline{x}_3=(1,1)'$, 试用感知器法则求感知器模型 $f(\underline{x})=\text{sign}(\underline{w} \cdot \underline{x})$ 。 $\underline{w}=(w_{(0)}, w_{(1)}, w_{(2)})'$, $\underline{x}=(1, x_{(1)}, x_{(2)})'$, ($\eta=1$)
- 解：
 - 取初值 $\underline{w}^{(0)}=0$
 - 对 $\underline{x}_1=(3,3)'$, $\underline{w}^{(0)} \cdot \underline{x}_1=0$, 未能被正确分类
 - 更新 $\underline{w}^{(0)}$, $\underline{w}^{(1)} = \underline{w}^{(0)} + \eta(t-o)\underline{x}_1 = (2, 6, 6)'$
 - 对 $\underline{x}_1=(3,3)'$, $\underline{x}_2=(4,3)'$, 显然 $\underline{w}^{(1)} \cdot \underline{x}_1 > 0$, $\underline{w}^{(1)} \cdot \underline{x}_2 > 0$, 被正确分类, 无需更新 $\underline{w}^{(1)}$
 - 对 $\underline{x}_3=(1,1)'$, $\underline{w}^{(1)} \cdot \underline{x}_3 > 0$, 被误分类
 - 更新 $\underline{w}^{(1)}$, $\underline{w}^{(2)} = \underline{w}^{(1)} + \eta(t-o)\underline{x}_3 = (0, 4, 4)'$
- 如此下去, 最终得感知器模型为 $f(\underline{x})=\text{sign}(2x_{(1)}+2x_{(2)}-6)$

正例为 $x_1=(3,3)'$, $x_2=(4,3)'$, 反例为
 $x_3=(1,1)'$

Artificial Neural Networks

迭代次数	误分类点	w	$w \cdot x$
0	—	$(0,0,0)'$	0
1	x_1	$(2,6,6)'$	$6x_{(1)}+6x_{(2)}+2$
2	x_3	$(0,4,4)'$	$4x_{(1)}+4x_{(2)}$
3	x_3	$(-2,2,2)'$	$2x_{(1)}+2x_{(2)}-2$
4	x_3	$(-4,0,0)'$	-4
5	x_1	$(-2,6,6)'$	$6x_{(1)}+6x_{(2)}-2$
6	x_3	$(-4,4,4)'$	$4x_{(1)}+4x_{(2)}-4$
7	x_3	$(-6,2,2)'$	$2x_{(1)}+2x_{(2)}-6$
8	无	$(-6,2,2)'$	$2x_{(1)}+2x_{(2)}-6$

Artificial Neural Networks

- 尽管当训练样例线性可分时，感知器法则可以成功地找到一个权向量，但如果样例不是线性可分时它将不能收敛。因此，人们设计了另一个训练法则来克服这个不足，称为delta法则（delta rule）。
- 如果训练样本不是线性可分的，那么delta法则会收敛到目标概念的最佳近似。
- delta法则的关键思想是使用梯度下降（gradient descent）来搜索可能权向量的假设空间，以找到最佳拟合训练样例的权向量。

Artificial Neural Networks

- 考虑简单的线性单元

$$O = w_0 + w_1x_1 + \dots + w_nx_n$$

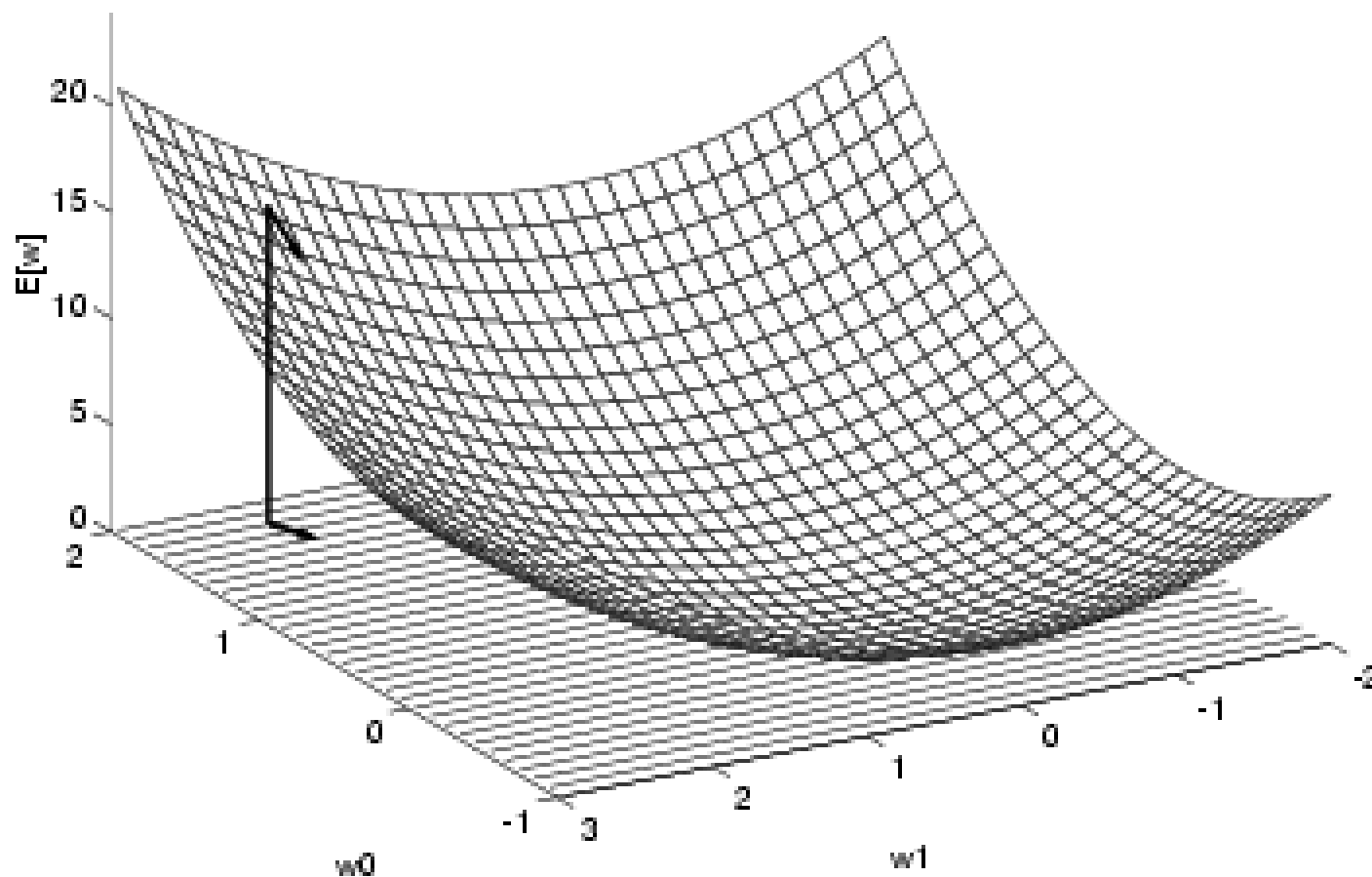
- 为推导权值学习法则，先要确定一个度量标准来衡量假设（权向量）相对于训练样例的训练误差（training error）。
一个常用的度量标准为：

$$E[\vec{w}] \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

- 其中D是训练样例集合， t_d 是训练样例d的目标输出， o_d 是线性单元对训练样例d的输出。

Artificial Neural Networks

- 包含可能权向量的整个假设空间和与它们相关联的E值。



Artificial Neural Networks

- 梯度下降法则的推导

- 可以通过计算E相对向量的每个分量的导数来得到沿误差曲面最陡峭下降的方向。这个向量导数被称为E对于w 的梯度，记作 $\nabla E[\vec{w}]$

$$\nabla E[\vec{w}] \equiv \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

Training rule:

$$\Delta \vec{w} = -\eta \nabla E[\vec{w}]$$

i.e.,

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

Artificial Neural Networks

$$E[\vec{w}] \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$$\begin{aligned} \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\ &= \sum_d (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \vec{w} \cdot \vec{x}_d) \\ \frac{\partial E}{\partial w_i} &= \sum_d (t_d - o_d) (-x_{i,d}) \end{aligned}$$

其中 $x_{i,d}$ 表示训练样例 d 的第 i 个输入分量， o_d 表示感知器输出， t_d 表示训练样例的目标值

梯度下降权值更新法则： $\Delta w_i = \eta \sum_{d \in D} (t_d - o_d) x_{i,d}$

Artificial Neural Networks

- 修改后的训练法则与之前给出的相似，只是在迭代计算每个训练样例时根据下面的公式来更新权值

$$\Delta w_i = \eta \sum_{d \in D} (t_d - o_d) x_{id}$$

$$\Delta w_i = \eta (t - o) x_i$$

与感知器训练
法则相同否？

其中t, o, 和xi分别是目标值、单元输出和第i个训练样例的输入。

训练法则被称为增量
法则（delta rule），
或有时叫LMS法则
（least-mean-square
最小均方）、Adaline
法则、或Windrow-
Hoff法则（以它的发
明者命名）。

Artificial Neural Networks

• 两种算法的区别

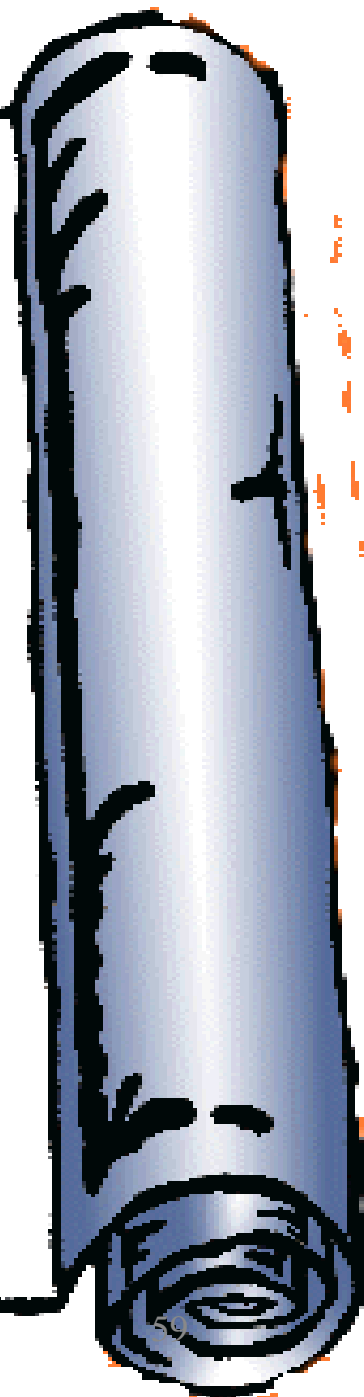
– 权值更新

- 感知器训练法则根据阈值化（thresholded）的感知器输出的误差更新权值
- 增量法则根据输入的非阈值化（unthresholded）线性组合的误差来更新权

– 收敛特性

- 感知器训练法则经过有限次的迭代收敛到一个能理想分类训练数据的假设，但条件是训练样例线性可分。
- 增量法则渐近收敛到最小误差假设，可能需要无限的时间，但无论训练样例是否线性可分都会收敛。

多层网络和反向传播算法

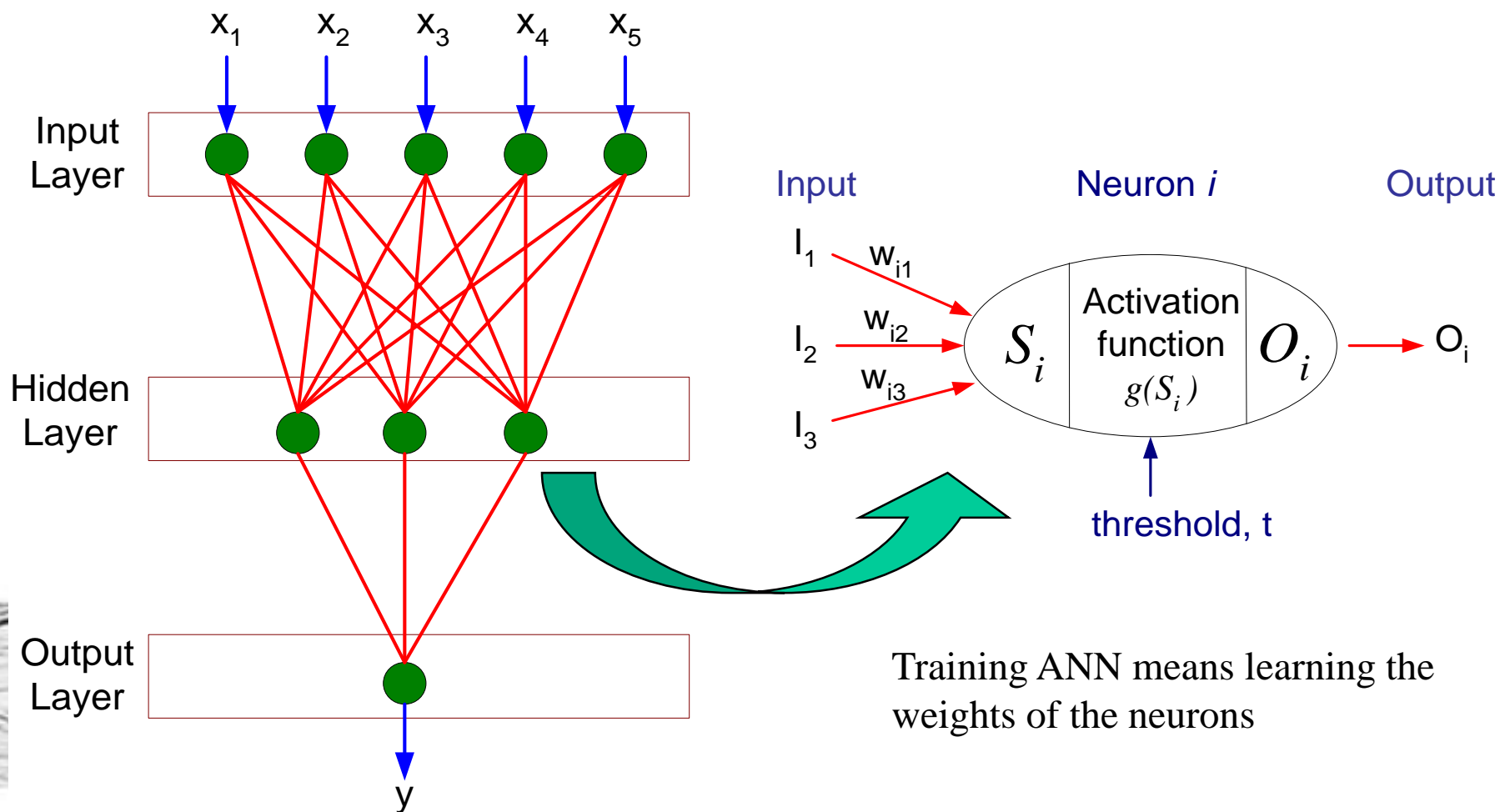


Artificial Neural Networks

- 多层感知机MLP和反向传播算法BP的区别？
- MLP是网络模型，BP是优化算法
- MLP可以很多层，但受限于激活函数，用BP优化时通常只有3层

Artificial Neural Networks

- 多层网络，多层感知器（MLP）

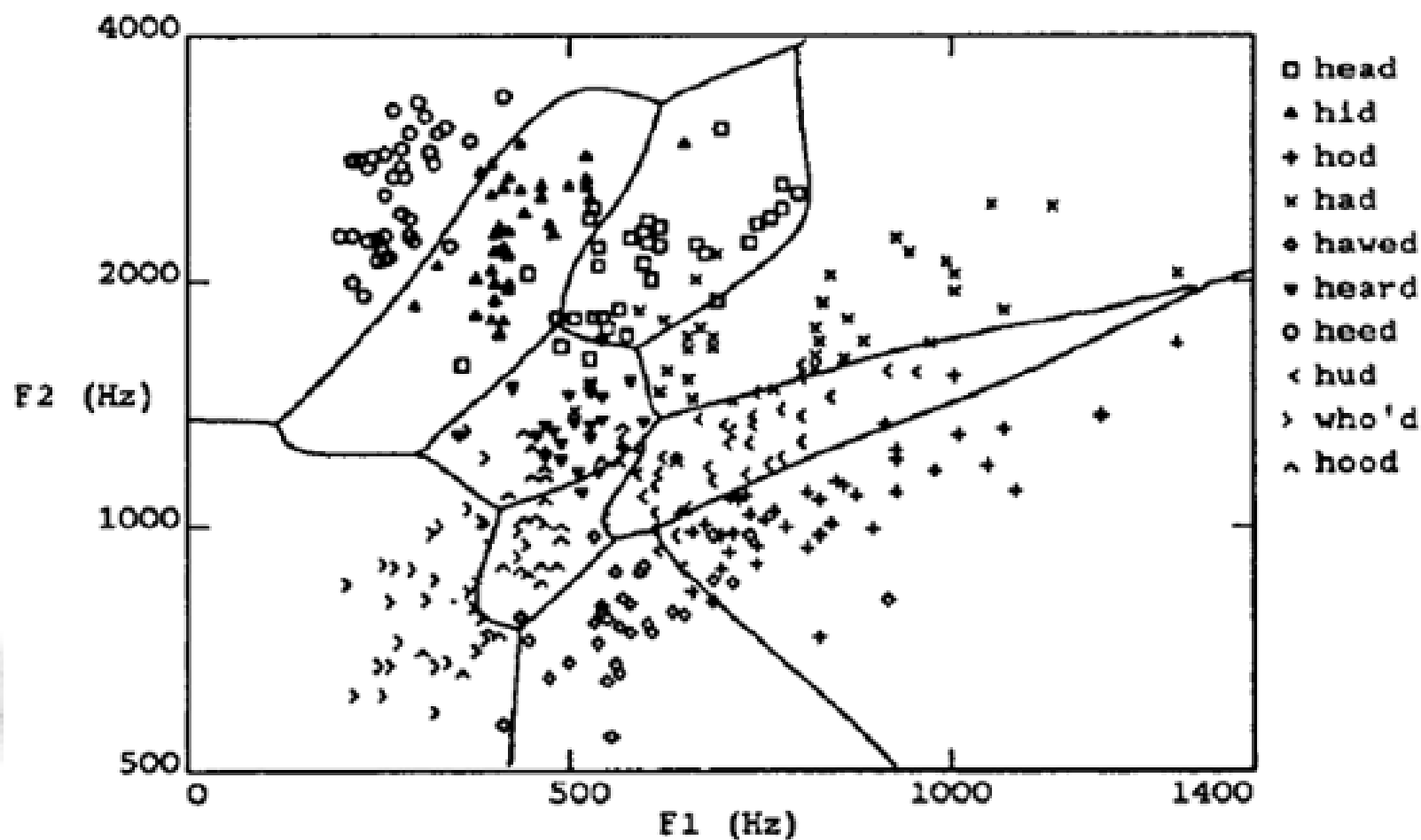


Artificial Neural Networks

- 单个感知器仅能表示线性决策面，而多层网络能够表示种类繁多的非线性曲面。
- 核心是可微阈值单元（即神经元）
 - 多个线性单元的连接仍产生线性函数，但我们希望构建表征非线性函数的网络
 - 修改感知器单元确实可以构建非线性函数，但它的不连续阈值使它不可微，不适合梯度下降算法
 - 神经元需要满足的条件
 - 输出是输入的非线性函数
 - 输出是输入的可微函数
 - 解决方案：引入Sigmoid单元，类似于感知器单元，但基于一个平滑的可微阈值函数

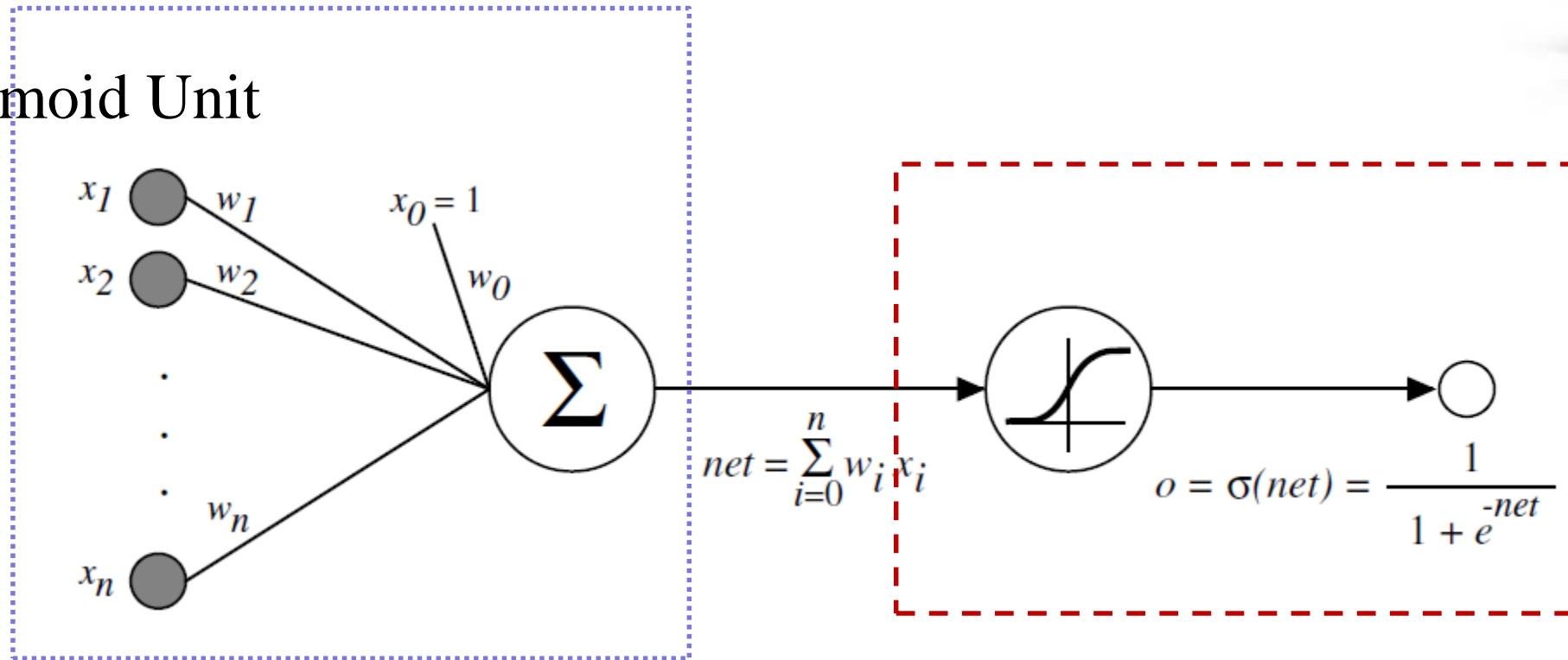
Artificial Neural Networks

- 语音识别任务：区分出现10种元音



Artificial Neural Networks

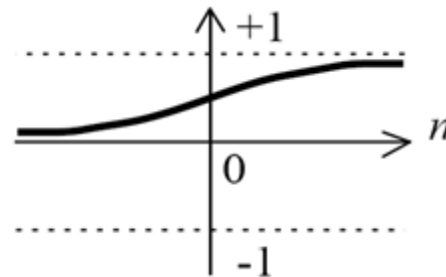
- Sigmoid Unit



$\sigma(x)$ is the sigmoid function

$$\frac{1}{1 + e^{-x}}$$

Artificial Neural Networks



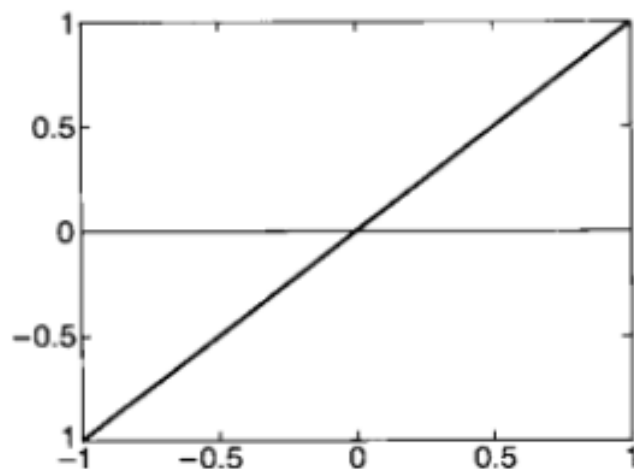
- sigmoid函数

- 也称logistic函数、sigmoid单元的挤压函数（squashing function）
- 输出范围是0到1
- 随输入单调递增
- 导数很容易以它的输出表示

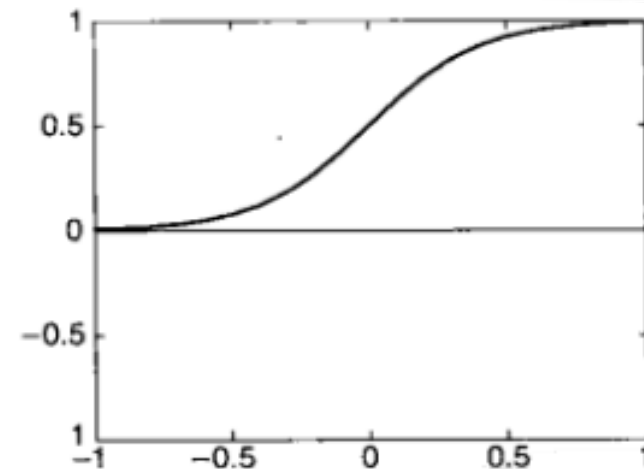
$$\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$$

Artificial Neural Networks

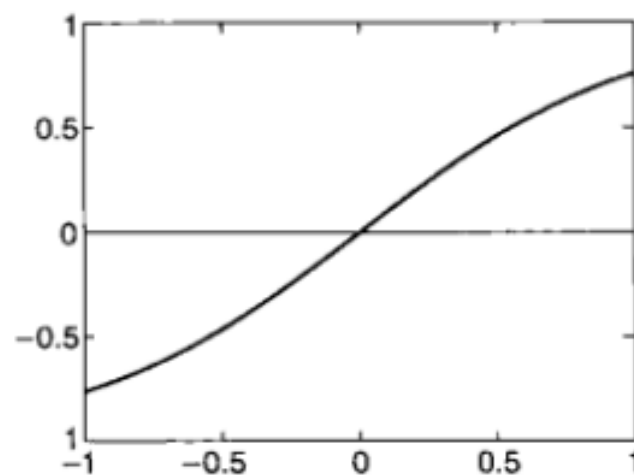
- Sigmoid函数被称为激活函数



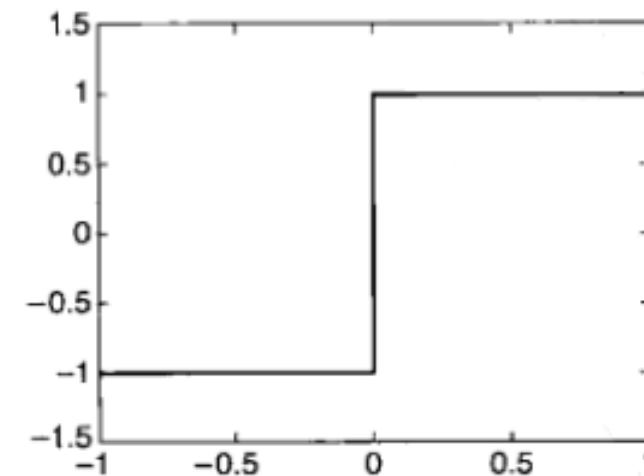
线性函数



S型函数



双曲正切函数



符号函数

人工神经网络中激活函数的类型

Artific

- 一个sigmoid单元的梯度下降法则

$$E[\vec{w}] \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

net_d : 网络输出

o_d : 激活函数输出

$$= \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2$$
$$= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d)$$

$$= \sum_d (t_d - o_d) \left(-\frac{\partial o_d}{\partial w_i} \right)$$

相比于感知机，这里增加了激活函数

$$= - \sum_d (t_d - o_d) \frac{\partial o_d}{\partial net_d} \frac{\partial net_d}{\partial w_i}$$

Artificial Neural Networks

But we know:

$$\frac{\partial o_d}{\partial net_d} = \frac{\partial \sigma(net_d)}{\partial net_d} = o_d(1 - o_d)$$

$$\frac{\partial net_d}{\partial w_i} = \frac{\partial (\vec{w} \cdot \vec{x}_d)}{\partial w_i} = x_{i,d}$$

So:

$$\frac{\partial E}{\partial w_i} = - \sum_{d \in D} (t_d - o_d) o_d (1 - o_d) x_{i,d}$$

$$\Delta w_i = -\eta \frac{\partial E_d}{\partial w_i} = \eta \sum_{d \in D} (t_d - o_d) o_d (1 - o_d) x_{i,d}$$

Artificial Neural Networks

- 多层网络的优化方法：反向传播（Back Propagation）算法

- 对于由一系列确定的单元互连形成的多层网络，反向传播算法可用来学习这个网络的权值。它采用梯度下降方法试图最小化网络输出值和目标值之间的误差平方。
- 网络的误差定义公式，多个输出单元的网络输出的误差：

$$E(\vec{w}) \equiv \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{outputs}} (t_{kd} - o_{kd})^2$$

- 其中outputs是网络输出单元的集合， t_{kd} 和 o_{kd} 是与训练样例d和第k个输出单元相关的输出值。

Artificial Neural Networks

- 反向传播算法面临的学习任务

- 搜索一个巨大的假设空间，这个空间由网络中所有的单元的所有可能的权值定义
- 搜索目标为使目标函数 E 最小化的权值。
- 在多层网络中，误差曲面可能有多个局部极小值，梯度下降不能保证收敛到全局极小值。

Artificial Neural Networks

- 反向传播法则的推导

- 随机梯度下降算法迭代处理训练样例，每次处理一个，对于每个训练样例 d ，利用关于这个样例的误差 E_d 的梯度修改权值

$$E_d(\vec{w}) = \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2$$

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}}$$

Artificial Neural Networks

• 符号说明

- x_{ji} : 单元j的第i个输入
- w_{ji} : 与 x_{ji} 相关联的权值
- net_j : 单元j的输入的加权和
- o_j : 单元j计算出的输出
- t_j : 单元j的目标输出
- σ : sigmoid函数
- outputs: 网络最后一层的输出单元的集合
- **Downstream(j)**: 单元j的输出所能到达单元的集合

Artificial Neural Networks

– 随机梯度下降法则的推导

- 注意权值 w_{ji} 仅能通过 net_j 影响网络的其他部分

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial net_j} \frac{\partial net_j}{\partial w_{ji}} = \frac{\partial E_d}{\partial net_j} x_{ji}$$

– 下面我们分情况讨论

- 单元 j 是网络的一个输出单元
- 单元 j 是网络的一个隐藏单元

Artificial Neural Networks

- 输出单元的权值训练法则
 - netj仅能通过oj影响其余的网络

$$\frac{\partial E_d}{\partial net_j} = \frac{\partial E_d}{\partial o_j} \frac{\partial o_j}{\partial net_j}$$

$$\begin{aligned}\frac{\partial E_d}{\partial o_j} &= \frac{\partial}{\partial o_j} \frac{1}{2} (t_j - o_j)^2 \\ &= \frac{1}{2} 2(t_j - o_j) \frac{\partial (t_j - o_j)}{\partial o_j} \\ &= -(t_j - o_j)\end{aligned}$$

Artificial Neural Networks

$$o_j = \sigma(\text{net}_j)$$

sigmoid函数的导数为 $\sigma(\text{net}_j)(1-\sigma(\text{net}_j))$

$$\frac{\partial o_j}{\partial \text{net}_j} = \frac{\partial \sigma(\text{net}_j)}{\partial \text{net}_j} = o_j(1 - o_j)$$

$$\begin{aligned} \frac{\partial E_d}{\partial \text{net}_j} &= \frac{\partial E_d}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} \\ &= -(t_j - o_j) o_j (1 - o_j) \end{aligned}$$

Artificial Neural Networks

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}}$$

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial net_j} x_{ji}$$

$$\frac{\partial E_d}{\partial net_j} = -(t_j - o_j) o_j (1 - o_j)$$

– 最终得到输出单元的随机梯度下降法则：

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}} = \eta x_{ji} (t_j - o_j) o_j (1 - o_j)$$

Artificial Neural Networks

- 隐藏单元的权值训练法则

- 对于网络中的内部单元或者说隐藏单元的情况，推导 w_{ji} 必须考虑 w_{ji} 间接地影响网络输出，从而影响 E_d 。
- net_j 只能通过 $Downstream(j)$ 中的单元影响网络输出（再影响 E_d ）

$$\frac{\partial E_d}{\partial net_j} = \sum_{k \in Downstream(j)} \frac{\partial E_d}{\partial net_k} \frac{\partial net_k}{\partial net_j}$$

Artificial Neural Networks

$$\begin{aligned}\frac{\partial E_d}{\partial net_j} &= \sum_{k \in \text{Downstream}(j)} \frac{\partial E_d}{\partial net_k} \frac{\partial net_k}{\partial net_j} \\&= \sum_{k \in \text{Downstream}(j)} -\delta_k \frac{\partial net_k}{\partial net_j} \\&= \sum_{k \in \text{Downstream}(j)} -\delta_k \frac{\partial net_k}{\partial o_j} \frac{\partial o_j}{\partial net_j} \\&= \sum_{k \in \text{Downstream}(j)} -\delta_k w_{kj} \frac{\partial o_j}{\partial net_j} \\&= \sum_{k \in \text{Downstream}(j)} -\delta_k w_{kj} o_j (1 - o_j) \\&= -o_j (1 - o_j) \sum_{k \in \text{Downstream}(j)} \delta_k w_{kj}\end{aligned}$$

$$\Delta w_{ji} = \eta x_{ji} o_j (1 - o_j) \sum_{k \in \text{Downstream}(j)} \delta_k w_{kj}$$

$$\delta_k = -\frac{\partial E_d}{\partial net_k}$$

$$\delta_j = o_j (1 - o_j) \sum_{k \in \text{Downstream}(j)} \delta_k w_{kj}$$

– 隐含单元的随机
梯度下降法则

Artificial

delta法则 $\Delta w_i = \eta \sum_{d \in D} (t_d - o_d) x_{id}$ BP算法:

输出单元 $\Delta w_{ki} = \eta x_{ki} (t_k - o_k) o_k (1 - o_k)$

隐藏单元 $\Delta w_{hi} = \eta x_{hi} o_h (1 - o_h) \sum_{k \in \text{Downstream}(j)} \delta_k w_{kh}$

- 反向传播算法与delta训练法则很相似，均是依照以下三者的乘积来更新每一个权值w：学习速率 η 、该权值应用的输入值 x_{ji} 、这个单元输出的误差。
- 区别：反向传播算法与delta法则不同点是误差项被替换成一个更复杂的误差项 δ_j 。
- 输出单元k的误差项
 - δ_k 与delta法则中的 $(t_k - o_k)$ 相似，但乘上了sigmoid挤压函数的导数 $o_k(1 - o_k)$ 。
- 隐藏单元h的误差项
 - 因为训练样例仅对网络的输出提供了目标值 t_k ，所以缺少直接的目标值来计算隐藏单元的误差值，采取以下的间接方法计算隐藏单元的误差项：对受隐藏单元h影响的每一个单元的误差 δ_k 进行加权求和，每个误差 δ_k 权值为 w_{kh} ， w_{kh} 就是从隐藏单元h到输出单元k的权值。这个权值刻画了隐藏单元h对于输出单元k的误差应负责的程度。

Artificial Neural Networks

- BP可以学习任意的无环网络

- 算法可以简单地推广到任意深度的无环前馈网络
- 第m层的单元r的 δ_r 值由更深的第m+1层 δ 值根据下式计算

$$\delta_r = o_r(1 - o_r) \sum_{s \in m+1 \text{ 层}} w_{sr} \delta_s$$

- 将这个算法推广到任何有向无环结构也同样简单，而不论网络中的单元是否被排列在统一的层上，计算任意内部单元的 δ 的法则是：

$$\delta_r = o_r(1 - o_r) \sum_{s \in \text{Downstream}(r)} w_{sr} \delta_s$$

- $\text{Downstream}(r)$ 是在网络中单元r的直接下游单元的集合，即输入中包括r的输出所有单元

Artificial Neural Networks

• 本章小结:

- 人工神经网络为学习实数值和向量值函数提供了一种实际的方法，对于连续值和离散值的属性都可以使用，并且对训练数据中的噪声具有很好的稳健性。
- 反向传播算法是最常见的网络学习算法
- 反向传播算法考虑的假设空间是固定连接的有权网络所能表示的所有函数的空间
- 包含3层单元的前馈网络能够以任意精度逼近任意函数，只要每一层有足够数量的单元。

Artificial Neural Networks

- 反向传播算法使用梯度下降方法搜索可能假设的空间，迭代减小网络的误差以拟合训练数据。
- 梯度下降收敛到训练误差相对网络权值的局部极小值。只要训练误差是假设参数的可微函数，梯度下降可用来搜索很多连续参数构成的假设空间。



THE END !