

2020.2.18 - 2020.6.2 9:50-12:15@6C102雨课堂+腾讯会议
《模式识别与机器学习》



第四章 人工神经网络

2020.3.3



Xuegong Zhang



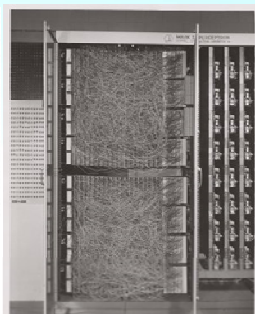
感知器

- 为什么把它叫做学习机器？

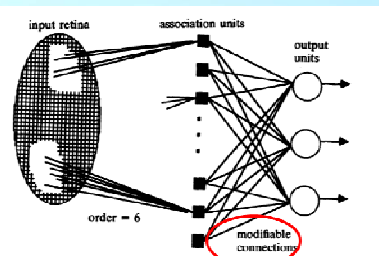
- ① 因为它是一台机器
- ② 因为它会学习！

- 它不是编好程序的冯诺依曼计算机，是一台根据训练数据自我调整的学习机器

$$y = \text{sgn}\left(\sum_{i=1}^n w_i x_i + w_0\right)$$



<https://en.wikipedia.org/wiki/Perceptron>



M. Clelland, A sociological study of the official history of the perceptron controversy, *Social Studies of Science*, 1996

Xuegong Zhang

2



ADALINE

Widrow & Hoff, Adaptive switching circuits, 1960 IRE Western Electric Show and Convention Record, Part 4, pp.96-104, Aug, 1960

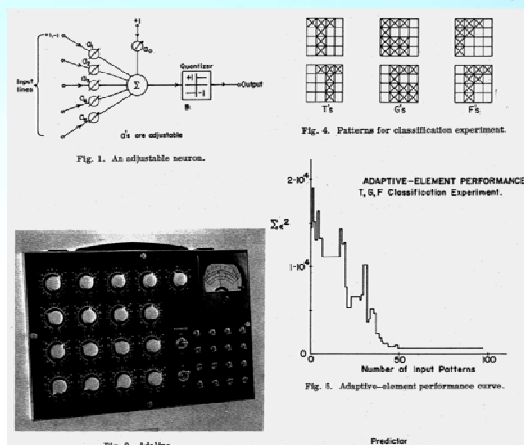


Fig. 2. Adaline. Widrow & Hoff, Adaptive switching circuits, 1960 IRE Western Electric Show and Convention Record, Part 4, pp.96-104, Aug, 1960

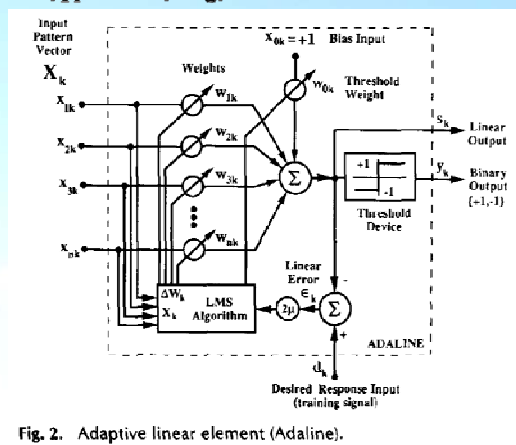


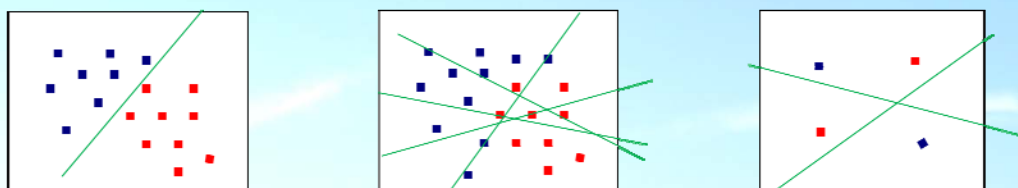
Fig. 2. Adaptive linear element (Adaline).

Widrow & Lehr, 30 years of adaptive neural networks: Perceptron, Madaline, and Backpropagation, *Proceedings of the IEEE*, 78(9): 1415-1442, 1990

Xuegong Zhang

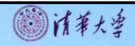
3

线性不可分情况怎么办？

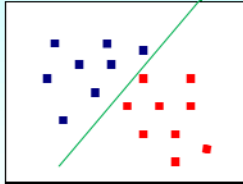


Xuegong Zhang

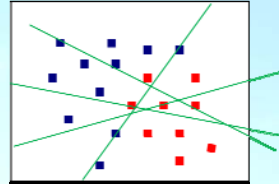
4



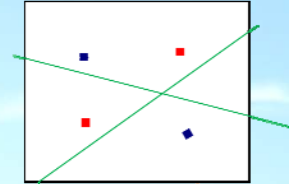
线性不可分情况怎么办？



线性可分



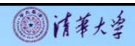
线性不可分
但错误可容忍



线性不可分
且无调和余地

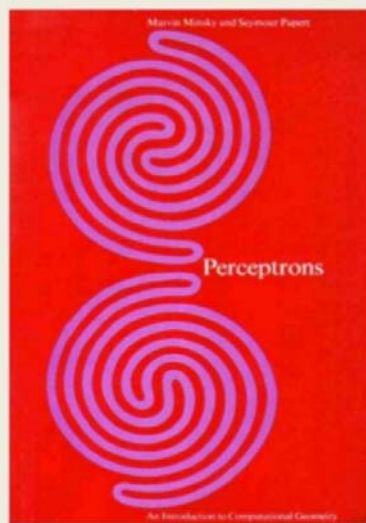
Xuegong Zhang

5

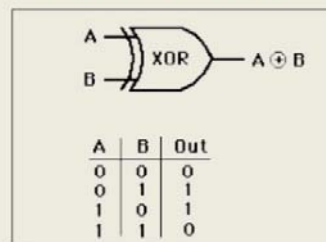


对感知器（乃至整个机器学习）的质疑

1969: Perceptrons can't do XOR!



<http://www.i-programmer.info/images/stories/Babag/AI/book.jpg>



<http://hyperphysics.phy-astr.gsu.edu/hbase/electronic/tetron/xor.gif>



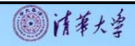
Minsky & Papert

<https://constructingkids.files.wordpress.com/2013/05/minsky-papert-71-csolumon-x640.jpg>

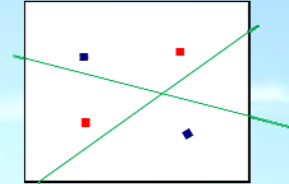
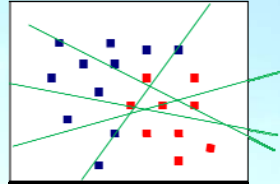
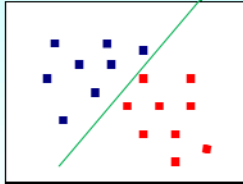
Xuegong Zh

<https://pmirla.github.io/2016/08/16/AI-Winter.html>

6



线性不可分情况怎么办？



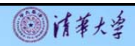
有什么办法？



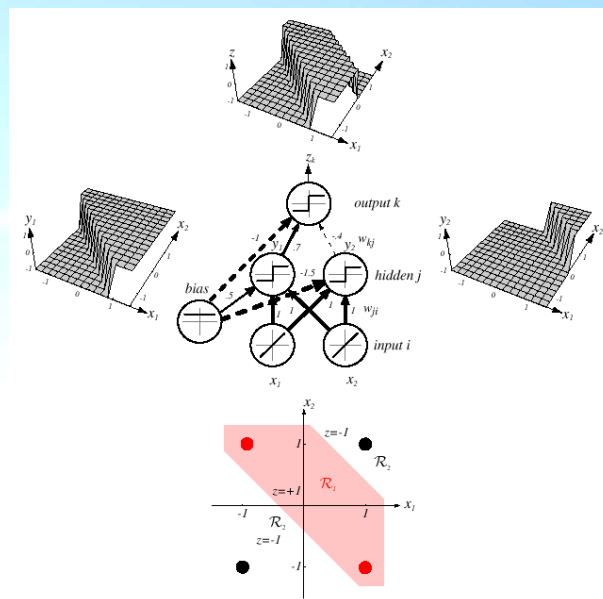
请弹幕回答

Xuegong Zhang

7



用多个感知器构造非线性分类器

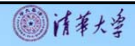


分段线性分类的思想

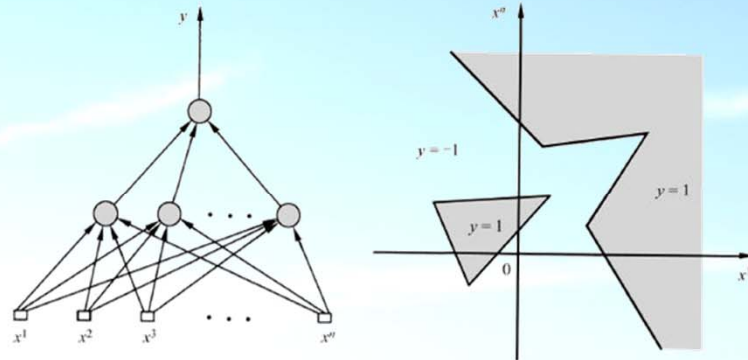
Duda, Hart & Stork, Pattern Classification (2nd edition), John Wiley & Sons, 2001

6

8

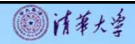


早期的多层感知器（神经网络）



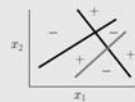
Xuegong Zhang

9



用感知器布尔分解非线性函数

Consider a target function f whose '+' and '-' regions are illustrated below.

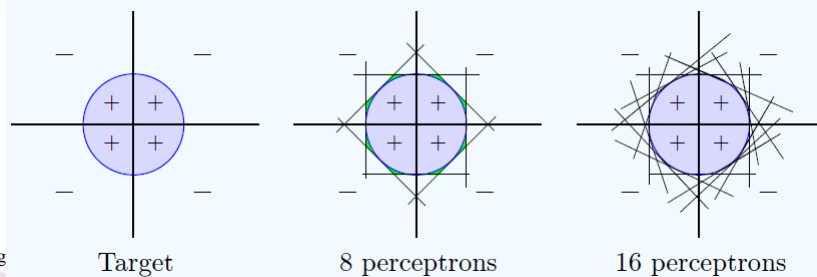
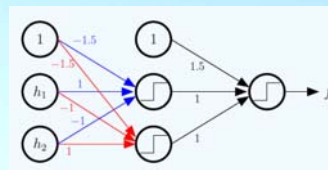


The target f has three perceptron components h_1, h_2, h_3 :



Show that

$$f = \overline{h_1}h_2h_3 + h_1\overline{h_2}h_3 + h_1h_2\overline{h_3}.$$



Xuegong Zhang

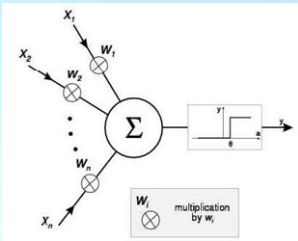
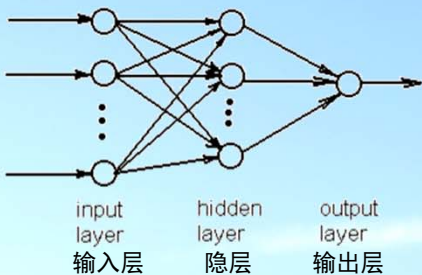
Target

8 perceptrons

16 perceptrons

10

清华大学

非线性判别函数

$$g_k(\mathbf{x}) \equiv y_k = \text{Sign} \left(\sum_j w_{kj} \text{Sign} \left(\sum_i w_{ji} x_i + w_{j0} \right) + w_{k0} \right)$$

Xuegang Zhang
11

清华大学

机器学习的基本要素：感知器版

- 怎样造一个学习机器？
 - 它需要老师
 - 我们设计它（特征和模型） $y = \text{sgn}(\sum_{i=1}^n w_i x_i + w_0)$
 - 它需要训练/学习材料
 - 训练数据 $\{(x_1, y_1), \dots, (x_n, y_n)\}$, $x_j \in \mathbb{R}^{d+1}$, $y_j \in \{-1, 1\}$
 - 我们需要为它树立学习的目标
 - 目标函数、学习准则 $\min J_p(\mathbf{w}) = \sum_{j \in \mathcal{P}} (-\alpha^j y_j)$
 - 我们需要告诉它怎样学
 - 学习/训练算法 $\mathbf{w}(k+1) = \mathbf{w}(k) - \rho_k \nabla J = \mathbf{w}(k) + \rho_k \sum_{j \in \mathcal{P}} y_j$

清华大学

机器学习的基本要素：罗杰斯特回归版

- 怎样造一个学习机器？
 - 它需要老师
 - 我们设计它（特征和模型） $h(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$
 - 它需要训练/学习材料
 - 训练数据 $\{(x_1, y_1), \dots, (x_n, y_n)\}$, $x_j \in \mathbb{R}^{d+1}$, $y_j \in \{-1, 1\}$
 - 我们需要为它树立学习的目标
 - 目标函数、学习准则 $\min E(\mathbf{w}) = -\frac{1}{n} \sum_{j=1}^n \ln(1 + e^{-y_j \mathbf{w}^T x_j})$
 - 我们需要告诉它怎样学
 - 学习/训练算法 $\mathbf{w}(k+1) = \mathbf{w}(k) - \eta \nabla E$

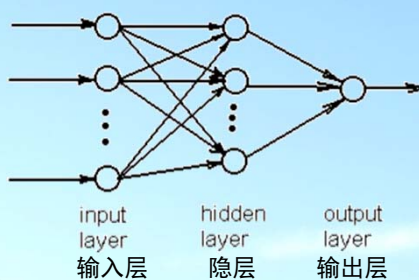
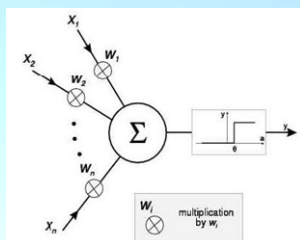
机器学习的基本要素：线性回归版

- 怎样造一个学习机器？
 - 它需要老师
 - 我们设计它（特征和模型） $f(\mathbf{x}) = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$
 - 它需要训练/学习材料
 - 训练数据 $\{(x_1, y_1), \dots, (x_n, y_n)\}$, $x_j \in \mathbb{R}^{d+1}$, $y_j \in \mathbb{R}$
 - 我们需要为它树立学习的目标
 - 目标函数、学习准则 $\min E = \frac{1}{n} \sum_{j=1}^n (f(x_j) - y_j)^2$
 - 我们需要告诉它怎样学
 - 学习/训练算法 $\mathbf{w}(k+1) = \mathbf{w}(k) - \rho_k \nabla E$

多层感知器版？

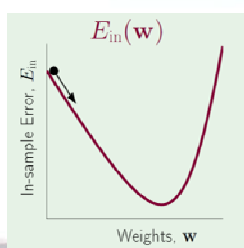
D

Xuegang Zhang
12



非线性判别函数

$$g_k(\mathbf{x}) \equiv y_k = \text{Sign} \left(\sum_j w_{kj} \text{Sign} \left(\sum_i w_{ji} x_i + w_{j0} \right) + w_{k0} \right)$$



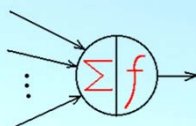
- 梯度下降，老套路还行吗？
 - 无法对中间权值求梯度



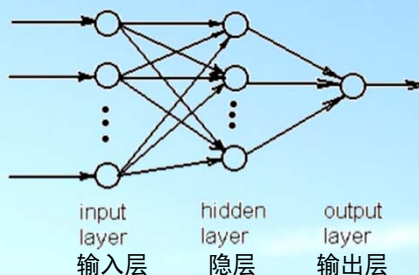
Xuegong Zhang

13

怎么办？



f : 激活函数/响应函数



非线性判别函数

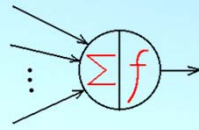
$$g_k(\mathbf{x}) \equiv y_k = f \left(\sum_j w_{kj} f \left(\sum_i w_{ji} x_i + w_{j0} \right) + w_{k0} \right)$$



Xuegong Zhang

14

怎么办？



f : 激活函数/响应函数

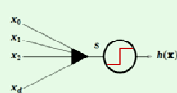
非线性判别函数

$$g_k(\mathbf{x}) \equiv y_k = f\left(\sum_j w_{kj} f\left(\sum_i w_{ji} x_i + w_{j0}\right) + w_{k0}\right)$$

$$s = \sum_{i=0}^d w_i x_i$$

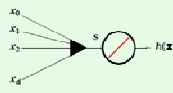
linear classification

$$h(\mathbf{x}) = \text{sign}(s)$$



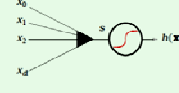
linear regression

$$h(\mathbf{x}) = s$$



logistic regression

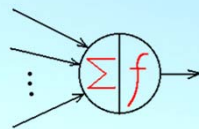
$$h(\mathbf{x}) = \theta(s)$$



Xuegong Zhang

15

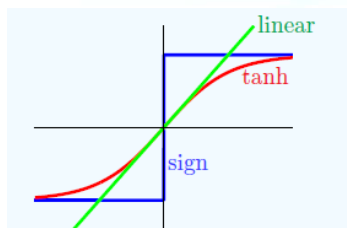
怎么办？



f : 激活函数/响应函数

非线性判别函数

$$g_k(\mathbf{x}) \equiv y_k = f\left(\sum_j w_{kj} f\left(\sum_i w_{ji} x_i + w_{j0}\right) + w_{k0}\right)$$



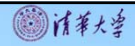
- 阶跃函数不行
- 线性激活函数行不行？



See 6.2 Problem 1, Duda et al, Pattern Classification, p. 335

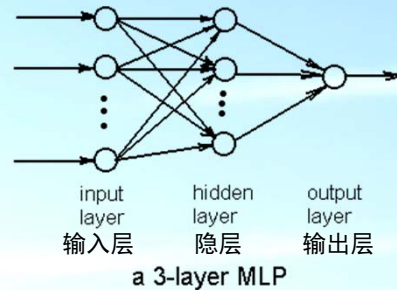
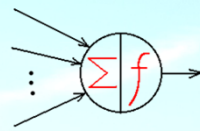
Xuegong Zhang

16



多层感知器MLP (Multi-layer Perceptron)

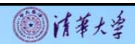
--- 最经典的人工神经网络Artificial Neural Network (ANN)



$$g_k(\mathbf{x}) \equiv y_k = f \left(\sum_j w_{kj} f \left(\sum_i w_{ji} x_i + w_{j0} \right) + w_{k0} \right)$$

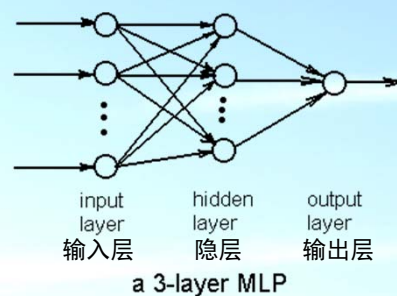
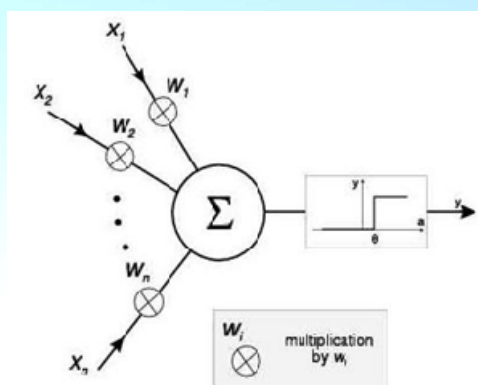
Xuegang Zhang

17



多层感知器MLP (Multi-layer Perceptron)

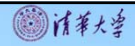
--- 最经典的人工神经网络Artificial Neural Network (ANN)



$$g_k(\mathbf{x}) \equiv y_k = f \left(\sum_j w_{kj} f \left(\sum_i w_{ji} x_i + w_{j0} \right) + w_{k0} \right)$$

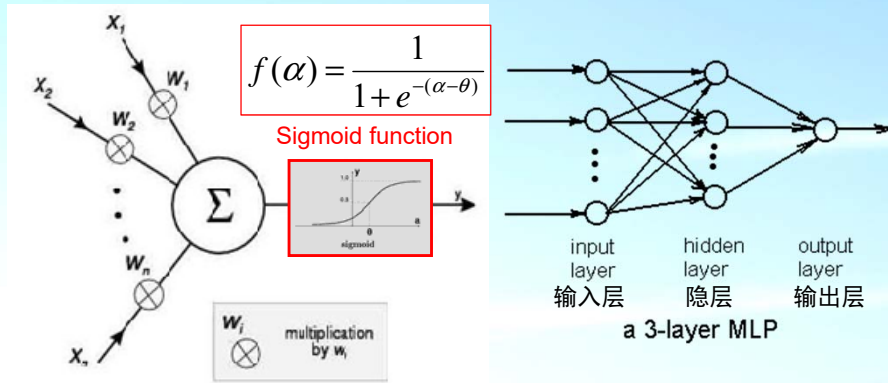
Xuegang Zhang

18



多层感知器MLP (Multi-layer Perceptron)

--- 最经典的人工神经网络Artificial Neural Network (ANN)



$$g_k(\mathbf{x}) \equiv y_k = f\left(\sum_j w_{kj} f\left(\sum_i w_{ji} x_i + w_{j0}\right) + w_{k0}\right)$$

Xuegong Zhang

19

多层感知器实现非线性的直观理解

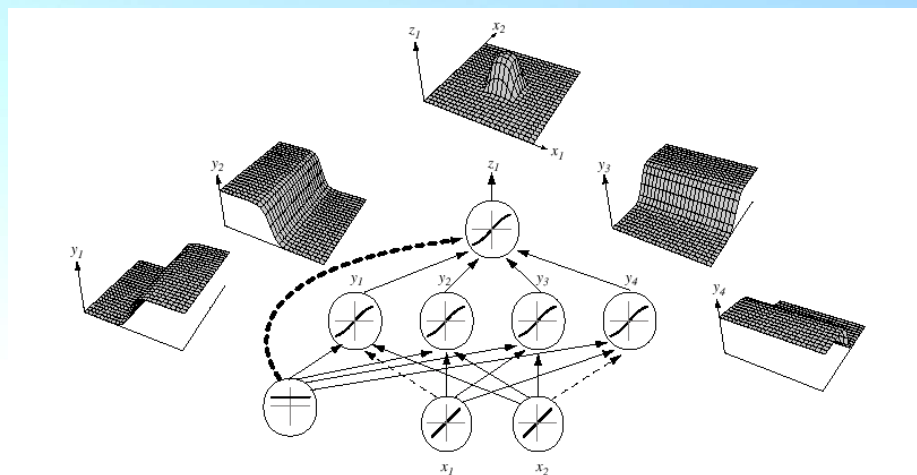
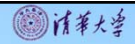


FIGURE 6.2. A 2-4-1 network (with bias) along with the response functions at different units; each hidden output unit has sigmoidal activation function $f(\cdot)$. In the case shown, the hidden unit outputs are paired in opposition thereby producing a “bump” at the output unit. Given a sufficiently large number of hidden units, any continuous function from input to output can be approximated arbitrarily well by such a network. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Xuegong Zha...

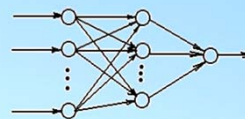
Duda, Hart & Stork, Pattern Classification (2nd edition), John Wiley & Sons, 2001

20



多层感知器的表示能力

$$g_k(\mathbf{x}) \equiv y_k = f \left(\sum_j w_{kj} f \left(\sum_i w_{ji} x_i + w_{j0} \right) + w_{k0} \right)$$



理论结果:

多层感知器能够实现任意连续的输入输出函数，只要有充分数目的隐层节点、合适的激活函数和权值。

- From the Kolmogorov theorem:

Any continuous function $g(\mathbf{x})$ defined on the unit hypercube I^n ($I=[0,1]$ and $n \geq 2$) can be represented in the form

$$g(\mathbf{x}) = \sum_{j=1}^{2n+1} \Xi_j \left(\sum_{i=1}^d \Psi_{ij}(x_i) \right)$$

for properly chosen functions Ξ_j and Ψ_{ij} .

Duda, Hart & Stork, Pattern Classification, John Wiley & Sons, 2001, p.287

Xuegong Zhang

21



怎样设计用作模式识别的多层感知器?

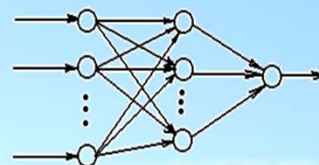
- 设计一个合适的网络

— 编码:

- 两类 0 vs. 1
- 多类 1-of-C (e.g., 1000, 0100, 0010, 0001) **one-hot编码**

— 结构:

- 三层网（一个隐层）可实现空间内任意的凸形区域划分
- 四层网（两个隐层）可实现任意形状（连续或不连续）区域划分



Xuegong Zhang

22

怎样设计用作模式识别的多层感知器？



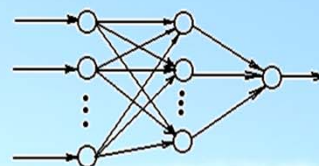
• 设计一个合适的网络

– 编码：

- 两类 0 vs. 1
- 多类 1-of-C (e.g., 1000, 0100, 0010, 0001)

– 结构：

- 三层网（一个隐层）可实现空间内任意的凸形成区域的划分
- 四层网（两个隐层）可实现任意形状（连续或不连续）区域划分



但是，
多少个隐节点？

Xuegang Zhang

23

怎样设计用作模式识别的多层感知器？



• 设计一个合适的网络

– 编码：

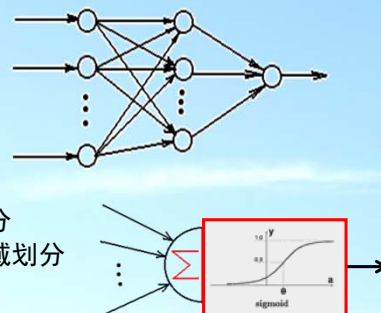
- 两类 0 vs. 1
- 多类 1-of-C (e.g., 1000, 0100, 0010, 0001)

– 结构：

- 三层网（一个隐层）可实现空间内任意的凸形成区域的划分
- 四层网（两个隐层）可实现任意形状（连续或不连续）区域划分

• 选择合适的结点函数

- Sigmoid函数 或 tanh函数 **双曲正切函数**
- ReLU ...



Xuegang Zhang

24

怎样设计用作模式识别的多层感知器？



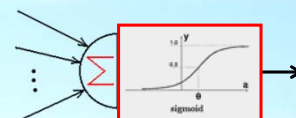
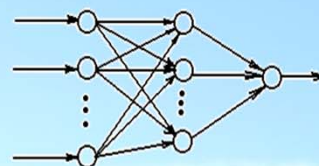
• 设计一个合适的网络

– 编码：

- 两类 0 vs. 1
- 多类 1-of-C (e.g., 1000, 0100, 0010, 0001)

– 结构：

- 三层网（一个隐层）可实现空间内任意的凸形成区域的划分
- 四层网（两个隐层）可实现任意形状（连续或不连续）区域划分

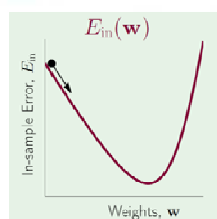


• 选择合适的结点函数

- Sigmoid函数 或tanh函数
- ReLU ...

• 训练神经网络参数

- 用什么训练？→训练样本
- 怎样训练？



→ 问题：当很多人集体完成一件事而出错的时候，怎样把错误追究到每个人？

Xuegong Zhang

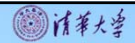
25

休息1分钟



Xuegong Zhang

26



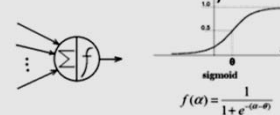
误差反向传播 Back-Propagation (BP) of the Error

LeCun, 1986; Rumelhart, Hinton & Williams, 1986; Parker, 1985



节点响应

for node j , $net_j = \sum_i w_{ij} O_i$, $O_j = f(net_j)$



实际输出
期望输出

for output node: $\hat{y}_j = O_j$ is the actual output, y_j is the desired output, error

误差

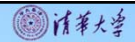
error: $E = \frac{1}{2} \sum_j (y_j - \hat{y}_j)^2$

误差梯度

gradients: $\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial net_j} \frac{\partial net_j}{\partial w_{ij}} = \delta_j O_i$, $\delta_j \equiv \frac{\partial E}{\partial net_j}$

Xuegong Zhang

27



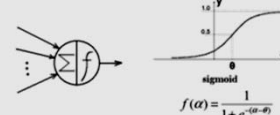
误差反向传播 Back-Propagation (BP) of the Error

LeCun, 1986; Rumelhart, Hinton & Williams, 1986; Parker, 1985



节点响应

for node j , $net_j = \sum_i w_{ij} O_i$, $O_j = f(net_j)$



实际输出
期望输出

for output node: $\hat{y}_j = O_j$ is the actual output, y_j is the desired output, error

误差


error: $E = \frac{1}{2} \sum_j (y_j - \hat{y}_j)^2$ 目标函数：平方误差

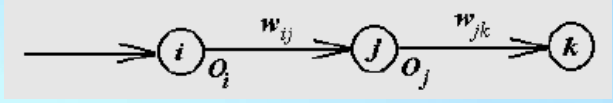
误差梯度

gradients: $\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial net_j} \frac{\partial net_j}{\partial w_{ij}} = \delta_j O_i$, $\delta_j \equiv \frac{\partial E}{\partial net_j}$

Xuegong Zhang

28





误差梯度

对输出节点

对隐层节点

权值学习

学习步长

gradients:

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial net_j} \frac{\partial net_j}{\partial w_{ij}} = \delta_j O_i, \quad \delta_j \equiv \frac{\partial E}{\partial net_j}$$

for output nodes: $O_j = \hat{y}_j, \quad \delta_j = \frac{\partial E}{\partial \hat{y}_j} \frac{\partial \hat{y}_j}{\partial net_j} = -(y_j - \hat{y}_j) f'(net_j)$


for hidden nodes (back-propagation):

$$\delta_j = \frac{\partial E}{\partial net_j} = \sum_k \frac{\partial E}{\partial net_k} \cdot \frac{\partial net_k}{\partial O_j} \cdot \frac{\partial O_j}{\partial net_j} = \sum_k \delta_k w_{jk} f'(net_j)$$


adaption (at iteration t):

$$w_{ij}(t+1) = w_{ij}(t) + \Delta w_{ij}(t), \quad \Delta w_{ij}(t) = -\eta \delta_j(t) O_i(t)$$

where: η - step length (learning rate)




Xuegong Zhang
29



误差反向传播 (BP)算法

LeCun, 1986; Rumelhart, Hinton & Williams, 1986; Parker, 1985

- Purpose: to train MLP weights with data
- Goal: **to minimize the error** (MSE)
- Algorithm: *gradient descent*
- Basic steps:
 - **forward**: compute the output with certain input
input nodes >>> hidden nodes >>> output nodes
 - **compare**: discrepancy between output and desired
 - **backward**: propagate the error back through the network
input nodes <<< hidden nodes <<< output nodes
 - **adaptation**: adapt the weights according to the errors



信号正向传播

计算误差

误差反向传播

梯度下降学习

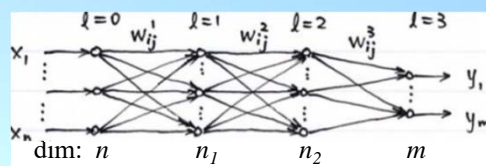
Xuegong Zhang
30

Sigmoid激活函数下BP算法的伪代码



1°. Initialize weights (with small random values), $t=0$

2°. Apply a training sample $\mathbf{x} = [x_1, \dots, x_n]^T \in R^n$ with desired output $D = [d_1, \dots, d_m]^T \in R^m$



3°. Forward calculation: $Y = [y_1, \dots, y_m]^T \in R^m$,

$$y_l = f\left(\sum_{j=1}^{n_2} w_{jl} f\left(\sum_{j=1}^{n_1} w_{jk} f\left(\sum_{i=1}^n w_{ij} x_i\right)\right)\right), l = 1, \dots, m$$

4°. Adjust weights from the output layer. For layer l ,

$$w_{ij}^l(t+1) = w_{ij}^l(t) + \eta \delta_j^l x_i^{l-1}, j = 1, \dots, n_l, i = 1, \dots, n_{l-1}$$

Where for the output layer

$$\delta_j^l = y_j(1 - y_j)(d_j - y_j), j = 1, \dots, m$$

$$f'(\alpha) = f(\alpha)(1 - f(\alpha))$$

for $f(\alpha) = 1/(1 + e^{-\alpha})$

and for hidden layers

$$\delta_j^l = x_j^l(1 - x_j^l) \sum_{k=1}^{n_{l+1}} \delta_k^{l+1} w_{jk}^{l+1}(t), j = 1, \dots, n_l$$

5°. Loop: if stop criterion not met, set $t=t+1$ and go to 2° with another sample.

Xuegong Zhang

31

机器学习的基本要素：多层感知器版



• 怎样造一个学习机器？

– 它需要老师

→ 我们设计它（特征和模型） $g_k(\mathbf{x}) = f(\sum_j w_{kj} f(\sum_i w_{ji} x_i + w_{j0}) + w_{k0})$

– 它需要训练/学习材料

→ 训练数据 $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, $\mathbf{x}_j \in R^{d+1}, y_j \in \{0,1\}$ or $[0,1]$

– 我们需要为它树立学习的目标

→ 目标函数、学习准则 $\min E(\mathbf{w}) = -\frac{1}{2} \sum_{j=1}^N (y - \hat{y})^2$

– 我们需要告诉它怎样学

→ 学习/训练算法 **BP算法** $w_{ij}^l(t+1) = w_{ij}^l(t) + \eta \delta_j^l x_i^{l-1}$

Xuegong Zhang

32



常见训练策略

- 逐一训练
 - 每个样本对网络逐一训练
- 随机训练
 - 从训练集中随机抽取样本进行训练，每次抽取一个或一组
- 整批训练
 - 所有样本一起训练（即把所有样本都输入一遍后再用累计误差进行训练）
- 查询式训练
 - 依据神经网络的输出挑选训练样本
 - “主动学习（Active Learning）”

Xuegong Zhang

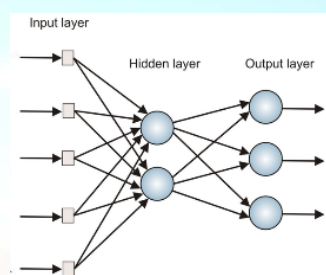
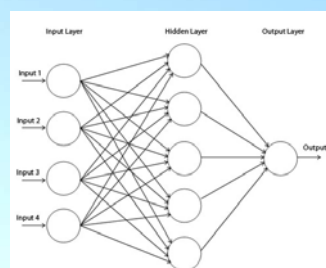
Duda, Hart & Stork, Pattern Classification, John Wiley & Sons, 2001, p.293-295, 480-481

33

多层感知器能学习干什么？

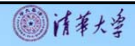


- Almost Everything!
 - 模式识别
 - 把类别编码称输出
 - 回归/函数估计/预测
 - 输出为实数
 - 数据压缩表示
 - 自己学自己：把输入当作输出



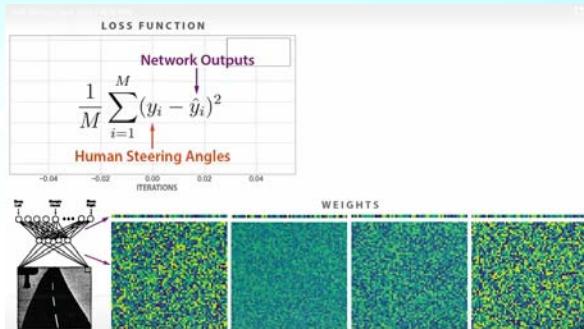
Xuegong Zhang

34



视频：早年自动驾驶实例 ALVINN and Autonomous Navigation in 1990s

从ALV到ALVNN



Self driving cars [S1E2 ALVINN].mp4 ~11:12

基于神经网络的自动驾驶



Stanford CS229 Lecture 2, 3:21-8:40

Xuegong Zhang

35



休息1分钟



Xuegong Zhang

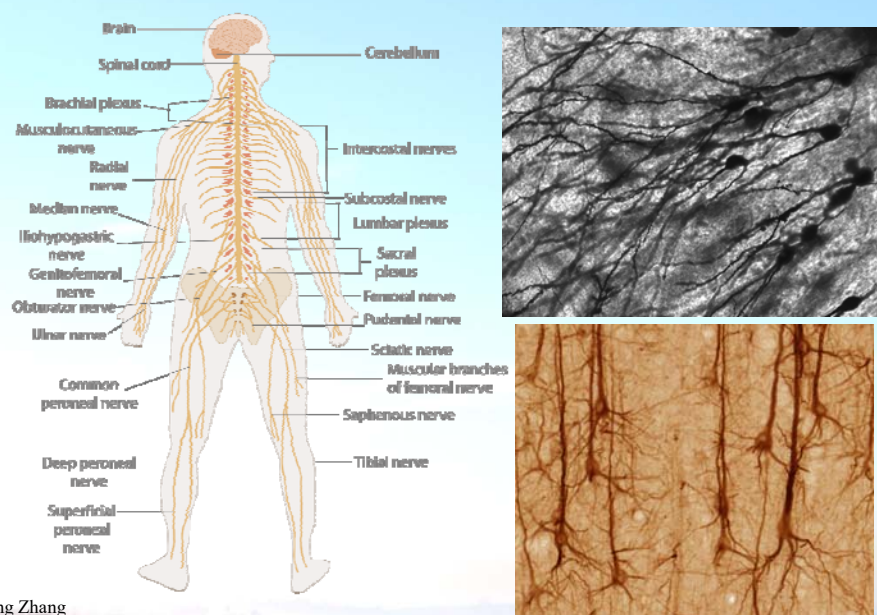
36

为什么称作神经网络？

Xuegong Zhang

37

Neurons 神经元

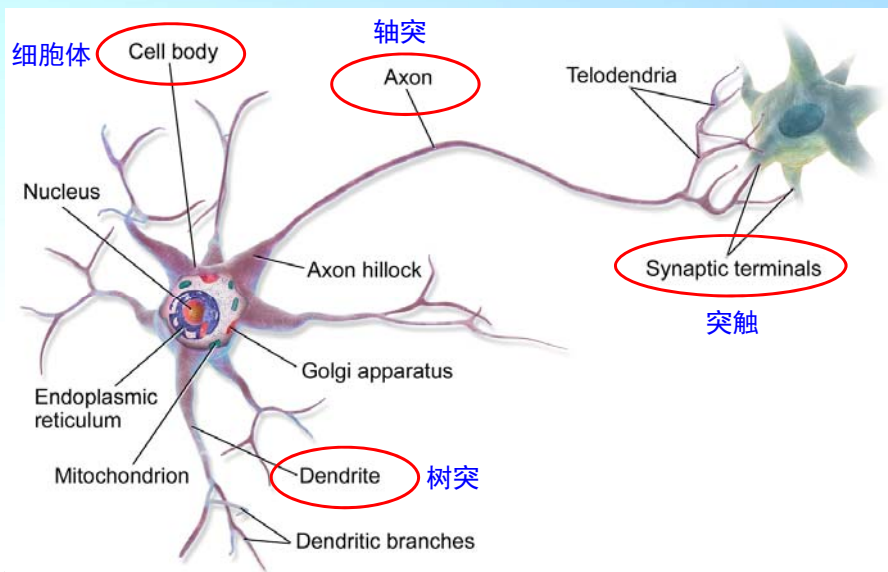
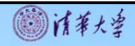


Xuegong Zhang

-- Wikipedia

38

Neurons 神经元

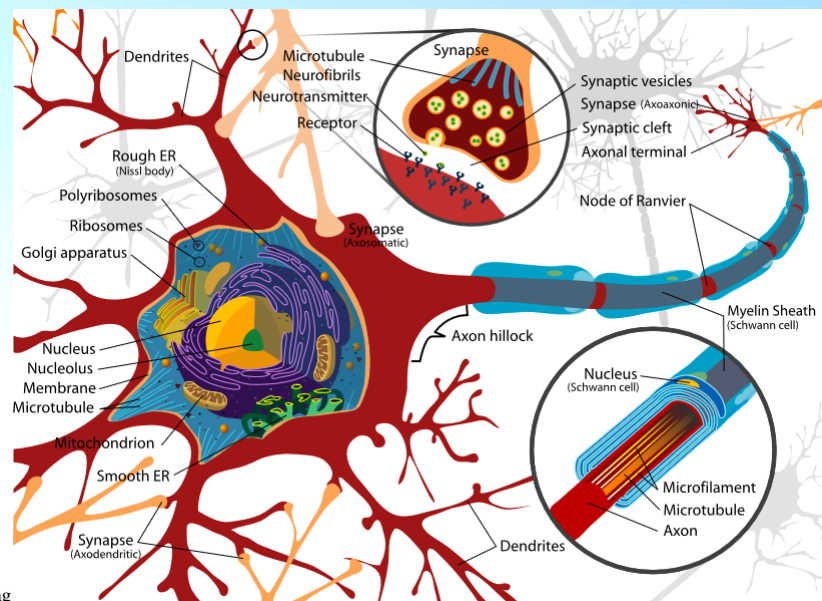
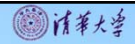


Xuegong Zi

39

<https://en.wikipedia.org/wiki/Neuron>

Neurons 神经元



Xuegong Zhang

40

<https://en.wikipedia.org/wiki/Neuron>

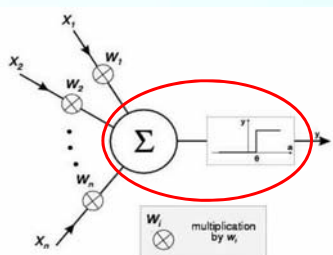


人工神经元

• 对自然神经元的数学模型

– 阈值逻辑单元(TLU- Threshold Logic Unit)模型 (McCulloch and Pitts in 1943)

$$y = \text{sgn}\left(\sum_{i=1}^n w_i x_i + w_0\right), \quad y = \begin{cases} +1 & \Rightarrow \text{class } A \\ -1 & \Rightarrow \text{class } B \end{cases}$$



BULLETIN OF
MATHEMATICAL BIOPHYSICS
VOLUME 5, 1943

A LOGICAL CALCULUS OF THE IDEAS IMMANENT IN NERVOUS ACTIVITY

WARREN S. MCCULLOCH AND WALTER PITTS

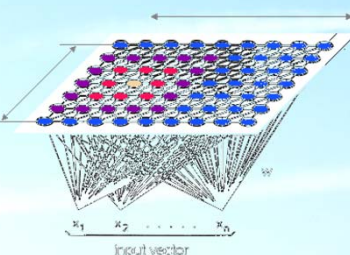
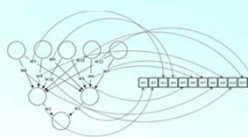
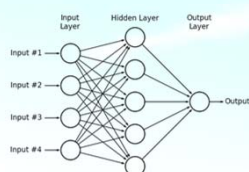
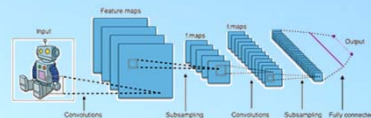
FROM THE UNIVERSITY OF ILLINOIS, COLLEGE OF MEDICINE,
DEPARTMENT OF PSYCHIATRY AT THE ILLINOIS NEUROPSYCHIATRIC INSTITUTE,
AND THE UNIVERSITY OF CHICAGO

Xuegong Zhang



• 人工神经网络 (ANN)

- 大量简单计算节点单元 (神经元)
- 经过一系列权值连接成一个网络
- 权值根据数据按某种算法进行学习



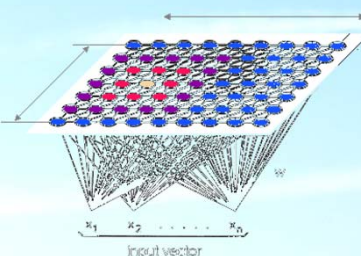
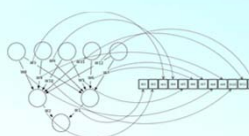
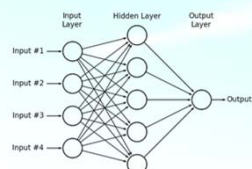
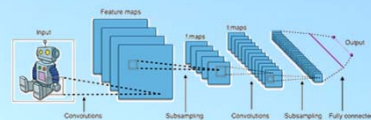
Xuegong Zhang

42



• 人工神经网络 (ANN)

- 大量简单计算节点单元 (神经元)
- 经过一系列权值连接成一个网络
- 权值根据数据按某种算法进行学习



• 神经网络的三要素：

- 网络结构
- 节点激活函数/计算 特性
- 学习算法

这些要素的不同构成了不同类型的神经网络

Xuegang Zhang

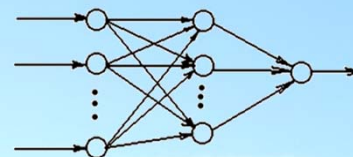
43

1980年代三种主要类型的神经网络



• 前馈型神经网络 Feedforward NN

- 代表性方法：多层感知器



Xuegang Zhang

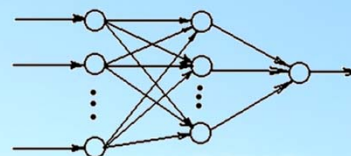
44

1980年代三种主要类型的神经网络



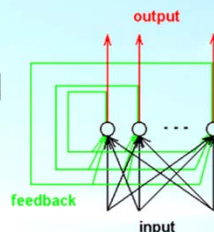
- 前馈型神经网络 Feedforward NN

- 代表性方法：多层感知器



- 反馈型神经网络 Feedback NN

- 代表性方法 Hopfield NN



Xuegang Zhang

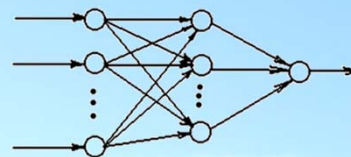
45

1980年代三种主要类型的神经网络



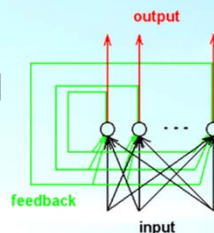
- 前馈型神经网络 Feedforward NN

- 代表性方法：多层感知器



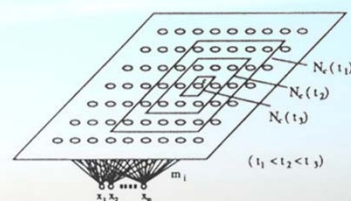
- 反馈型神经网络 Feedback NN

- 代表性方法 Hopfield NN



- 竞争学习神经网络 Competitive Learning NN

- 代表性方法：自组织映射 Self-organizing map



Xuegang Zhang

46



再看神经网络的训练

Xuegong Zhang

47

学习曲线

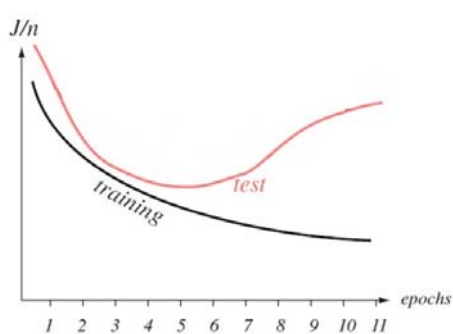


FIGURE 6.6. A learning curve shows the criterion function as a function of the amount of training, typically indicated by the number of epochs or presentations of the full training set. We plot the average error per pattern, that is, $1/n \sum_{p=1}^n J_p$. The validation error and the test or generalization error per pattern are virtually always higher than the training error. In some protocols, training is stopped at the first minimum of the validation set. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Xuegong Zhang

48

学习曲线

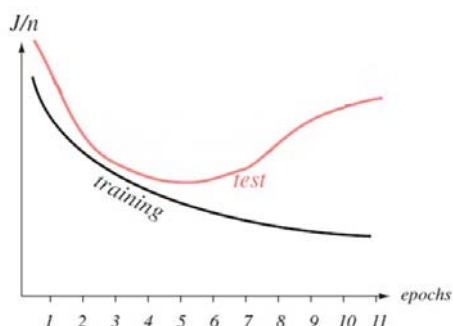


FIGURE 6.6. A learning curve shows the criterion function as a function of the amount of training, typically indicated by the number of epochs or presentations of the full training set. We plot the average error per pattern, that is, $1/n \sum_{p=1}^n J_p$. The validation error and the test or generalization error per pattern are virtually always higher than the training error. In some protocols, training is stopped at the first minimum of the validation set. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

观察：

- 训练错误率并非总是单调下降
- 并非训练错误率越小测试错误率也越小
 - 过学习现象：
训练错误率减小，测试错误率反而增大

Xuegong Zhang

49

学习曲线

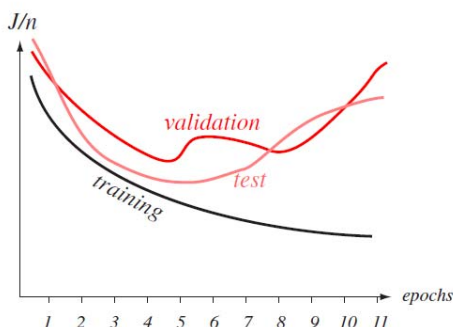


FIGURE 6.6. A learning curve shows the criterion function as a function of the amount of training, typically indicated by the number of epochs or presentations of the full training set. We plot the average error per pattern, that is, $1/n \sum_{p=1}^n J_p$. The validation error and the test or generalization error per pattern are virtually always higher than the training error. In some protocols, training is stopped at the first minimum of the validation set. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

建议：

- 密切关注训练过程中的误差变化
- 从训练集中分出一个内部验证集
 - 用验证集上的表现推断是否过学习

Xuegong Zhang

50



BP算法训练中可能遇到的问题

- 收敛慢
- 学习曲线震荡
- 过学习：推广能力差
- 可能的原因：
 - 网络结构或激活函数问题
 - 训练样本不充分或不合适
 - 初值不合适，学习率不合适，...
 - 缺乏恰当的预处理（比如归一化）
 - 需要更好的训练策略
 - 运气问题
 - 问题提出的有问题，比如 x 和 y 没有关系
 - ...

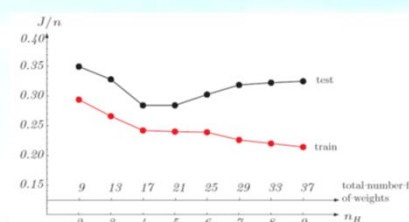
Xuegong Zhang

51



改进BP算法的可能策略

- 激活函数
 - Sigmoid, tanh, ReLU, ...
- 特征尺度调整或归一化
 - 有效值范围，相对重要性等
- 目标值调整
 - 比如在两类分类问题中，用0.1对0.9替代0对1
- 引入“伪样本”增大训练样本集
 - 如加噪、引入不变性特征的样本，平衡两类样本等
- 隐层节点数、层数
 - 与样本规模和/或问题复杂度相适应
 - $\sim n/10$?
 - 网络剪枝



Xuegong Zhang

Duda, Hart & Stork, Pattern Classification, John Wiley & Sons, 2001, p.306-318

52



改进BP算法的可能策略

- 初始化
 - 探索不同的初始化策略，以是算法有机会搜索到整个空间
- 学习率（步长）
 - 比如对平方误差采用 $\eta_{opt} = (\partial^2 J / \partial w^2)^{-1}$
 - 根据二阶导数动态调整
 - 对不同权值采用不同的学习率

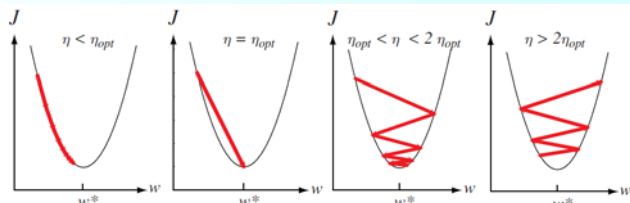


FIGURE 6.16. Gradient descent in a one-dimensional quadratic criterion with different learning rates. If $\eta < \eta_{opt}$, convergence is assured, but training can be needlessly slow. If $\eta = \eta_{opt}$, a single learning step suffices to find the error minimum. If $\eta_{opt} < \eta < 2\eta_{opt}$, the system will oscillate but nevertheless converge, but training is needlessly slow. If $\eta > 2\eta_{opt}$, the system diverges. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

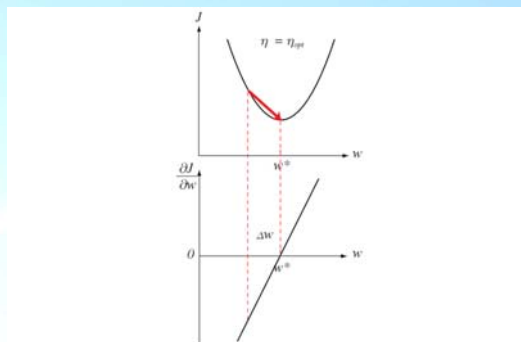


FIGURE 6.17. If the criterion function is quadratic (above), its derivative is linear (below). The optimal learning rate η_{opt} ensures that the weight value yielding minimum error, w^* , is found in a single learning step. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Xuegong Zhang

Duda, Hart & Stork, *Pattern Classification*, John Wiley & Sons, 2001, p.306-318

53

改进BP算法的可能策略

- 引入动量（Momentum）

$$\mathbf{w}(t+1) = \mathbf{w}(t) + (1 - \alpha)\Delta\mathbf{w}_{bp}(t) + \alpha(\mathbf{w}(t) - \mathbf{w}(t-1))$$
- 权值衰减（Weight decay） $\mathbf{w}^{new} = \mathbf{w}^{old}(1 - \epsilon)$
 - 等价于定义了新的正则化误差 $J_{ef}(\mathbf{w}) = J(\mathbf{w}) + \frac{2\epsilon}{\eta} \mathbf{w}^T \mathbf{w}$
- 引入提示输出（Hints）
 - 加入比较容易学习的扩增目标
- 终止条件
 - 用验证集决定是否提前终止

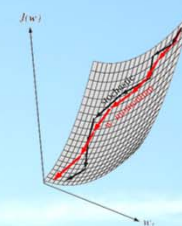


FIGURE 6.18. The incorporation of momentum into stochastic gradient descent by Eq. 37 (red arrows) reduces the variation in overall gradient directions and speeds learning. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

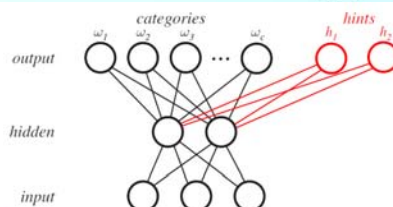


FIGURE 6.19. In learning with hints, the output layer of a standard network having c category units is augmented with hint units. During training, the target vectors are also augmented with signals for the hint units. In this way the input-to-hidden weights learn improved feature groupings. The hint units are not used during classification, and thus they and their hidden-to-output weights are removed from the trained network. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Xuegong Zhang

Duda, Hart & Stork, *Pattern Classification*, John Wiley & Sons, 2001, p.306-318

54



本章知识点

- 多层感知器原理
- BP算法
- 人工神经网络的基本概念

Xuegong Zhang

55



作业

- 第二、四章
 - 2.5 MSE与Fisher线性判别
 - 2.6 罗杰斯特函数
 - 2.7 罗杰斯特回归的梯度下降算法
 - 4.1 线性多层感知器实现的函数映射
 - 4.2 采用tanh()函数的BP算法
- 截止日期：3月9日
- 第三、四章
 - 3.3 Softmax小实验
 - 4.3 多层感知器小实验
- 截止日期：3月16日

Xuegong Zhang

56

