# 神经网络

陈飞宇

fchen@cqu.edu.cn

办公室：软件学院529
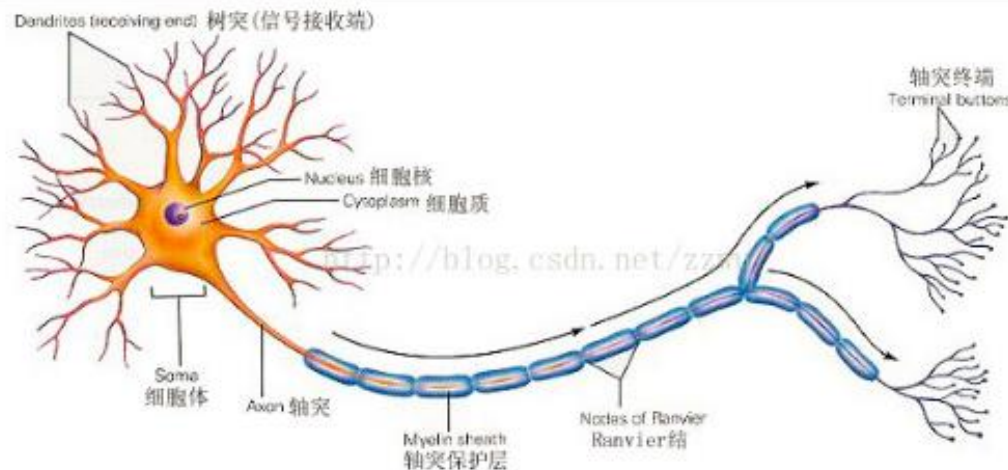
# How our brain works?
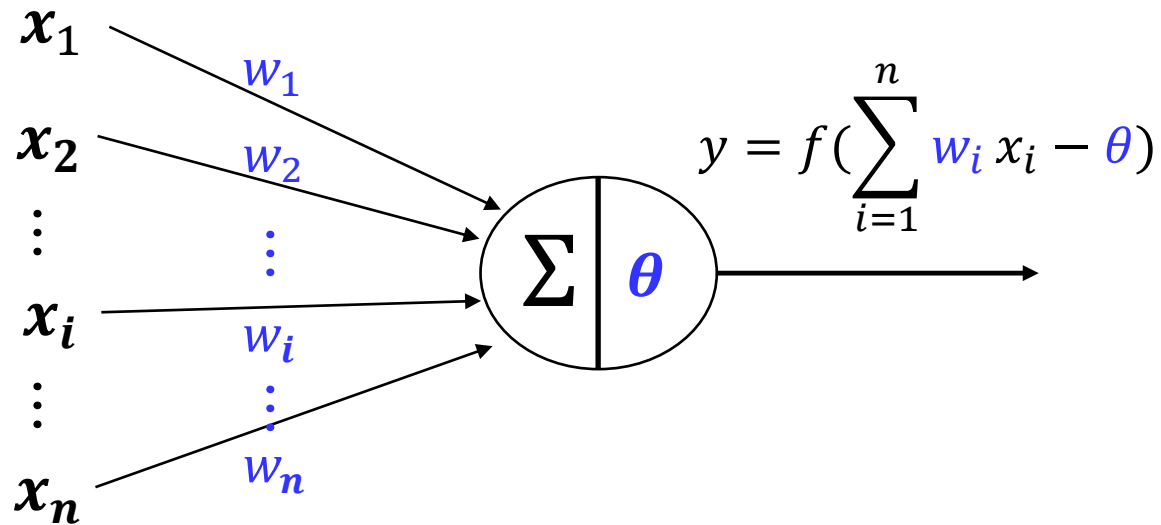


1 大脑半球像半个核桃

2 大脑皮层由灰质白质组成

Dendrites (receiving end) 树突(信号接收端)

Nucleus 细胞核
Cytoplasm 细胞质

Some 细胞体
Axon 轴突
Myelin sheath 轴突保护层

轴突终端
Terminal buttons

Nodes of Ranvier
Ranvier结

2

# 神经元模型

$$x_1$$

$$x_2$$

$$\vdots$$

$$x_i$$

$$\vdots$$

$$x_n$$

$$w_1$$

$$w_2$$

$$\vdots$$

$$w_i$$

$$\vdots$$

$$w_n$$

$$\sum \quad \theta$$

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

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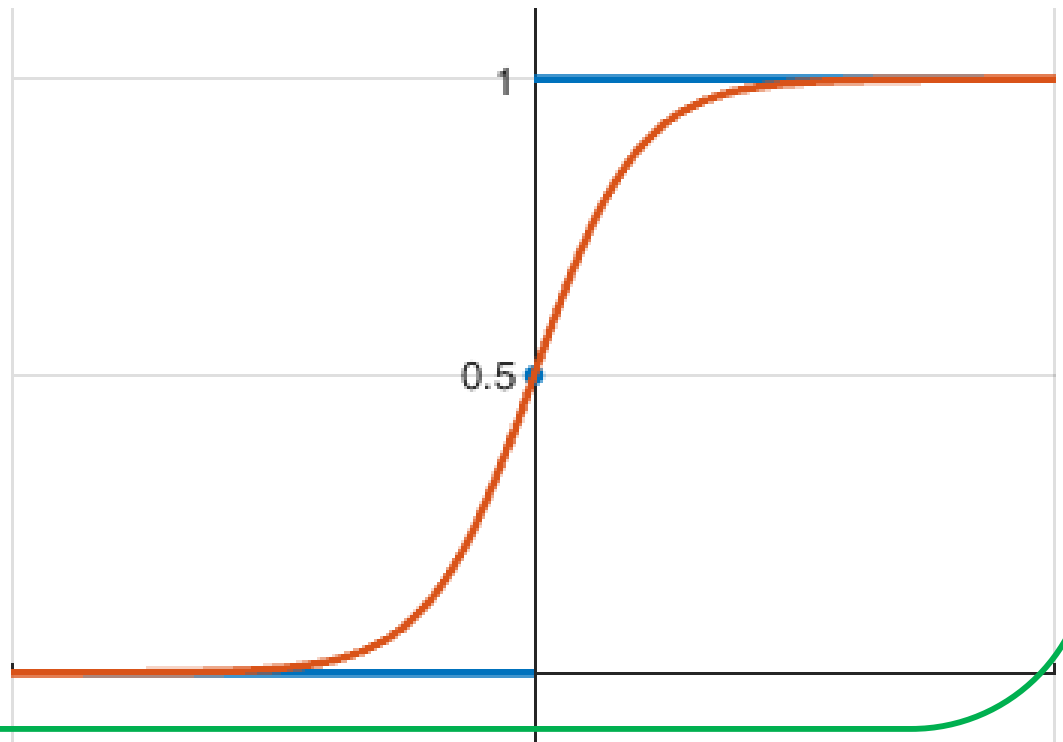- $\theta$ ：阈值(threshold)或称为偏置（bias）
- $y$ ：神经元状态（兴奋或抑制）

# 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。
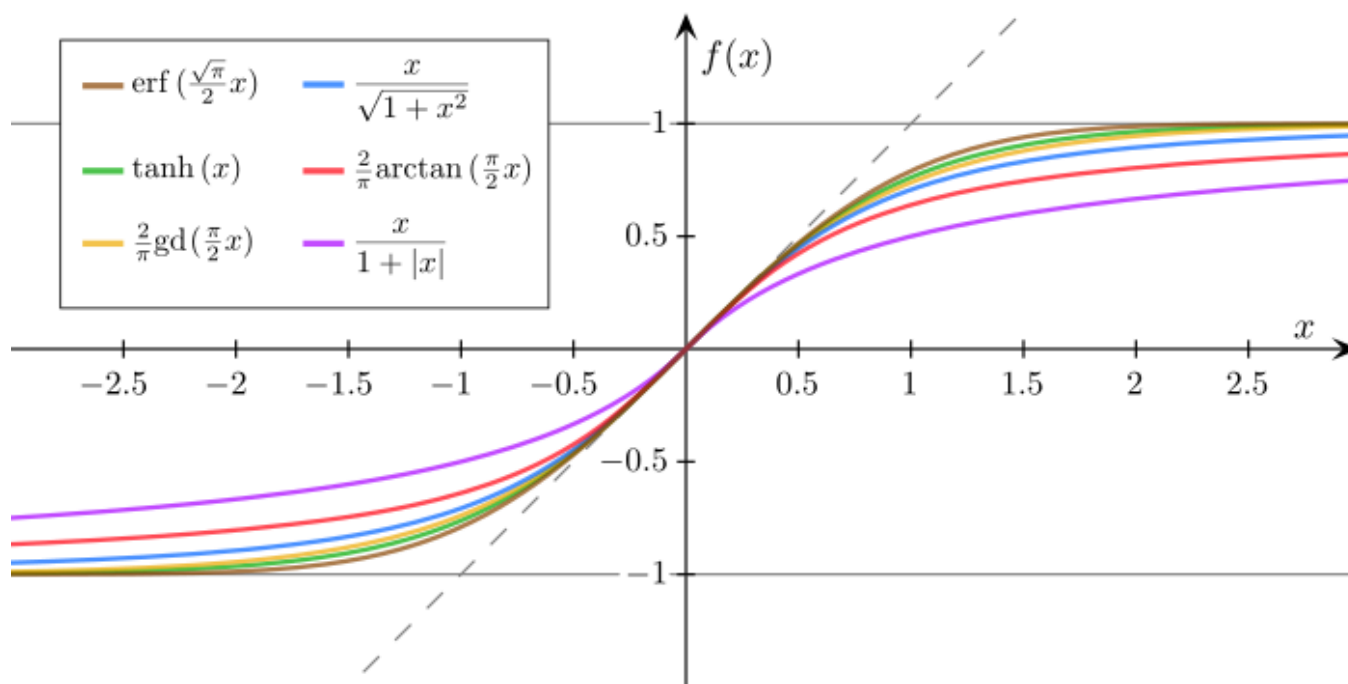
**Logistic function:**

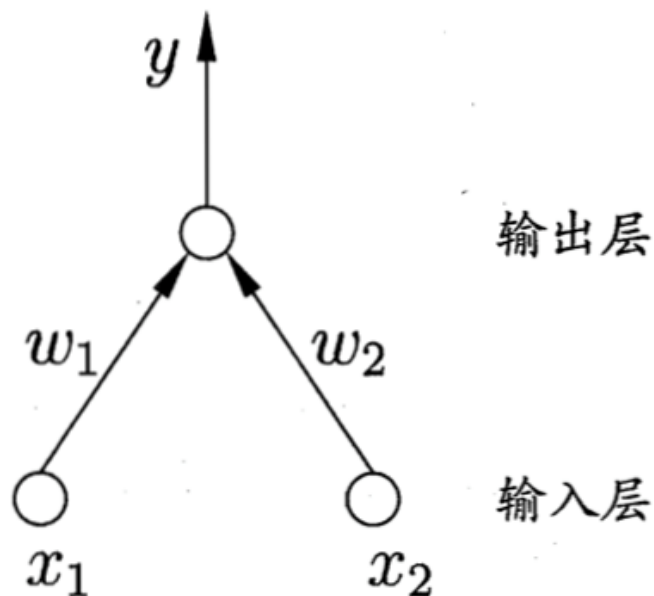$$y = \frac{1}{1 + e^{-z}}$$

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）.

# 感知机（**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（ Perceptron）

对于上述感知机模型，给定训练点$(x, y)$, 令$\boldsymbol{w_0} = \boldsymbol{\theta}, \boldsymbol{x_0} = -\boldsymbol{1}$(哑结点), 则$\boldsymbol{\widehat{w}^T\widehat{x}} = w_1 x_1 + w_2 x_2 + w_0 x_0$。

感知机模型为：

$$\min_{\boldsymbol{\widehat{w}}} \ (\boldsymbol{\widehat{y}} - \boldsymbol{y})\boldsymbol{\widehat{w}^T\widehat{x}}$$

其中$\hat{y} = f(\boldsymbol{\widehat{w}^T\widehat{x}})$为感知机当前的输出结果。

感知机参数$\boldsymbol{\widehat{w}}$的更新公式：

$$\boldsymbol{\widehat{w}} \leftarrow \boldsymbol{\widehat{w}} - \boldsymbol{\eta}\frac{\boldsymbol{\partial((\widehat{y} - y)\widehat{w}^T\widehat{x})}}{\boldsymbol{\partial\widehat{w}}} \leftarrow \boldsymbol{\widehat{w}} - \boldsymbol{\eta}(\boldsymbol{\widehat{y}} - \boldsymbol{y})\boldsymbol{\widehat{x}}$$

若样例$(x, y)$预测正确 ➡ $\boldsymbol{\widehat{y}} = \boldsymbol{y}$ ➡ $\boldsymbol{\widehat{w}}$ 不再更新

若样例$(x, y)$预测错误 ➡ $\boldsymbol{\widehat{y}} \neq \boldsymbol{y}$ ➡ $\boldsymbol{\widehat{w}}$ 继续更新

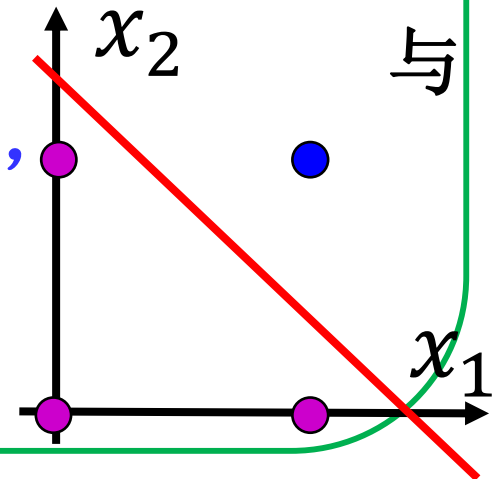# 逻辑运算 I

"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 $y = f(\sum_{i=1}^{n} w_i x_i - \theta) = f(\sum_{i=1}^{n} x_i - n)$
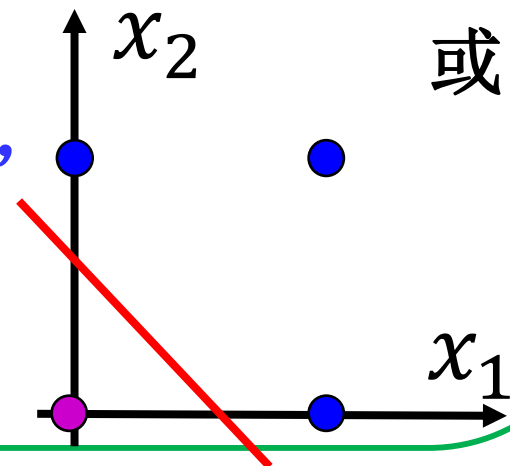
当且仅当 $x_1 = x_2 = \cdots = x_n = 1$ 时，$y = 1$。

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

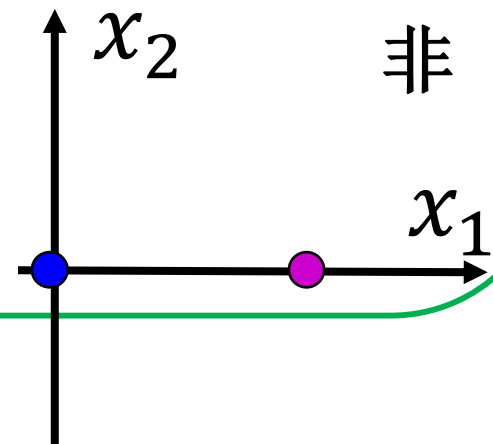当且仅当至少有一个 $x_i = 1$ 时，$y = 1$。

或

$x_2$

$x_1$

# 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

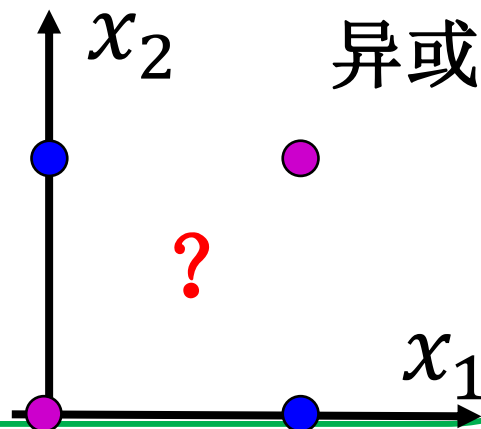"异或" 运算$(x_i \oplus x_j)$，其中$x_i \in \{0, 1\}$。

若感知机可以解决异或运算，则

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

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

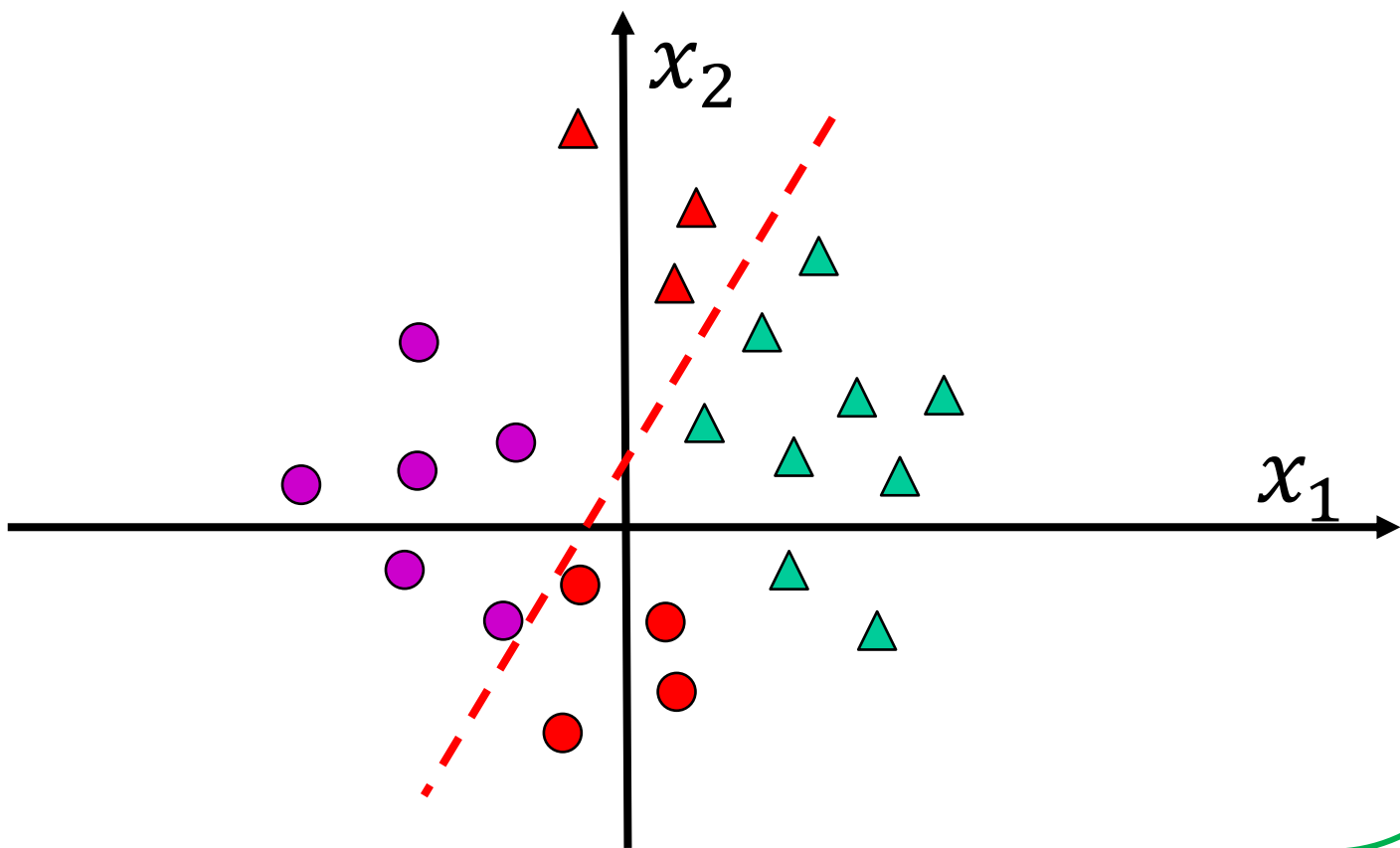当 $x_i \neq x_j$ 时，$y = f(\sum_{i=1}^{n} w_i x_i - \theta) = 1$

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

不存在这样的 $w_i$和$w_j$！！！

$x_2$  异或

?

$x_1$

# 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

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

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

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

因此 $(\hat{y}_i - y_i)(w^{\mathrm{T}}x_i - \theta)$ 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 $\theta$ 看作是 "哑结点" $x_0$ 所对应的权重, 则感知机最小化<span style="color:red">误分类点</span>到分类平面的距离和:

$$\min_{\widehat{w}} \quad \sum_{i=1}^{n} (\widehat{y}_i - y_i)\widehat{w}^{\mathrm{T}}\widehat{x}_i$$

$\widehat{w} = [w_0; w_1; \cdots; w_d]$, $\widehat{x}_i = [x_0; x_{i1}; \cdots; x_{id}]$。

<span style="color:red">注意, 上述求和仅当 $\widehat{y}_i \neq y_i$ 时起作用!</span>

# 感知机模型

感知机模型：

$$\min_{\widehat{w}} \quad \sum_{i=1}^{n}(\widehat{y}_i - y_i)\widehat{w}^{\mathrm{T}}\widehat{x}_i$$

梯度下降法：

$$\widehat{w} \longleftarrow \widehat{w} - \eta\Delta\widehat{w}$$

其中

$$\Delta\widehat{w} = \frac{\partial(\sum_{i=1}^{n}(\widehat{y}_i - y_i)\widehat{w}^{\mathrm{T}}\widehat{x}_i)}{\partial\widehat{w}} = \sum_{i=1}^{n}(\widehat{y}_i - y_i)\widehat{x}_i = X^{\mathrm{T}}(\widehat{y} - y)$$
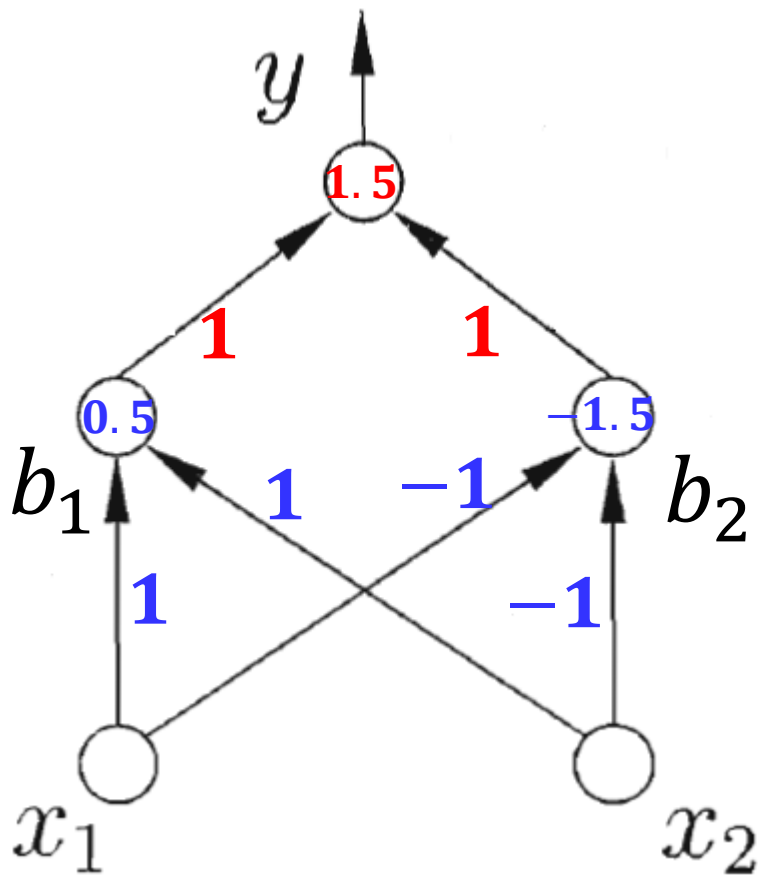
# 感知机原理

# 单隐层神经网络



第一层：
$$b_1 = f(v_{11}x_1 + v_{12}x_2 - \theta_1)$$
$$b_2 = f(v_{21}x_1 + v_{22}x_2 - \theta_2)$$

第二层：
$$y = g(w_1 b_1 + w_2 b_2 - \theta)$$

# 异或问题



$$b_1 = f(\ x_1 + x_2 - 0.5)$$
$$b_2 = f(-x_1 - x_2 + 1.5)$$
$$y = g(b_1 + b_2 - 1.5)$$

$$x_1 = 0, \ x_2 = 0;$$
$$b_1 = 0, b_2 = 1, y = 0$$
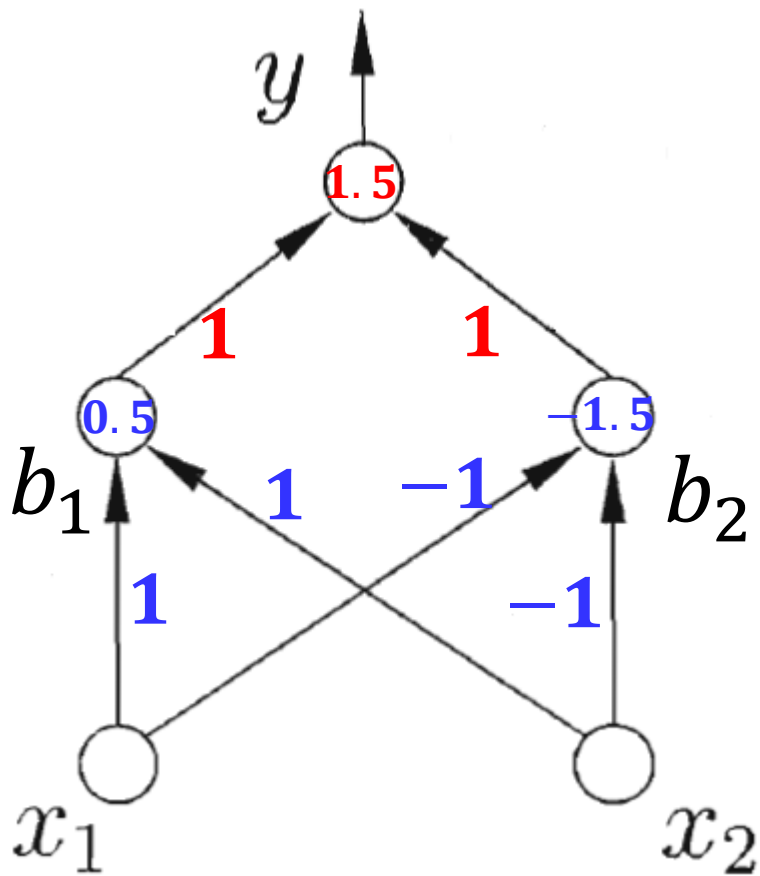$$x_1 = 0, \ x_2 = 1;$$
$$b_1 = 1, b_2 = 1, y = 1$$
$$x_1 = 1, \ x_2 = 0;$$
$$b_1 = 1, b_2 = 1, y = 1$$
$$x_1 = 1, \ x_2 = 1;$$
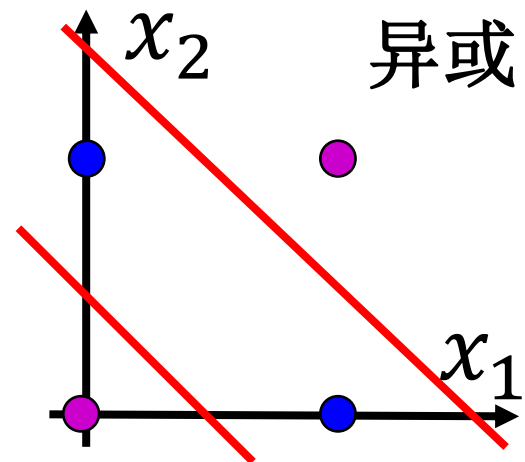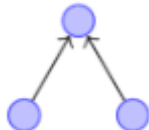$$b_1 = 1, b_2 = 0, y = 0$$

# 异或问题



$$x_1 = 0, \quad x_2 = 0, y = 0$$
$$x_1 = 0, \quad x_2 = 1, y = 1$$
$$x_1 = 1, \quad x_2 = 0, y = 1$$
$$x_1 = 1, \quad x_2 = 1, y = 0$$

# 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。

输入层

隐藏层

输出层

# 误差逆传播算法

$l$ 个阈值 $\theta_j$

$ql$ 个权重 $w$

$q$ 个阈值 $\gamma_h$

$dq$ 个权重 $v$

输出层

$y_1$   $y_j$   $y_l$

$w_{1j}$   $w_{2j}$   $w_{hj}$   $w_{qj}$

隐层   $b_1$   $b_2$   $b_h$   $b_q$

$v_{1h}$   $v_{ih}$   $v_{dh}$

输入层

$x_1$   $x_i$   $x_d$

第 $j$ 个输出神经元的输入

$$\beta_j = \sum_{h=1}^{q} w_{hj} b_h$$

第 $h$ 个隐层神经元的输入

$$\alpha_h = \sum_{i=1}^{d} v_{ih} x_i$$

**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



$$\hat{y}_j^k = f(\beta_j - {\color{red}\theta_j});$$

$$\beta_j = \sum_{h=1}^{q} {\color{red}\omega_{hj}} b_h$$

$$b_h = f(\alpha_h - {\color{red}\gamma_h})$$

$$\alpha_h = \sum_{i=1}^{d} {\color{red}v_{ih}} x_i$$

其中 $f$ 是对数几率函数。

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

**BP**网络均方误差： $E_k = \frac{1}{2}\sum_{j=1}^{l}(\hat{y}_j^k - y_j^k)^2$

# 误差逆传播算法

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

$$\beta_j = \sum_{h=1}^{q} \omega_{hj} b_h$$

$$b_h = f(\alpha_h - \gamma_h)$$

$$\alpha_h = \sum_{i=1}^{d} v_{ih} x_i$$

其中 $f$ 是对数几率函数。

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

$$\frac{\partial E_k}{\partial \omega_{hj}} = \frac{\partial E_k}{\partial \hat{y}_j^k} \frac{\partial \hat{y}_j^k}{\partial \beta_j} \frac{\partial \beta_j}{\partial \omega_{hj}}$$
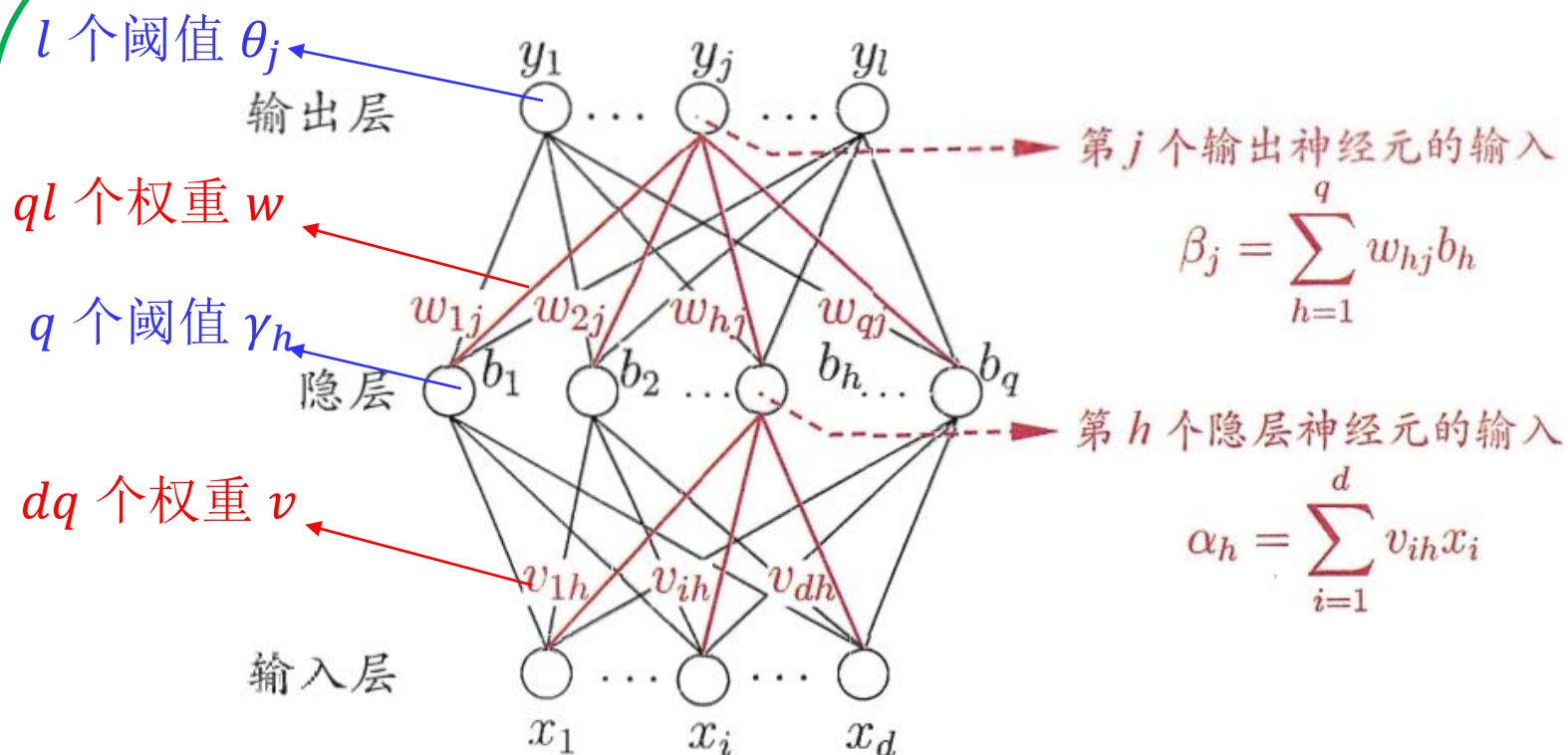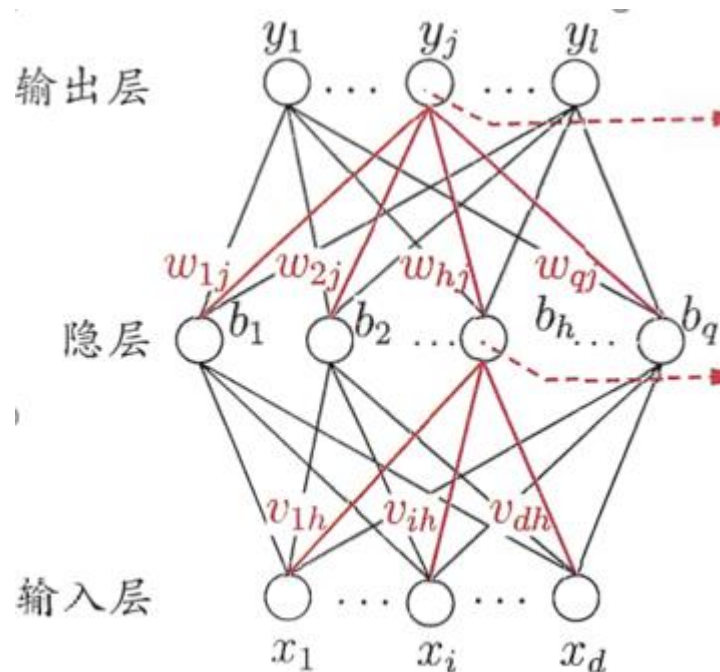
$$\frac{\partial E_k}{\partial \theta_j} = -\frac{\partial E_k}{\partial \hat{y}_j^k} \frac{\partial \hat{y}_j^k}{\partial \beta_j}$$

$$\frac{\partial E_k}{\partial v_{ih}} = \sum_{j=1}^{l} \frac{\partial E_k}{\partial \hat{y}_j^k} \frac{\partial \hat{y}_j^k}{\partial \beta_j} \frac{\partial \beta_j}{\partial b_h} \frac{\partial b_h}{\partial \alpha_h} \frac{\partial \alpha_h}{\partial v_{ih}}$$

$$\frac{\partial E_k}{\partial \gamma_h} = -\sum_{j=1}^{l} \frac{\partial E_k}{\partial \hat{y}_j^k} \frac{\partial \hat{y}_j^k}{\partial \beta_j} \frac{\partial \beta_j}{\partial b_h} \frac{\partial b_h}{\partial \alpha_h}$$

$$\frac{\partial E_k}{\partial \hat{y}_j^k} \ \frac{\partial \hat{y}_j^k}{\partial \beta_j}$$

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

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

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

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

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

$$\sum_{j=1}^{l} \frac{\partial E_k}{\partial \hat{y}_j^k} \frac{\partial \hat{y}_j^k}{\partial \beta_j} \frac{\partial \beta_j}{\partial b_h} \frac{\partial b_h}{\partial \alpha_h}$$

$$\beta_j = \sum_{h=1}^{q} \omega_{hj} b_h; \quad b_h = f(\alpha_h - \gamma_h)$$

$$\frac{\partial \beta_j}{\partial b_h} = \omega_{hj}$$

$$\frac{\partial b_h}{\partial \alpha_h} = \frac{\partial f(\alpha_h - \gamma_h)}{\partial \alpha_h} = f(1 - f) = b_h(1 - b_h)$$

$$\frac{\partial \beta_j}{\partial b_h} \frac{\partial b_h}{\partial \alpha_h} = b_h(1 - b_h)\omega_{hj}$$

$$e_h = \sum_{j=1}^{l} \frac{\partial E_k}{\partial \hat{y}_j^k} \frac{\partial \hat{y}_j^k}{\partial \beta_j} \frac{\partial \beta_j}{\partial b_h} \frac{\partial b_h}{\partial \alpha_h} = b_h(1 - b_h) \sum_{j=1}^{l} \omega_{hj} g_j$$

# 误差逆传播算法

$$\frac{\partial E_k}{\partial \omega_{hj}} = \frac{\partial E_k}{\partial \hat{y}_j^k} \frac{\partial \hat{y}_j^k}{\partial \beta_j} \frac{\partial \beta_j}{\partial \omega_{hj}} = -g_j b_h$$

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

$$\frac{\partial E_k}{\partial v_{ih}} = \sum_{j=1}^{l} \frac{\partial E_k}{\partial \hat{y}_j^k} \frac{\partial \hat{y}_j^k}{\partial \beta_j} \frac{\partial \beta_j}{\partial b_h} \frac{\partial b_h}{\partial \alpha_h} \frac{\partial \alpha_h}{\partial v_{ih}} = -e_h x_i$$

$$\frac{\partial E_k}{\partial \gamma_h} = -\sum_{j=1}^{l} \frac{\partial E_k}{\partial \hat{y}_j^k} \frac{\partial \hat{y}_j^k}{\partial \beta_j} \frac{\partial \beta_j}{\partial b_h} \frac{\partial b_h}{\partial \alpha_h} = e_h$$

迭代公式：

$$\omega_{hj} \leftarrow \omega_{hj} + \eta_\omega g_j b_h$$

$$\theta_j \leftarrow \theta_j - \eta_\theta g_j$$

$$v_{ih} \leftarrow v_{ih} + \eta_v e_h x_i$$

$$\gamma_h \leftarrow \gamma_h - \eta_\gamma e_h$$

其中$\eta \in (0, 1)$为学习率。

# 误差逆传播算法

BP算法优化准则：

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

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

梯度下降法：

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

# 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 $n-1, n-2, \cdots, 2$ 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