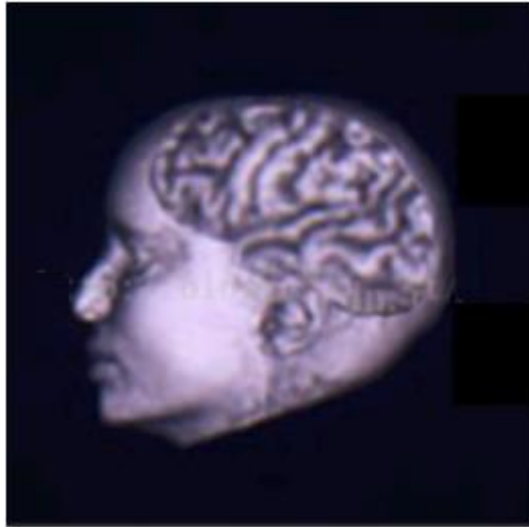

神经网络

陈飞宇

fchen@cqu.edu.cn

办公室：软件学院529

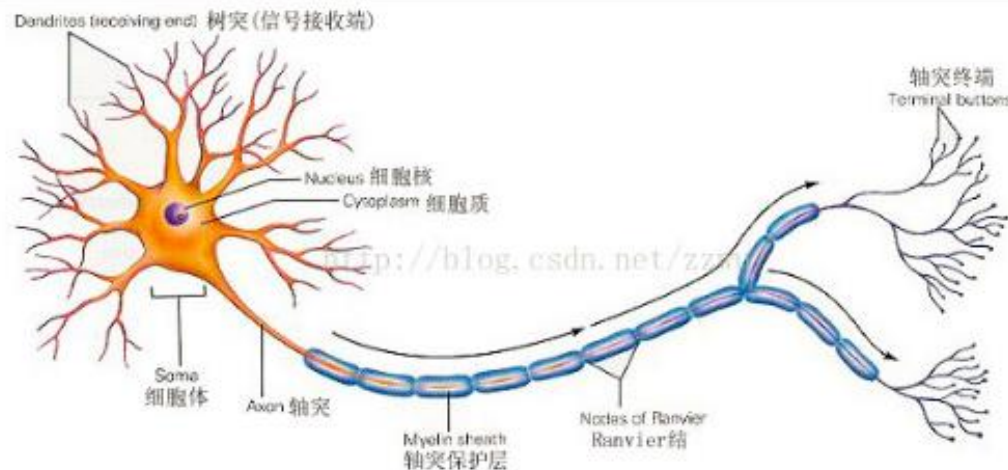
How our brain works?



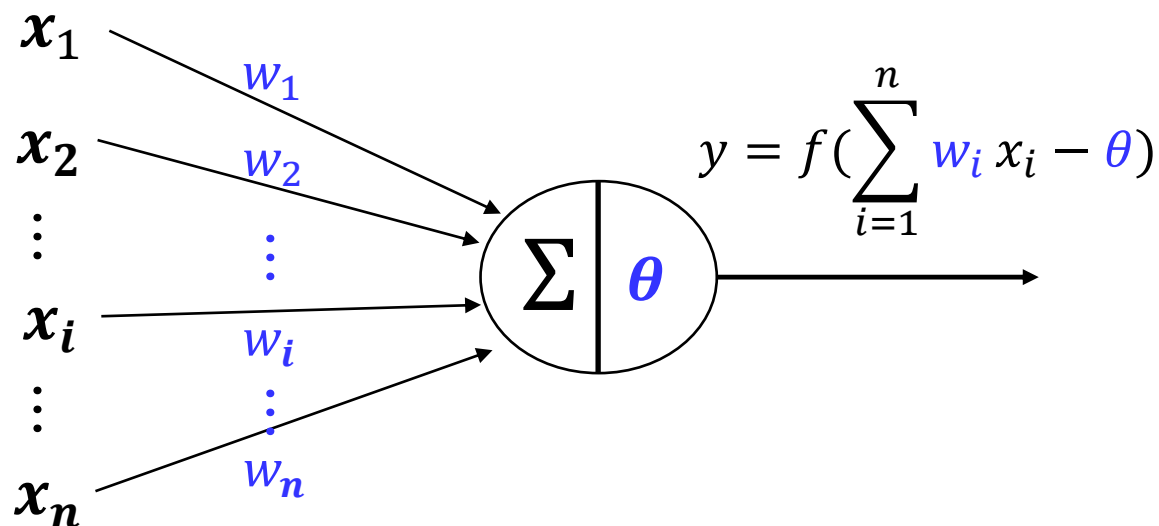
1 大脑半球像半个核桃



2 大脑皮层由灰质白质组成



神经元模型



- x_i : 第 i 个神经元的输入
- w_i : 第 i 个神经元的权重
- θ : 阈值(threshold)或称为偏置 (bias)
- y : 神经元状态 (兴奋或抑制)

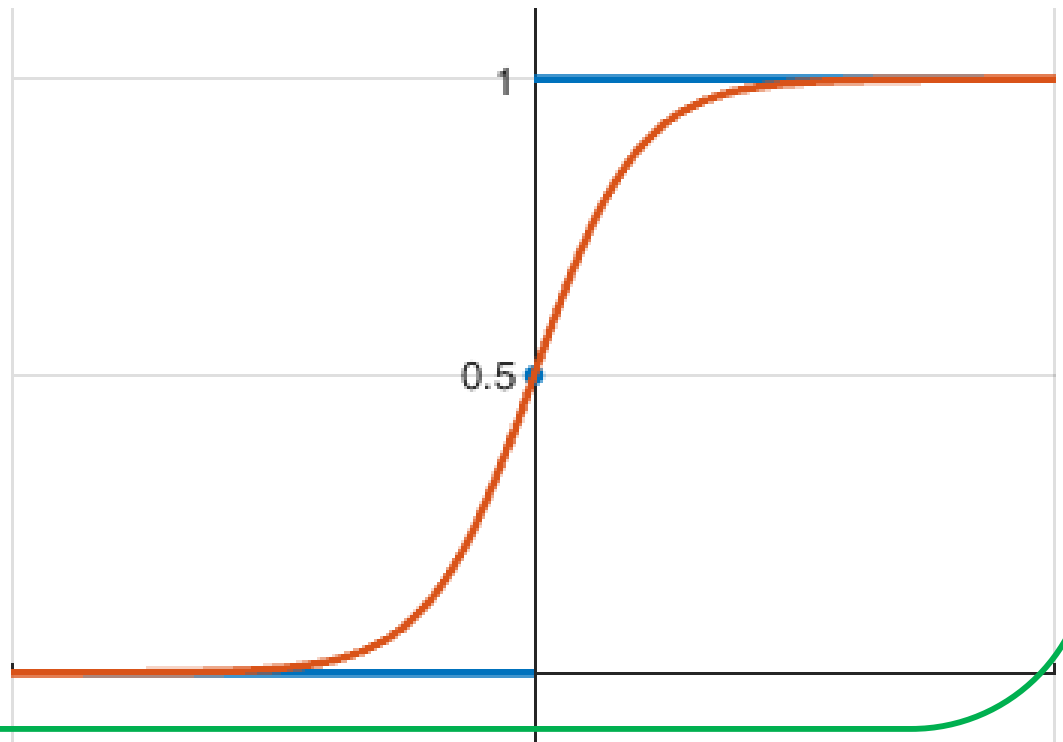
激活函数（ Activation function ）

由于单位阶跃函数不是一个连续函数，我们通常选择一些性质好的函数作为替代函数。

Logistic function:

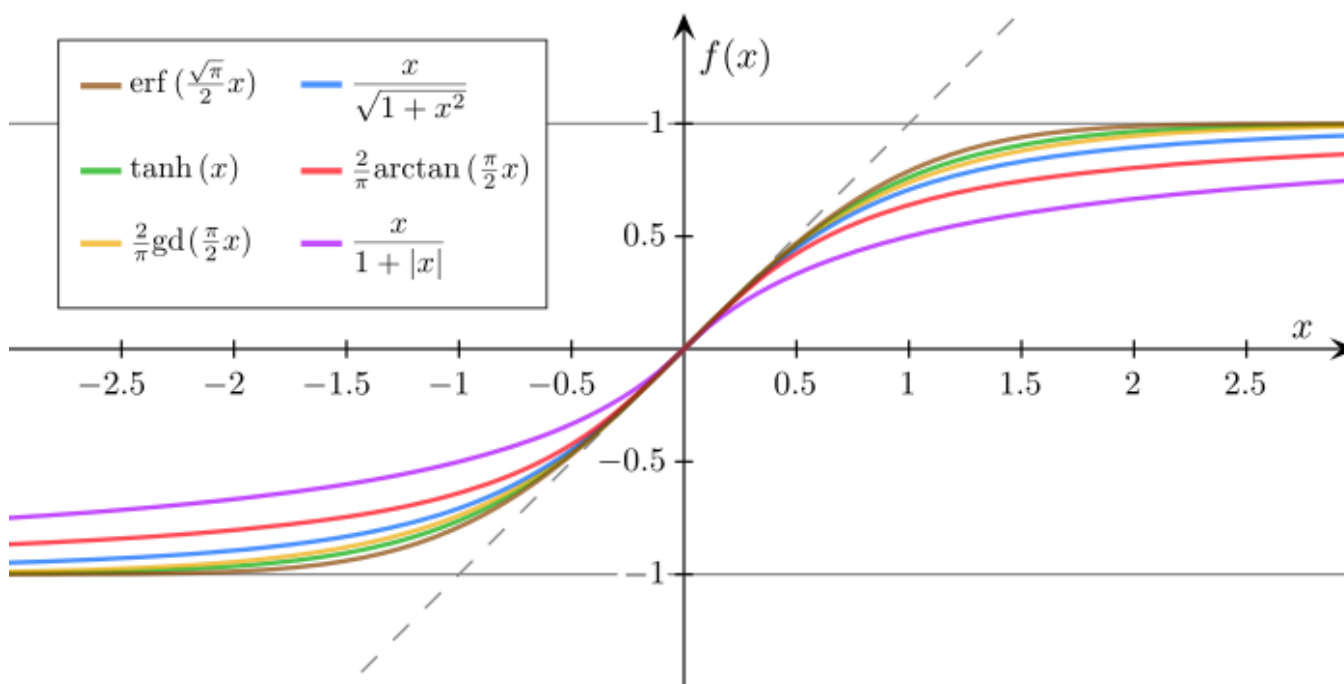
$$y = \frac{1}{1 + e^{-z}}$$

Unit-step function and logistic function



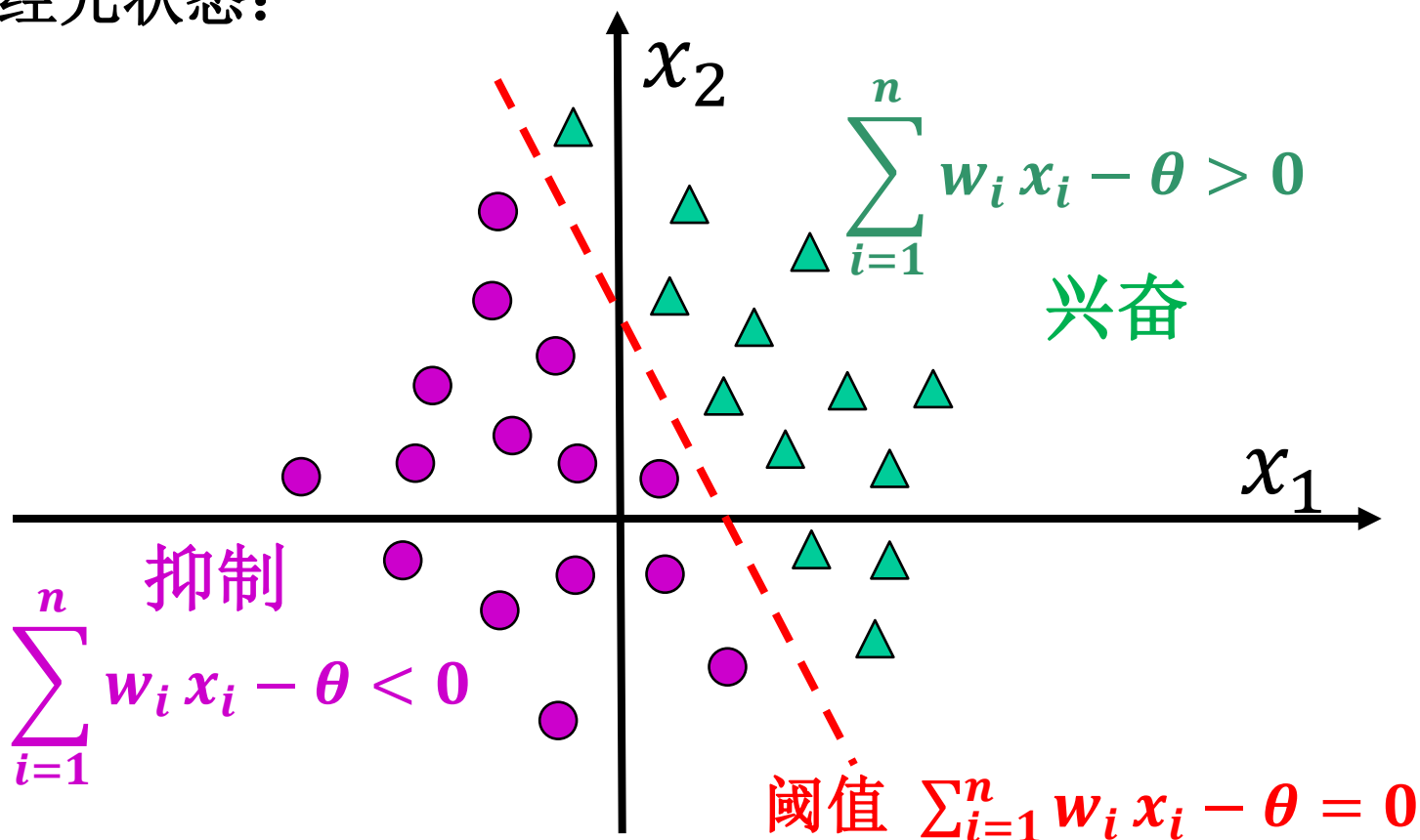
激活函数 (Activation function)

$f(\cdot)$ 函数称为激活函数(activation function)或挤压函数(Squashing function)。



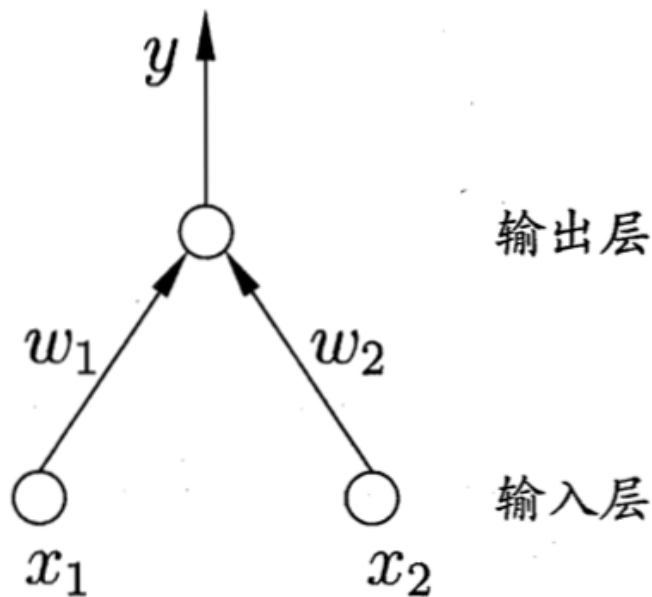
二分类问题

神经元状态:



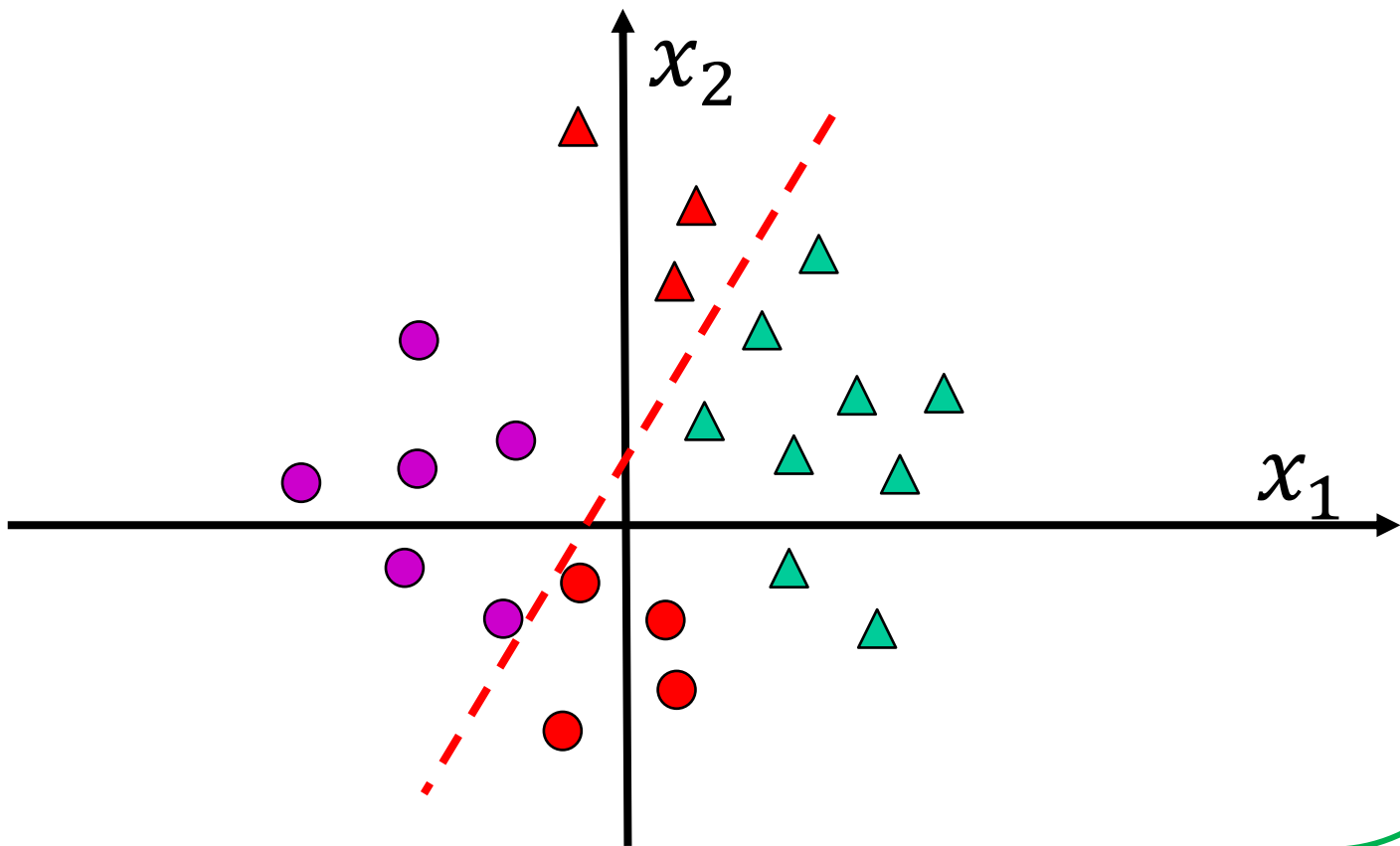
多层感知机（Multilayer Perceptron）

感知机（Perceptron）由两层神经元组成，是最简单的神经网络。



感知机原理

考虑二分类问题，**红色点**代表着误分类点。



误分类点

考虑输出变量 $y_i \in \{0, 1\}$, 定义

$$\gamma_i = (\hat{y}_i - y_i) \frac{w^T x_i - \theta}{\|w\|} \geq 0$$

其中 $\hat{y}_i = f(w^T x_i - \theta)$ 。

若 $\hat{y}_i = y_i$, 则 $\gamma_i = 0$; 若 $\hat{y}_i \neq y_i$, 则 $\gamma_i = \frac{|w^T x_i - \theta|}{\|w\|}$

因此 $(\hat{y}_i - y_i)(w^T x_i - \theta)$ 衡量着误分类点 x_i 到分类超平面的距离。

感知机模型

令 $w_0 = \theta$, $x_0 = -1$, 即将阈值 θ 看作是“哑结点” x_0 所对应的权重, 则感知机最小化误分类点到分类平面的距离和:

$$\min_{\hat{w}} \sum_{i=0}^n (\hat{y}_i - y_i) \hat{w}^T \hat{x}_i$$

$$\hat{w} = [w_0; w_1; \cdots; w_n], \quad \hat{x} = [x_0; x_1; \cdots; x_n]。$$

注意, 上述求和仅当 $\hat{y}_i \neq y_i$ 时起作用!

感知机模型

感知机模型：

$$\min_{\hat{\mathbf{w}}} \sum_{i=0}^n (\hat{y}_i - y_i) \hat{\mathbf{w}}^T \hat{\mathbf{x}}_i$$

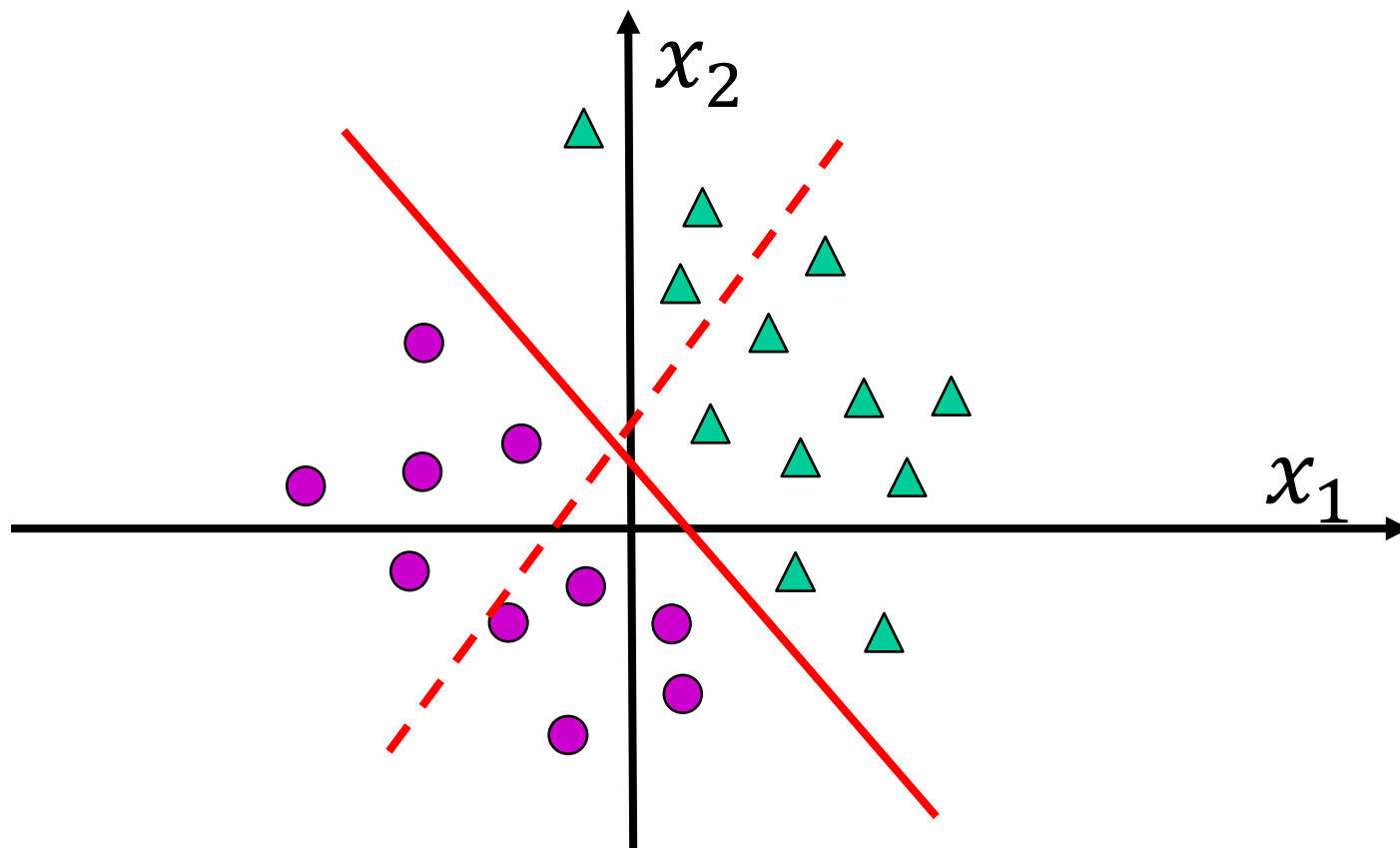
梯度下降法：

$$\hat{\mathbf{w}} \leftarrow \hat{\mathbf{w}} - \eta \Delta \hat{\mathbf{w}}$$

其中

$$\Delta \hat{\mathbf{w}} = \frac{\partial (\sum_{i=0}^n (\hat{y}_i - y_i) \hat{\mathbf{w}}^T \hat{\mathbf{x}}_i)}{\partial \hat{\mathbf{w}}} = \sum_{i=0}^n (\hat{y}_i - y_i) \hat{\mathbf{x}}_i$$

感知机原理



逻辑运算 I

“与”运算 $(x_1 \wedge x_2 \wedge \cdots \wedge x_n)$, 其中 $x_i \in \{0, 1\}$ 。

令权重 $w_1 = w_2 = \cdots = w_n = 1$, 阈值 $\theta = n$ 。

则输出 $y = f(\sum_{i=1}^n w_i x_i - \theta) = f(\sum_{i=1}^n x_i - n)$

当且仅当 $x_1 = x_2 = \cdots = x_n = 1$ 时, $y = 1$ 。

逻辑运算 II

“或”运算 $(x_1 \vee x_2 \vee \cdots \vee x_n)$, 其中 $x_i \in \{0, 1\}$ 。

令权重 $w_1 = w_2 = \cdots = w_n = 1$, 阈值 $\theta = 1/2$ 。

输出 $y = f(\sum_{i=1}^n w_i x_i - \theta) = f(\sum_{i=1}^n x_i - 1/2)$

当且仅当至少有一个 $x_i = 1$ 时, $y = 1$ 。

逻辑运算 III

“非”运算($\sim x_i$), 其中 $x_i \in \{0, 1\}$ 。

令 $w_i = -2$, 其他 $w_j = 0$, 阈值 $\theta = -1$.

输出 $y = f(\sum_{i=1}^n w_i x_i - \theta) = f(-2x_i + 1)$

当 $x_i = 1$ 时, $y = f(-1) = 0$;

当 $x_i = 0$ 时, $y = f(1) = 1$ 。

逻辑运算 IV

“异或”运算($x_i \oplus x_j$), 其中 $x_i \in \{0, 1\}$ 。

若感知机可以解决异或运算, 则

当 $x_i = x_j$ 时, $y = f(\sum_{i=1}^n w_i x_i - \theta) = 0$

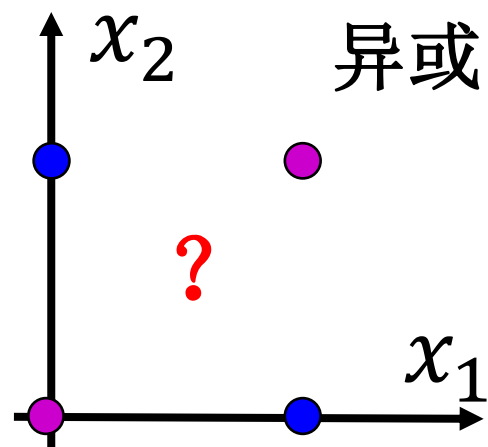
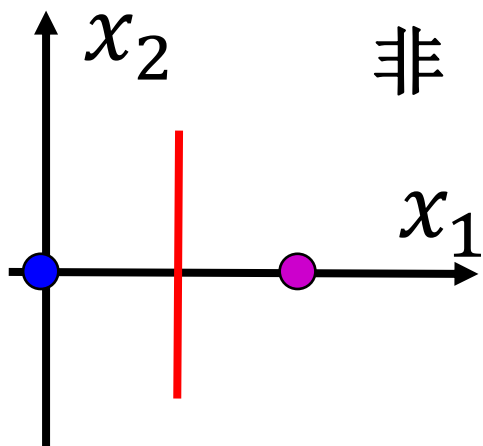
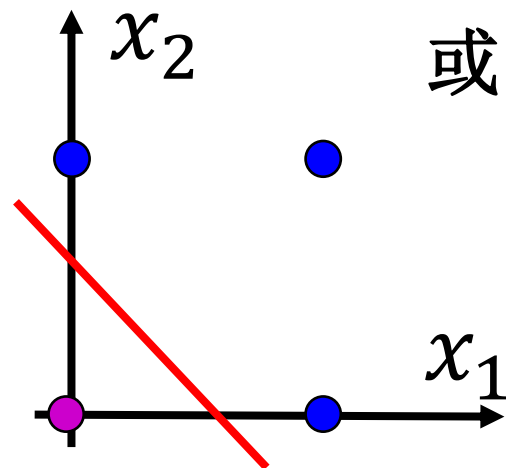
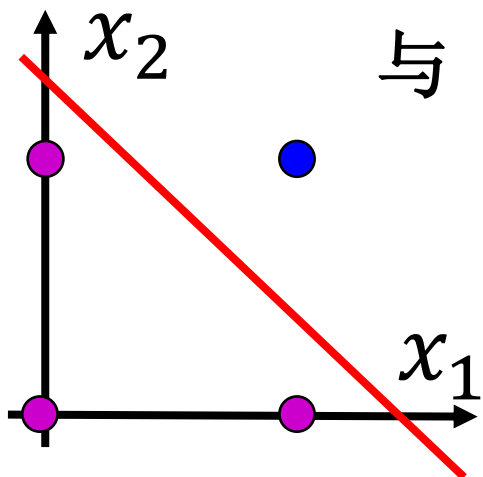
$$w_i + w_j - \theta < 0, -\theta < 0$$

当 $x_i \neq x_j$ 时, $y = f(\sum_{i=1}^n w_i x_i - \theta) = 1$

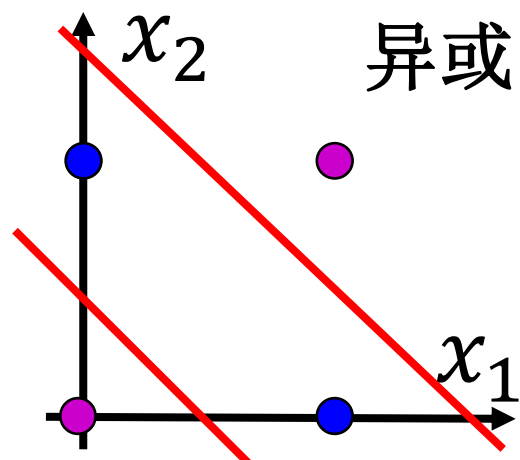
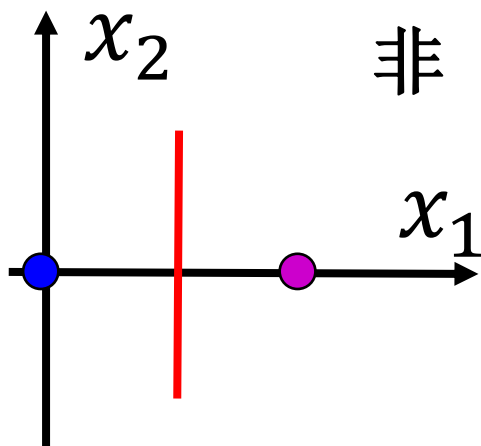
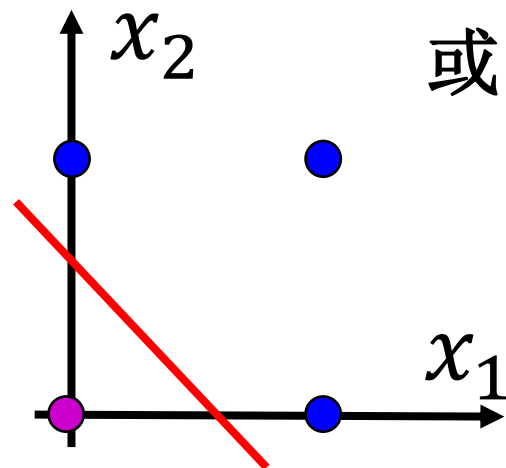
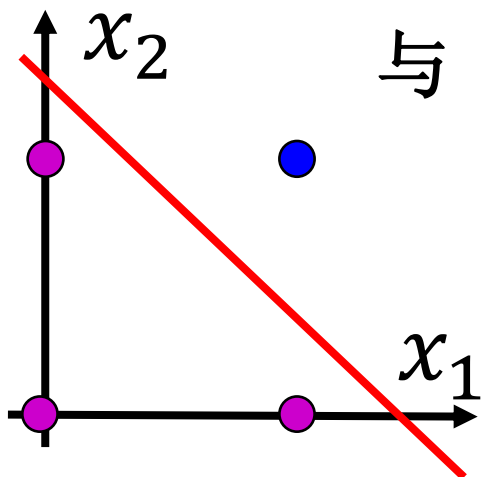
$$w_i - \theta \geq 0, w_j - \theta \geq 0$$

不存在这样的 w_i 和 w_j ! ! !

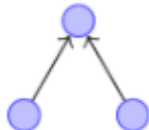

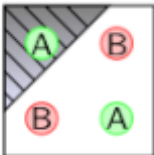
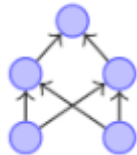

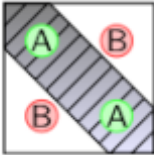
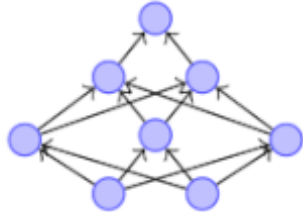


逻辑运算



逻辑运算

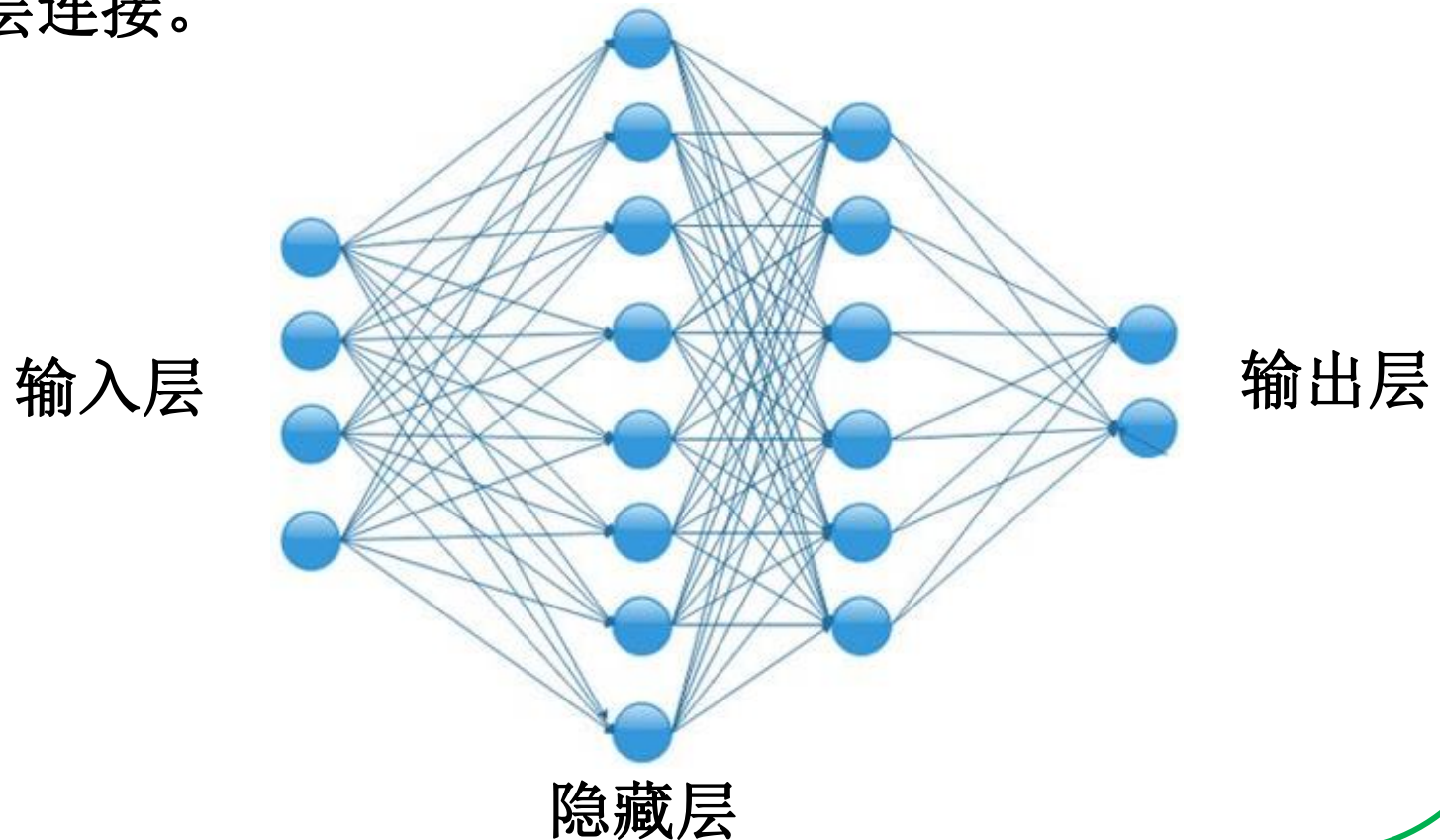


多层神经网络

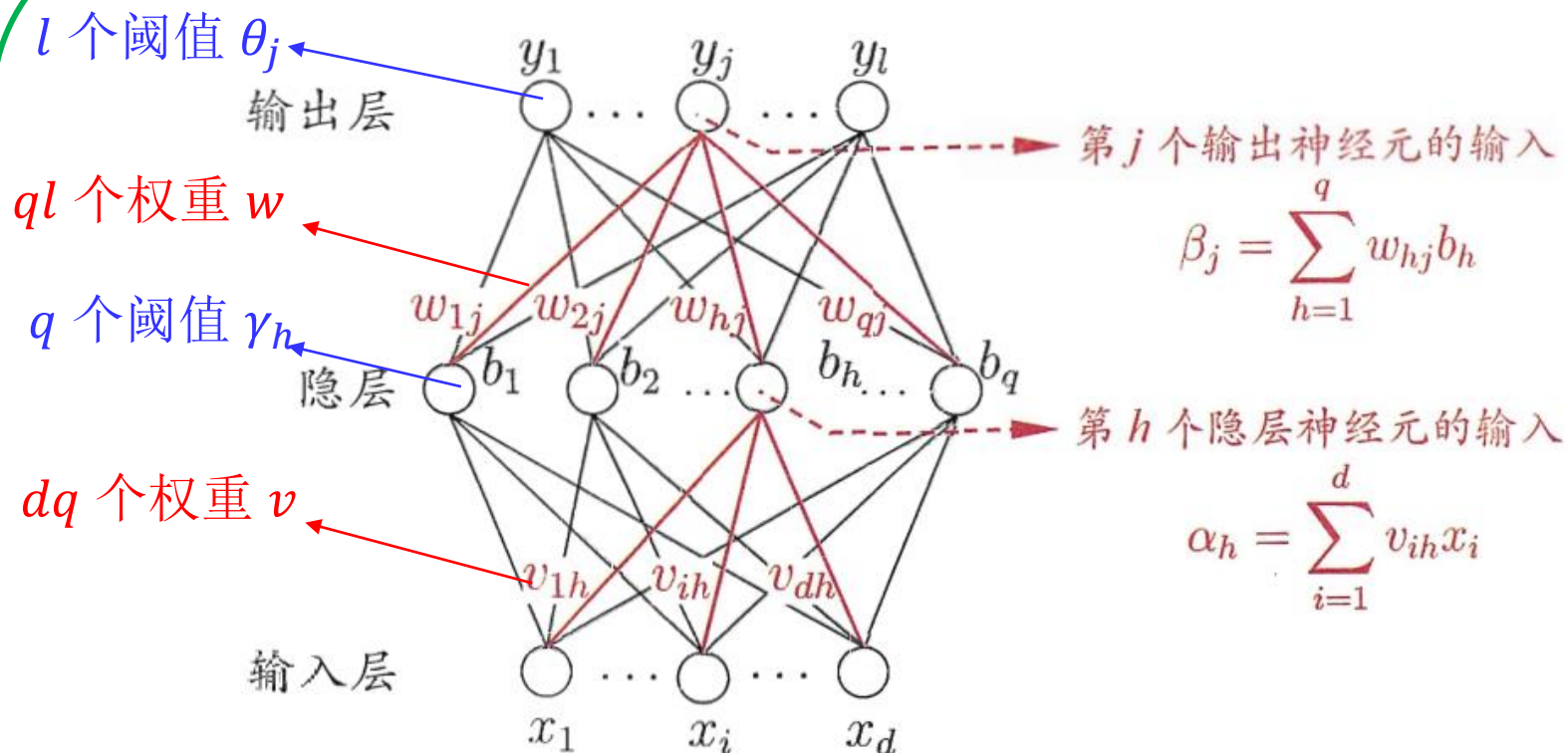
结构	决策区域类型	区域形状	异或问题
无隐层 	由一超平面分成两个		
单隐层 	开凸区域或闭凸区域		
双隐层 	任意形状（其复杂度由单元数目确定）		

多层前馈神经网络

多层前馈神经网络：每层神经元与下一层神经元完全相连，神经元之间不存在同层连接，也不存在跨层连接。



误差逆传播算法



总参数: $l + ql + q + dq$

图 5.7 BP 网络及算法中的变量符号

误差逆传播算法

BP算法优化准则:

$$\min_{\theta, \omega, \gamma, \nu} E_k = \frac{1}{2} \sum_{j=1} (\hat{y}_j^k - y_j^k)^2$$

其中预测值 \hat{y}_j^k 与权重 $\theta, \omega, \gamma, \nu$ 相关。

梯度下降法:

$$\frac{\partial E_k}{\partial \theta}, \quad \frac{\partial E_k}{\partial \omega}, \quad \frac{\partial E_k}{\partial \gamma}, \quad \frac{\partial E_k}{\partial \nu}$$

更新 ω

$$E_k = \frac{1}{2} \sum_{j=1}^l (\hat{y}_j^k - y_j^k)^2;$$

$$\hat{y}_j^k = f(\beta_j - \theta_j); \quad \beta_j = \sum_{h=1}^q \omega_{hj} b_h$$

$$\text{链式法则: } \frac{\partial E_k}{\partial \omega_{hj}} = \frac{\partial E_k}{\partial \hat{y}_j^k} \frac{\partial \hat{y}_j^k}{\partial \beta_j} \frac{\partial \beta_j}{\partial \omega_{hj}}$$

$$\frac{\partial E_k}{\partial \hat{y}_j^k} = \frac{\partial (\frac{1}{2} \sum_{j=1}^l (\hat{y}_j^k - y_j^k)^2)}{\partial \hat{y}_j^k} = \hat{y}_j^k - y_j^k$$

$$\text{因为 } f'(x) = f(x)(1 - f(x))$$

$$\frac{\partial \hat{y}_j^k}{\partial \beta_j} = f(\beta_j - \theta_j)(1 - f(\beta_j - \theta_j)) = \hat{y}_j^k(1 - \hat{y}_j^k)$$

更新 ω

$$E_k = \frac{1}{2} \sum_{j=1}^l (\hat{y}_j^k - y_j^k)^2;$$

$$\hat{y}_j^k = f(\beta_j - \theta_j); \quad \beta_j = \sum_{h=1}^q \omega_{hj} b_h$$

$$\text{链式法则: } \frac{\partial E_k}{\partial \omega_{hj}} = \frac{\partial E_k}{\partial \hat{y}_j^k} \frac{\partial \hat{y}_j^k}{\partial \beta_j} \frac{\partial \beta_j}{\partial \omega_{hj}}$$

$$\text{令 } g_j = -\frac{\partial E_k}{\partial \hat{y}_j^k} \frac{\partial \hat{y}_j^k}{\partial \beta_j} = (y_j^k - \hat{y}_j^k) \hat{y}_j^k (1 - \hat{y}_j^k)$$

$$\frac{\partial \beta_j}{\partial \omega_{hj}} = \frac{\partial (\sum_{h=1}^q \omega_{hj} b_h)}{\partial \omega_{hj}} = b_h$$

更新 ω

$$E_k = \frac{1}{2} \sum_{j=1}^l (\hat{y}_j^k - y_j^k)^2;$$

$$\hat{y}_j^k = f(\beta_j - \theta_j); \quad \beta_j = \sum_{h=1}^q \omega_{hj} b_h$$

$$\text{链式法则: } \frac{\partial E_k}{\partial \omega_{hj}} = \frac{\partial E_k}{\partial \hat{y}_j^k} \frac{\partial \hat{y}_j^k}{\partial \beta_j} \frac{\partial \beta_j}{\partial \omega_{hj}} = -g_j b_h$$

$$\omega_{hj} \leftarrow \omega_{hj} + \eta g_j b_h$$

更新 θ

$$E_k = \frac{1}{2} \sum_{j=1}^l (\hat{y}_j^k - y_j^k)^2; \quad \hat{y}_j^k = f(\beta_j - \theta_j);$$

链式法则: $\frac{\partial E_k}{\partial \theta_j} = \frac{\partial E_k}{\partial \hat{y}_j^k} \frac{\partial \hat{y}_j^k}{\partial \theta_j} = - \frac{\partial E_k}{\partial \hat{y}_j^k} \frac{\partial \hat{y}_j^k}{\partial \beta_j} = g_j$

$$\theta_j \leftarrow \theta_j - \eta g_j$$

同理可得

$$\frac{\partial E_k}{\partial \gamma}, \quad \frac{\partial E_k}{\partial \nu}$$

BP算法小结

核心思想：利用前向传播，计算第 n 层输出值

优化目标：输出值和实际值的残差。

计算方法：将残差按影响逐步传递回第 $n - 1, n - 2, \dots, 2$ 层，以修正各层参数。（即所谓的误差逆传播）

主要工具：链式法则（复合函数求偏导）。

BP算法局限性

- 容易过拟合！

早停、正则化

- 容易陷入局部最优！

选取多次初值、随机梯度下降法

- 难以设置隐层个数！

试错法