

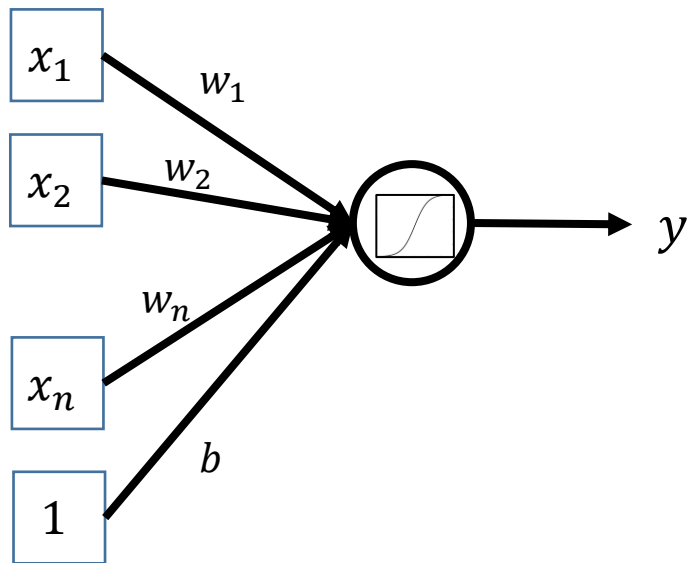
深度学习 第四讲

多层感知机与BP算法

王文中

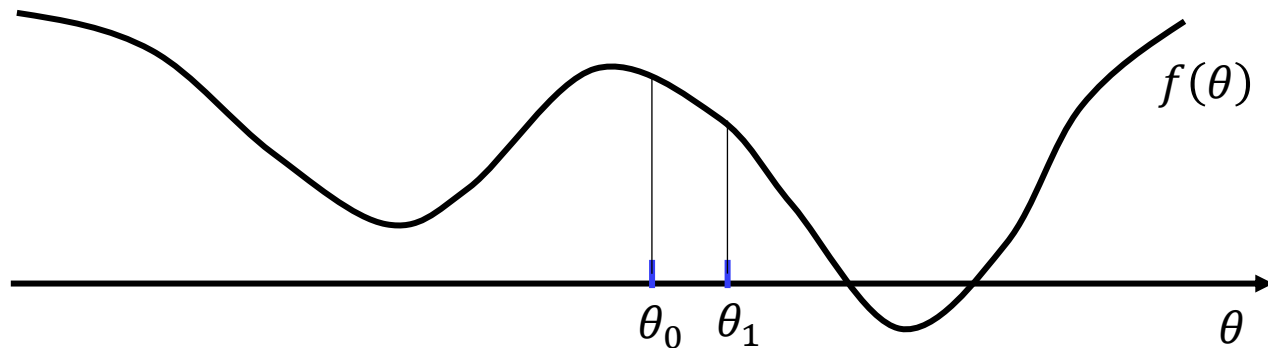
安徽大学计算机学院

回顾：神经元模型



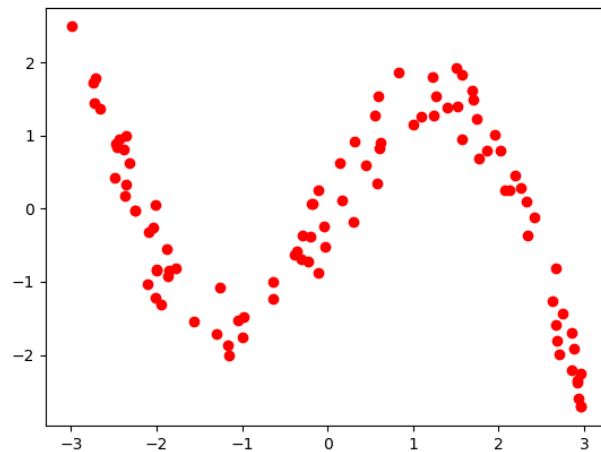
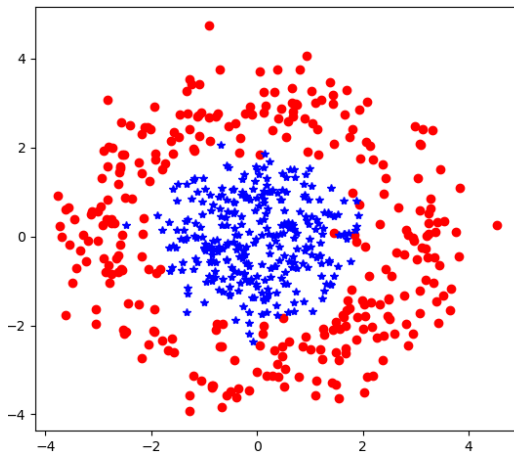
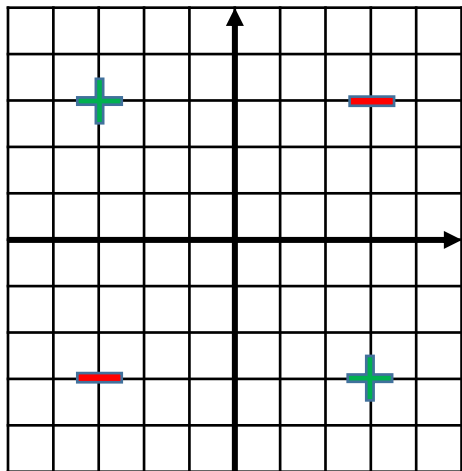
$$z = \sum_{i=1}^n w_i x_i + b = W^T X + b, y = f(z)$$

回顾：梯度下降法

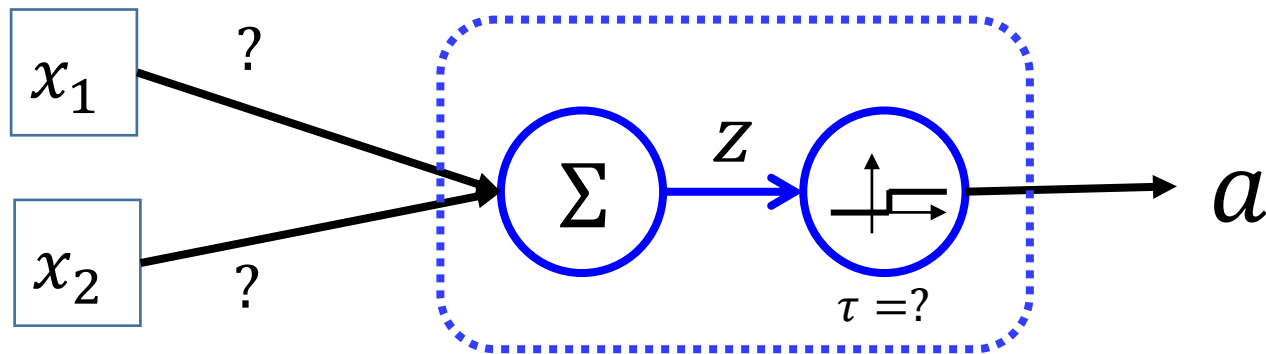


- 0: 猜测一个初始值 θ_0
- 1: 如果 θ_0 是局部极小值，退出
- 2: 更新 θ_0 为 θ_1 : $\theta_1 = \theta_0 - \alpha \cdot \frac{df}{d\theta_0}$
- 3: $\theta_0 = \theta_1$ ，转1

回顾：非线性问题



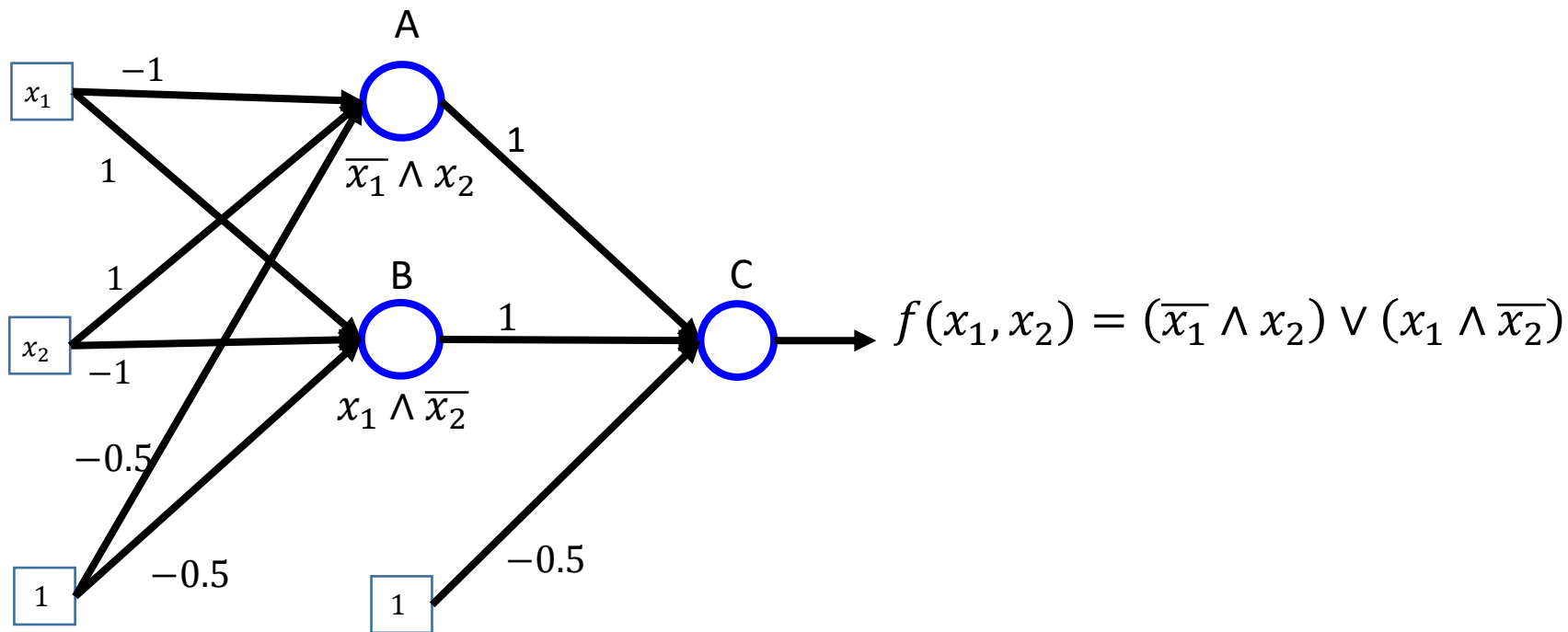
异或问题



$$y = x_1 XOR x_2 = (\overline{x_1} \wedge x_2) \vee (x_1 \wedge \overline{x_2})$$

x_1	x_2	y
0	0	0
1	0	1
0	1	1
1	1	0

用MCP神经网络实现异或函数



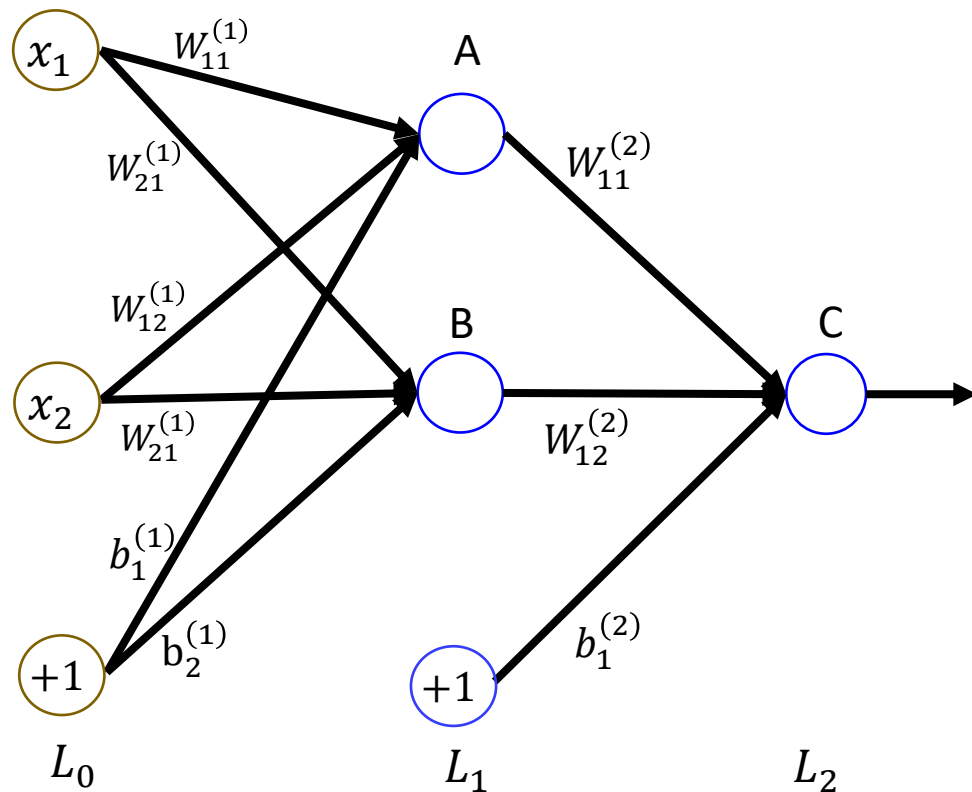
利用神经网络处理非线性问题

- 单个神经元只能实现线性分类/回归
- 使用三个MCP神经元可以实现异或运算
 - 三个MCP神经元构成一个三层神经网络：输入层、中间层以及输出层
 - 中间层神经元对输入层信号做非线性变换
 - 输出层神经元接收中间层神经元的输出结果
 - 三个神经元的参数依赖人工设计

多层感知机

Multi-Layer Perceptron

多层感知机



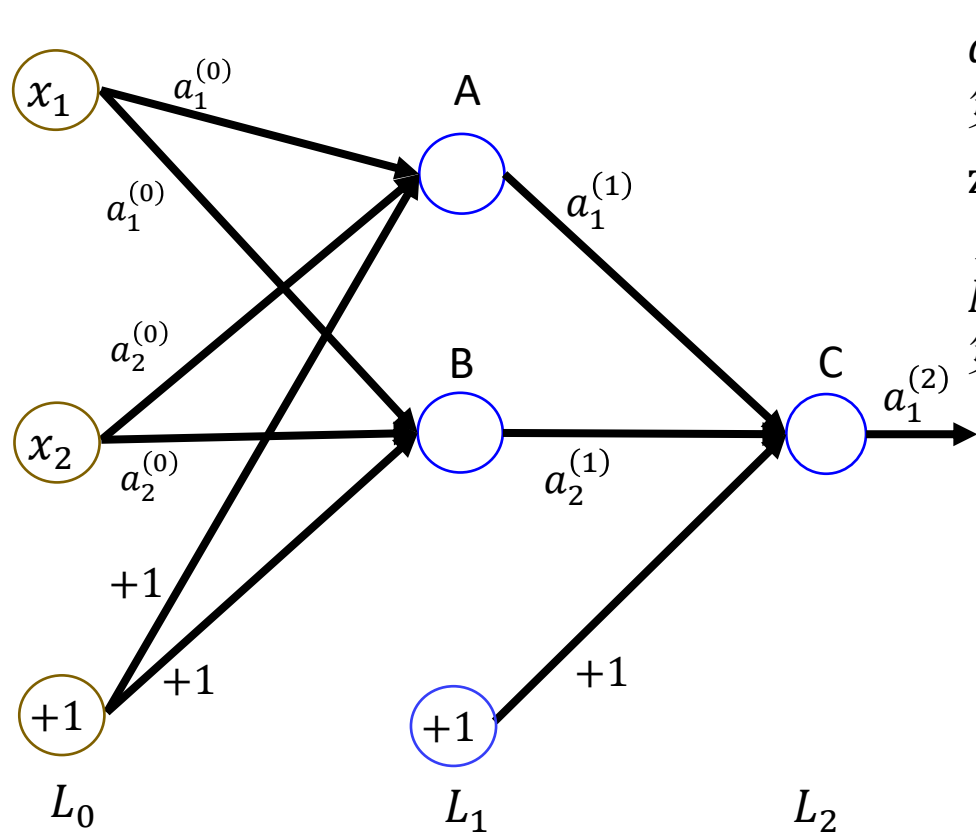
多层感知机：由多个神经元构成的一个神经网络，这些神经元分成若干层，相邻两层的神经元之间由突触连接。

$W_{ij}^{(l)}$:
第 l 层第 i 个神经元与第 $l - 1$ 层第 j 个神经元之间的连接权重参数

$b_i^{(l)}$:
第 l 层第 i 个神经元的偏置参数

第一层称为输入层(Input Layer)
最后一层称为输出层(Output Layer)
中间层称为隐层(Hidden Layer)

多层感知机



$a_i^{(l)}$:
第 l 层第 i 个神经元的响应(Activation)

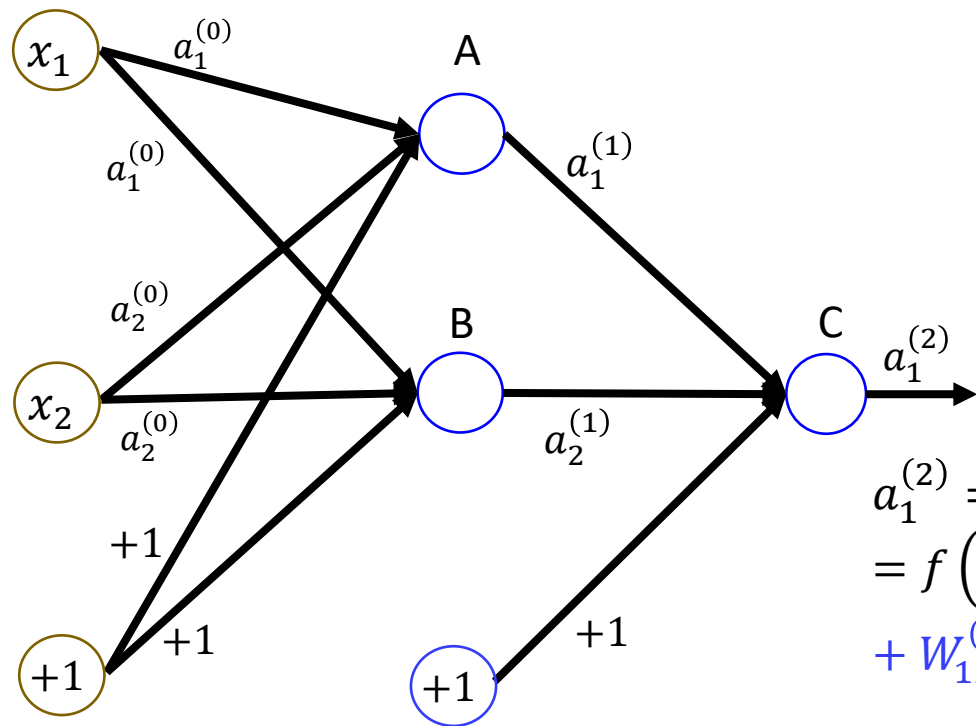
$z_i^{(l)}$:
第 l 层第 i 个神经元的净响应(Activation)

N_l :
第 l 层神经元数目

$$a_i^{(0)} = x_i, i = 1, \dots, N_0$$
$$z_i^{(l)} = \sum_{j=1}^{N_{l-1}} W_{ij}^{(l)} a_j^{(l-1)} + b_i^{(l)}$$
$$a_i^{(l)} = f(z_i^{(l)})$$

f : 非线性响应函数
(Activation Function)

多层感知机



$$a_1^{(2)} = f(W_{11}^{(2)} a_1^{(1)} + W_{12}^{(2)} a_2^{(1)} + b_1^{(2)})$$

$$a_1^{(1)} = f(W_{11}^{(1)} a_1^{(0)} + W_{12}^{(1)} a_2^{(0)} + b_1^{(1)})$$

$$= f(W_{11}^{(1)} x_1 + W_{12}^{(1)} x_2 + b_1^{(1)})$$

$$a_2^{(1)} = f(W_{21}^{(1)} a_1^{(0)} + W_{22}^{(1)} a_2^{(0)} + b_2^{(1)})$$

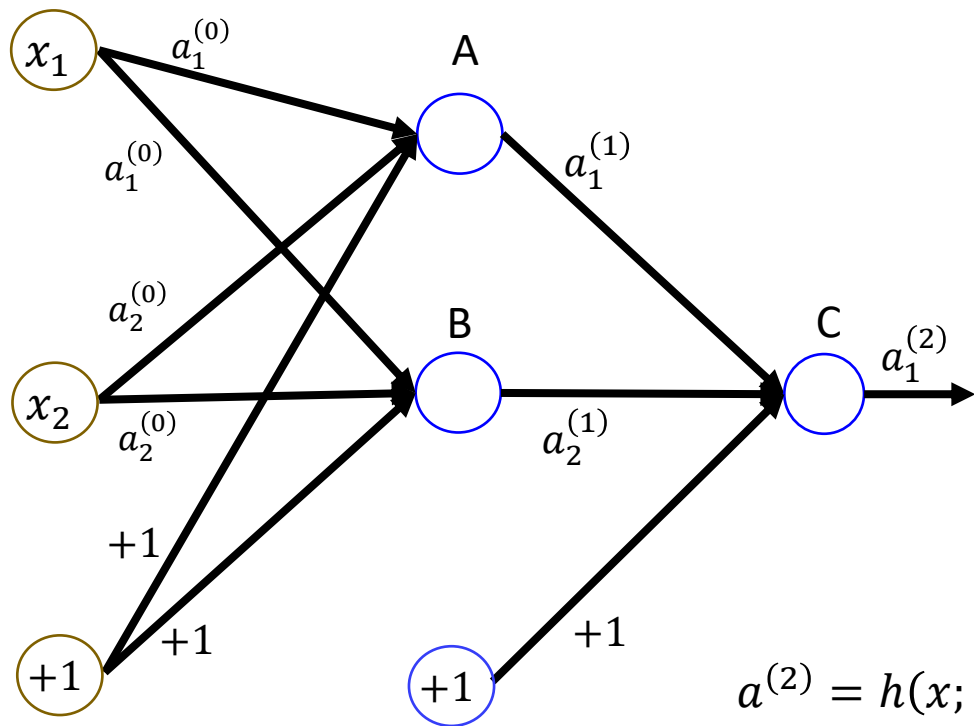
$$= f(W_{21}^{(1)} x_1 + W_{22}^{(1)} x_2 + b_2^{(1)})$$

$$a_1^{(2)} = h(x_1, x_2; \Theta)$$

$$= f(W_{11}^{(2)} f(W_{11}^{(1)} x_1 + W_{12}^{(1)} x_2 + b_1^{(1)}) + W_{12}^{(2)} f(W_{21}^{(1)} x_1 + W_{22}^{(1)} x_2 + b_2^{(1)}) + b_1^{(2)})$$

$$\Theta = \{W_{ij}^{(l)}, b_k^{(l)}\}$$

多层感知机



$$\begin{aligned}
 z^{(1)} &= (z_1^{(1)}, z_2^{(1)})^T \\
 &= \begin{pmatrix} W_{11}^{(1)} a_1^{(0)} + W_{12}^{(1)} a_2^{(0)} + b_1^{(1)} \\ W_{21}^{(1)} a_1^{(0)} + W_{22}^{(1)} a_2^{(0)} + b_2^{(1)} \end{pmatrix} \\
 &= W^{(1)} a^{(0)} + b^{(1)}
 \end{aligned}$$

$$a^{(1)} = f(z^{(1)}) = \begin{pmatrix} f(z_1^{(1)}) \\ f(z_2^{(1)}) \end{pmatrix}$$

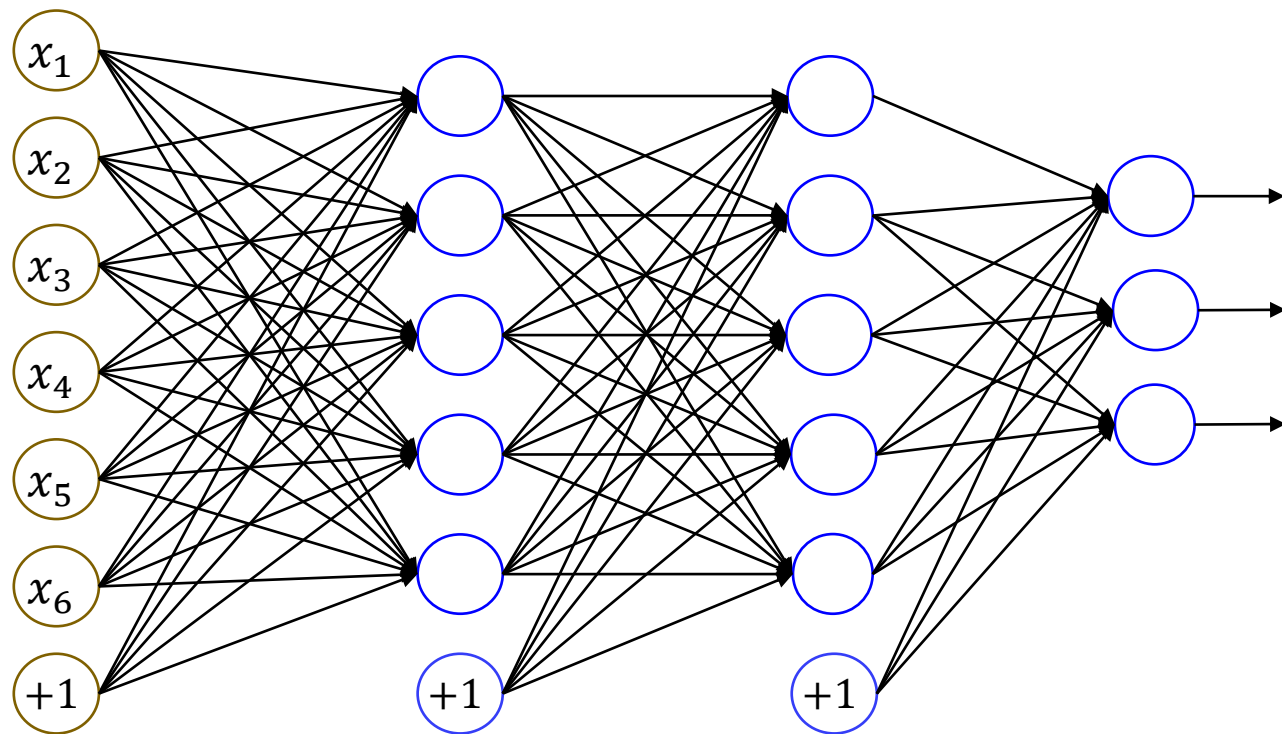
$$z^{(2)} = W^{(2)} a^{(1)} + b^{(2)}$$

$$a^{(2)} = f(z^{(2)})$$

$$a^{(2)} = h(x; \Theta) = f(W^{(2)} f(W^{(1)} x + b^{(1)}) + b^{(2)})$$

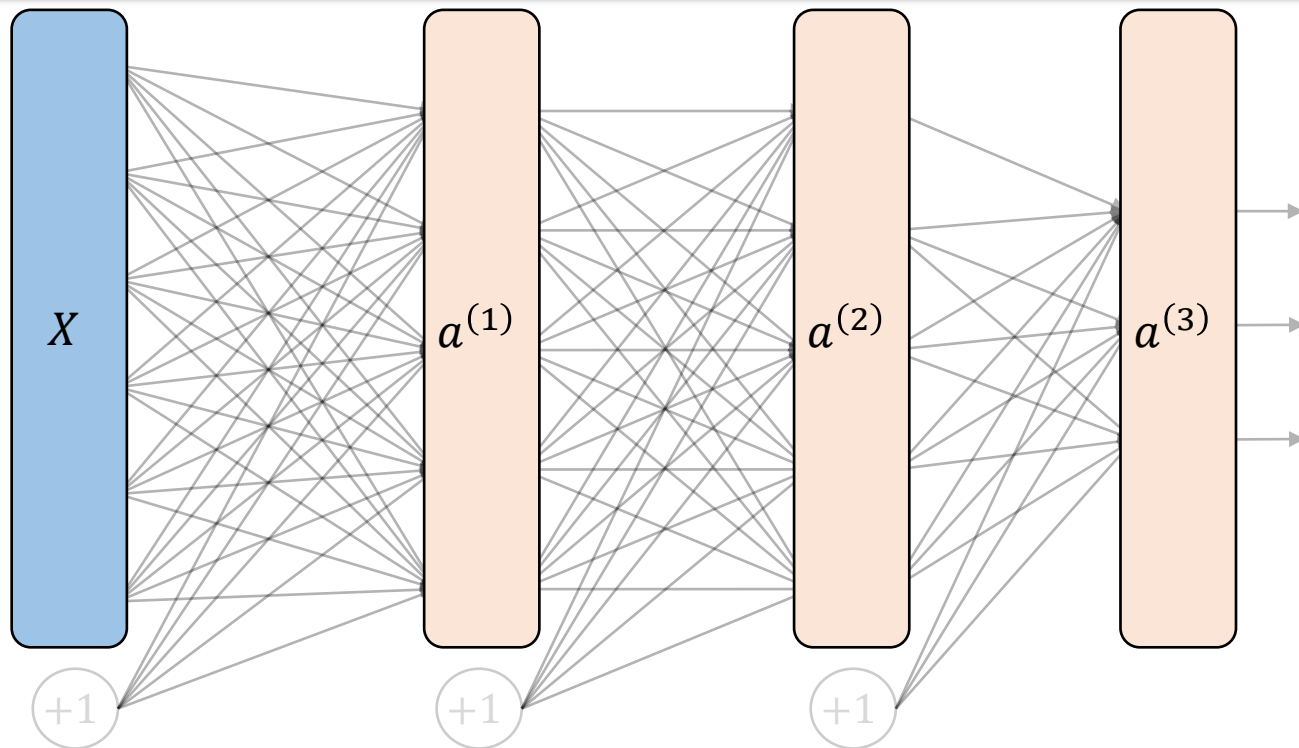
$$\Theta = \{W_{ij}^{(l)}, b_k^{(l)}\}$$

多层感知机



$$y = h(X; \Theta) = f(W^{(3)} f(W^{(2)} f(W^{(1)} X + b^{(1)}) + b^{(2)}) + b^{(3)})$$

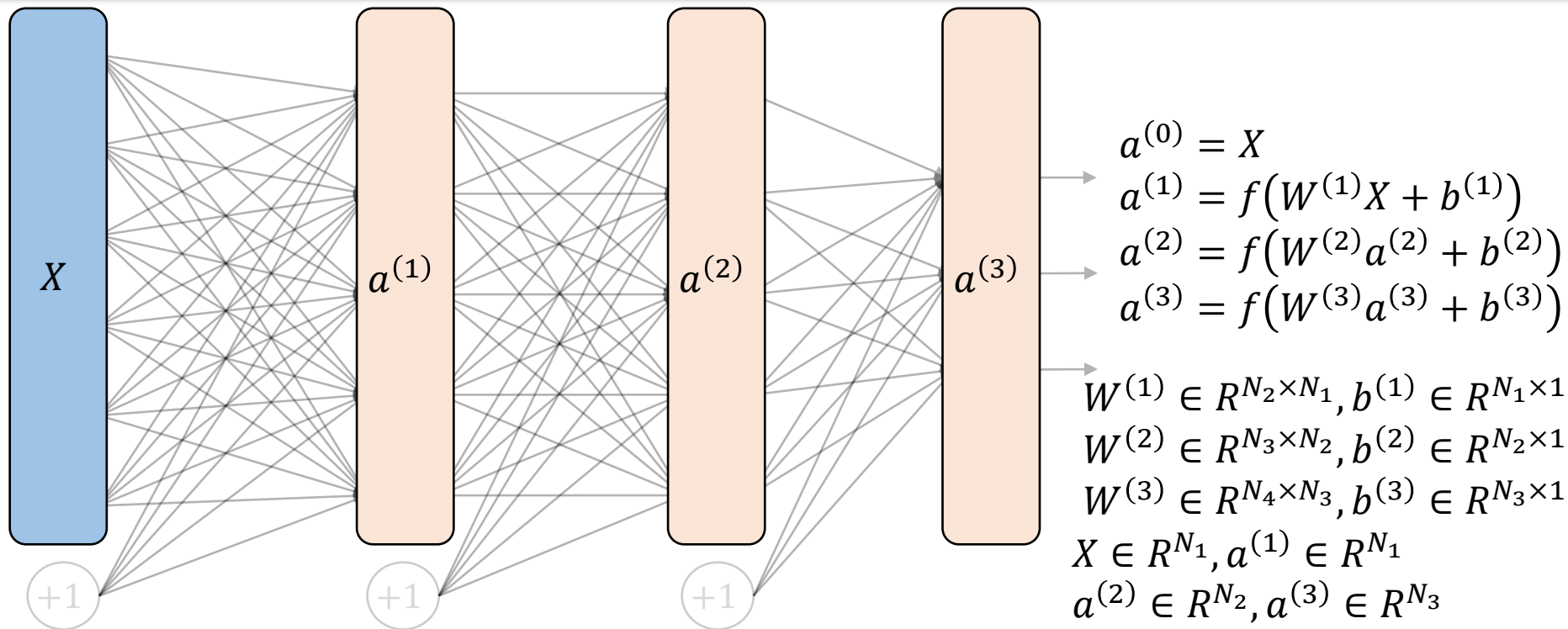
多层感知机



$$\begin{aligned}a^{(0)} &= X \\a^{(1)} &= f(W^{(1)}X + b^{(1)}) \\a^{(2)} &= f(W^{(2)}a^{(1)} + b^{(2)}) \\a^{(3)} &= f(W^{(3)}a^{(2)} + b^{(3)})\end{aligned}$$

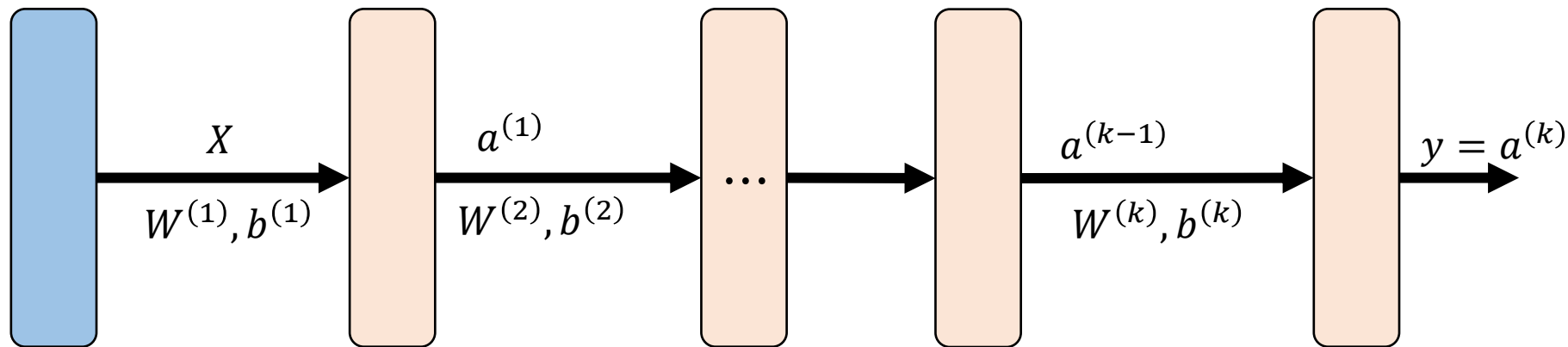
$$y = h(X; \Theta) = f(W^{(3)}f(W^{(2)}f(W^{(1)}X + b^{(1)}) + b^{(2)}) + b^{(3)})$$

多层感知机



$$y = h(X; \Theta) = f(W^{(3)}f(W^{(2)}f(W^{(1)}X + b^{(1)}) + b^{(2)}) + b^{(3)})$$

多层感知机=复合函数



$$a^{(l)} = f_l(a^{(l)}; W^{(l)}, b^{(l)}) = f_l(W^{(l)} a^{(l-1)} + b^{(l)})$$

$$y = h(X; \Theta) = f_k(f_{k-1}(\dots f_1(X; W^{(1)}, b^{(1)}) \dots; W^{(k-1)}, b^{(k-1)}); W^{(k)}, b^{(k)})$$

多层感知机学习算法

用梯度下降法训练多层感知器

- 输入: $D = \{X^{(i)}, Y^{(i)}\}_{i=1}^n$
- 输出: 网络参数 Θ
- 1. 初始化: 用随机数初始化 $W^{(l)}, b^{(l)}$, $l_0 = Inf$
- 2. *While True*:
 - 2.1 计算损失 $l_1 = l(\Theta; D)$
 - 2.2 *if* $|l_1 - l_0| < \epsilon$ *Break*, *else* $l_0 = l_1$
 - 2.3 计算梯度: $\Delta W^{(l)}, \Delta b^{(l)}$
 - 2.4 更新参数: $W^{(l)} = W^{(l)} - \alpha \Delta W^{(l)}, b^{(l)} = b^{(l)} - \alpha \Delta b^{(l)}$
- 3. 返回 Θ

损失函数

多层感知器:

$$y = h(X; \Theta) = f_k(f_{k-1}(\cdots f_2(X; W^{(1)}, b^{(1)}) \cdots; W^{(k-2)}, b^{(k-2)}); W^{(k-1)}, b^{(k-1)})$$

➤ 用损失函数评估参数 Θ 的"好坏"

➤ 根据预测的 \hat{Y} 与真实 Y 之间的差异计算参数 Θ 的损失:

➤ $l(\Theta; X, Y) = l(\hat{Y}, Y), \hat{Y} = h(X; \Theta)$

➤ 回归问题: $Y \in \mathbb{R}^m, \hat{Y} \in \mathbb{R}^m$

➤ 均方误差损失(MSE loss): $l(\Theta; X, Y) = \frac{1}{2} \|\hat{Y} - Y\|^2 = \frac{1}{2} \|h(X; \Theta) - Y\|^2$

➤ 分类问题: $Y \in \{0, 1, \dots, K\}, \hat{\rho} \in (0, 1)^K$

➤ 对数损失(log-loss)/交叉熵损失(Cross-Entropy Loss)

➤ $l(\Theta; X, Y) = -\sum_{j=1}^K 1(Y = j) \ln(\hat{\rho}_j), \rho = h(X; \Theta)$

计算梯度 $\partial l / \partial \Theta$

$$l(\Theta; X, Y) = l(\hat{Y}, Y)$$

$$\hat{Y} = h(X; \Theta) = f_k(f_{k-1}(\cdots f_1(X; W^{(1)}, b^{(1)}) \cdots; W^{(k-1)}, b^{(k-1)}); W^{(k)}, b^{(k)})$$

$$\Delta W^{(l)} = \frac{\partial l}{\partial W^{(l)}} = ?$$

$$\Delta b^{(l)} = \frac{\partial l}{\partial b^{(l)}} = ?$$

计算梯度 $\partial l / \partial \Theta$

$$l(\Theta; X, Y) = l(\hat{Y}, Y)$$

$$\hat{Y} = h(X; \Theta) = f_k(f_{k-1}(\cdots f_l(a^{(l)}; \mathbf{W}^{(l)}, \mathbf{b}^{(l)}) \cdots); \cdots)$$

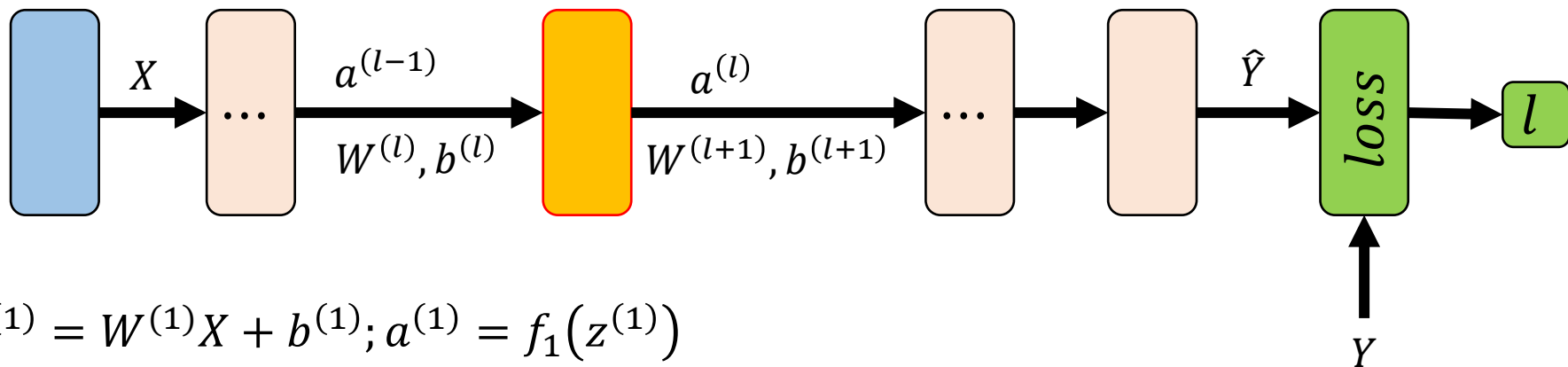
$$\Delta W^{(l)} = \frac{\partial l}{\partial W^{(l)}} = ?$$

$$\Delta b^{(l)} = \frac{\partial l}{\partial b^{(l)}} = ?$$

$$z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)}$$

$$a^{(l)} = f_l(z^{(l)})$$

计算梯度 $\partial l / \partial \Theta$



$$z^{(1)} = W^{(1)}X + b^{(1)}; a^{(1)} = f_1(z^{(1)})$$

\vdots

$$z^{(l)} = W^{(l)}a^{(l-1)} + b^{(l)}; a^{(l)} = f_l(z^{(l)})$$

$$z^{(l+1)} = W^{(l+1)}a^{(l)} + b^{(l+1)}; a^{(l+1)} = f_{l+1}(z^{(l+1)})$$

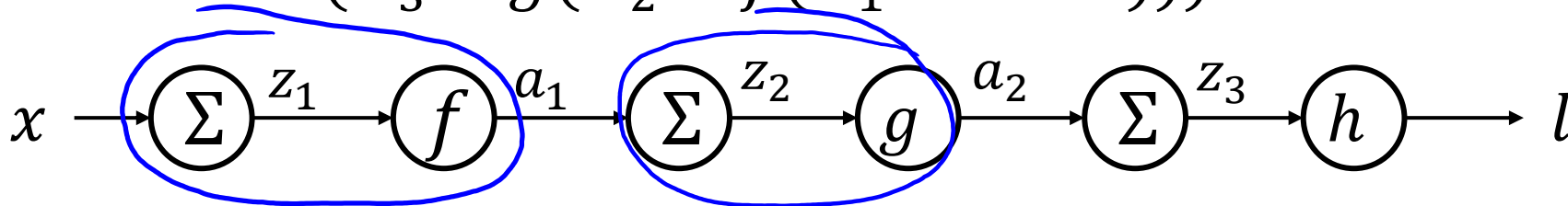
\vdots

$$\hat{Y} = f_K(z^{(K)}); l = l(\hat{Y}, Y)$$

$$\frac{\partial l}{\partial W^{(l)}} = \frac{\partial l}{\partial z^{(l)}} \cdot \frac{\partial z^{(l)}}{\partial W^{(l)}}$$

复合函数求导

$$l = h(w_3 \times g(w_2 \times f(w_1 \times x + b)))$$



$$l = h(z_3);$$

$$z_3 = w_3 \times a_2;$$

$$a_2 = g(z_2);$$

$$z_2 = w_2 \times a_1;$$

$$a_1 = f(z_1);$$

$$z_1 = w_1 \times x + b;$$

复合函数求导

$$l = h(w_3 \times g(w_2 \times f(w_1 \times x + b)))$$

$$l = h(z_3);$$

$$z_3 = w_3 \times a_2;$$

$$a_2 = g(z_2);$$

$$z_2 = w_2 \times a_1;$$

$$a_1 = f(z_1);$$

$$z_1 = w_1 \times x + b;$$

$$\frac{\partial l}{\partial z_3} = h'(z_3); \frac{\partial z_3}{\partial w_3} = a_2, \frac{\partial z_3}{\partial a_2} = w_3;$$

$$\frac{\partial a_2}{\partial z_2} = g'(z_2); \frac{\partial z_2}{\partial w_2} = a_1; \frac{\partial z_2}{\partial a_1} = w_2;$$

$$\frac{\partial a_1}{\partial z_1} = f'(z_1); \frac{\partial z_1}{\partial w_1} = x; \frac{\partial z_1}{\partial b} = 1$$

复合函数求导

$$l = h(w_3 \times g(w_2 \times f(w_1 \times x + b)))$$

$$\frac{\partial l}{\partial z_3} = h'(z_3); \frac{\partial z_3}{\partial w_3} = a_2, \frac{\partial z_3}{\partial a_2} = w_3;$$

$$\frac{\partial a_2}{\partial z_2} = g'(z_2); \frac{\partial z_2}{\partial w_2} = a_1; \frac{\partial z_2}{\partial a_1} = w_2;$$

$$\frac{\partial a_1}{\partial z_1} = f'(z_1); \frac{\partial z_1}{\partial w_1} = x; \frac{\partial z_1}{\partial b} = 1$$

$$\frac{\partial l}{\partial w_3} = \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial w_3} = h'(z_3) \times a_2,$$

$$\frac{\partial l}{\partial w_2} = \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial w_2} = h'(z_3) \times w_3 \times g'(z_2) \times a_1;$$

$$\begin{aligned} \frac{\partial l}{\partial w_1} &= \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \times \frac{\partial z_1}{\partial w_1} \\ &= h'(z_3) \times w_3 \times g'(z_2) \times w_2 \times f'(z_1) \times x; \end{aligned}$$

$$\begin{aligned} \frac{\partial l}{\partial b} &= \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \times \frac{\partial z_1}{\partial b} \\ &= h'(z_3) \times w_3 \times g'(z_2) \times w_2 \times f'(z_1) \times 1 \end{aligned}$$

复合函数求导

$$l = h(w_3 \times g(w_2 \times f(w_1 \times x + b)))$$

$$\frac{\partial l}{\partial w_3} = \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial w_3} = h'(z_3) \times a_2,$$

$$\frac{\partial l}{\partial w_2} = \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial w_2} = h'(z_3) \times w_3 \times g'(z_2) \times a_1;$$

$$\frac{\partial l}{\partial w_1} = \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \times \frac{\partial z_1}{\partial w_1} = h'(z_3) \times w_3 \times g'(z_2) \times w_2 \times f'(z_1) \times x;$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \times \frac{\partial z_1}{\partial b} = h'(z_3) \times w_3 \times g'(z_2) \times w_2 \times f'(z_1) \times 1$$

复合函数求导

$$\frac{\partial l}{\partial w_3} = \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial w_3},$$

$$\frac{\partial l}{\partial w_2} = \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial w_2};$$

$$\frac{\partial l}{\partial w_1} = \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \times \frac{\partial z_1}{\partial w_1};$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \times \frac{\partial z_1}{\partial b}$$

$$\delta_3 \equiv \frac{\partial l}{\partial z_3}$$

$$\begin{aligned}\delta_2 &\equiv \frac{\partial l}{\partial z_2} = \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \\ &= \delta_3 \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2}\end{aligned}$$

$$\begin{aligned}\delta_1 &\equiv \frac{\partial l}{\partial z_1} \\ &= \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \\ &= \delta_2 \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1}\end{aligned}$$

复合函数求导

$$\frac{\partial l}{\partial w_3} = \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial w_3},$$

$$\frac{\partial l}{\partial w_2} = \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial w_2};$$

$$\frac{\partial l}{\partial w_1} = \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \times \frac{\partial z_1}{\partial w_1};$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial z_3} \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \times \frac{\partial z_1}{\partial b}$$

$$\delta_3 \equiv \frac{\partial l}{\partial z_3}$$

$$\delta_2 \equiv \frac{\partial l}{\partial z_2} = \delta_3 \times \frac{\partial z_3}{\partial a_2} \times \frac{\partial a_2}{\partial z_2}$$

$$\delta_1 \equiv \frac{\partial l}{\partial z_1} = \delta_2 \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1}$$

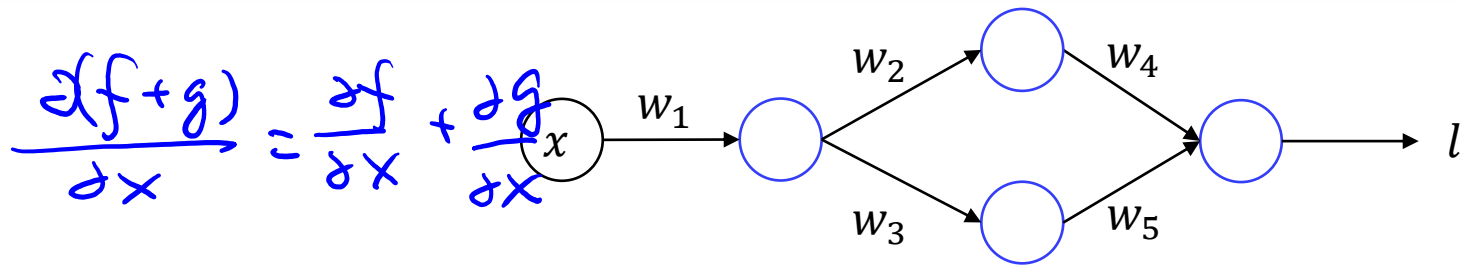
$$\frac{\partial l}{\partial w_3} = \delta_3 \times \frac{\partial z_3}{\partial w_3} = \delta_3 \times a_2,$$

$$\frac{\partial l}{\partial w_2} = \delta_2 \times \frac{\partial z_2}{\partial w_2} = \delta_2 \times a_1;$$

$$\frac{\partial l}{\partial w_1} = \delta_1 \times \frac{\partial z_1}{\partial w_1} = \delta_1 \times x;$$

$$\frac{\partial l}{\partial b} = \delta_1 \times \frac{\partial z_1}{\partial b} = \delta_1$$

复合函数求导



$$l = f(w_4 f(w_2 f(w_1 x + b_1) + b_2) + w_5 f(w_3 f(w_1 x + b_1) + b_3))$$

$$l = f(z_4);$$

$$z_4 = w_4 a_2 + w_5 a_3;$$

$$a_2 = f(z_2); z_2 = w_2 a_1 + b_2$$

$$a_3 = f(z_3); z_3 = w_3 a_1 + b_3;$$

$$a_1 = f(z_1); z_1 = w_1 x + b_1;$$

$$\frac{\partial l}{\partial w_1} = ? \quad \frac{\partial l}{\partial b_1} = ?$$

复合函数求导

$$l = f(w_4 f(w_2 f(w_1 x + b_1) + b_2) + w_5 f(w_3 f(w_1 x + b_1) + b_3) + b)$$

$$l = f(z_4);$$

$$z_4 = w_4 a_2 + w_5 a_3 + b;$$

$$a_2 = f(z_2); z_2 = w_2 a_1 + b_2$$

$$a_3 = f(z_3); z_3 = w_3 a_1 + b_3;$$

$$a_1 = f(z_1); z_1 = w_1 x + b_1;$$

$$\begin{aligned} \frac{\partial l}{\partial w_1} &= \frac{\partial l}{\partial z_4} \frac{\partial z_4}{\partial w_1} = \frac{\partial l}{\partial z_4} \frac{\partial z_4}{\partial a_2} \frac{\partial a_2}{\partial w_1} + \frac{\partial l}{\partial z_4} \frac{\partial z_4}{\partial a_3} \frac{\partial a_3}{\partial w_1} \\ &= \frac{\partial l}{\partial z_4} \frac{\partial z_4}{\partial a_2} \frac{\partial a_2}{\partial z_2} \frac{\partial z_2}{\partial w_1} + \frac{\partial l}{\partial z_4} \frac{\partial z_4}{\partial a_3} \frac{\partial a_3}{\partial z_3} \frac{\partial z_3}{\partial w_1} \\ &= \frac{\partial l}{\partial z_2} \frac{\partial z_2}{\partial w_1} + \frac{\partial l}{\partial z_3} \frac{\partial z_3}{\partial w_1} = \delta_2 \frac{\partial z_2}{\partial w_1} + \delta_3 \frac{\partial z_3}{\partial w_1} \end{aligned}$$

$$\delta_4 \equiv \frac{\partial l}{\partial z_4}, \delta_2 \equiv \frac{\partial l}{\partial z_2} = \frac{\partial l}{\partial z_4} \frac{\partial z_4}{\partial a_2} \frac{\partial a_2}{\partial z_2} = \delta_4 w_4 f'(z_2), \delta_3 \equiv \frac{\partial l}{\partial z_3} = \frac{\partial l}{\partial z_4} \frac{\partial z_4}{\partial a_3} \frac{\partial a_3}{\partial z_3} = \delta_4 w_5 f'(z_3)$$

复合函数求导

$$l = f(w_4 f(w_2 f(w_1 x + b_1) + b_2) + w_5 f(w_3 f(w_1 x + b_1) + b_3) + b)$$

$$l = f(z_4);$$

$$z_4 = w_4 a_2 + w_5 a_3 + b;$$

$$a_2 = f(z_2); z_2 = w_2 a_1 + b_2$$

$$a_3 = f(z_3); z_3 = w_3 a_1 + b_3;$$

$$a_1 = f(z_1); z_1 = w_1 x + b_1;$$

$$\frac{\partial l}{\partial w_1} = \delta_2 \frac{\partial z_2}{\partial w_1} + \delta_3 \frac{\partial z_3}{\partial w_1}$$

$$\frac{\partial z_2}{\partial w_1} = \frac{\partial z_2}{\partial a_1} \frac{\partial a_1}{\partial z_1} \frac{\partial z_1}{\partial w_1}; \quad \frac{\partial z_3}{\partial w_1} = \frac{\partial z_3}{\partial a_1} \frac{\partial a_1}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial l}{\partial w_1} = \delta_2 \frac{\partial z_2}{\partial a_1} \frac{\partial a_1}{\partial z_1} \frac{\partial z_1}{\partial w_1} + \delta_3 \frac{\partial z_3}{\partial a_1} \frac{\partial a_1}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$

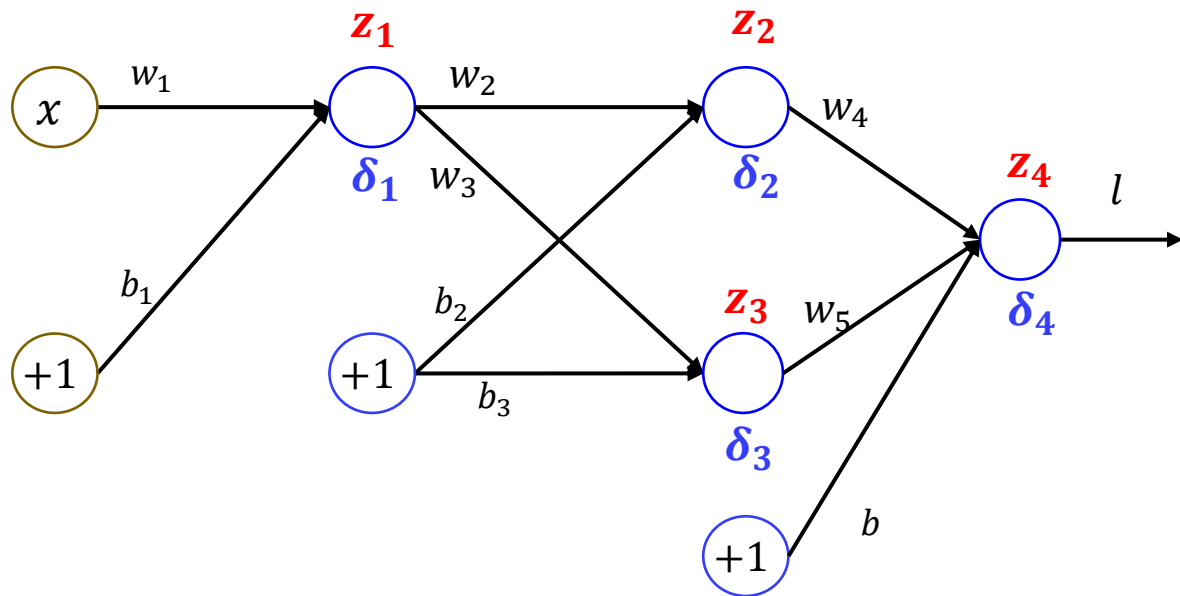
$$= \left(\delta_2 \frac{\partial z_2}{\partial a_1} \frac{\partial a_1}{\partial z_1} + \delta_3 \frac{\partial z_3}{\partial a_1} \frac{\partial a_1}{\partial z_1} \right) \frac{\partial z_1}{\partial w_1} = \delta_1 \frac{\partial z_1}{\partial w_1} = \delta_1 x$$

$$\frac{\partial l}{\partial z_1} \equiv \delta_1$$

$$\delta_1 \equiv \delta_2 \frac{\partial z_2}{\partial a_1} \frac{\partial a_1}{\partial z_1} + \delta_3 \frac{\partial z_3}{\partial a_1} \frac{\partial a_1}{\partial z_1} = \left(\delta_2 \frac{\partial z_2}{\partial a_1} + \delta_3 \frac{\partial z_3}{\partial a_1} \right) \frac{\partial a_1}{\partial z_1} = (\delta_2 w_2 + \delta_3 w_3) f'(z_1)$$

复合函数求导

$$l = f(w_4 f(w_2 f(w_1 x + b_1) + b_2) + w_5 f(w_3 f(w_1 x + b_1) + b_2) + b)$$



$$\delta_4 \equiv \frac{\partial l}{\partial z_4}$$

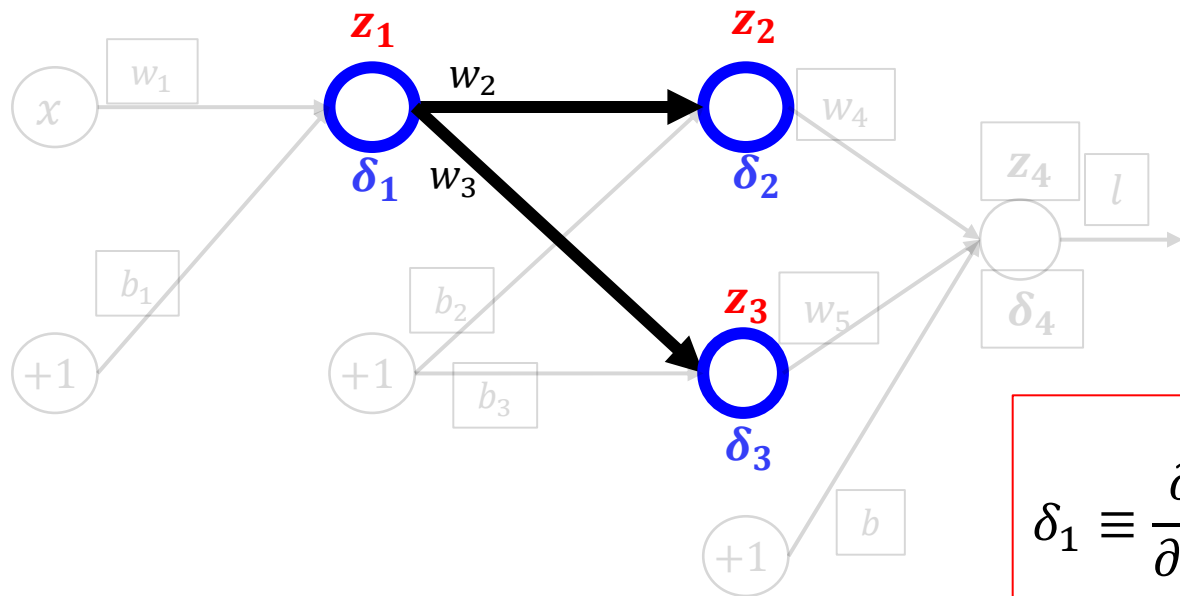
$$\delta_2 \equiv \frac{\partial l}{\partial z_2} = \delta_4 w_4 f'(z_2)$$

$$\delta_3 \equiv \frac{\partial l}{\partial z_3} = \delta_4 w_5 f'(z_3)$$

$$\delta_1 \equiv \frac{\partial l}{\partial z_1} = (\delta_2 w_2 + \delta_3 w_3) f'(z_1)$$

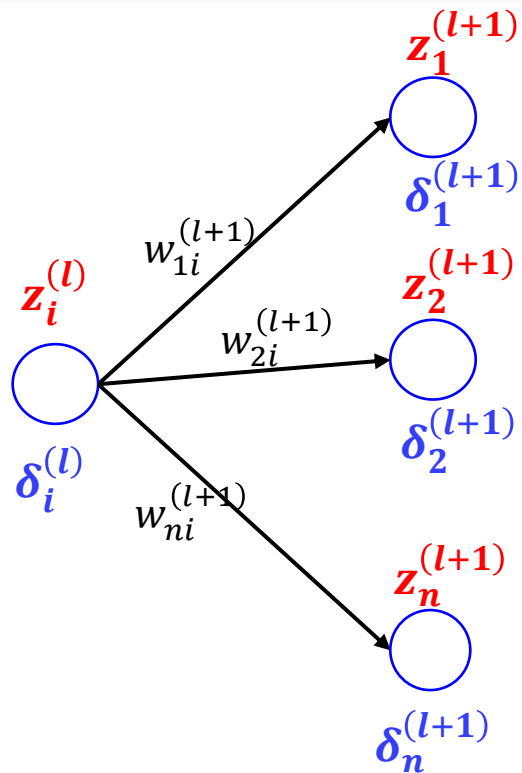
复合函数求导

$$l = f(w_4 f(w_2 f(w_1 x + b_1) + b_2) + w_5 f(w_3 f(w_1 x + b_1) + b_3) + b)$$



$$\delta_1 \equiv \frac{\partial l}{\partial z_1} = (\delta_2 w_2 + \delta_3 w_3) f'(z_1)$$

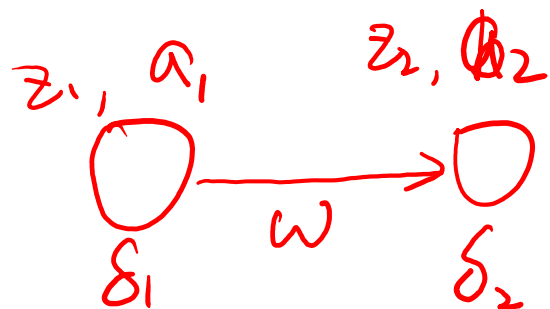
响应前向传播、误差后向传播(Error Back Propagation)



$$z_j^{(l+1)} = \dots + w_{ji}^{(l+1)} f(z_i^{(l)}) + \dots$$
$$a_j^{(l+1)} = f(z_j^{(l+1)})$$

$$\delta_i^{(l)} = \left[\sum_{j=1}^n w_{ji}^{(l+1)} \delta_j^{(l+1)} \right] f'(z_i^{(l)})$$

$$\Delta W_{ji}^{(l)} = \delta_j^{(l)} a_i^{(l)}$$



$$\underline{\Delta W = a_1 \delta_2} \quad \checkmark \text{ Update}$$

$$\delta_1 = \delta_2 \cdot w \cdot f'(z_1)$$

B.P.

$$a_2 = f(w a_1)$$

F.P.

误差后向传播

$$l(\Theta; X, Y) = l(\hat{Y}, Y)$$

$$\hat{Y} = h(X; \Theta) = f_k(f_{k-1}(\cdots f_1(X; W^{(1)}, b^{(1)}) \cdots; W^{(k-1)}, b^{(k-1)}); W^{(k)}, b^{(k)})$$

$$z^{(l)} = W^{(l)}a^{(l-1)} + b^{(l)}; a^{(l)} = f_l(z^{(l)})$$

$$\delta^{(l)} \equiv \frac{\partial l}{\partial z^{(l)}} \equiv \begin{pmatrix} \delta_1^{(l)} \\ \delta_2^{(l)} \\ \vdots \\ \delta_{N_l}^{(l)} \end{pmatrix}$$

$$\Delta W^{(l)} = \frac{\partial l}{\partial W^{(l)}} = \frac{\partial l}{\partial z^{(l)}} a^{(l-1)T} = \delta^{(l)} a^{(l-1)T}$$

$$\Delta b^{(l)} = \frac{\partial l}{\partial b^{(l)}} = \frac{\partial l}{\partial z^{(l)}} = \delta^{(l)}$$

$$\delta_i^{(l)} = \left[\sum_{j=1}^{N_{l+1}} w_{ji}^{(l+1)} \delta_j^{(l+1)} \right] f' \left(z_i^{(l)} \right) = \begin{bmatrix} w_{1i}^{(l+1)} \\ w_{2i}^{(l+1)} \\ \vdots \\ w_{N_{l+1}i}^{(l+1)} \end{bmatrix}^T \begin{bmatrix} \delta_1^{(l+1)} \\ \delta_2^{(l+1)} \\ \vdots \\ \delta_{N_{l+1}}^{(l+1)} \end{bmatrix} f' \left(z_i^{(l)} \right) = \begin{bmatrix} w_{1i}^{(l+1)} \\ w_{2i}^{(l+1)} \\ \vdots \\ w_{N_{l+1}i}^{(l+1)} \end{bmatrix}^T \delta^{(l+1)} f' \left(z_i^{(l)} \right)$$

误差后向传播

$$\delta^{(l)} \equiv \begin{pmatrix} \delta_1^{(l)} \\ \delta_2^{(l)} \\ \vdots \\ \delta_{N_l}^{(l)} \end{pmatrix}$$

$$\delta_i^{(l)} = \begin{bmatrix} w_{1i}^{(l+1)} \\ w_{2i}^{(l+1)} \\ \vdots \\ w_{N_{l+1}i}^{(l+1)} \end{bmatrix}^T \delta^{(l+1)} f'(z_i^{(l)})$$

$$A \odot B = \begin{bmatrix} A_1 B_1 \\ A_2 B_2 \\ \vdots \\ A_n B_n \end{bmatrix}$$

$$\begin{aligned} \delta^{(l)} \equiv \frac{\partial l}{\partial z^{(l)}} &\equiv \begin{pmatrix} \delta_1^{(l)} \\ \delta_2^{(l)} \\ \vdots \\ \delta_{N_l}^{(l)} \end{pmatrix} = \left\{ \begin{bmatrix} w_{1,1}^{(l+1)} & w_{2,1}^{(l+1)} & \cdots & w_{N_{l+1},1}^{(l+1)} \\ w_{1,2}^{(l+1)} & w_{2,2}^{(l+1)} & & w_{N_{l+1},2}^{(l+1)} \\ \vdots & \vdots & & \vdots \\ w_{1,N_l}^{(l+1)} & w_{2,N_l}^{(l+1)} & \cdots & w_{N_{l+1},N_l}^{(l+1)} \end{bmatrix} \delta^{(l+1)} \right\} \odot f'(z^{(l)}) \\ &= W^{(l+1)T} \delta^{(l+1)} \odot f'(z^{(l)}) \end{aligned}$$

$$f'(z^{(l)}) = \begin{bmatrix} f'(z_1^{(l)}) & f'(z_2^{(l)}) & \cdots & f'(z_{N_l}^{(l)}) \end{bmatrix}^T$$

使用BP算法计算梯度

$$l(\Theta; X, Y) = l(\hat{Y}, Y)$$

$$\hat{Y} = h(X; \Theta) = f(f(\dots f(X; W^{(1)}, b^{(1)}) \dots; W^{(k-1)}, b^{(k-1)}); W^{(k)}, b^{(k)})$$

➤ *Forward Propagation(FP)*

$$a^{(0)} = X; z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)}; a^{(l)} = f(z^{(l)}); \hat{Y} = a^{(L)}$$

$$l(\Theta; X, Y) = l(\hat{Y}, Y)$$

➤ *error Backward Propagation(BP)*

$$\delta^{(L)} = \frac{\partial l}{\partial z^{(L)}}; \delta^{(l)} = W^{(l+1)T} \delta^{(l+1)} \odot f'(z^{(l)})$$

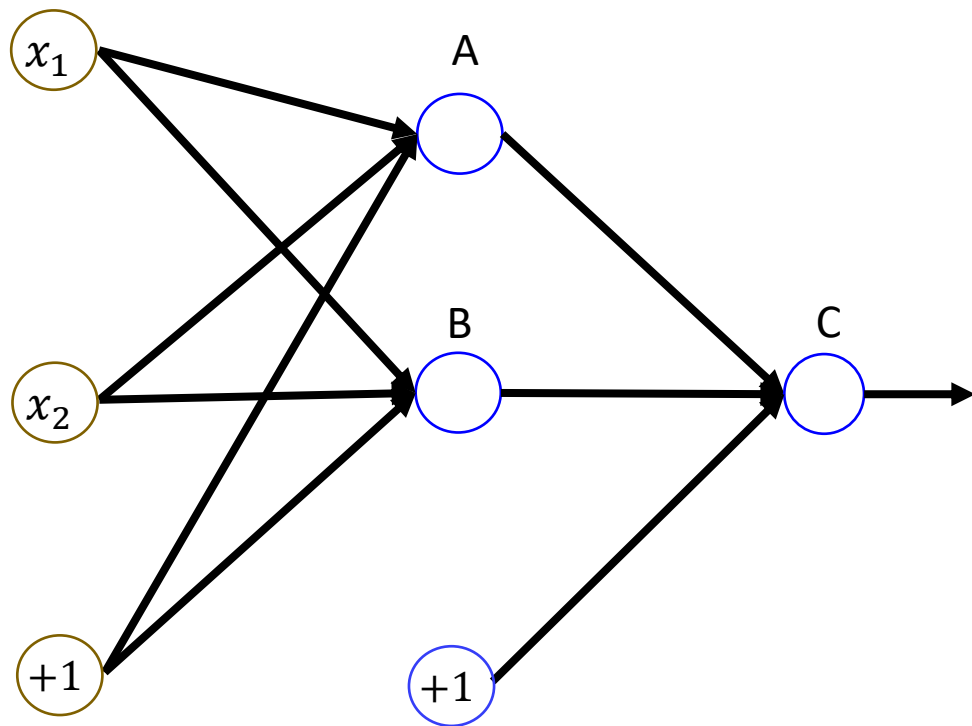
$$\Delta W^{(l)} = \delta^{(l+1)} a^{(l)T}; \Delta b^{(l)} = \delta^{(l+1)}$$

使用梯度下降法训练神经网络

- 输入: 训练样本 $D = \{(X^{(i)} \in R^d, Y^{(i)})\}_{i=1}^n$, 学习速率 α , 收敛条件 ϵ
- 输出: $\Theta = \{W^{(l)}, b^{(l)}\}_{l=1}^L$
- 1. 初始化 $\Theta, l_0 = \text{Inf}$
- 2. While True:
 - 2.1 FP: $\text{loss} = 0$; for $i = 1:n$ $\{a^{(0)} = X^{(i)}$ for $l = 1:L$ $\{z^{(l)} = W^{(l)}a^{(l-1)} + b^{(l)}; a^{(l)} = f(z^{(l)})\}$ $l_i = \text{loss}(a^{(L)}, Y^{(i)})$; $\text{loss} += l_i$; $\}$
 - 2.2 $l = \frac{l}{n}$; if $|l - l_0| < \epsilon$, break; else $l_0 = l$
 - 2.3 BP: $\delta^{(L)} = \frac{\partial l}{\partial z^{(L)}}$; for $l = L - 1:1$ $\{\Delta W^{(l)} = \delta^{(l)} a^{(l)T}, \Delta b = \delta^{(l)}; \delta^{(l)} = W^{(l+1)T} \delta^{(l+1)} \odot f'(z^{(l)})\}$
 - 2.4 GD: for $l = L:1$ $\{W^{(l)} = W^{(l)} - \alpha \Delta W^{(l)}, b^{(l)} = b^{(l)} - \alpha \Delta b^{(l)}\};$
- 3. 返回 Θ

批梯度下降(Batch Gradient Descent)

权值初始化



$$w_{ij}^{(l)} \sim N(0, \sigma)$$

梯度下降法(Gradient Descend)

- Batch GD:

- $l(\theta) = \frac{1}{n} \sum_{i=1}^n l(\theta; \mathbf{x}^{(i)}, y^{(i)})$

- $\theta^{(t+1)} = \theta^{(t)} - \eta \nabla l(\theta^{(t)}) = \theta^{(t)} - \eta \frac{1}{n} \sum_{i=1}^n \nabla l(\theta^{(t)}; \mathbf{x}^{(i)}, y^{(i)})$

- Stochastic GD:

- $\theta^{(t+1)} = \theta^{(t)} - \eta \nabla l(\theta^{(t)}; \mathbf{x}^{(i)}, y^{(i)})$

- Mini-Batch GD:

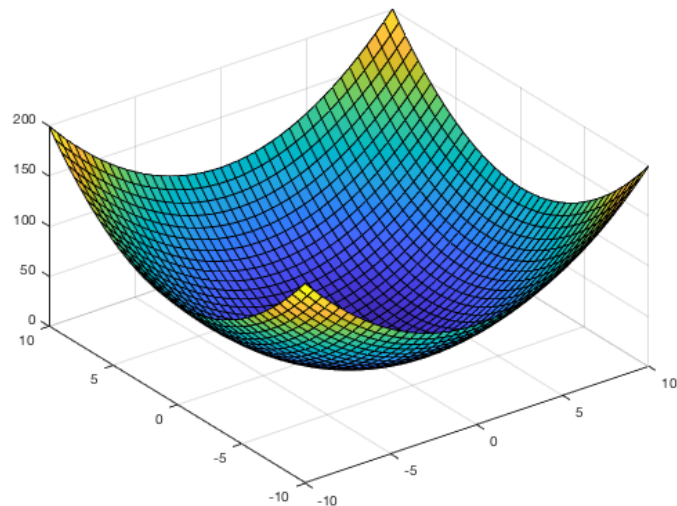
- $\theta^{(t+1)} = \theta^{(t)} - \eta \nabla l(\theta^{(t)}; B^{(k)})$

- $B^{(k)} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=(k-1)*batch_size+1}^{k*batch_size}$

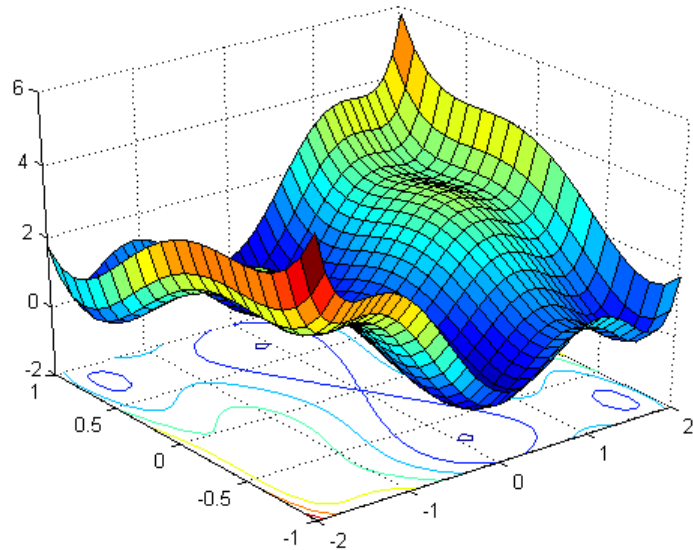
- Epoch: 遍历整个样本集合 \mathcal{D} , 每一个Epoch随机排列 \mathcal{D} 中元素

1次迭代(iteration):
前向+后向传播
权值更新

损失函数 $l(\theta)$ 的性质

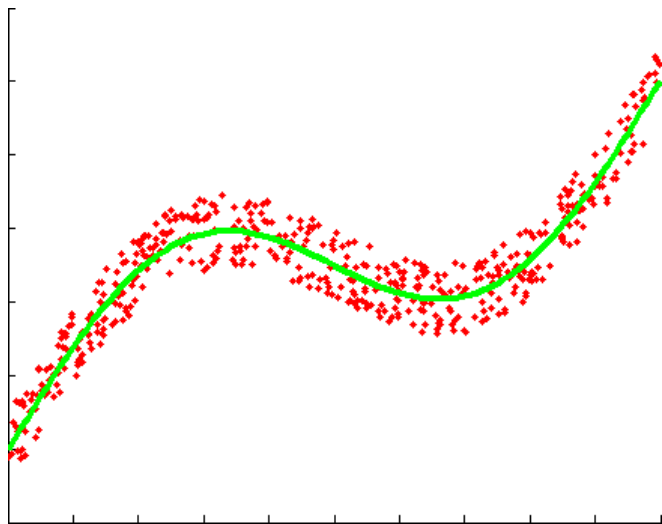


1个神经元, Logistic Regression

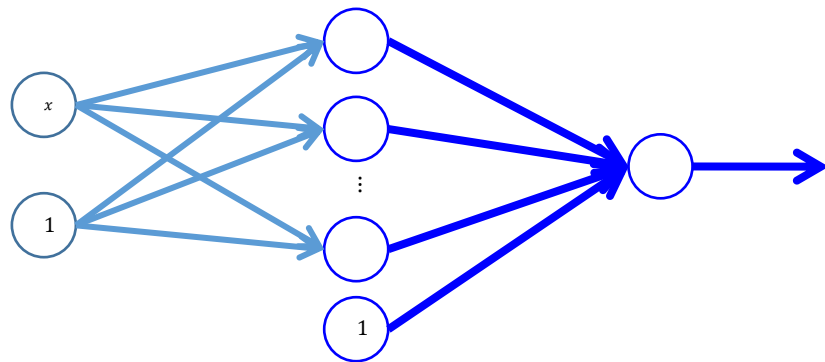


1个或多个隐层, 非线性响应函数, MLP

实例：曲线拟合

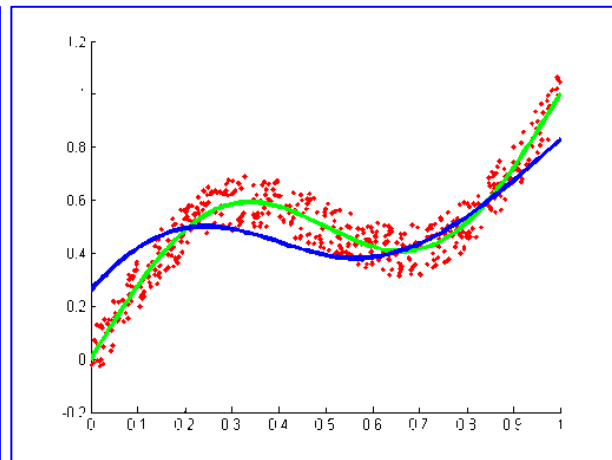
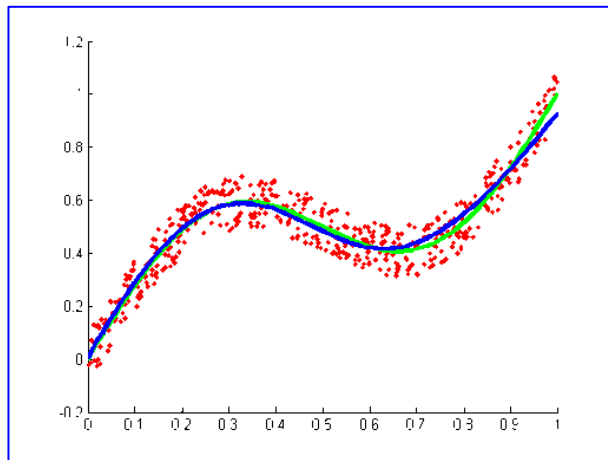
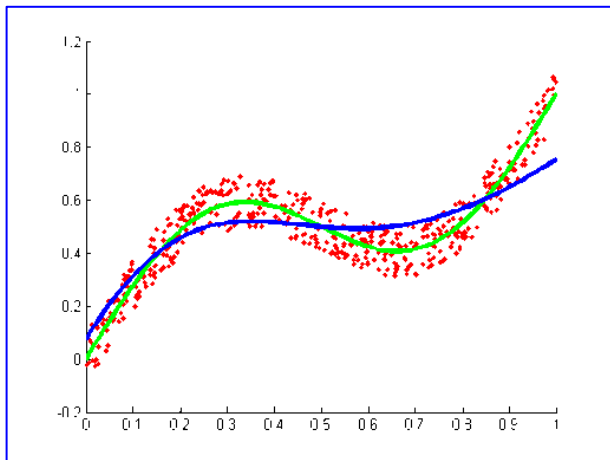


$$y = x + 0.3 \sin(2\pi x) + \varepsilon \quad x \in [0, 1] \quad \varepsilon \in [-0.1, 0.1]$$



$$l(h; x, y) = \frac{1}{2} (h(x) - y)^2$$

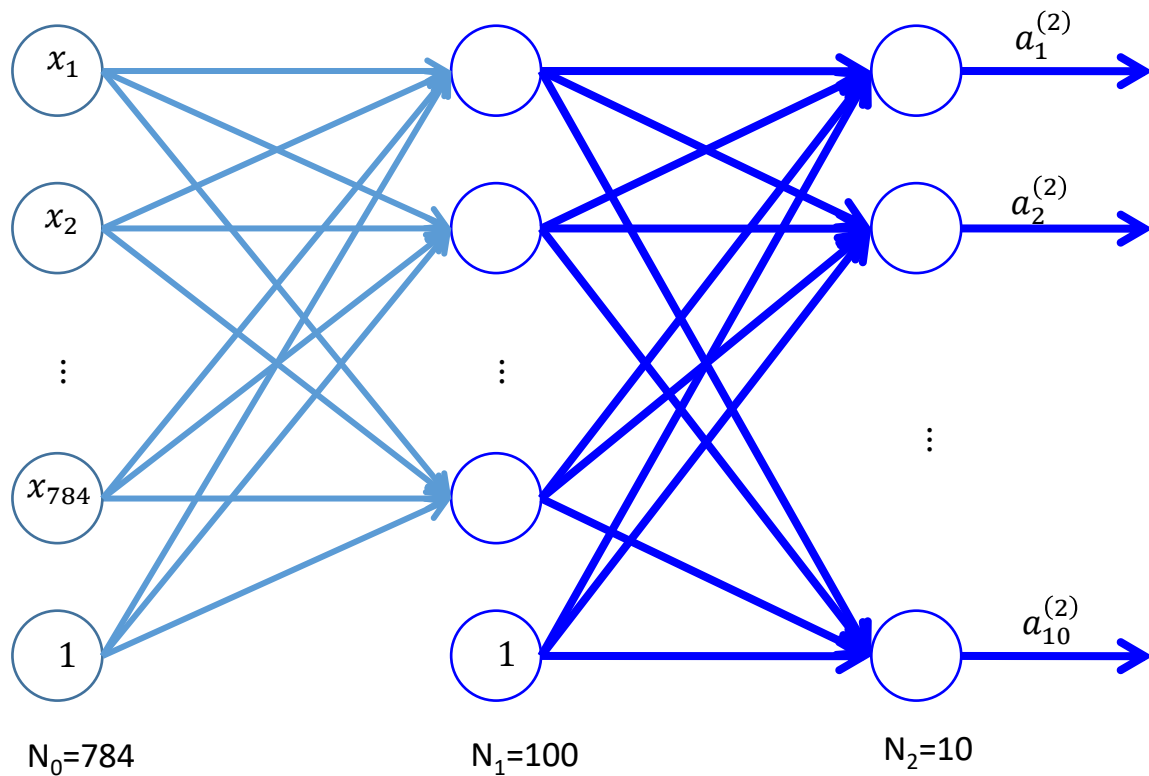
实例：曲线拟合



示例：MNIST手写体数字识别



实例:MNIST手写数字识别



实例:MNIST手写数字识别

$$P(y = k|x) = a_k^{(2)}(x) = \frac{\exp(z_k^{(2)}(x))}{\sum_{i=1}^{10} \exp(z_i^{(2)}(x))}$$

Softmax

$$P(y|x) = \prod_{j=1}^{10} [P(y = j|x)]^{y=j}$$

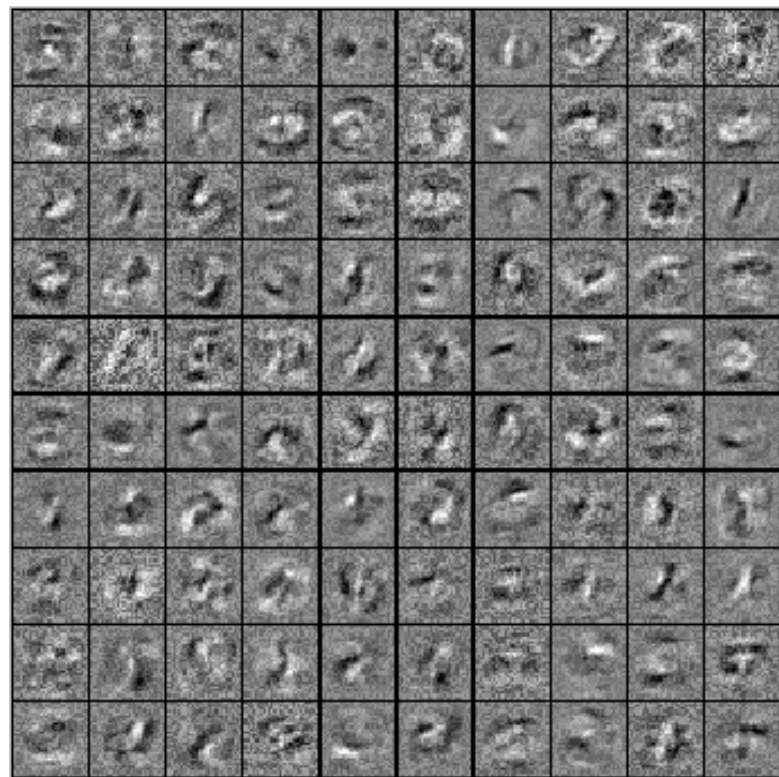
$$l(\theta; x, y) = -\log P(y|x) = -\sum_{j=1}^{10} \{y_j \log [a_j^{(2)}(x)]\}$$

$$l_{reg}(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^{10} \{y_j^{(i)} \log [a_j^{(2)}(x^{(i)})]\} + \frac{\lambda}{2m} \left\{ \sum_{i=1}^{784} \sum_{j=1}^{100} (w_{j,i}^{(1)})^2 + \sum_{i=1}^{100} \sum_{j=1}^{10} (w_{j,i}^{(2)})^2 \right\}$$

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_{10} \end{pmatrix}; y_k = \begin{cases} 1 & x \in k \\ 0 & x \notin k \end{cases}$$

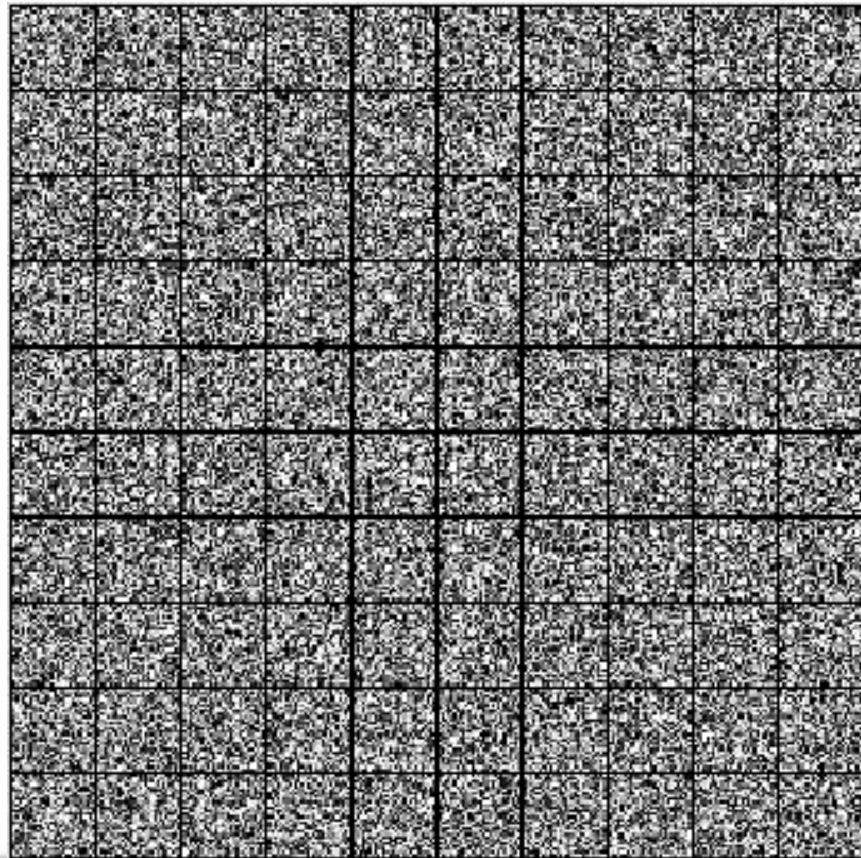
One-hot vector

实例:MNIST手写数字识别



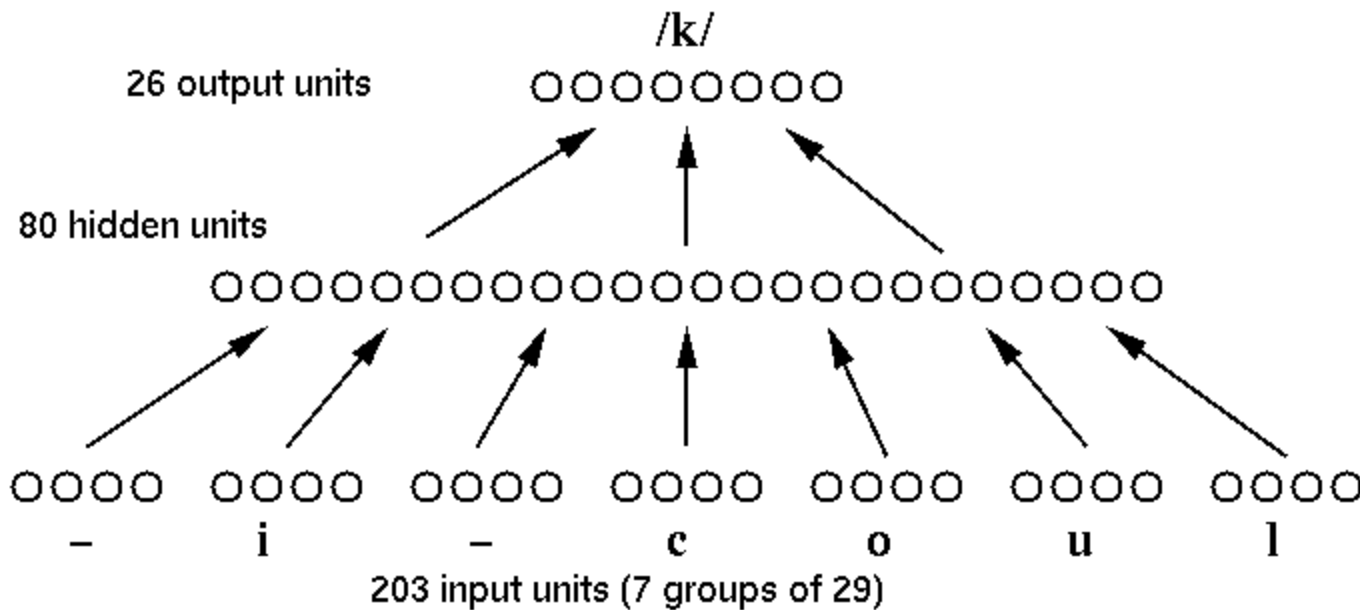
准确率=97.3%

实例:MNIST手写数字识别



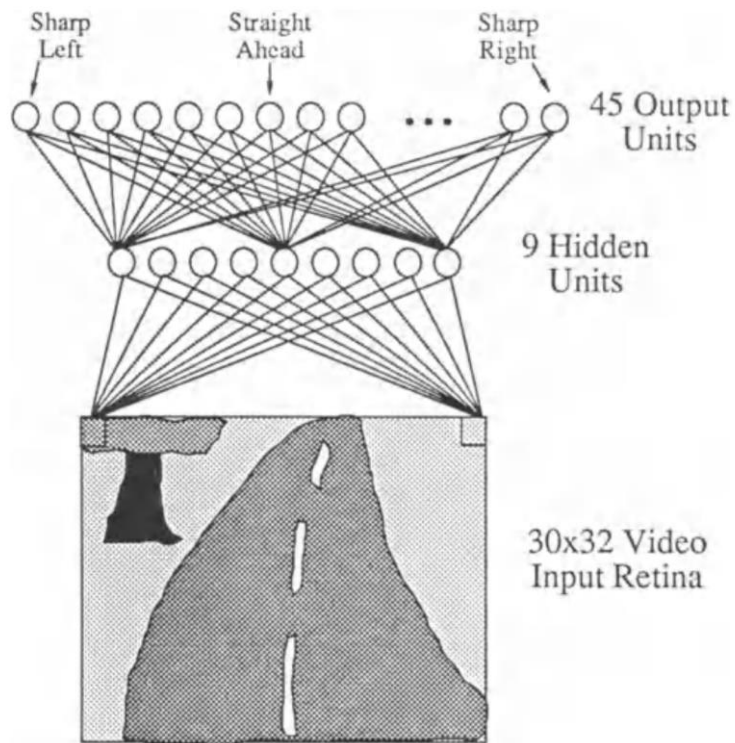
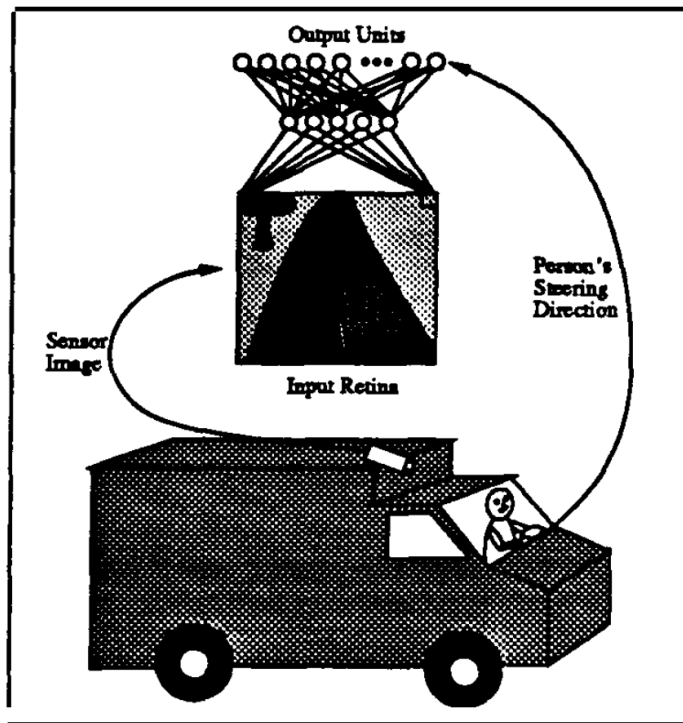
应用实例

NETalk, 1987



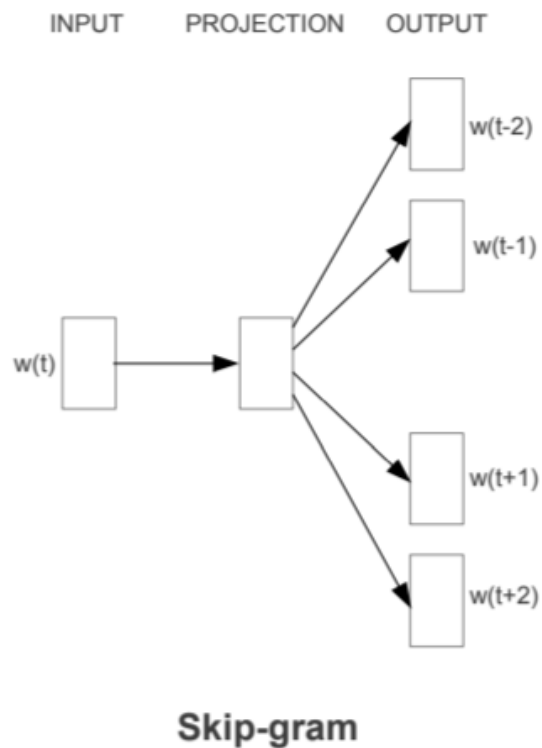
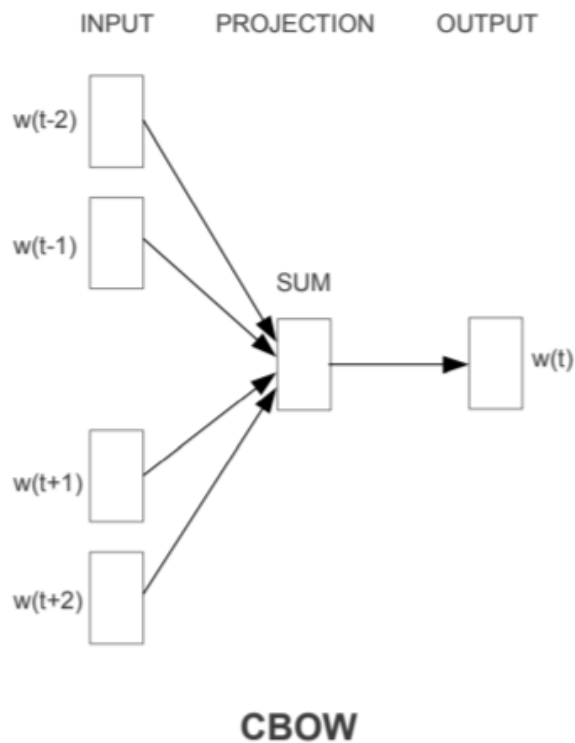
Sejnowski, T.J., and Rosenberg, C.R. (1987). "Parallel networks that learn to pronounce English text" in *Complex Systems*, 1, 145-168

ALVINN, CMU 1988



Dean Pomerleau : Autonomous Land Vehicle in a Neural Network, NIPS, 1988

词嵌入(Word Embedding)



Skip-gram

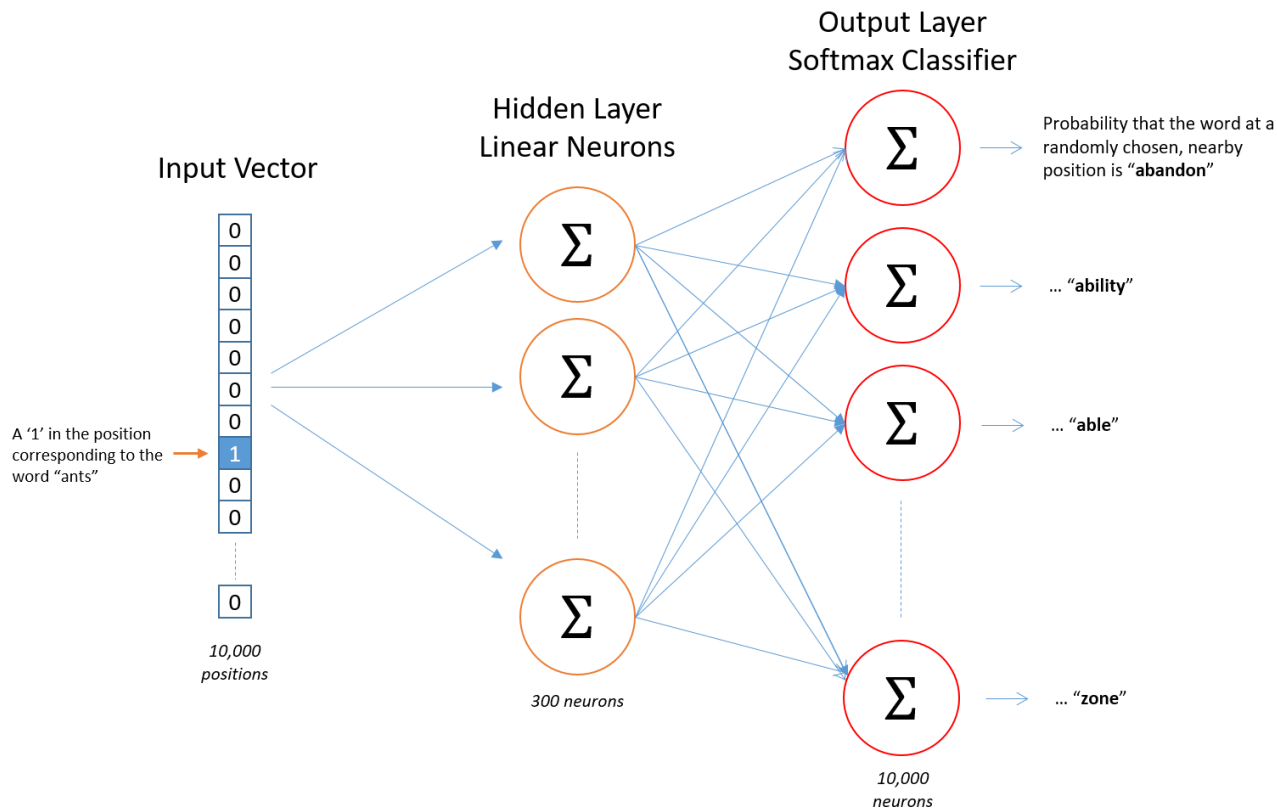
Source Text

Training Samples

<div>The quick brown fox jumps over the lazy dog.</div>	→	(the, quick) (the, brown)
<div>The quick brown fox jumps over the lazy dog.</div>	→	(quick, the) (quick, brown) (quick, fox)
<div>The quick brown fox jumps over the lazy dog.</div>	→	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
<div>The quick brown fox jumps over the lazy dog.</div>	→	(fox, quick) (fox, brown) (fox, jumps) (fox, over)

<http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/>

Skip-gram



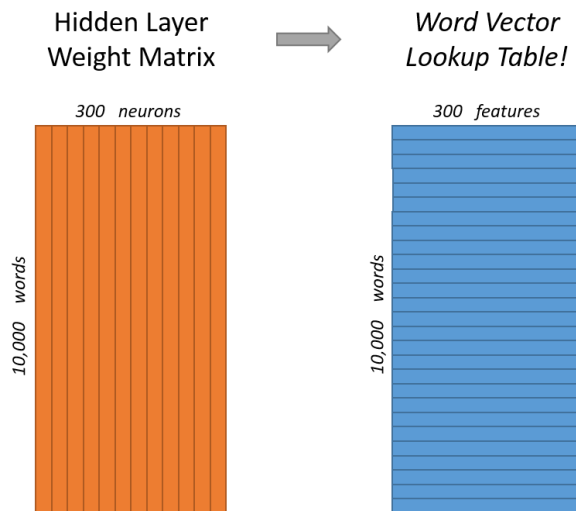
Skip-gram

$$|V| = 10,000 \quad x \equiv \text{ont_hot}('the') = [0, 0, \dots, 0, 1, 0, \dots, 0] \in \{0, 1\}^{10,000}$$

$$\text{index}('the') = 5793: x_{5793} = 1$$

$$E \in R^{10,000 \times 300}: E = [e_1, e_2, \dots, e_{10,000}]$$
$$z = xE = e_{5793} \leftarrow \text{Embedding of word 'the'}$$

$$\begin{bmatrix} 0 & 0 & 0 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ 10 & 12 & 19 \\ 11 & 18 & 25 \end{bmatrix} = \begin{bmatrix} 10 & 12 & 19 \end{bmatrix}$$



Skip-gram

$$|V| = 10,000 \quad x \equiv \text{ont_hot}('the') = [0, 0, \dots, 0, 1, 0, \dots, 0] \in \{0, 1\}^{10,000}$$

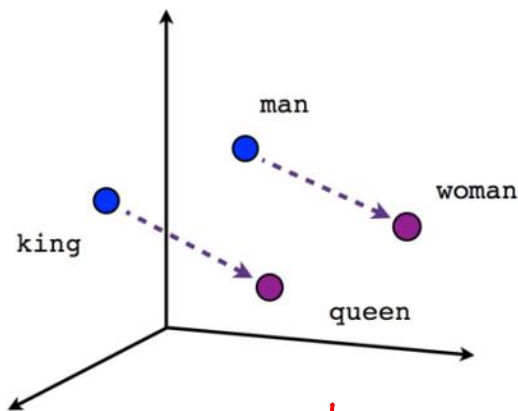
$$\text{index}('the') = 5793: x_{5793} = 1$$

$$E \in R^{10,000 \times 300}: E = [e_1, e_2, \dots, e_{10,000}]$$
$$z = xE = e_{5793} \leftarrow \text{Embedding of word 'the'}$$

$$\text{输出层: } y = Wz + b, W \in R^{10,000 \times 300}, b \in R^{10,000}$$

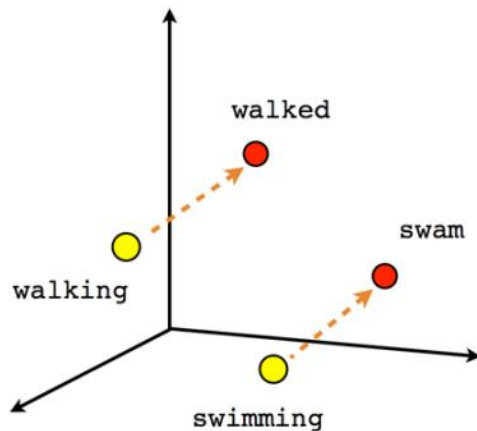
$$\rho = \text{softmax}(y) \equiv [\rho_1, \rho_2, \dots, \rho_{10,000}], \rho_i = P(\text{word}_i \text{ 是 } the \text{ 的近邻})$$

词向量

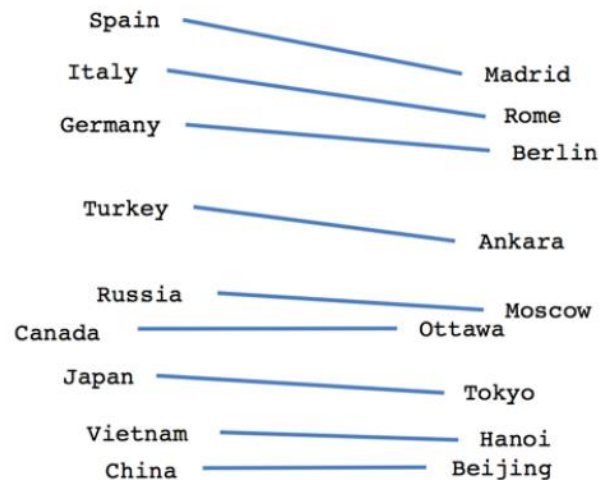


Queen - king + Man
= Woman

Male-Female



Verb tense



Country-Capital

神经语言模型(Neural Language Model)

“The quick brown fox jumps over the lazy dog.”

句子=单词序列: $S = (w_1, w_2, \dots, w_T)$

语言模型: 句子的概率分布模型 $P(S) = ?$

$$\begin{aligned} P(S) &= P(w_1, w_2, \dots, w_T) \\ &= P(w_1)P(w_2|w_1)P(w_3|w_2, w_1) \cdots P(w_T|w_{T-1}, \dots, w_1) \end{aligned}$$

神经语言模型(Neural Language Model)

$$\begin{aligned} P(S) &= P(w_1, w_2, \dots, w_T) \\ &= P(w_1)P(w_2|w_1)P(w_3|w_2, w_1) \cdots P(w_T|w_{T-1}, \dots, w_1) \end{aligned}$$

$$n - gram: P(w_t|w_{t-1}, w_{t-2}, \dots, w_1) = \prod_{i=1}^T P(w_t|w_{t-1}, w_{t-2}, \dots, w_{t-n+1})$$

$$n = 1: unigram: P(w_1, w_2, \dots, w_T) = \prod_{i=1}^T P(w_i)$$

$$n = 2: bigram: P(w_1, w_2, \dots, w_T) = \prod_{i=1}^T P(w_i|w_{i-1})$$

$$n = 3: trigram: P(w_1, w_2, \dots, w_T) = \prod_{i=1}^T P(w_i|w_{i-1}, w_{i-2})$$

神经语言模型(Neural Language Model)

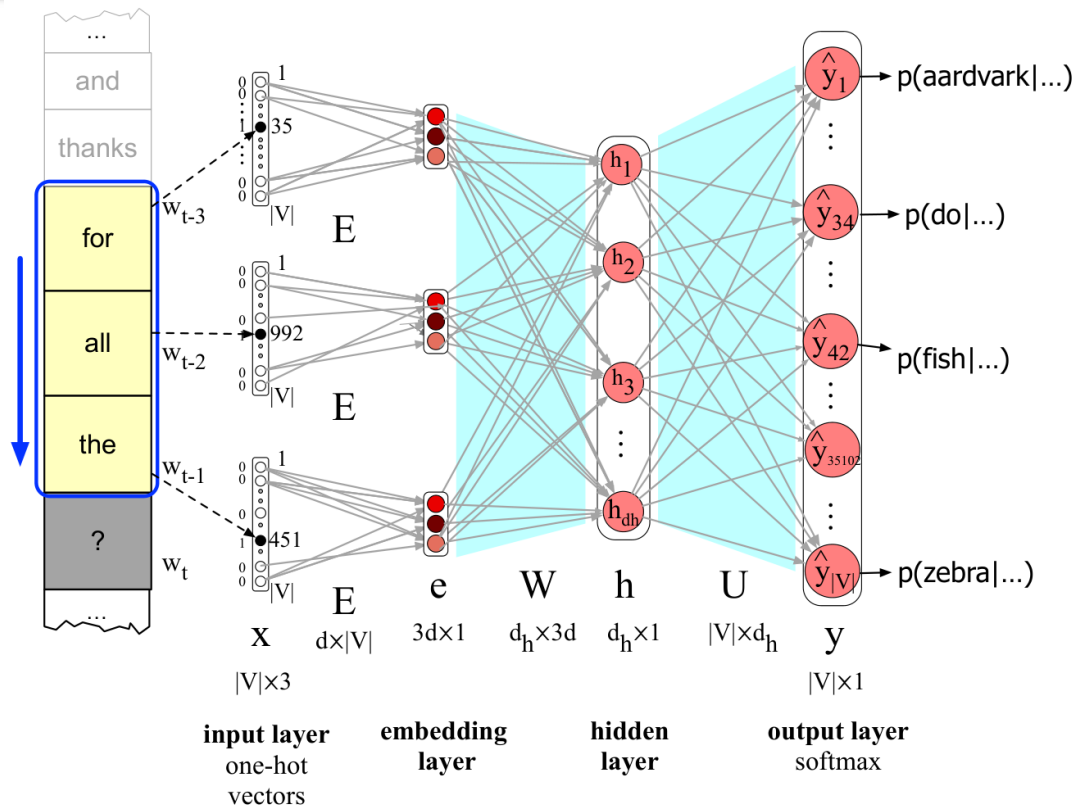
$$n - gram: P(w_t | w_{t-1}, w_{t-2}, \dots, w_1) = \prod_{i=1}^T P(w_t | w_{t-1}, w_{t-2}, \dots, w_{t-n+1})$$

如何表示 $P(w_t | w_{t-1}, w_{t-2}, \dots, w_{t-n+1})$?

输入 $w_{t-1}, w_{t-2}, \dots, w_{t-n+1}$, 输出 $P(w_t | w_{t-1}, w_{t-2}, \dots, w_{t-n+1})$

神经语言模型：用一个神经网络表示 $P(w_t | w_{t-1}, w_{t-2}, \dots, w_{t-n+1})$ ；输入层为 $n-1$ 个单词 $w_{t-1}, w_{t-2}, \dots, w_{t-n+1}$ ，输出层为下一个单词的概率分布 $P(w_t | w_{t-1}, w_{t-2}, \dots, w_{t-n+1})$

神经语言模型(Bengio,2003)



Bengio, Y. et al: A neural probabilistic language model. JMLR, 3:1137–1155. 2003

神经语言模型(Bengio,2003)

$$|V| = 10,000 \quad P(w| 'for', 'all', 'the')$$

$$x_1 \equiv \text{ont_hot}('for') = [0, 0, \dots, 1, \dots, 0, \dots, 0] \in \{0, 1\}^{10,000}$$

$$x_2 \equiv \text{ont_hot}('all') = [0, 0, \dots, 0, 1, 0, \dots, 0] \in \{0, 1\}^{10,000}$$

$$x_3 \equiv \text{ont_hot}('the') = [0, 0, \dots, 0, 0, 0, \dots, 1, \dots, 0] \in \{0, 1\}^{10,000}$$

$$E \in R^{10,000 \times 300}: E = [e_1, e_2, \dots, e_{10,000}]^T, e_i \in R^{300}$$

$$v_1 = x_1 E = e_{3079} \Leftarrow \text{Embedding of word 'for'}$$

$$v_2 = x_2 E = e_{1908} \Leftarrow \text{Embedding of word 'all'}$$

$$v_3 = x_3 E = e_{5793} \Leftarrow \text{Embedding of word 'the'}$$

$$v = (v_1, v_2, v_3) \in R^{900}$$

$$z_h = W_h v + b_h \in R^{1,000}, a_h = \tanh(z_h)$$

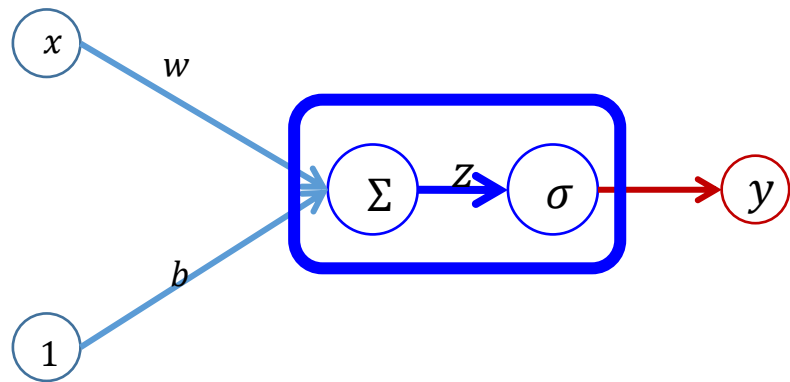
$$W_h \in R^{1,000 \times 900}, b_h \in R^{1,000}$$

$$P(w| 'for', 'all', 'the') = \text{softmax}(W_o a_h + b_o), W_o \in R^{10,000 \times 1,000}, b_o \in R^{10,000}$$

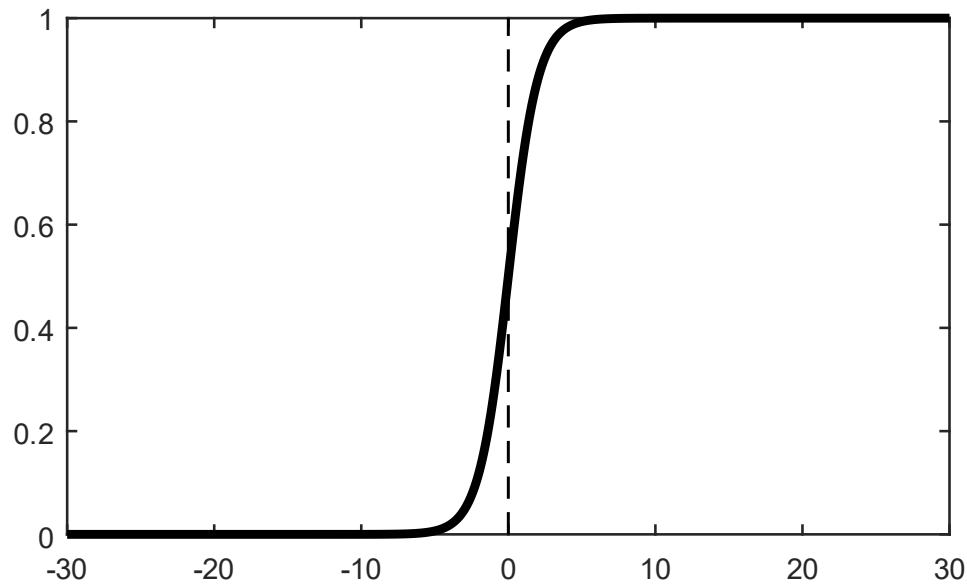
Bengio, Y. et al: A neural probabilistic language model. JMLR, 3:1137–1155. 2003

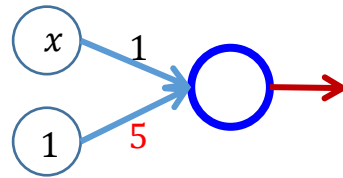
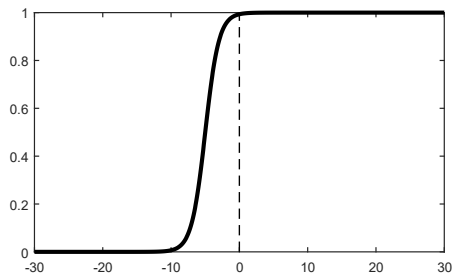
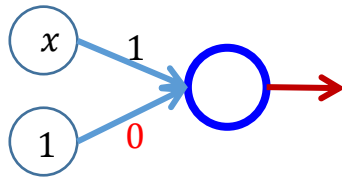
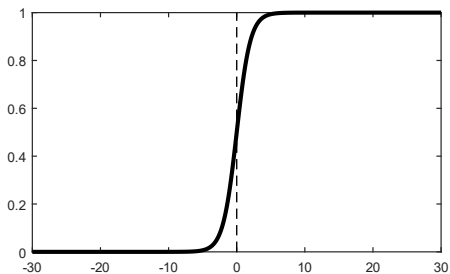
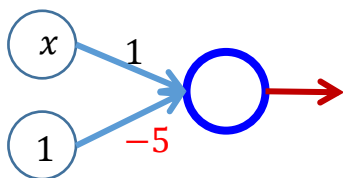
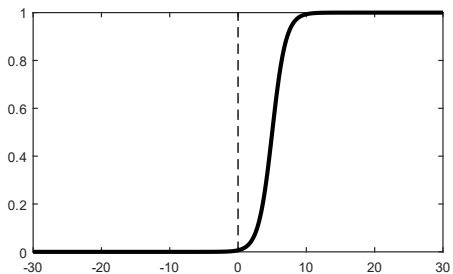
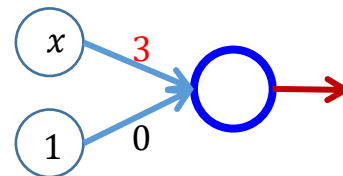
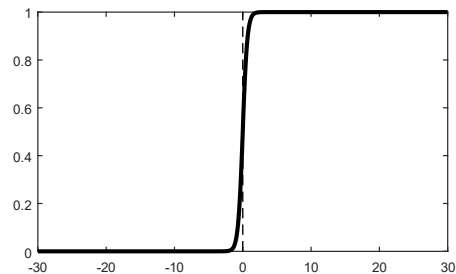
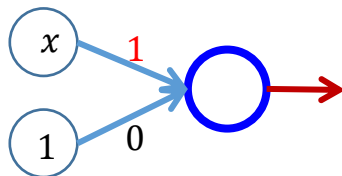
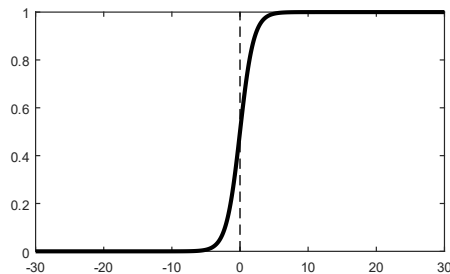
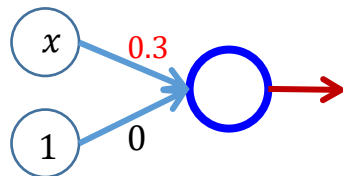
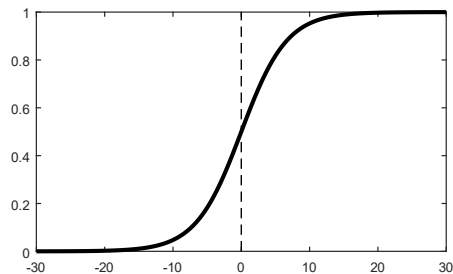
通用逼近定理

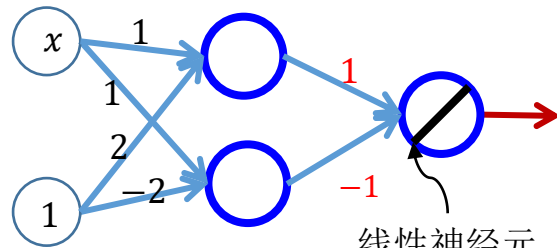
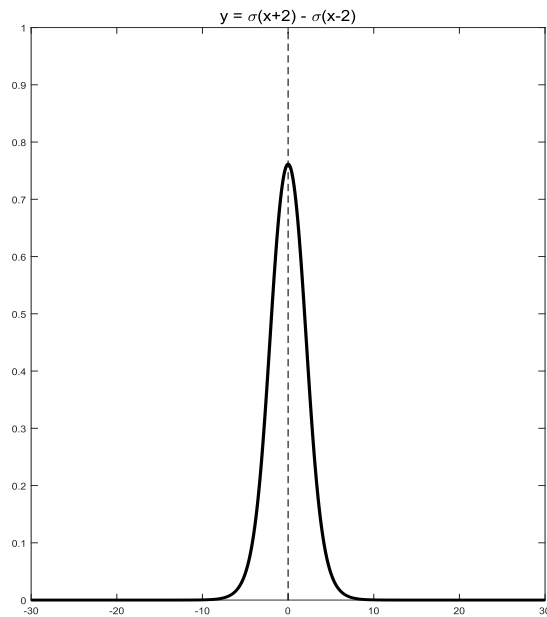
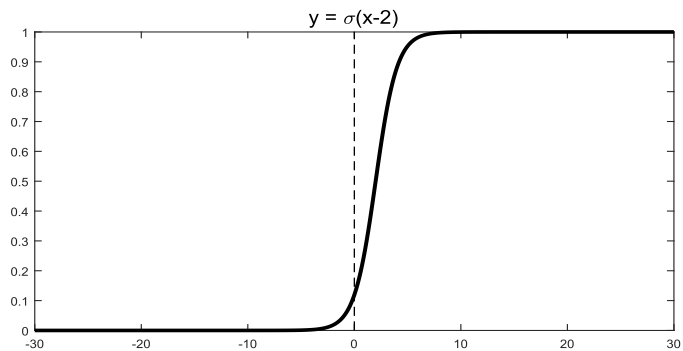
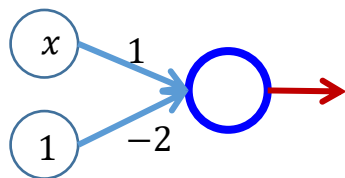
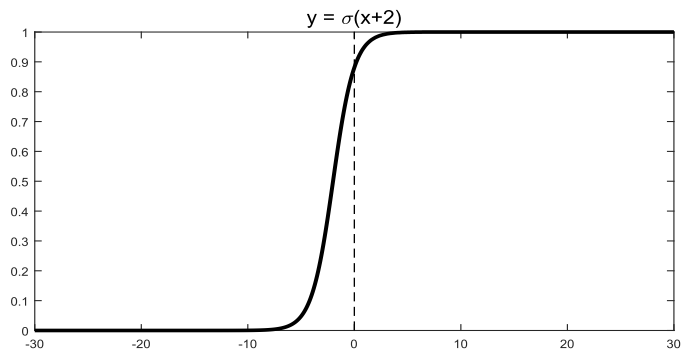
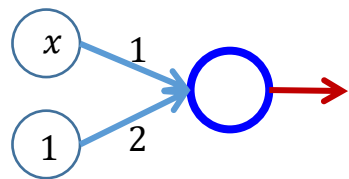
单隐层神经网络的表示能力



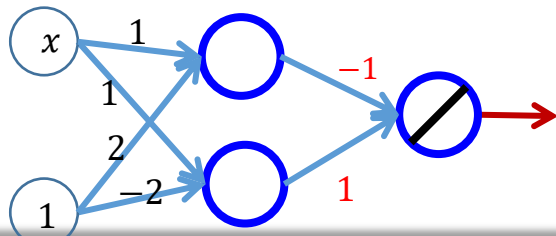
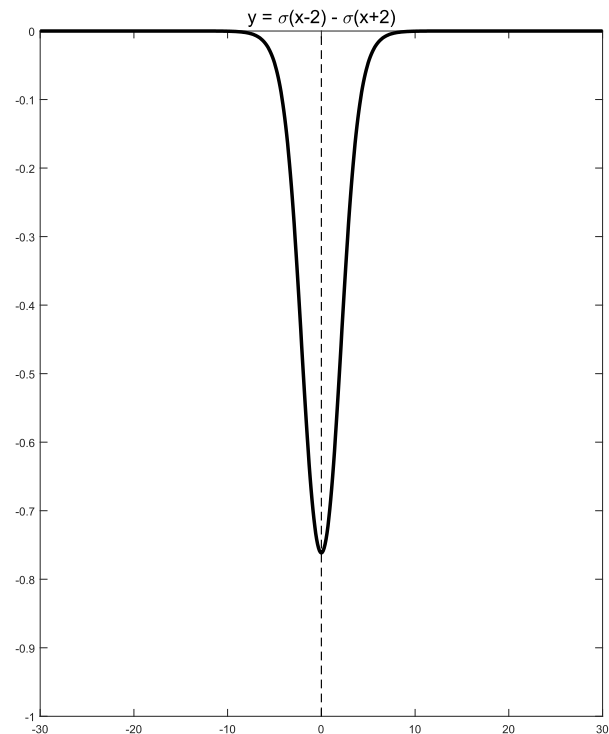
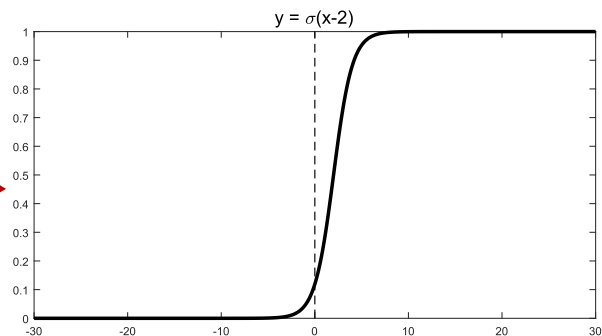
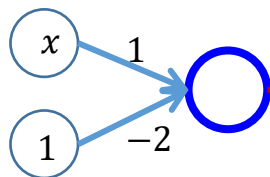
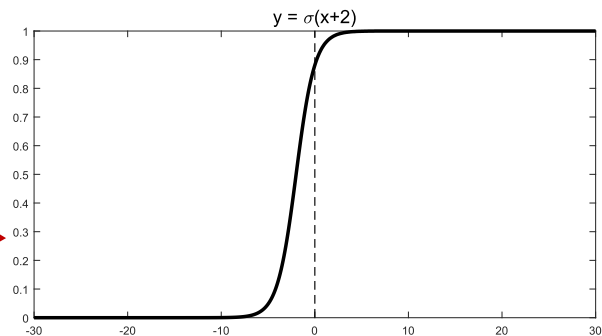
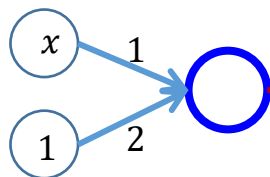
$$y = \frac{1}{1 + e^{-(wx+b)}}$$

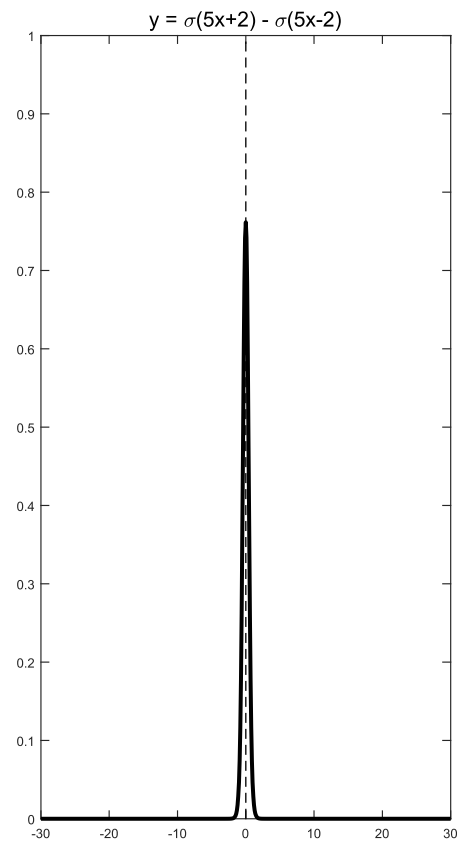
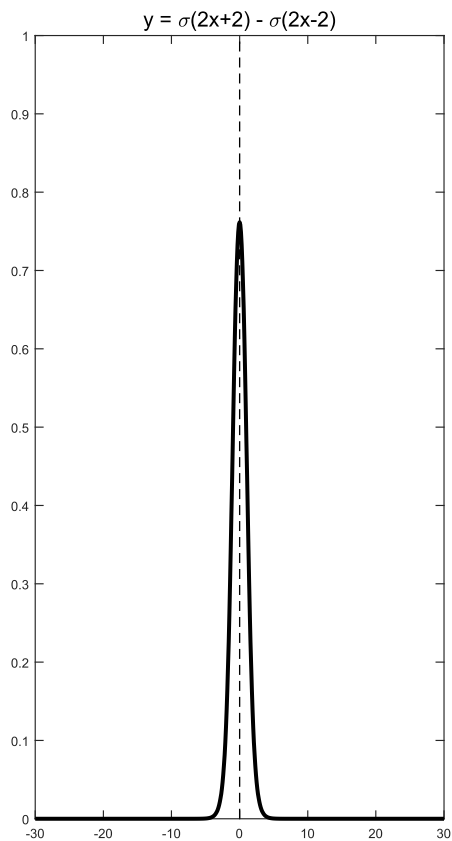
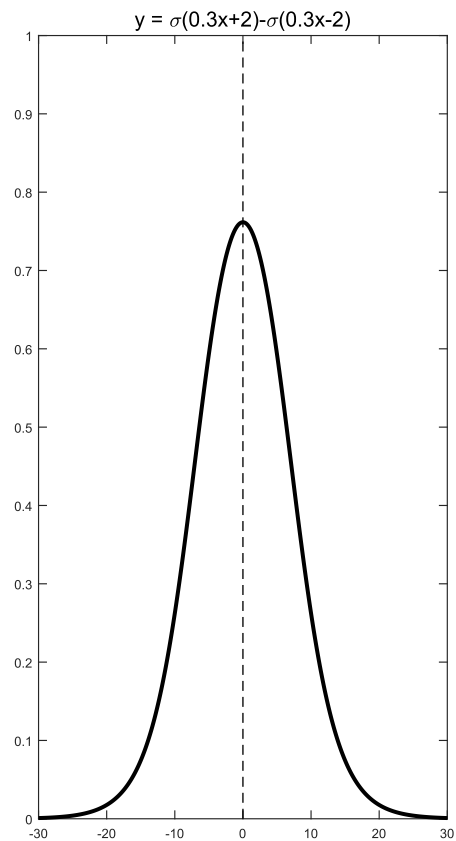


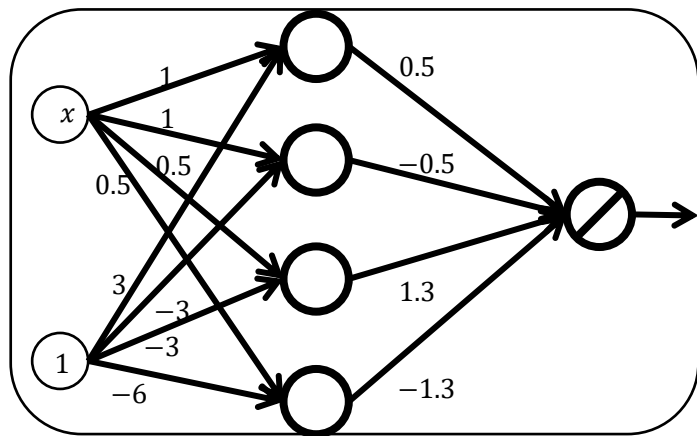




线性神经元
 $f(z) = z$

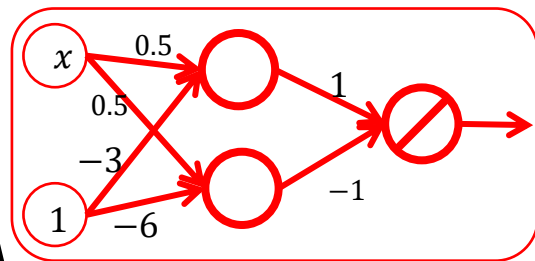




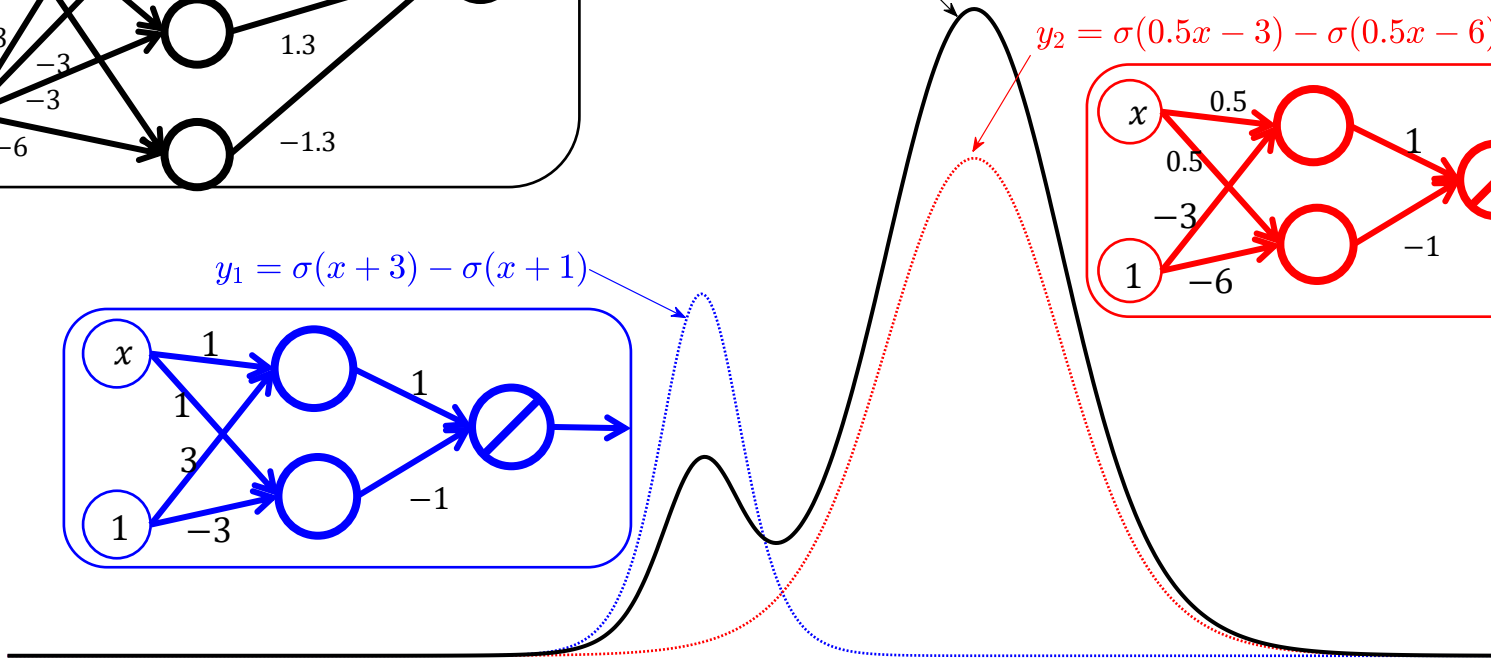
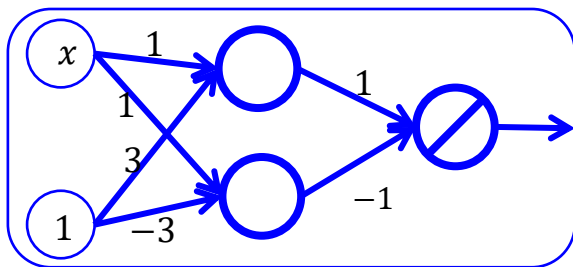


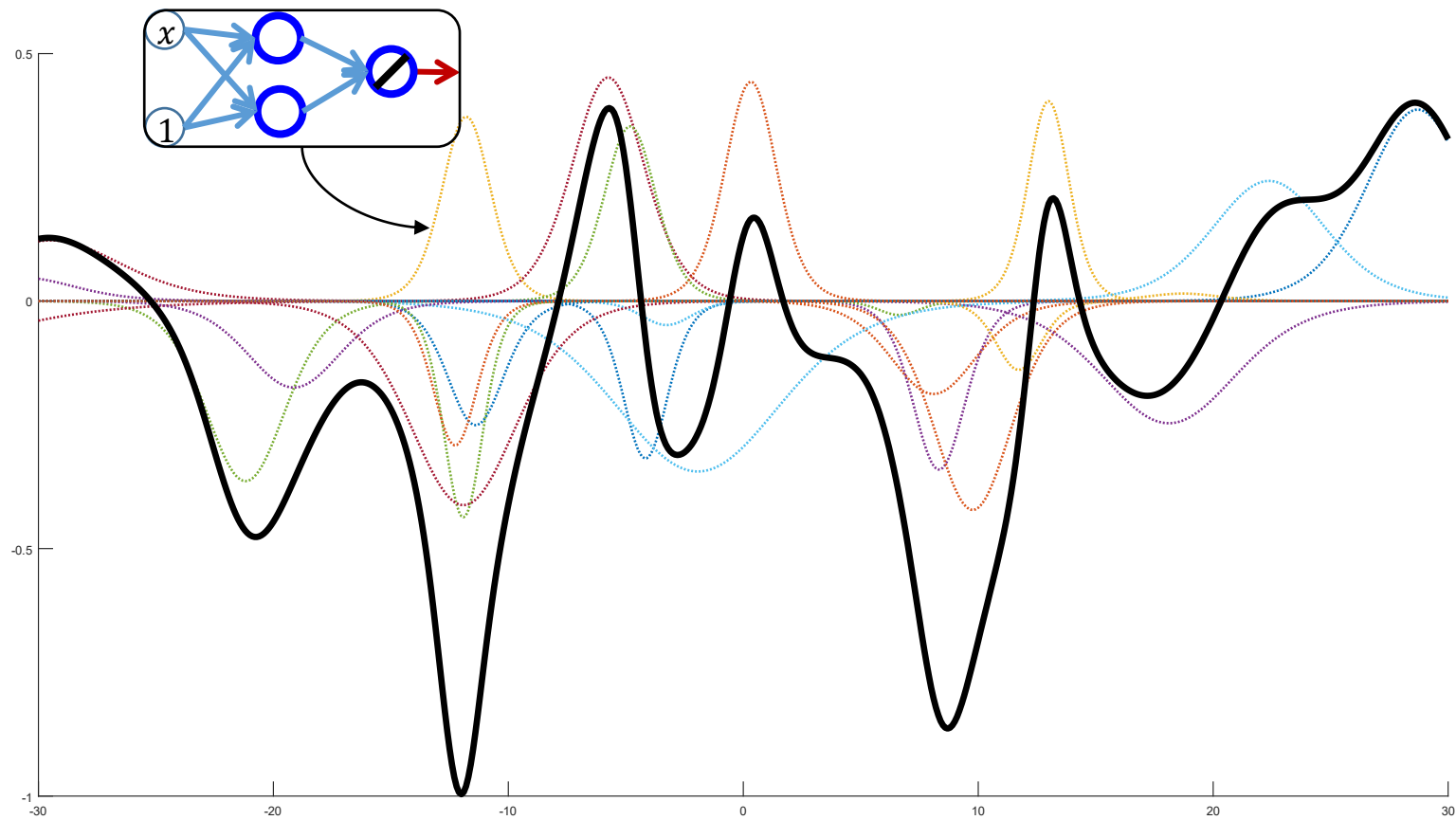
$$y = 0.5y_1 + 1.3y_2$$

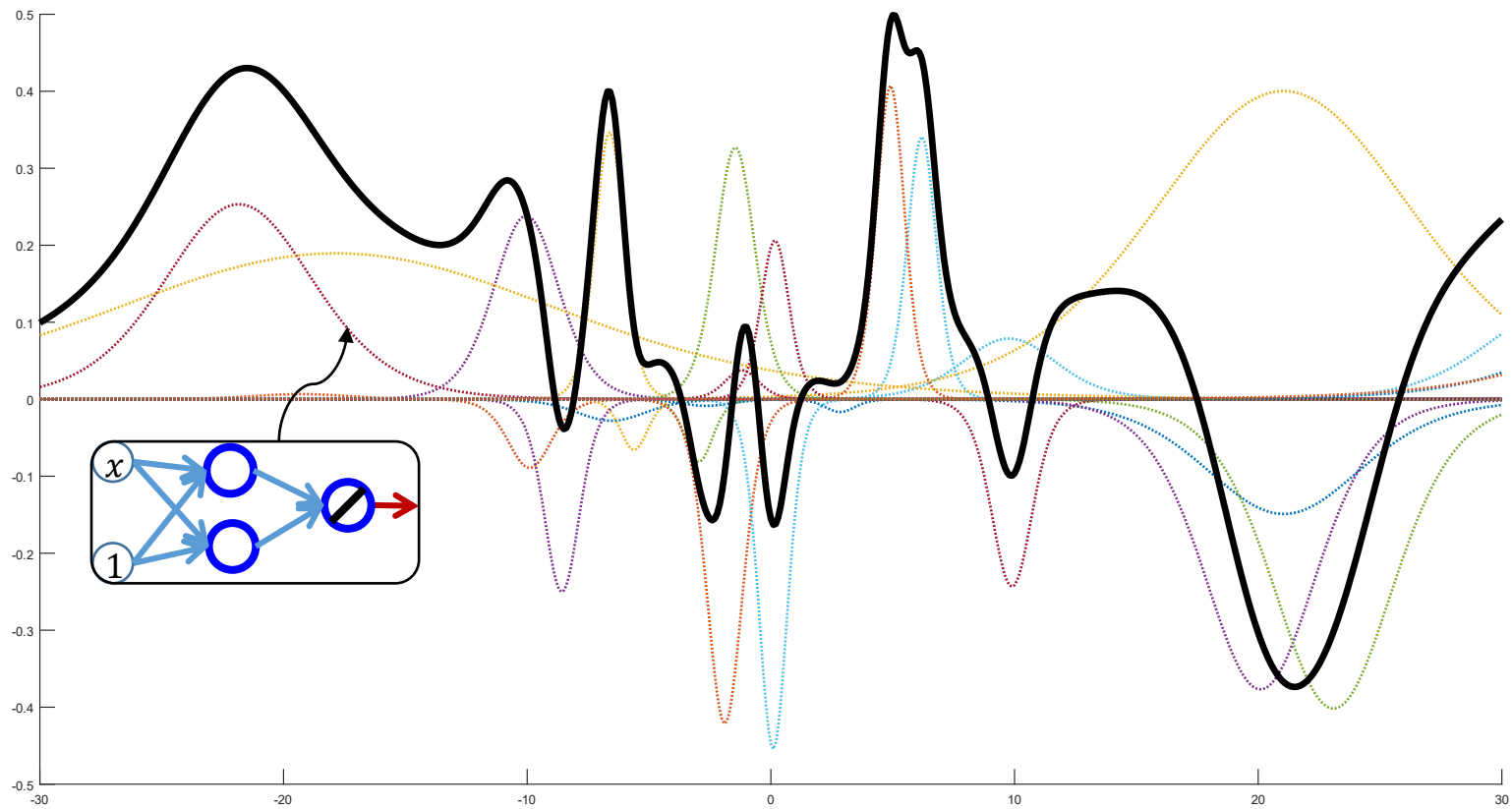
$$y_2 = \sigma(0.5x - 3) - \sigma(0.5x - 6)$$

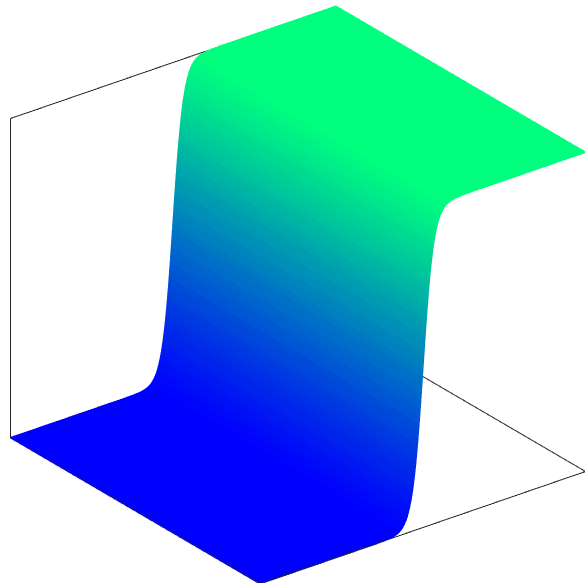
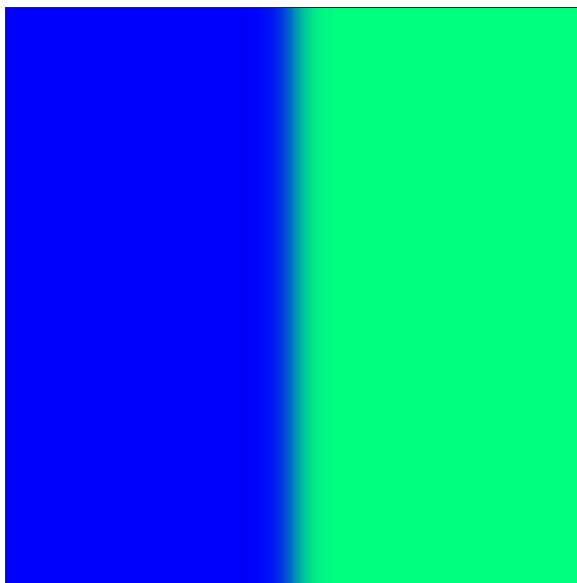
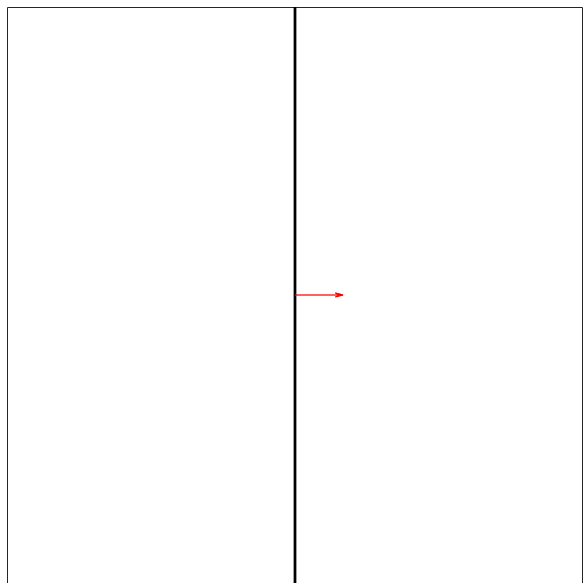


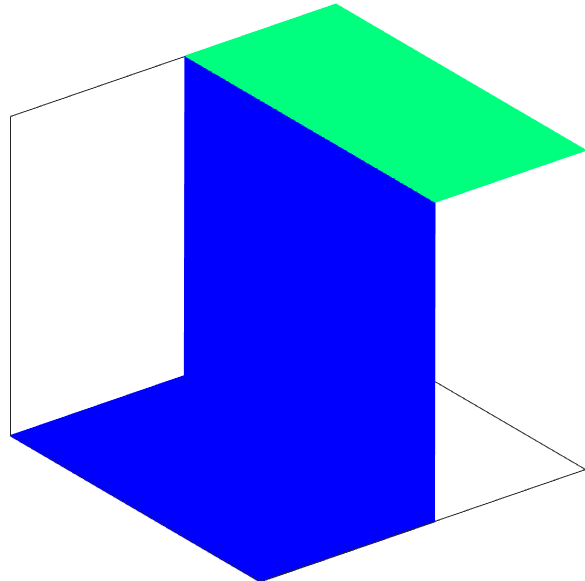
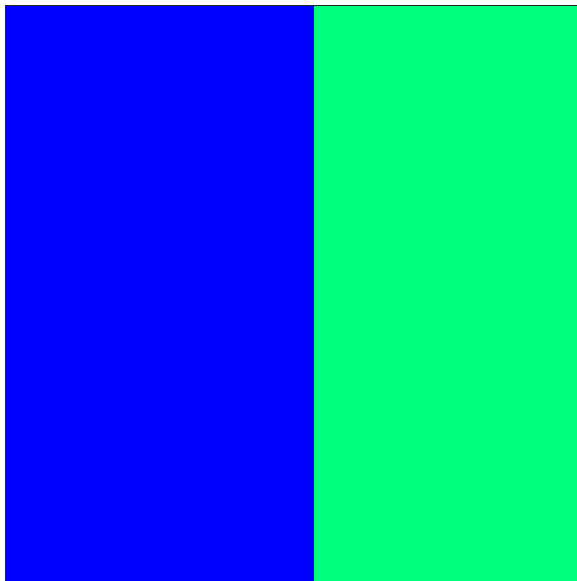
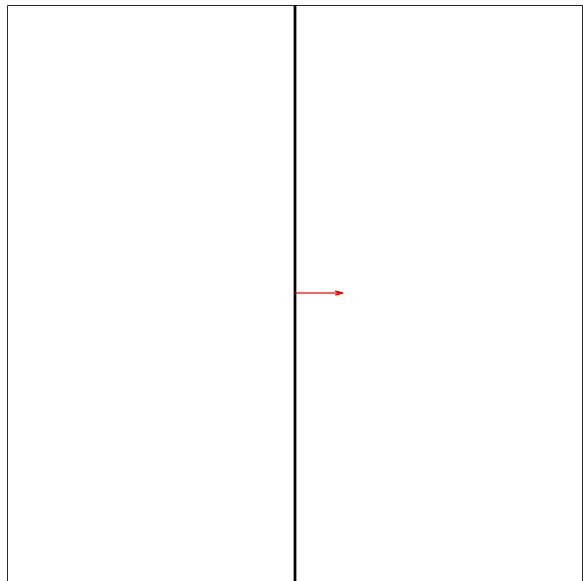
$$y_1 = \sigma(x + 3) - \sigma(x + 1)$$

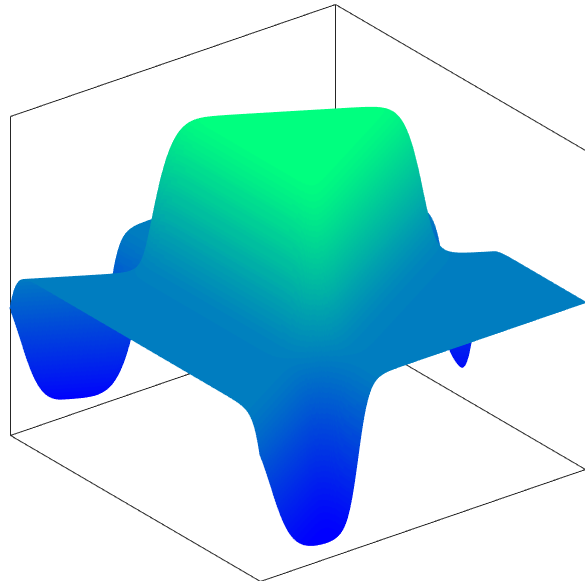
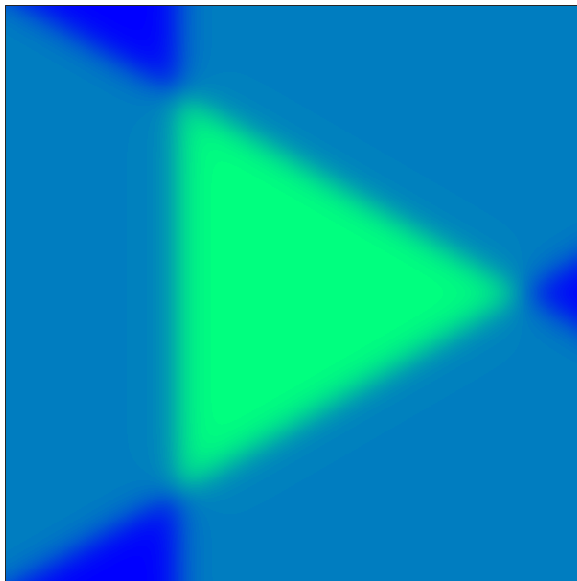
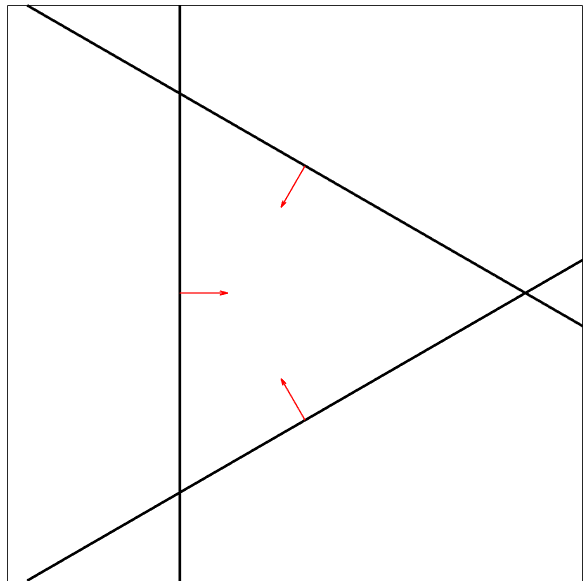


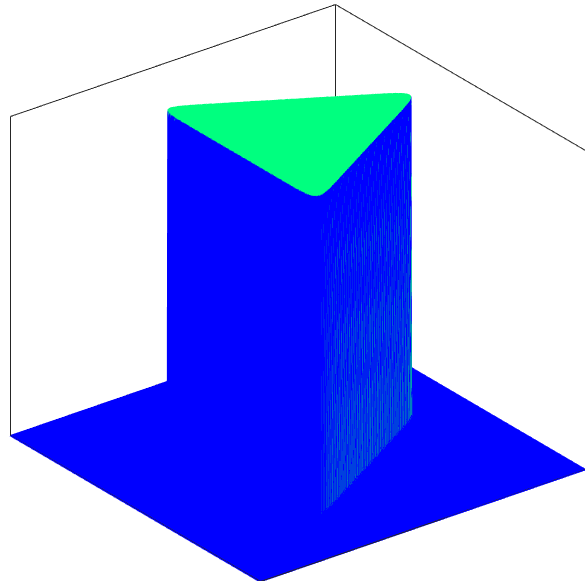
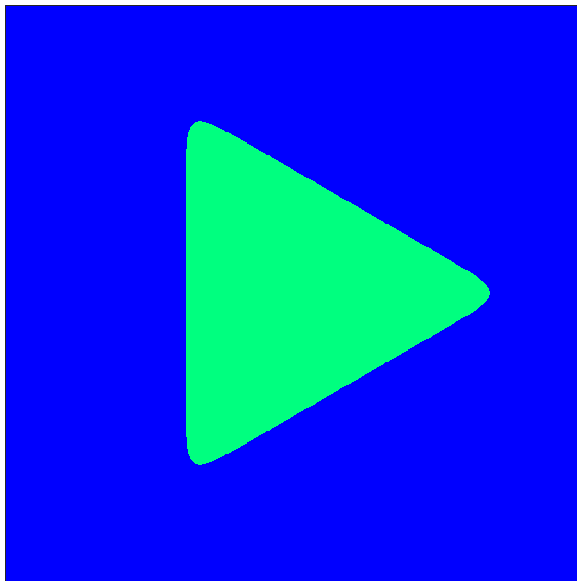
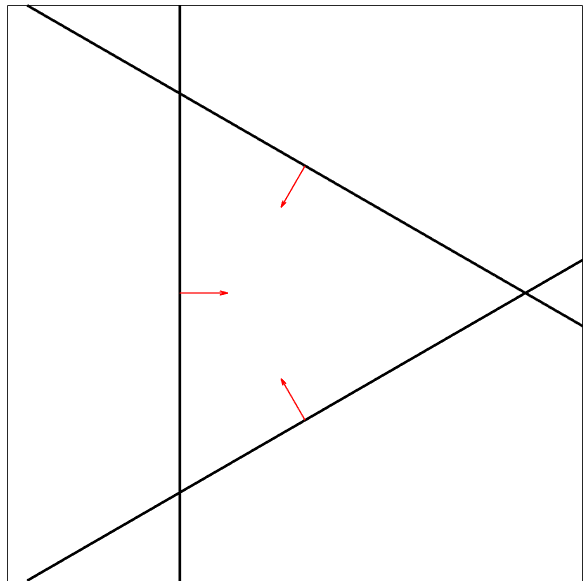


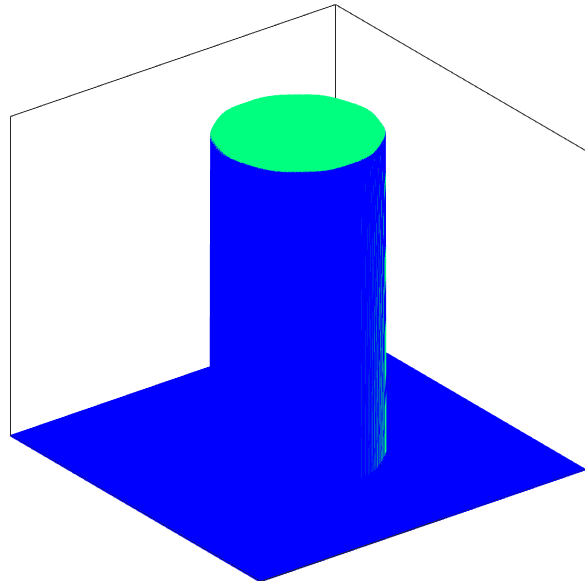
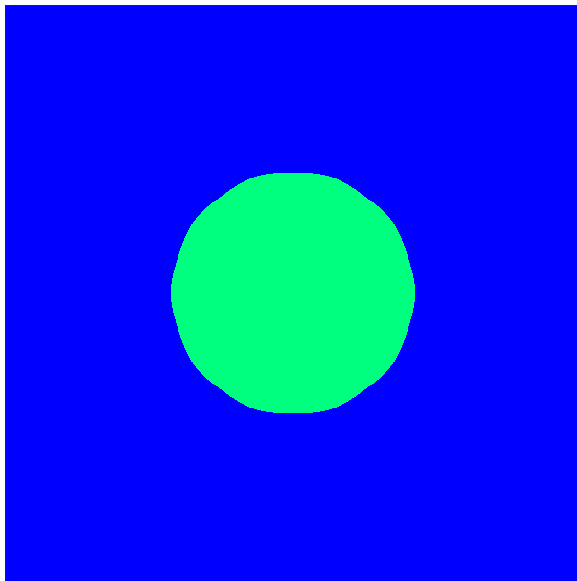
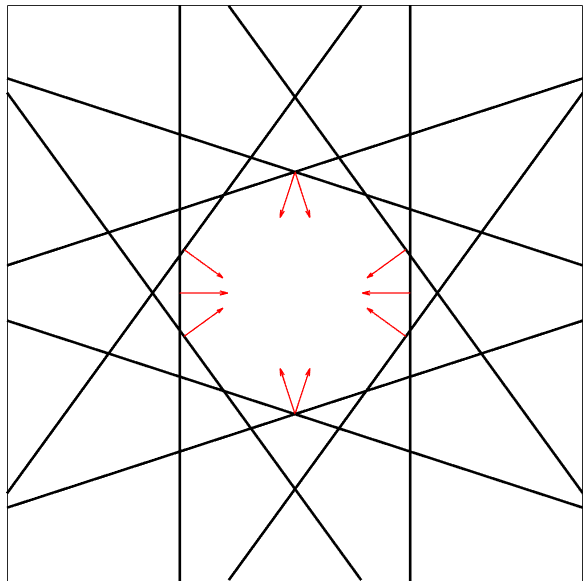


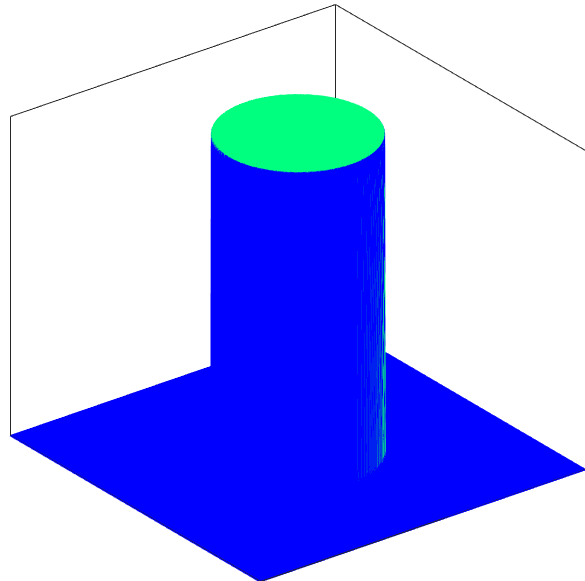
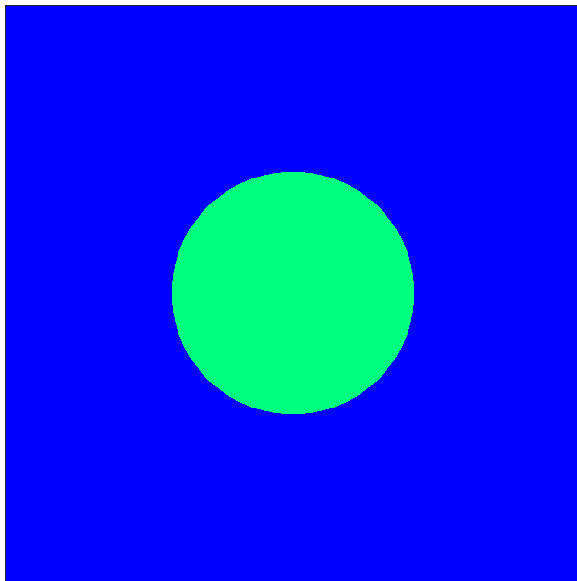
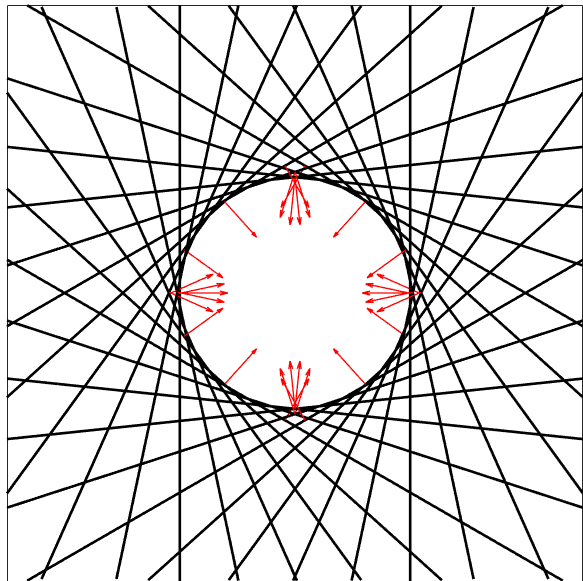


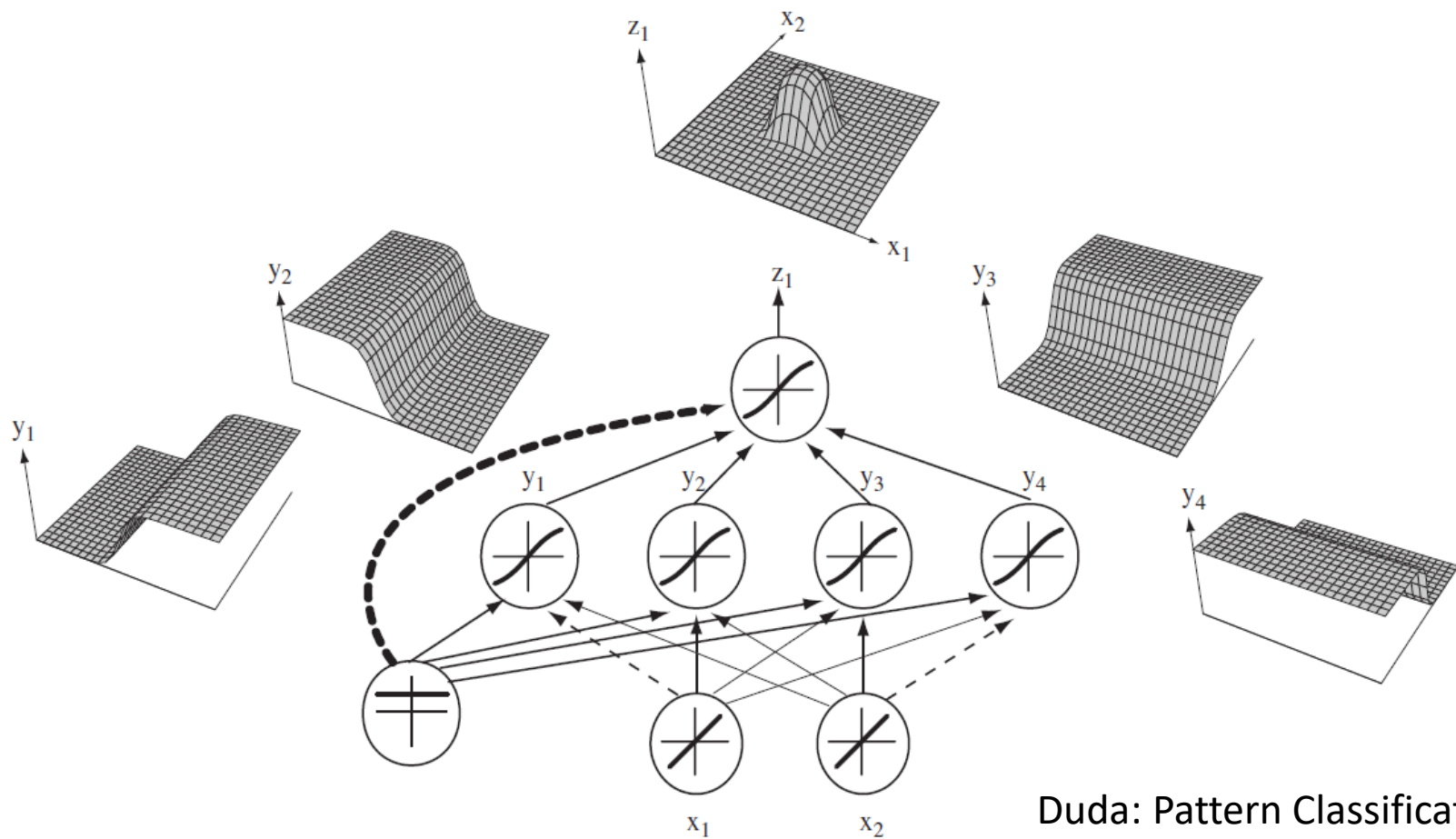




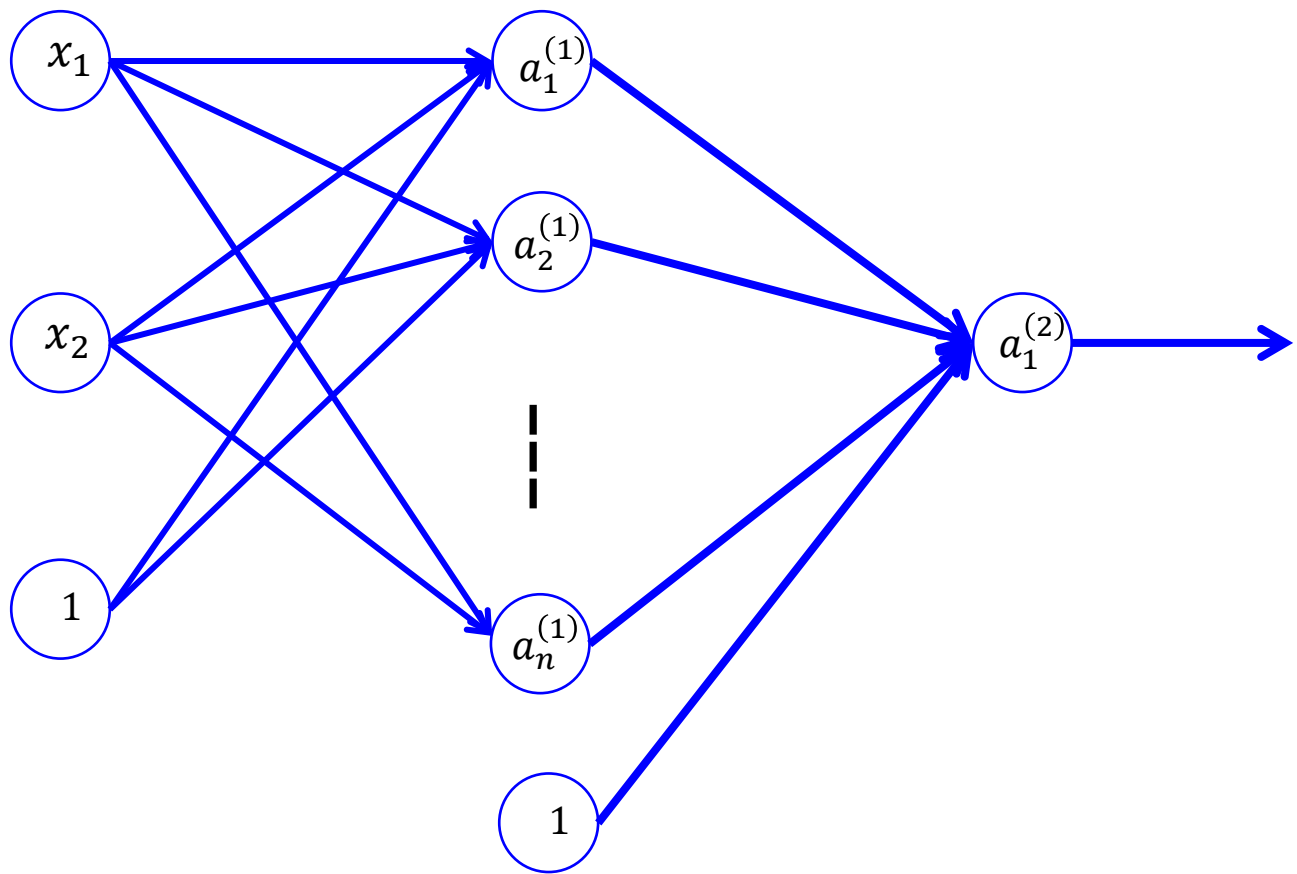


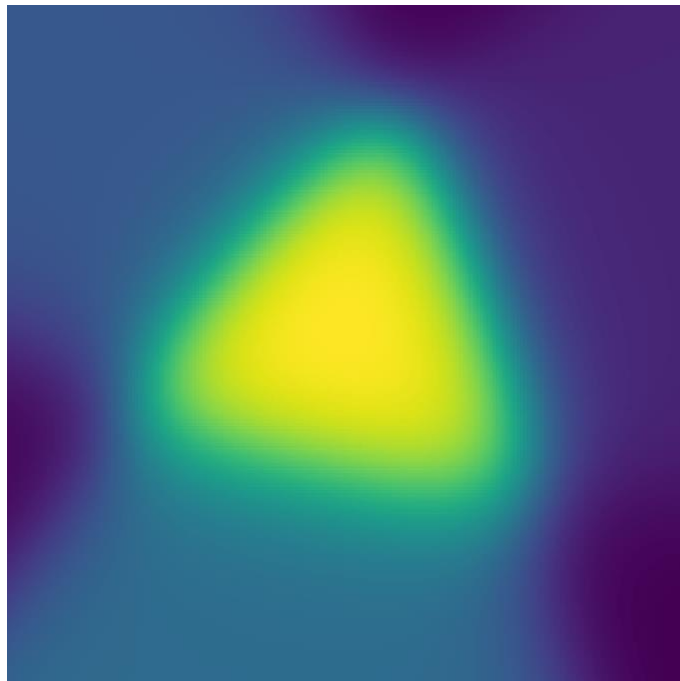
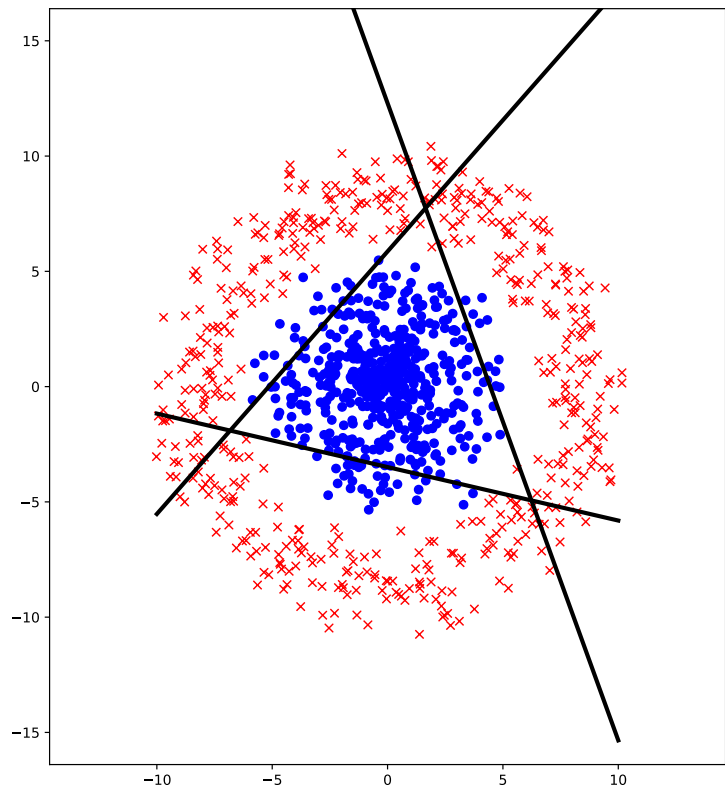




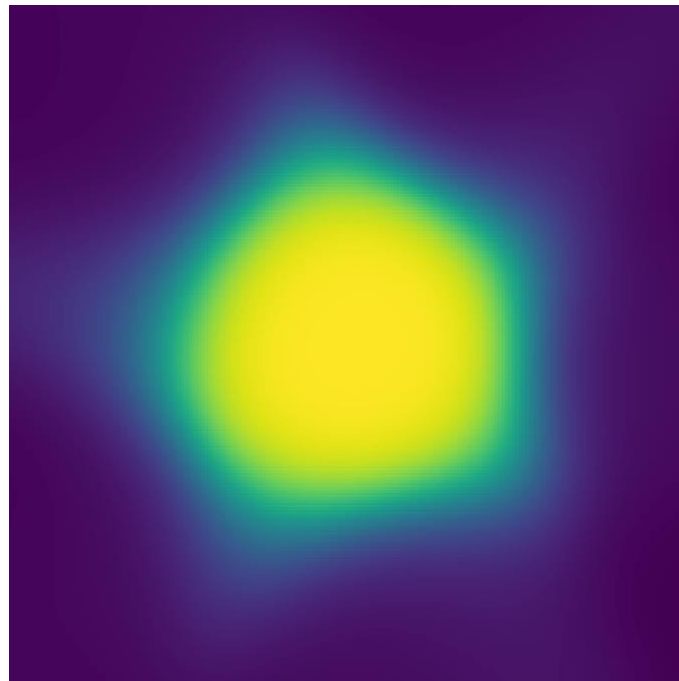
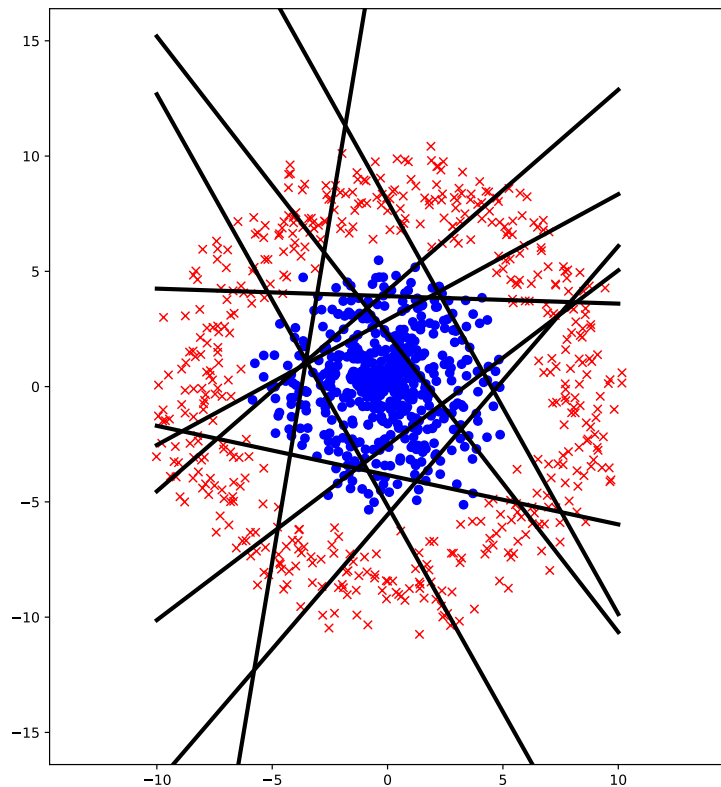


Duda: Pattern Classification

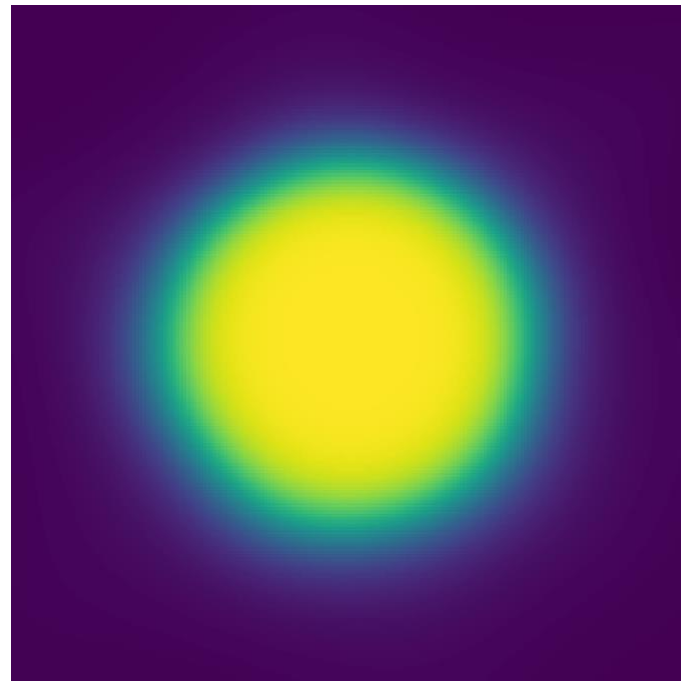
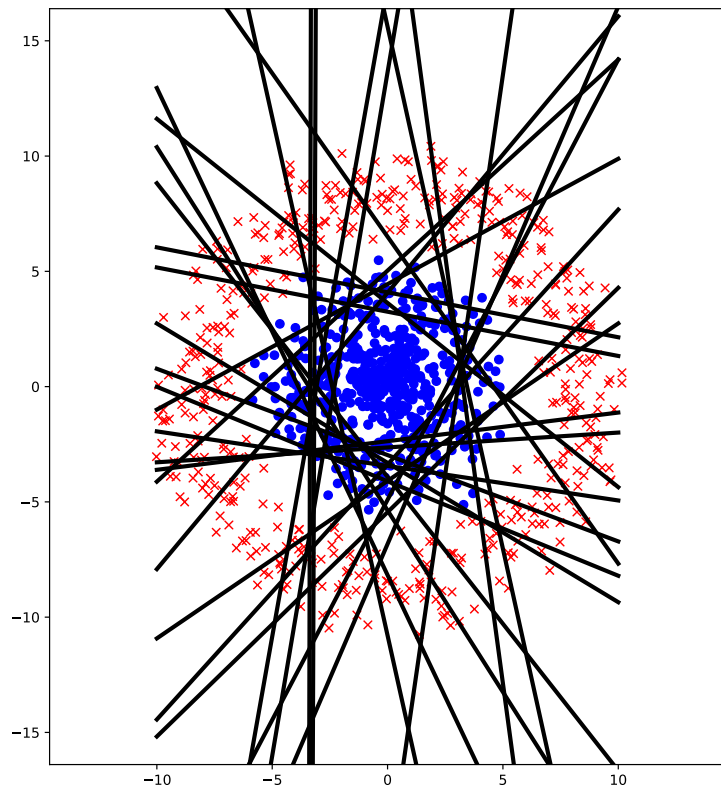




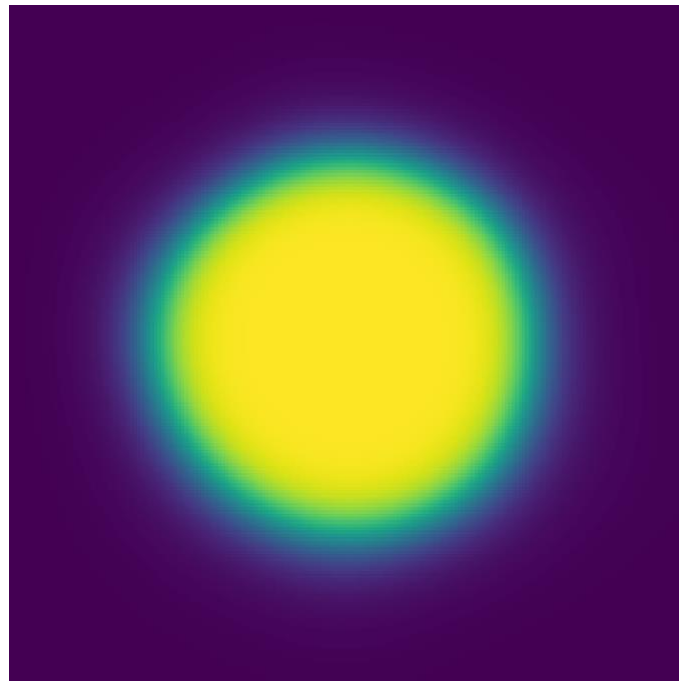
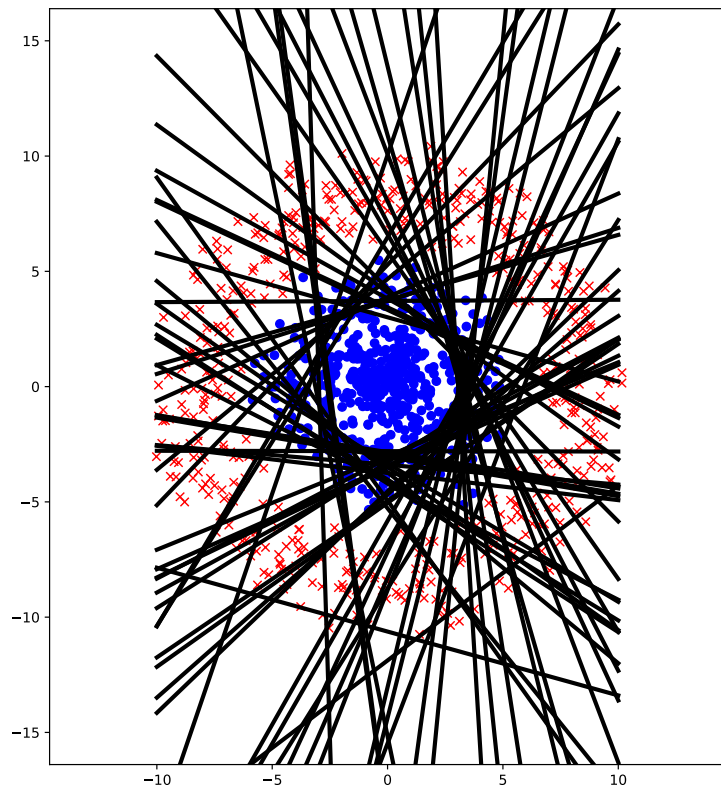
Hidden = 10

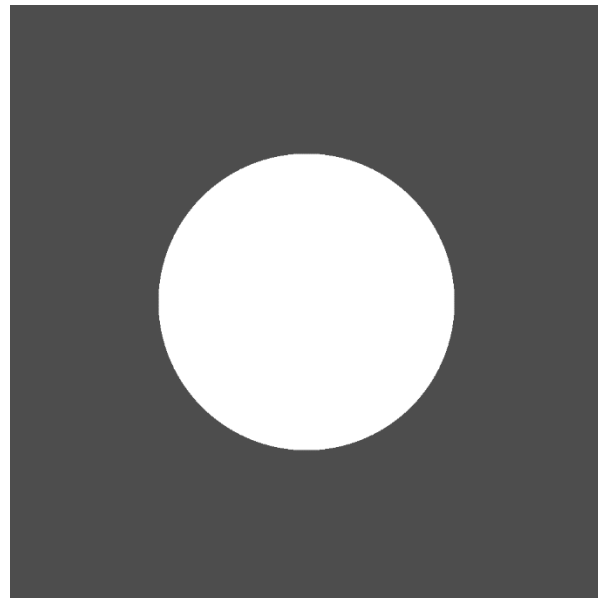
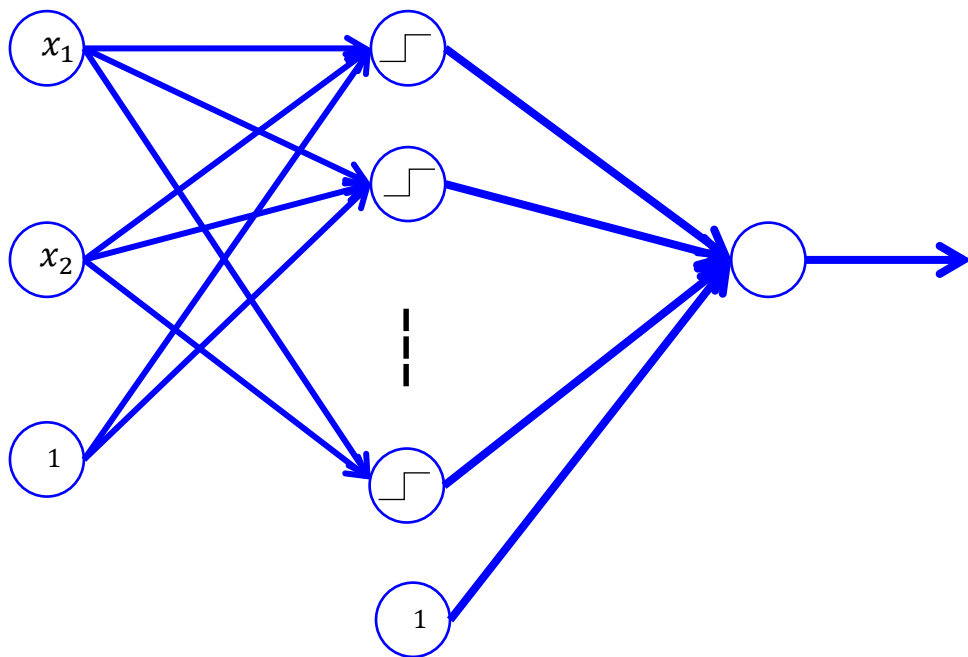


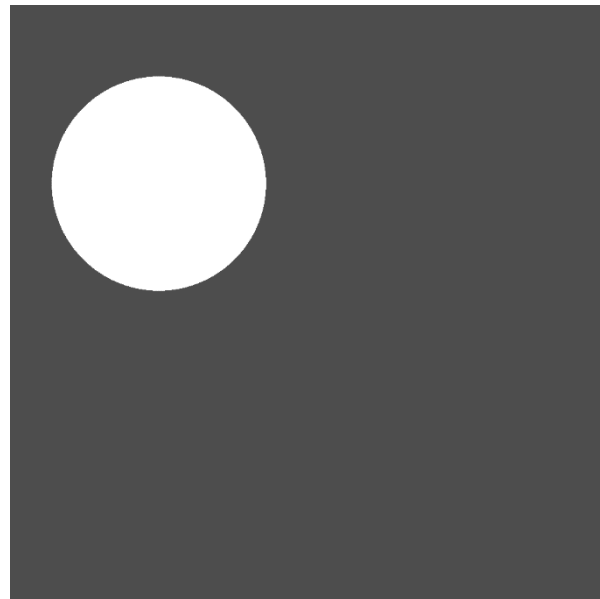
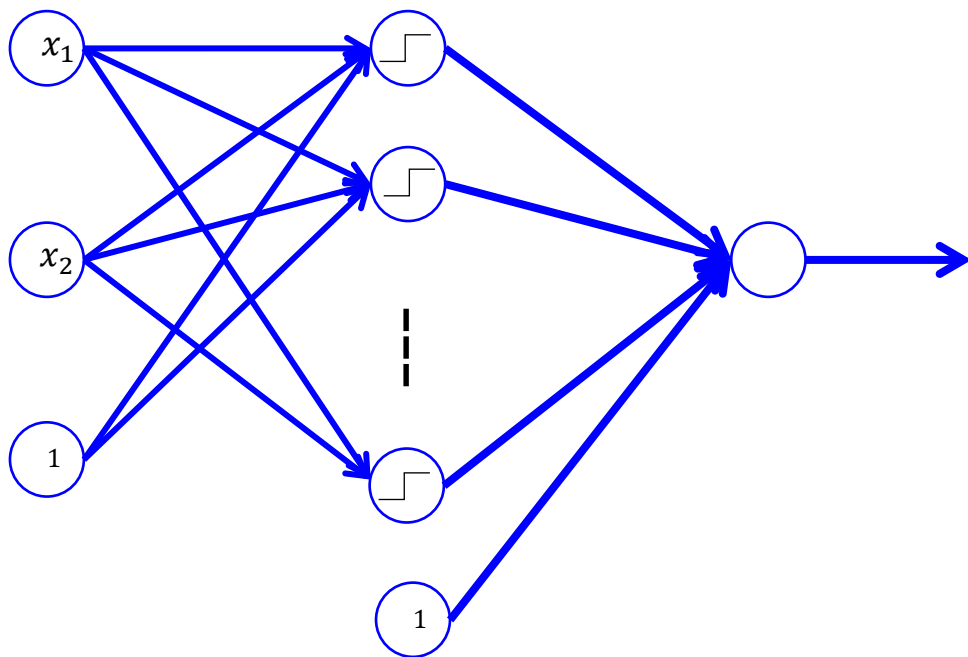
Hidden = 100

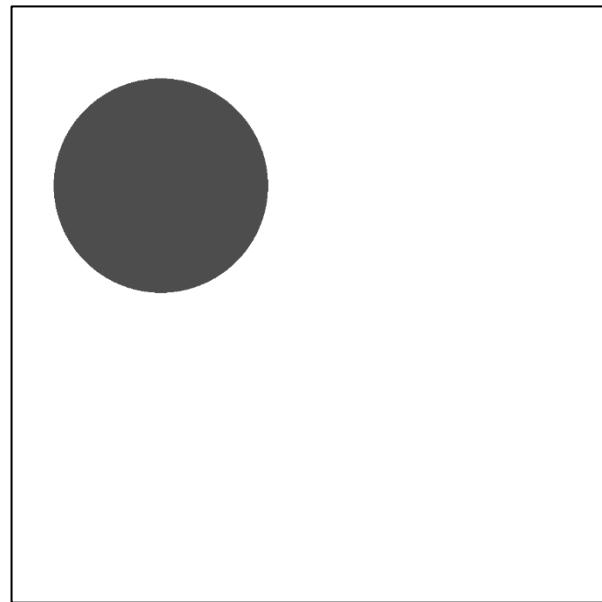
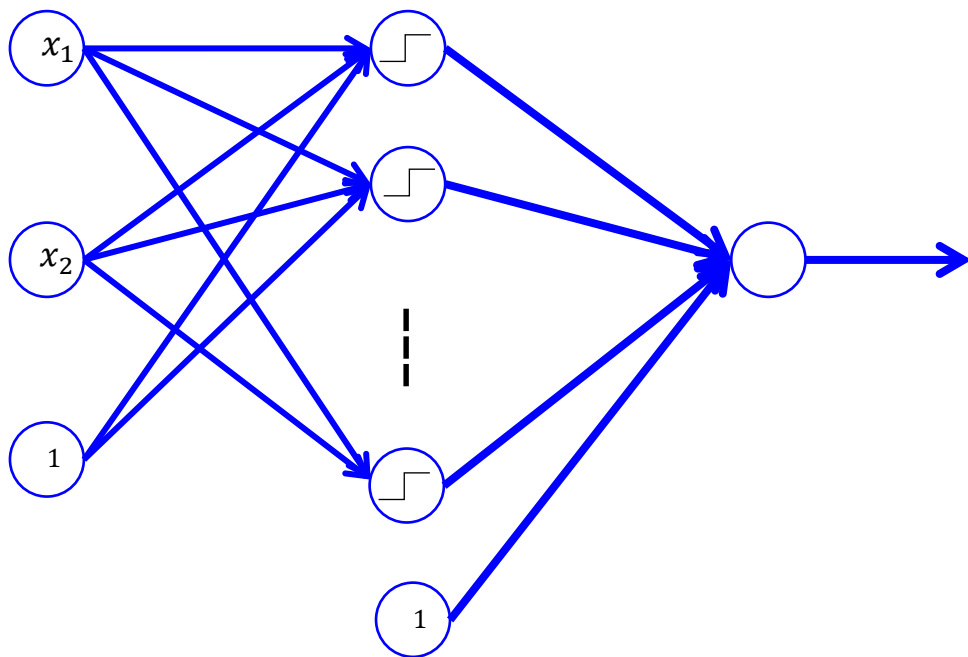


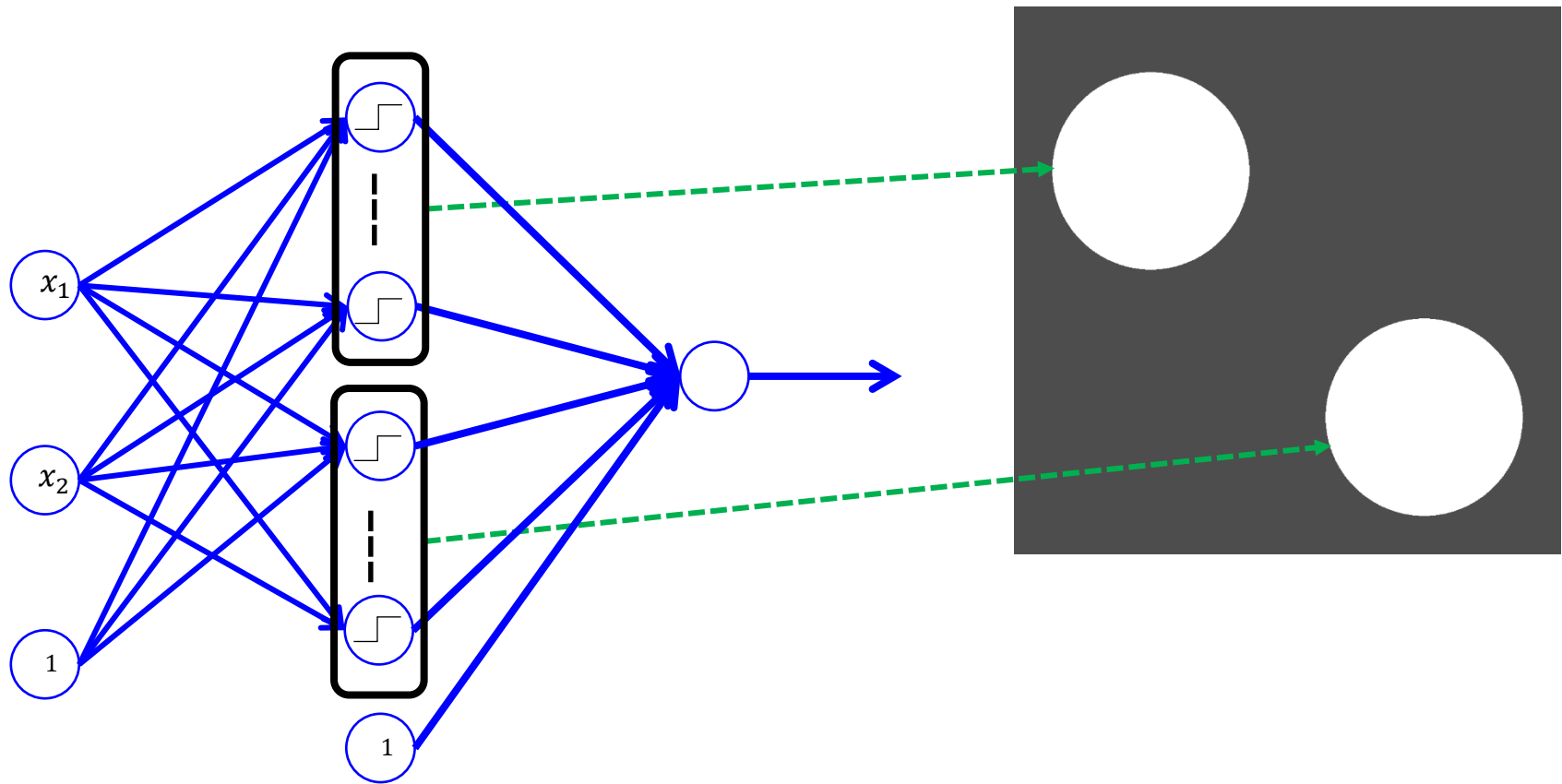
Hidden = 300

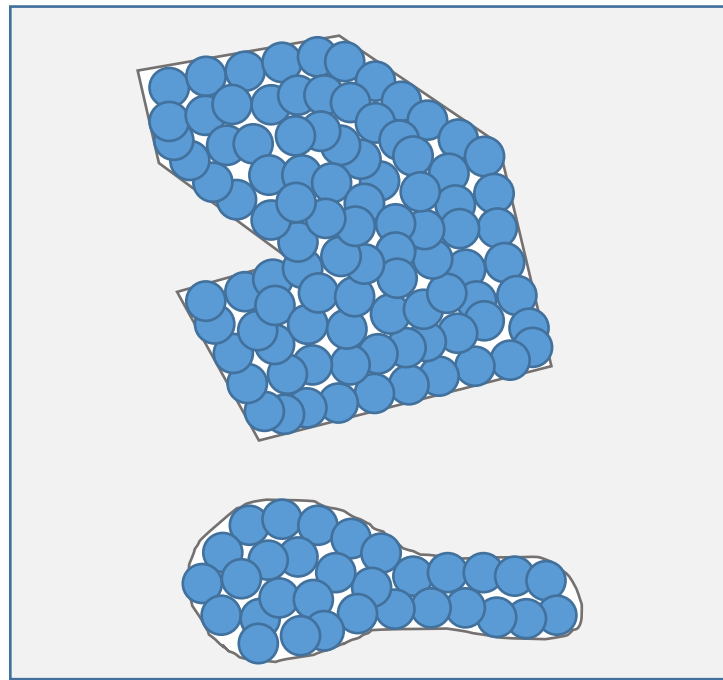
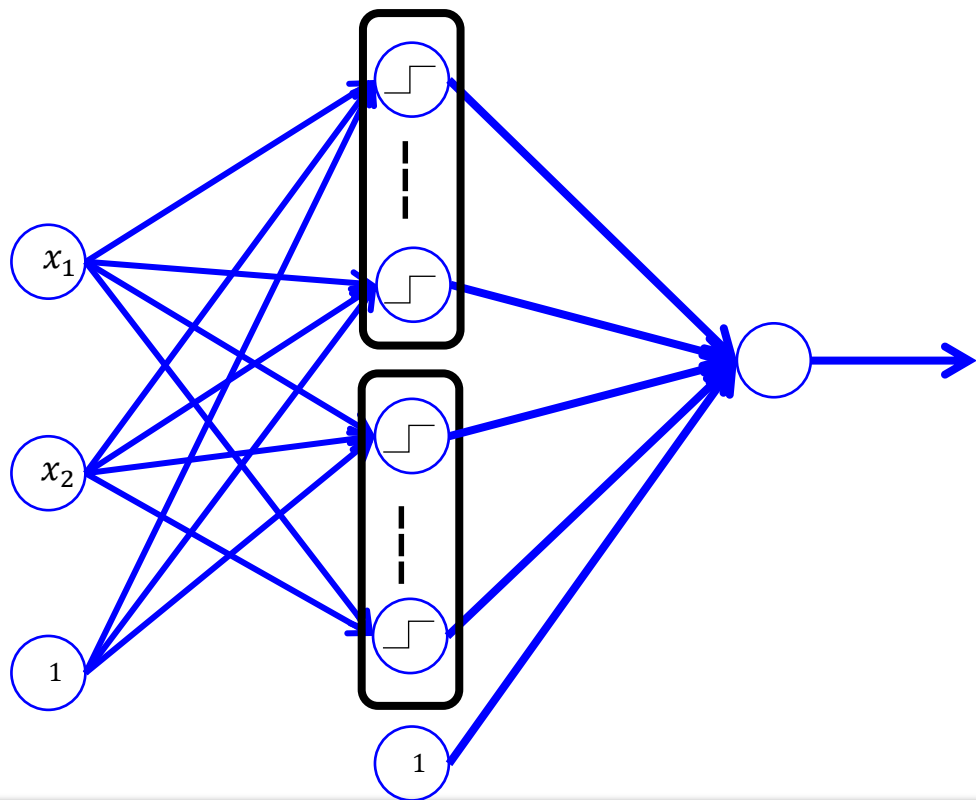


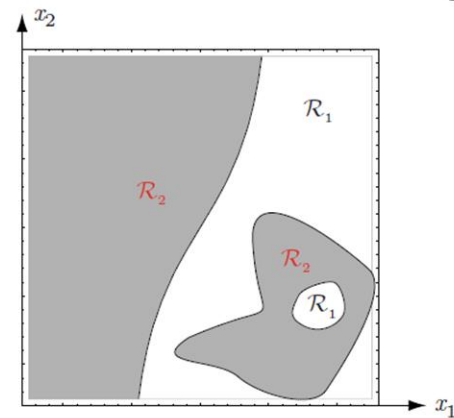
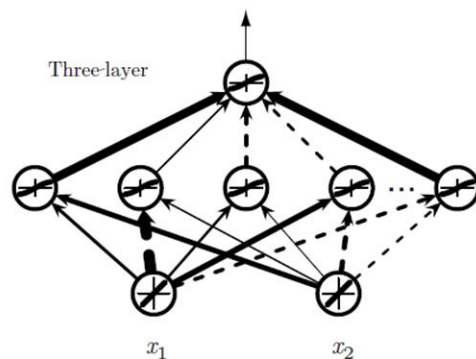
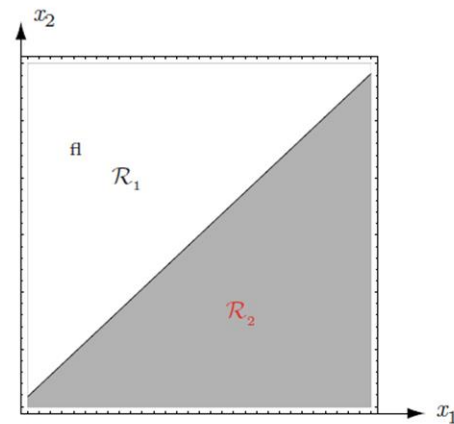
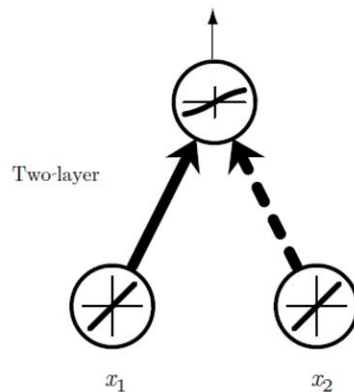












通用逼近定理(Universal Approximation Theorem)

- 单隐层神经网络（输出单元为线性响应、隐层单元为非线性（sigmoid, tanh, Relu...））
 - 可以逼近任意连续函数
 - 可以拟合任意决策边界
 - 可以表示任意布尔函数
- 困难：
 - 可能需要巨大数目（甚至无穷个）的神经元才能达到足够的逼近精度

总结

- 多层感知机(Multi-Layer Perceptron, MLP)
 - 多层神经元叠加构成的一个神经网络
 - 通过复合(Composition)的方式表示一个函数
- 多层感知机的学习算法
 - 梯度下降法
 - 误差后向传播算法(error Back Propagation, BP算法)
 - 本质是复合函数求导的链式法则
- 通用逼近定理(Universal Approximation Theorem)
 - 单隐层神经网络可以表示任意连续函数
 - 不确定是否能真正学习出任意连续函数
 - 可能需要巨大数目的隐层神经元才能达到合适的表示精度