

5.1 人工神经网络简介

AlphaGo中的机器学习

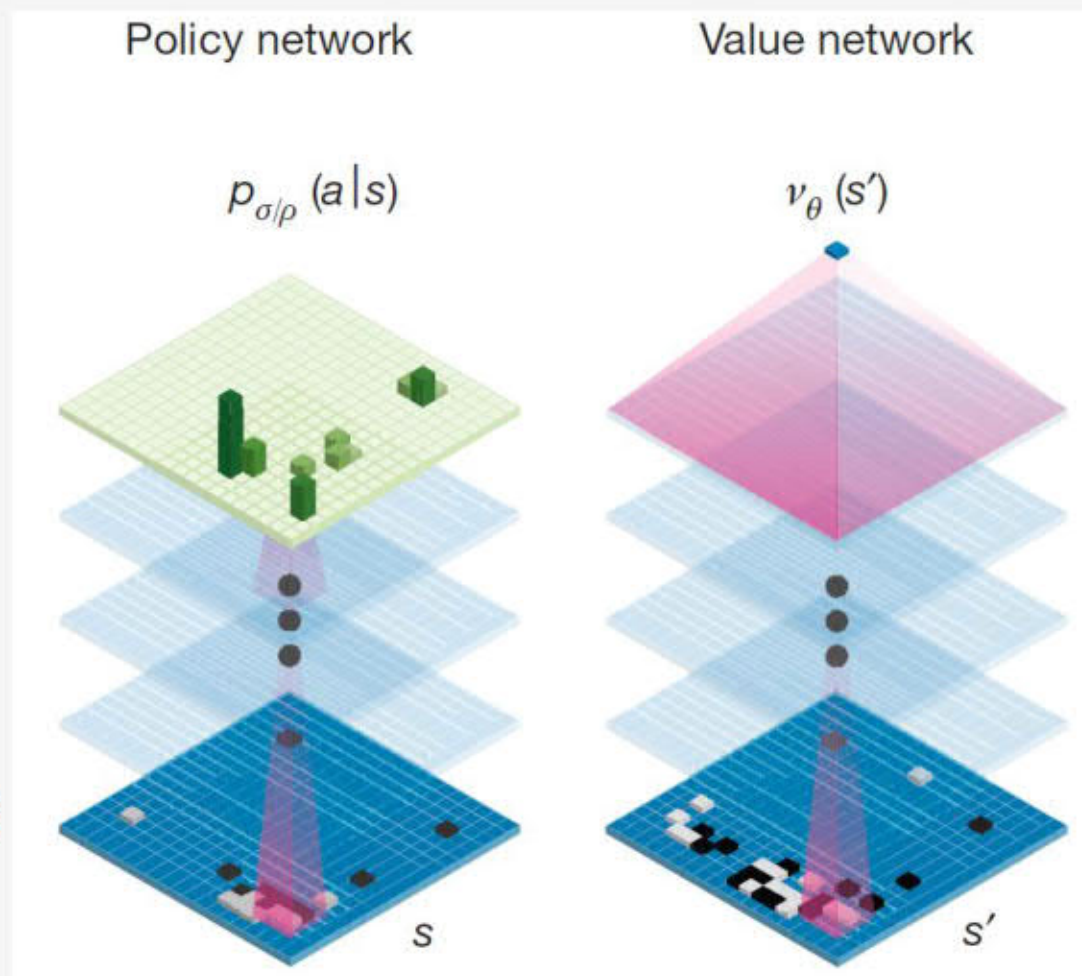
□ 策略网络

- 监督学习
 - 预测下一步人类移动的最佳结果
- 强化学习
 - 学习去选择下一步移动去最大化获胜率

□ 价值网络

- 在给定当前状态的情况下获胜的期望

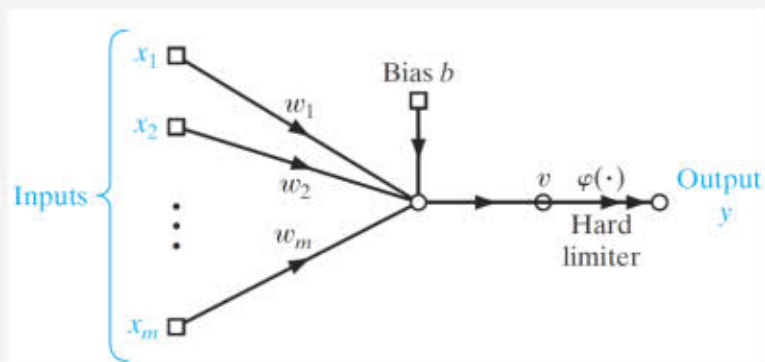
□ 通过（深度）神经网络实现



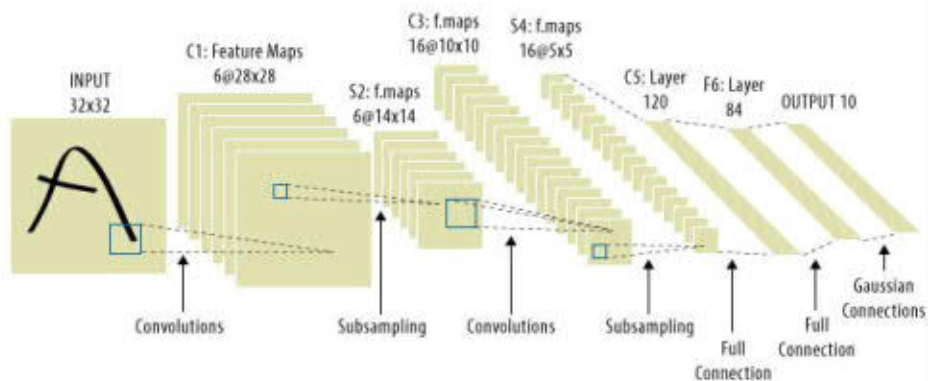
5.1 人工神经网络简介

人工神经网络

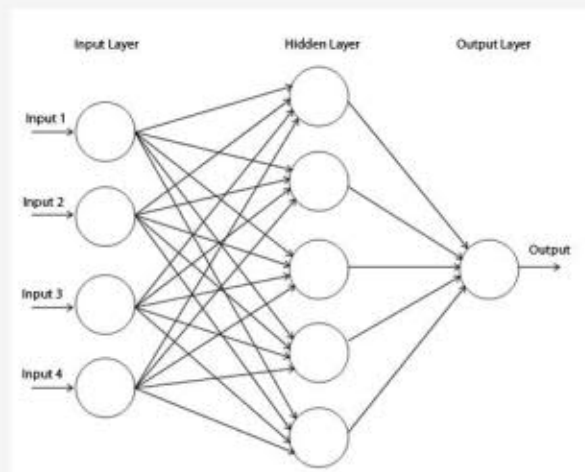
神经网络是深度学习的基础



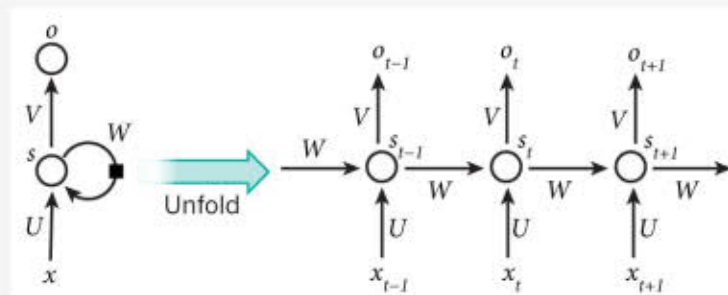
感知机



卷积神经网络



多层感知机



循环神经网络

5.1 人工神经网络简介

真正的神经元学习

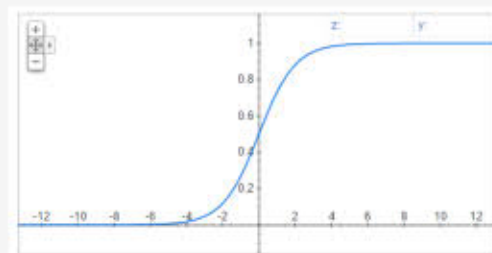
- ▣ 突触通过经验改变尺寸和连接强度
- ▣ **Hebbian学习**: 如果两个连接在一起的神经元同时被激发, 那么他们之间的连接会变强
- ▣ “Neurons that fire together, wire together.”
- ▣ 这些都推动了人工神经网络的研究

5.2 激活函数与损失函数

非线性激活函数

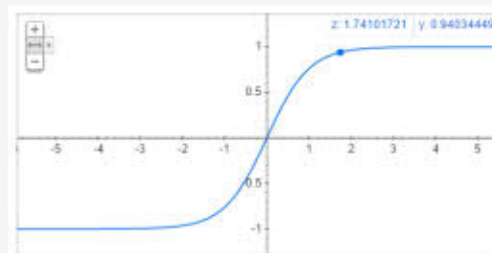
□ Sigmoid激活函数

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



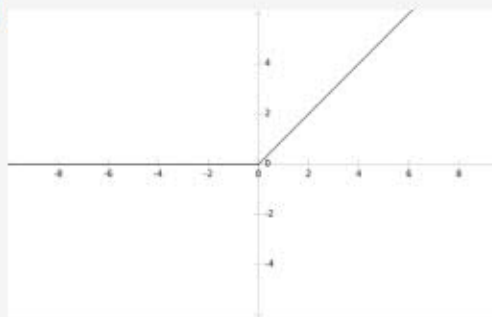
□ Tanh激活函数

$$\tanh(z) = \frac{1 - e^{-2z}}{1 + e^{-2z}}$$



□ 线性整流函数 (Rectified Linear Unit, ReLU)

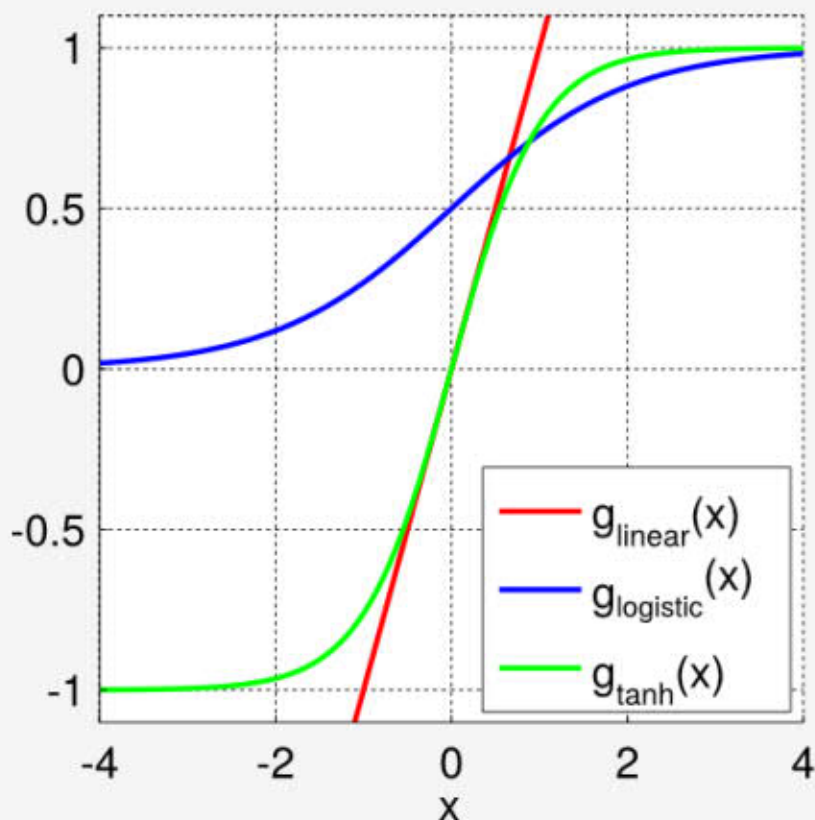
$$\text{ReLU}(z) = \max(0, z)$$



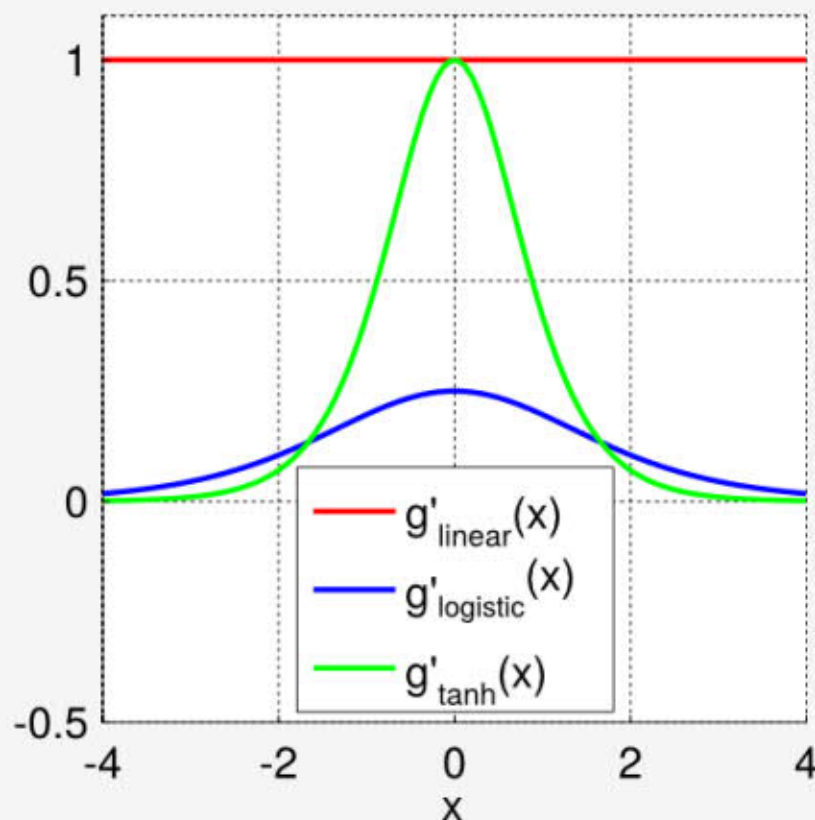
5.2 激活函数与损失函数

激活函数

Some Common Activation Functions



Activation Function Derivatives

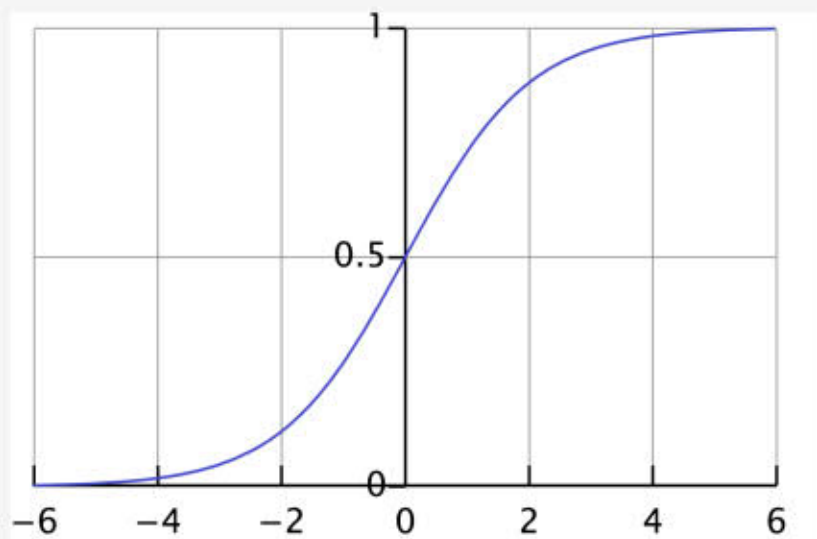


5.2 激活函数与损失函数

激活函数

逻辑斯谛Sigmoid函数

$$f_{\text{Sigmoid}}(x) = \frac{1}{1 + e^{-x}}$$



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

□ 导数:

$$f'_{\text{Sigmoid}}(x) = f_{\text{Sigmoid}}(x)(1 - f_{\text{Sigmoid}}(x))$$

□ 输出范围 [0,1]

□ 受生物神经元的启发产生，可以被看做一个人工神经元在当前输入下被“激活”的概率

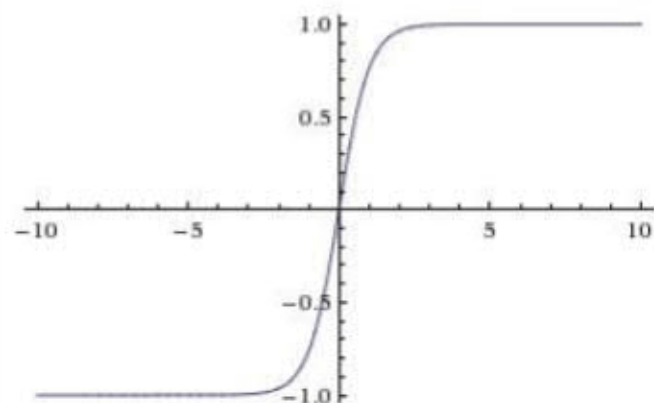
□ 边界值会使梯度消失 (为什么?)

5.2 激活函数与损失函数

激活函数

Tanh函数

$$f_{\tanh}(x) = \frac{\sinh(x)}{\cosh(x)} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



□ 导数:

$$f'_{\tanh}(x) = 1 - f_{\tanh}(x)^2$$

- 负值较大的输入经过tanh函数会映射为负输出
- 只有接近零的输入会被映射为接近零的输出
- 使网络训练时被“卡住”的可能性降低

5.2 激活函数与损失函数

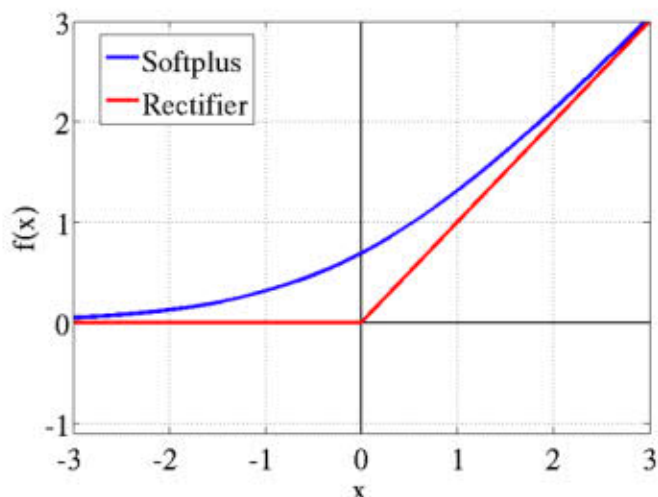
激活函数

ReLU函数(rectified linear unit)

$$f_{ReLU}(x) = \max(0, x)$$

ReLU可以被近似为softplus函数

$$f_{softplus}(x) = \log(1 + e^x)$$



其他激活函数:

Leaky ReLU, Exponential LU, Maxout 等

□ 导数:

$$f'_{ReLU}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

□ 另一个版本Noise ReLU:

$$f_{NoisyReLU}(x) = \max(0, x + N(0, \delta(x)))$$

□ ReLU可以被近似为softplus函数

$$f_{softplus}(x) = \log(1 + e^x)$$

□ x 增加时ReLU的梯度不会消失

□ 可以用来对正值输入进行建模

□ 由于无需计算指数函数, 所以它的计算速度很快

□ 使用它可以不再需要“预训练”过程

5.2 激活函数与损失函数

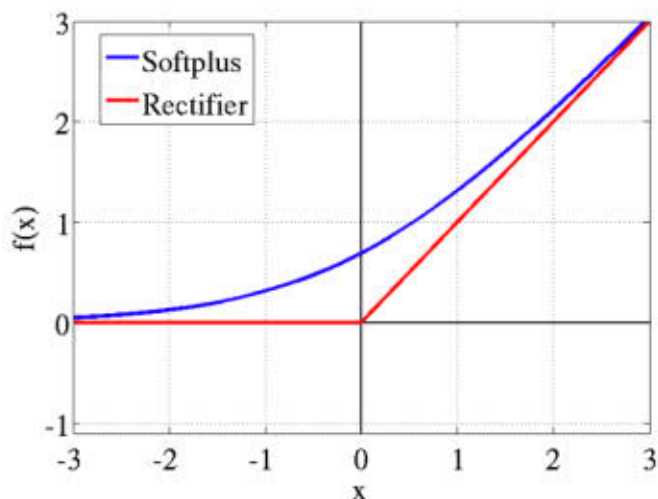
激活函数

ReLU函数(rectified linear unit)

$$f_{ReLU}(x) = \max(0, x)$$

ReLU可以被近似为softplus函数

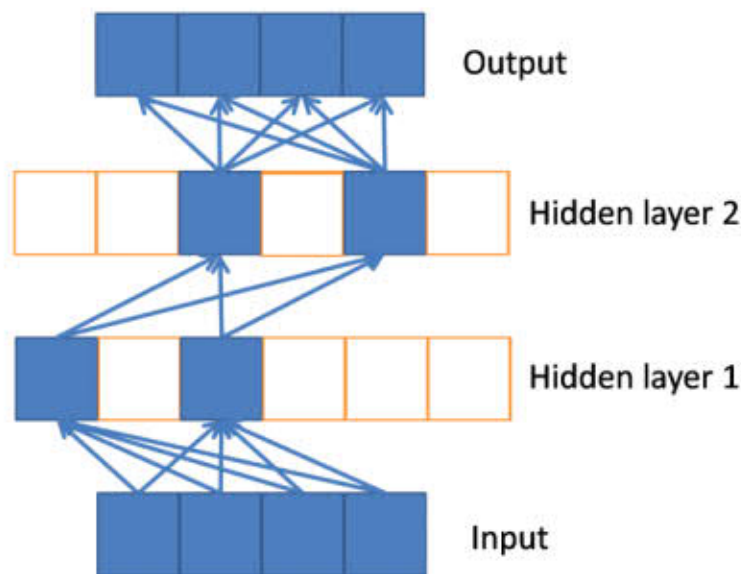
$$f_{softplus}(x) = \log(1 + e^x)$$



其他激活函数:

Leaky ReLU, Exponential LU, Maxout 等

- 非线性取决于激活或者未激活的神经元构成的传播路径的选择
- 允许稀疏表示:
 - 对于特定的输入, 只有一部分神经元是活跃的



激活信号和梯度的稀疏传播

5.2 激活函数与损失函数

误差/损失函数

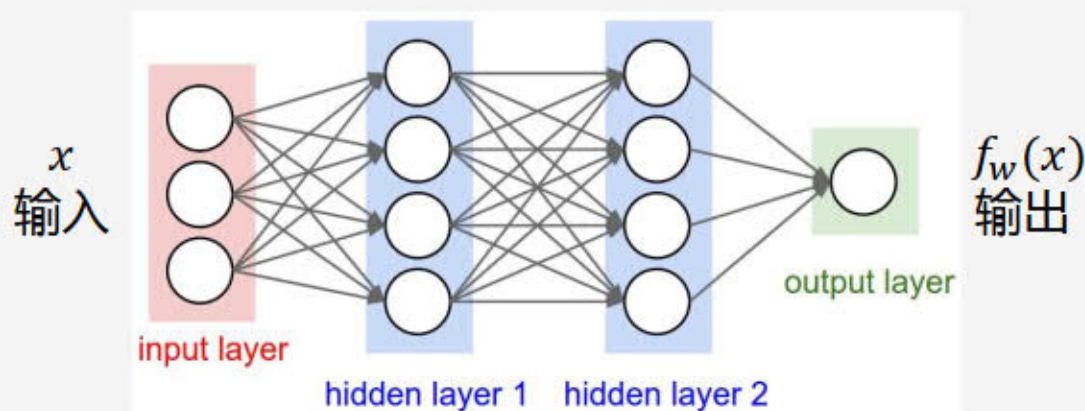
□ 随机梯度下降

- 从随机抽取的样本中进行更新（实际上执行批处理更新）：

$$w = w - \eta \frac{\partial \mathcal{L}(w)}{\partial w}$$

□ 一个二值输入的平方误差损失：

$$\mathcal{L}(w) = \frac{1}{2} (y - f_w(x))^2$$



5.2 激活函数与损失函数

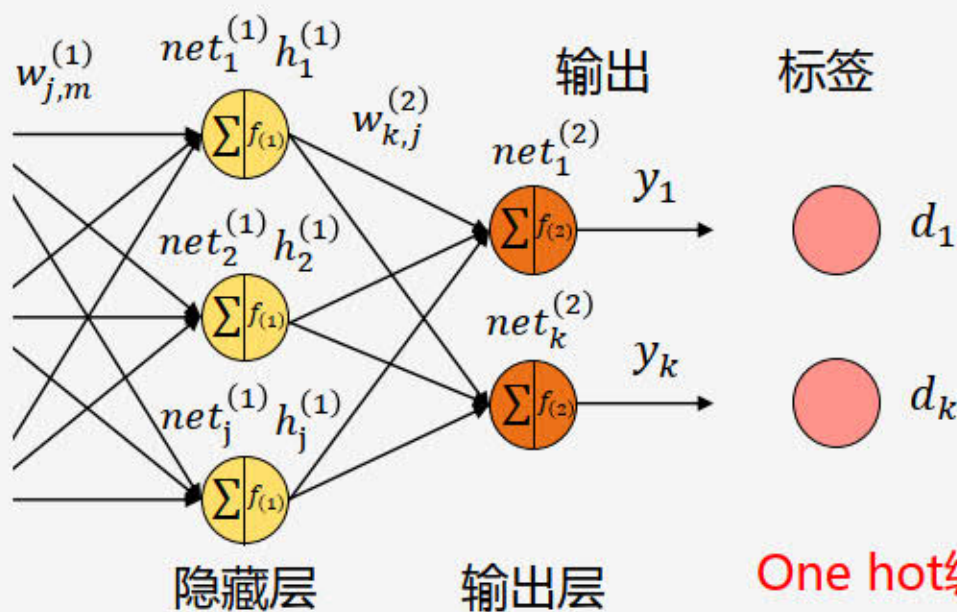
误差/损失函数

多分类问题的Softmax损失 (交叉熵损失)

(类标签服从多项分布)

$$\mathcal{L}(w) = - \sum_k d_k \log \hat{y}_k$$

$$\text{其中 } \hat{y}_k = \frac{\exp(\sum_j w_{k,j}^{(2)} h_j^{(1)})}{\sum_{k'} \exp(\sum_j w_{k',j}^{(2)} h_j^{(1)})}$$



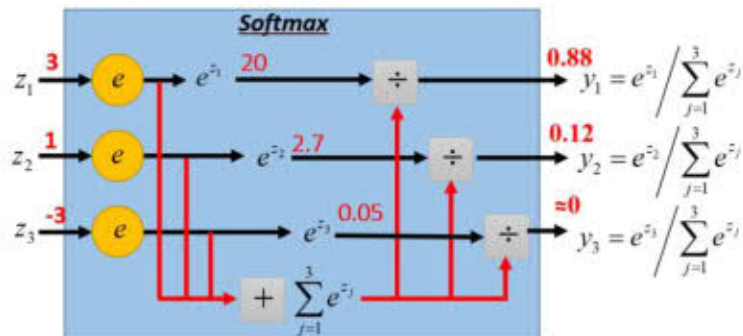
Multi-class Classification (3 classes as example)

$$\begin{aligned} C_1: w^1, b_1 & \quad z_1 = w^1 \cdot x + b_1 \\ C_2: w^2, b_2 & \quad z_2 = w^2 \cdot x + b_2 \\ C_3: w^3, b_3 & \quad z_3 = w^3 \cdot x + b_3 \end{aligned}$$

Probability:

- $1 > y_i > 0$
- $\sum_i y_i = 1$

$$y_i = P(C_i | x)$$

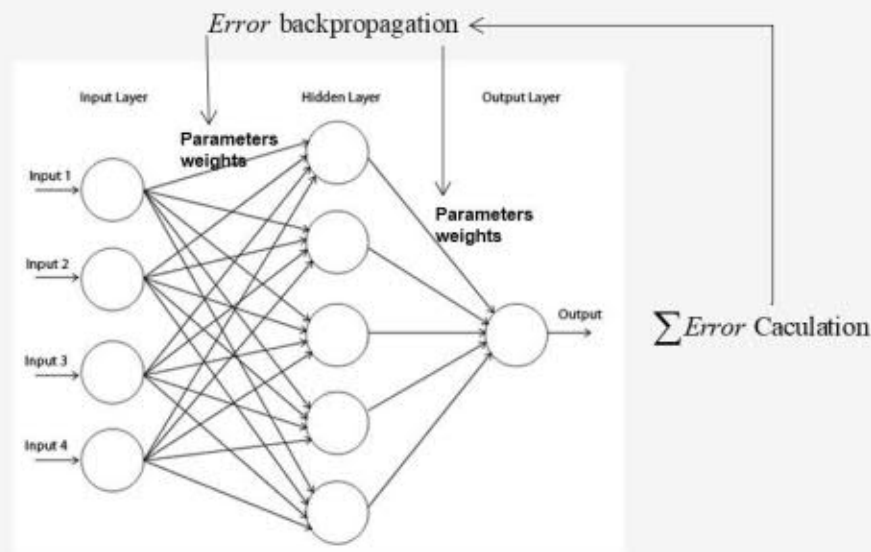


One hot编码的类标签

5.3 反向传播算法介绍

隐含层和反向传播 (1986~)

- 多层神经网络的有效算法之一是反向传播算法
- 它在1986年被重新提出，神经网络重新变得流行



注意：反向传播算法最开始在1974年被Werbos^[1]发现，之后在1985年左右被Rumelhart, Hinton, Williams^[2] 以及Parker^[3] 分别独立重新发现

[1] Werbos, Paul John (1975). Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences. Harvard University.

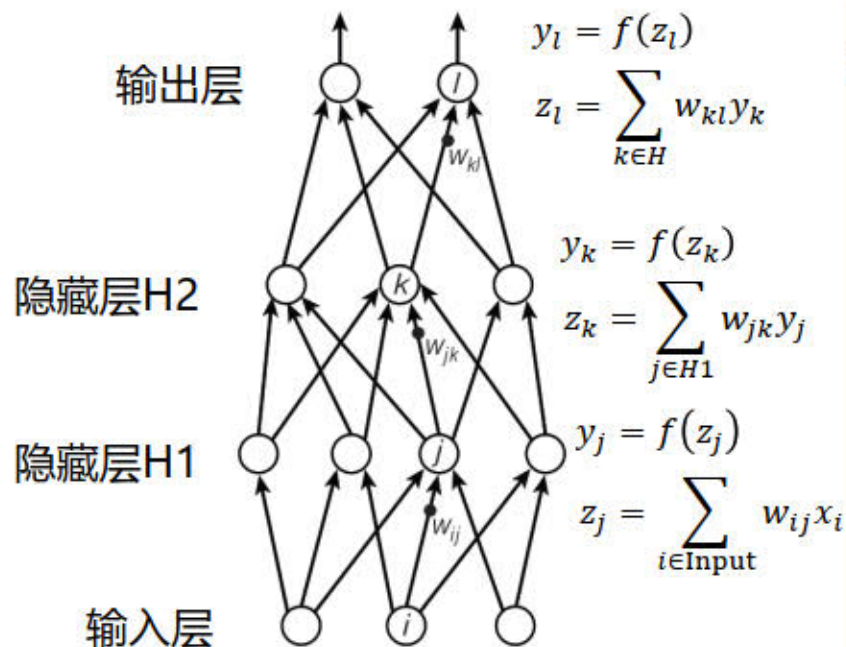
[2] Rumelhart, David E.; Hinton, Geoffrey E.; Williams, Ronald J. (8 October 1986). "Learning representations by back-propagating errors". Nature. 323 (6088): 533–536.

[3] Parker, D. B. (1985). "Learning Logic," Technical Report TR-47, Center for Computational Research in Economics and Management Science, Massachusetts Institute of Technology, Cambridge, MA.

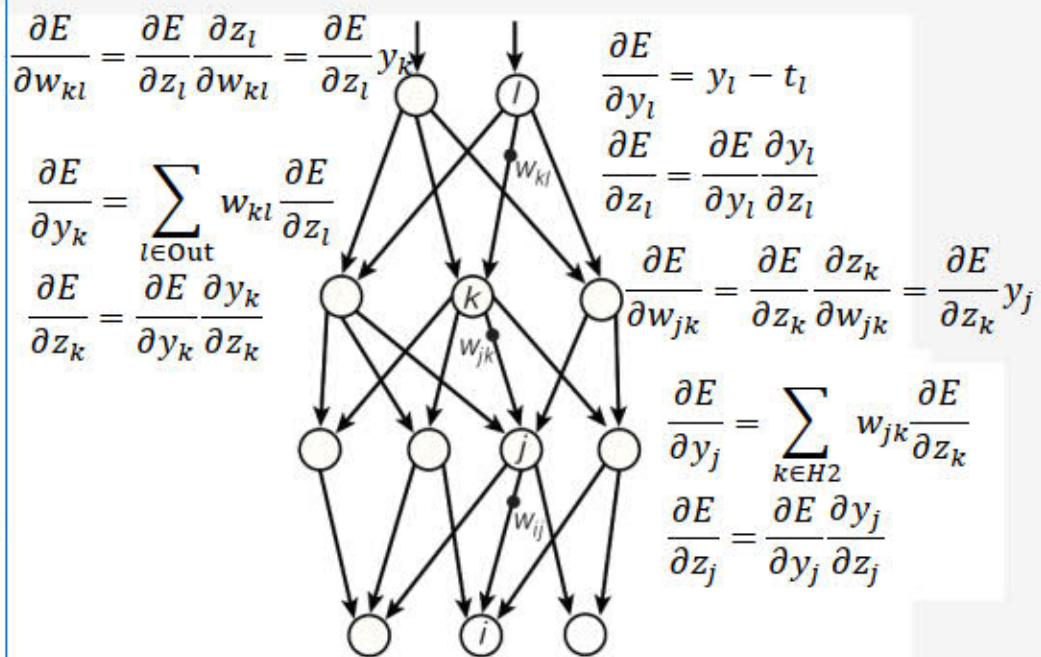
5.3 反向传播算法介绍

通过反向传播学习神经网络

输入数据逐层计算得到输出

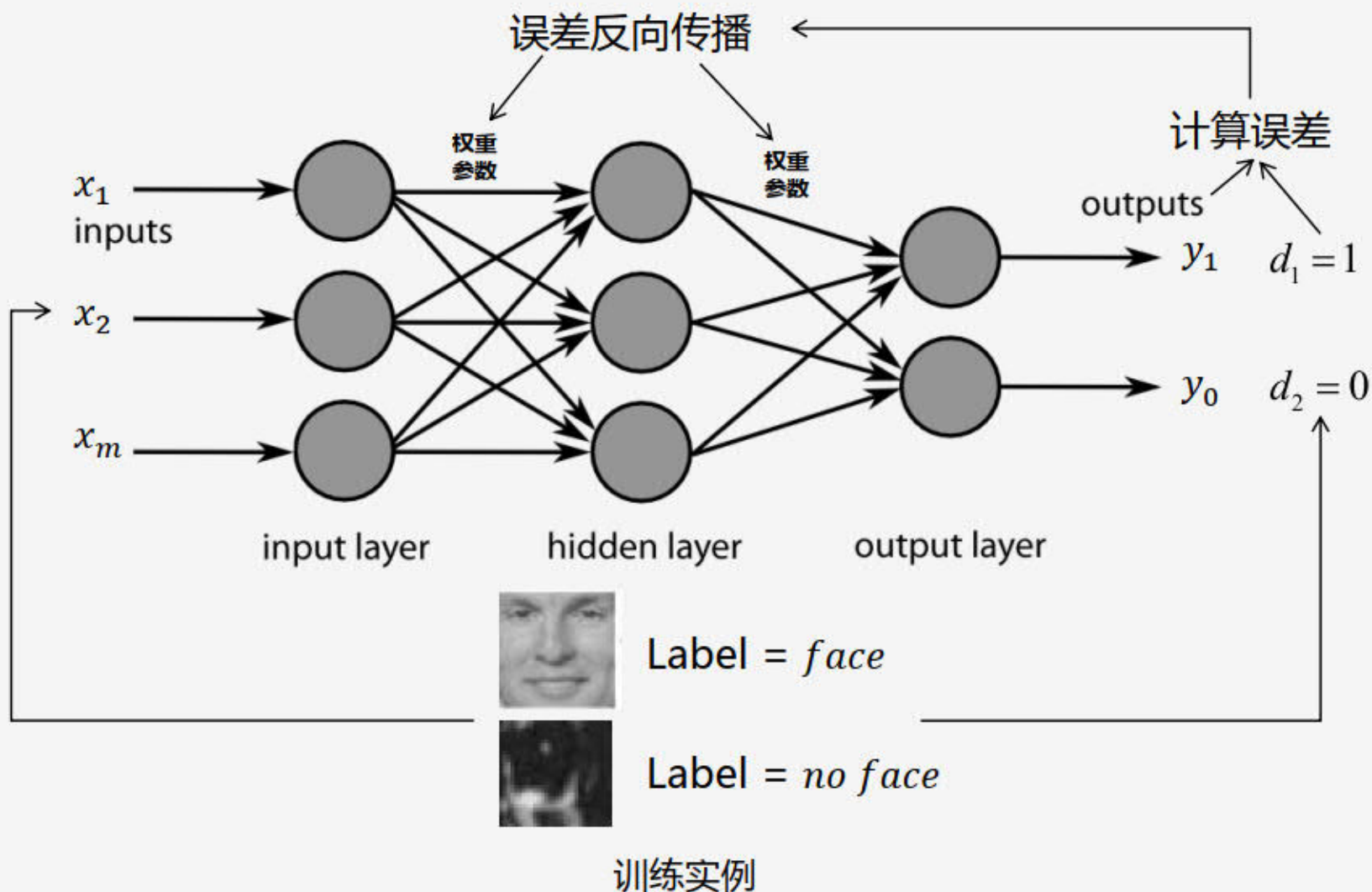


将输出与正确答案比较得到
误差，然后反向求导数



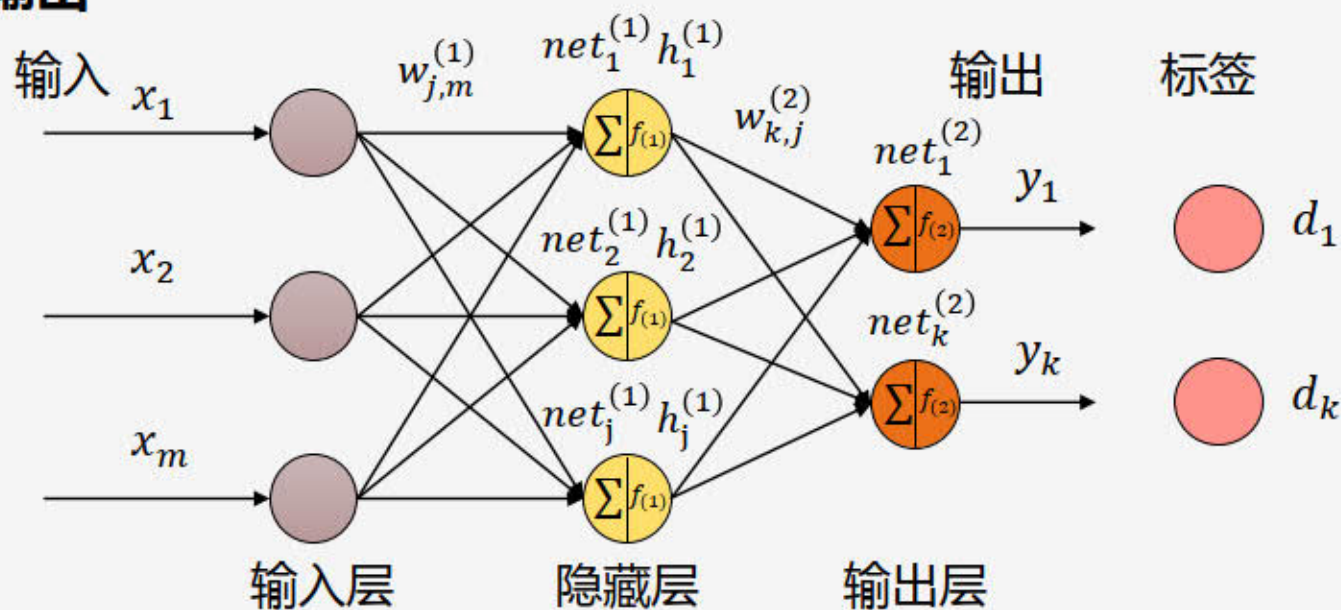
5.3 反向传播算法介绍

通过反向传播学习神经网络



5.3 反向传播算法介绍

预测输出



双层前馈神经网络

前馈预测

$$x = (x_1, \dots, x_m) \xrightarrow{h_j^{(1)} = f_{(1)}(net_j^{(1)}) = f_{(1)}\left(\sum_m w_{j,m}^{(1)} x_m\right)} h_j^{(1)} \xrightarrow{y_k = f_{(2)}(net_k^{(2)}) = f_{(2)}\left(\sum_j w_{k,j}^{(2)} h_j^{(1)}\right)} y_k$$

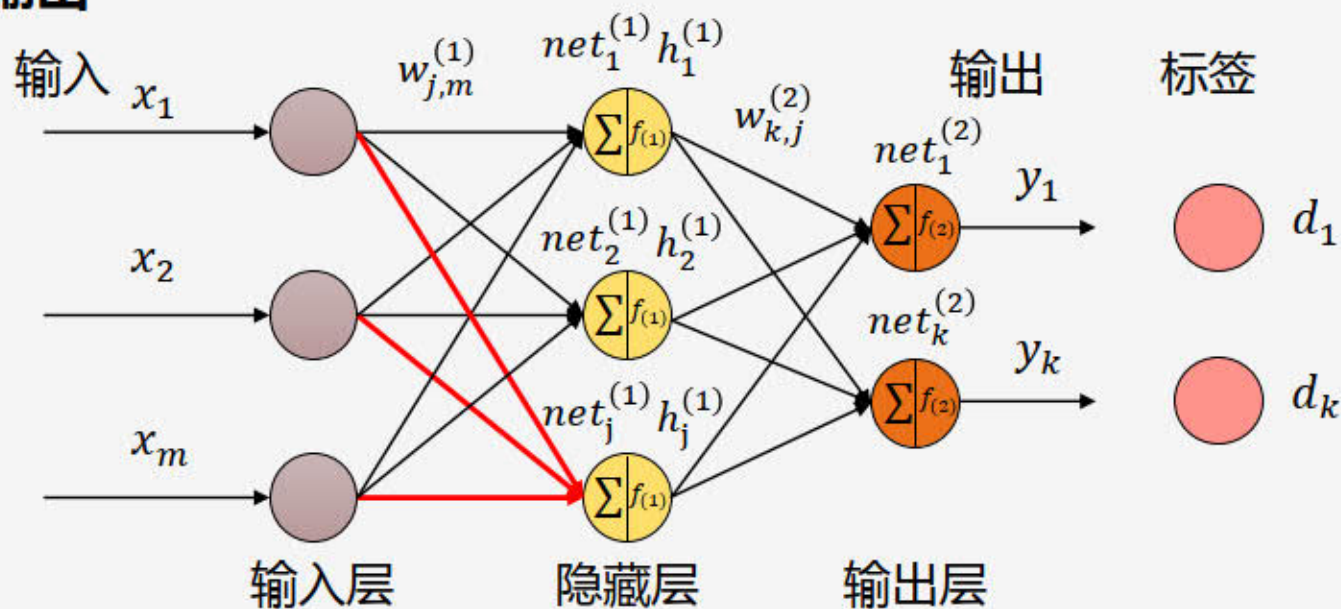
其中,

$$net_j^{(1)} = \sum_m w_{j,m}^{(1)} x_m$$

$$net_k^{(2)} = \sum_j w_{k,j}^{(2)} h_j^{(1)}$$

5.3 反向传播算法介绍

预测输出



双层前馈神经网络

前馈预测

$$x = (x_1, \dots, x_m) \xrightarrow{h_j^{(1)}} h_j^{(1)} \xrightarrow{y_k} y_k$$
$$h_j^{(1)} = f_{(1)}(net_j^{(1)}) = f_{(1)}\left(\sum_m w_{j,m}^{(1)} x_m\right) \quad y_k = f_{(2)}(net_k^{(2)}) = f_{(2)}\left(\sum_j w_{k,j}^{(2)} h_j^{(1)}\right)$$

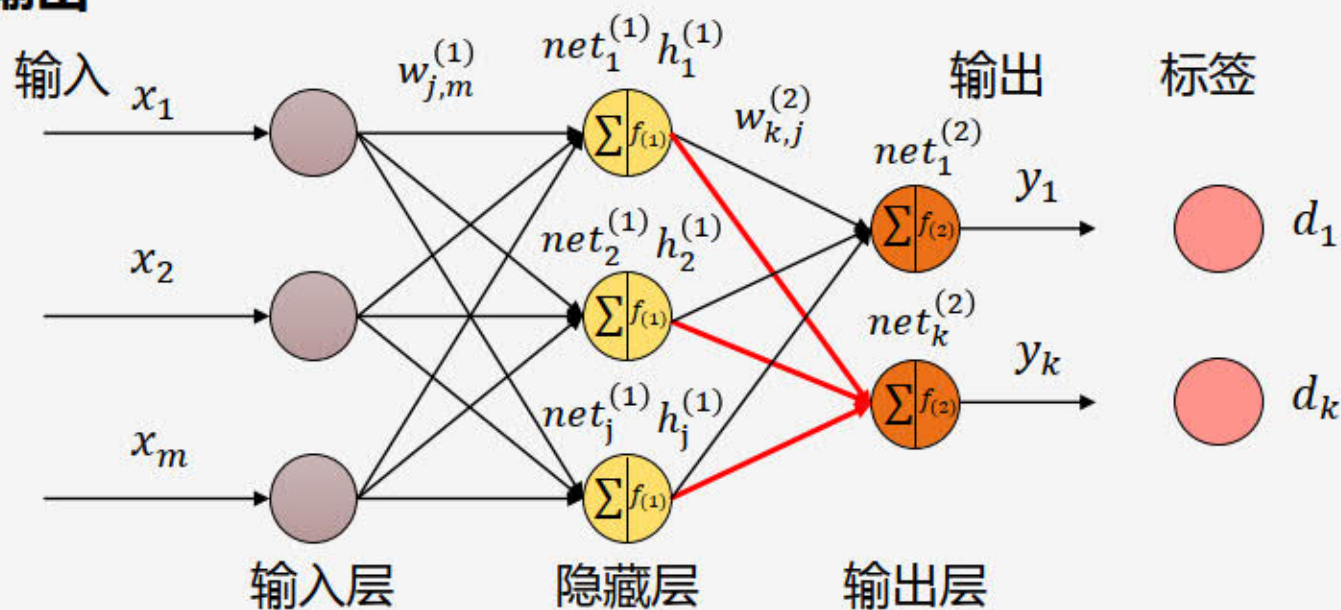
其中,

$$net_j^{(1)} = \sum_m w_{j,m}^{(1)} x_m$$

$$net_k^{(2)} = \sum_j w_{k,j}^{(2)} h_j^{(1)}$$

5.3 反向传播算法介绍

预测输出



双层前馈神经网络

前馈预测

$$x = (x_1, \dots, x_m) \longrightarrow h_j^{(1)} \longrightarrow y_k$$
$$h_j^{(1)} = f_{(1)}(net_j^{(1)}) = f_{(1)}\left(\sum_m w_{j,m}^{(1)} x_m\right) \quad y_k = f_{(2)}(net_k^{(2)}) = f_{(2)}\left(\sum_j w_{k,j}^{(2)} h_j^{(1)}\right)$$

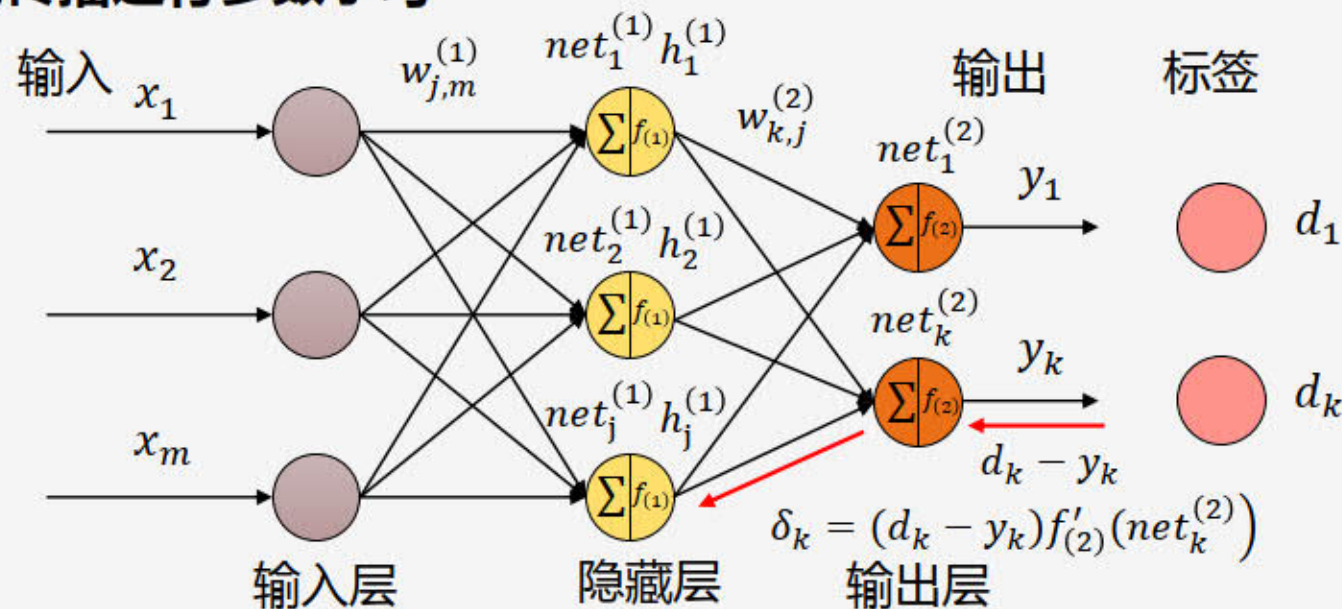
其中,

$$net_j^{(1)} = \sum_m w_{j,m}^{(1)} x_m$$

$$net_k^{(2)} = \sum_j w_{k,j}^{(2)} h_j^{(1)}$$

5.3 反向传播算法介绍

反向传播进行参数学习



$$net_j^{(1)} = \sum_m w_{j,m}^{(1)} x_m$$

双层前馈神经网络

$$net_k^{(2)} = \sum_j w_{k,j}^{(2)} h_j^{(1)}$$

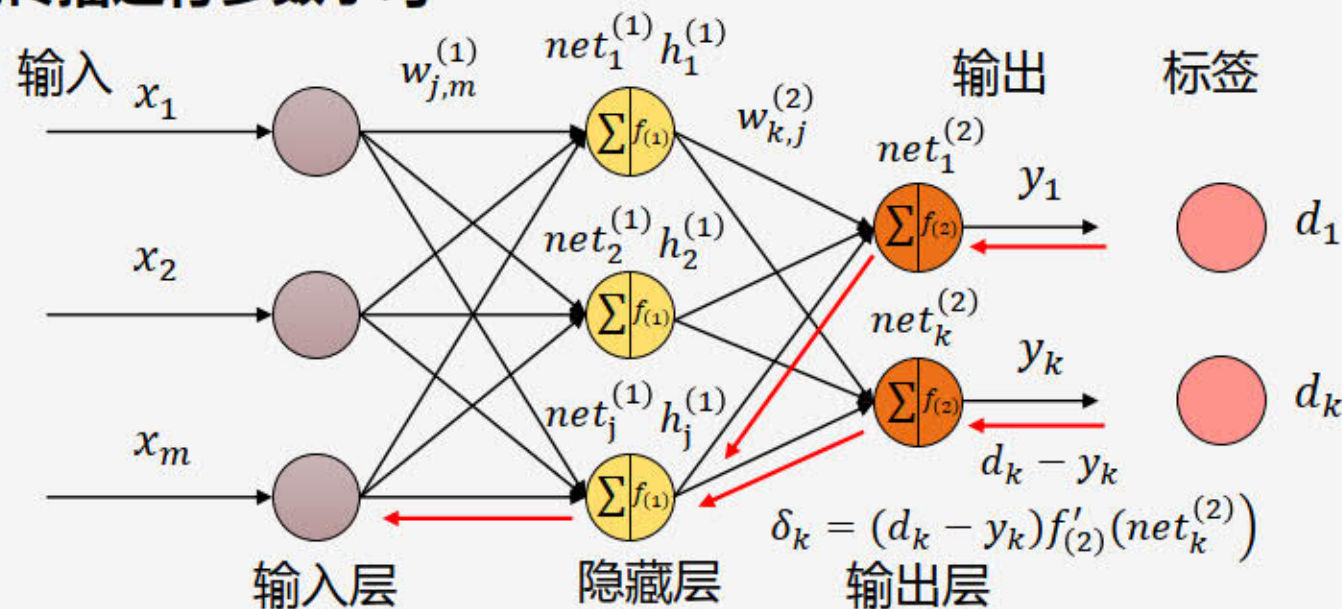
$$w_{k,j}^{(2)} = w_{k,j}^{(2)} + \Delta w_{k,j}^{(2)}$$

$$\Delta w_{k,j}^{(2)} = \eta \text{Error}_k \text{Output}_j = \eta \delta_k h_j^{(1)} \quad E(W) = \frac{1}{2} \sum_k (y_k - d_k)^2$$

$$\Delta w_{k,j}^{(2)} = -\eta \frac{\partial E(W)}{\partial w_{k,j}^{(2)}} = -\eta (y_k - d_k) \frac{\partial y_k}{\partial net_k^{(2)}} \frac{\partial net_k^{(2)}}{\partial w_{k,j}^{(2)}} = \eta (d_k - y_k) f'_{(2)}(net_k^{(2)}) h_j^{(1)} = \eta \delta_k h_j^{(1)}$$

5.3 反向传播算法介绍

反向传播进行参数学习



$$net_j^{(1)} = \sum_m w_{j,m}^{(1)} x_m$$

双层前馈神经网络

$$net_k^{(2)} = \sum_j w_{k,j}^{(2)} h_j^{(1)}$$

$$w_{j,m}^{(1)} = w_{j,m}^{(1)} + \Delta w_{j,m}^{(1)}$$

$$\Delta w_{j,m}^{(1)} = \eta \text{Error}_j \quad \text{Out}_m = \eta \delta_j x_m$$

$$E(W) = \frac{1}{2} \sum_k (y_k - d_k)^2$$

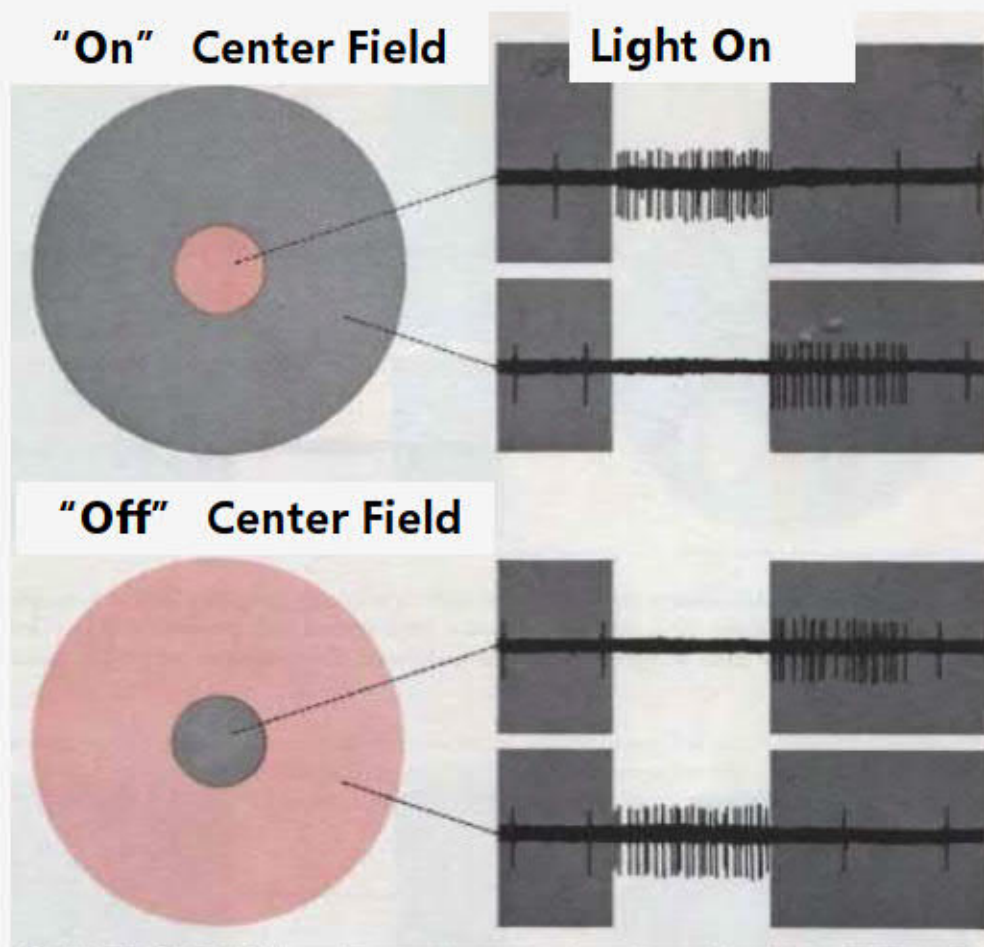
$$\Delta w_{j,m}^{(1)} = -\eta \frac{\partial E(W)}{\partial w_{j,m}^{(1)}} = -\eta \frac{\partial E(W)}{\partial h_j^{(1)}} \frac{\partial h_j^{(1)}}{\partial w_{j,m}^{(1)}} = \eta \sum_k (d_k - y_k) f'_{(2)}(net_k^{(2)}) w_{k,j}^{(2)} x_m f'_{(1)}(net_j^{(1)}) = \eta \delta_j x_m$$

5.4 卷积神经网络

卷积神经网络：感受野

感受野

- 视网膜中的神经元只响应视野的受限区域中的光刺激
- 两个视网膜神经节细胞感受野的动物实验
 - 视野是视网膜的圆形区域
 - 当中心点亮周围变暗时，细胞（上部）响应
 - 当中心变暗并且周围照亮时，细胞（下部）响应
 - 当中心和周围都被照亮时，两个单元都给出响应和关闭响应，但都没有响应像只有中心或周围被照亮时那样强烈



5.4 卷积神经网络

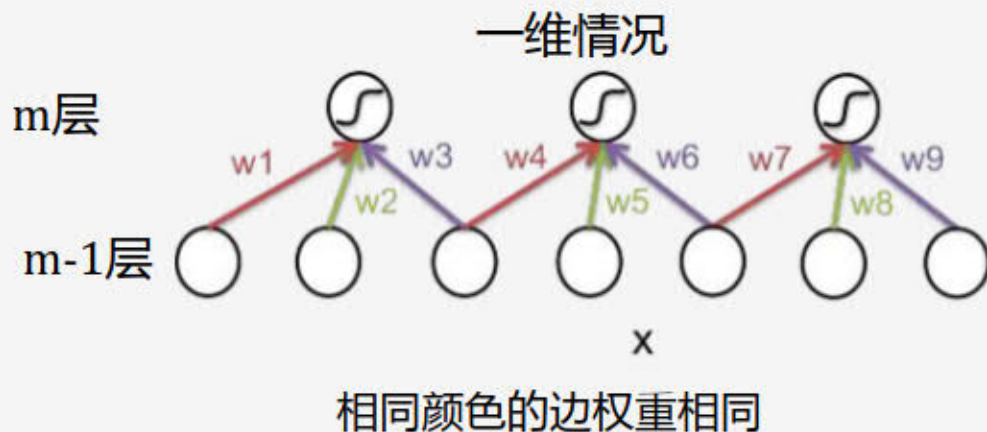
卷积神经网络 (Convolutional Neural Network)

□ 通过局部相关性稀疏连接

- 过滤器：m层中隐含单元的输入来自m-1层中具有空间连接的感受野的单元的子集

□ 共享权重

- 每个过滤器都在整个视野中复制。这些复制的单元共享相同的权重并形成特征映射。



二维情况(下标是权重)

1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

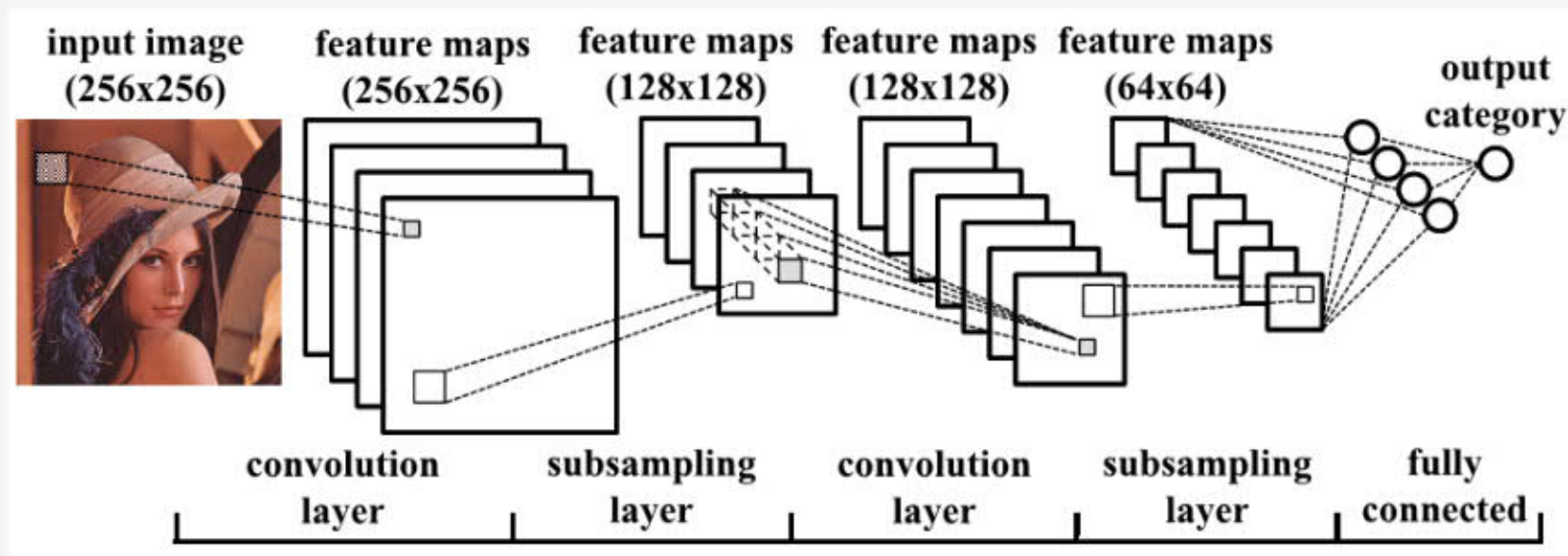
m-1层

4		

m层的一个过滤器

5.4 卷积神经网络

卷积神经网络

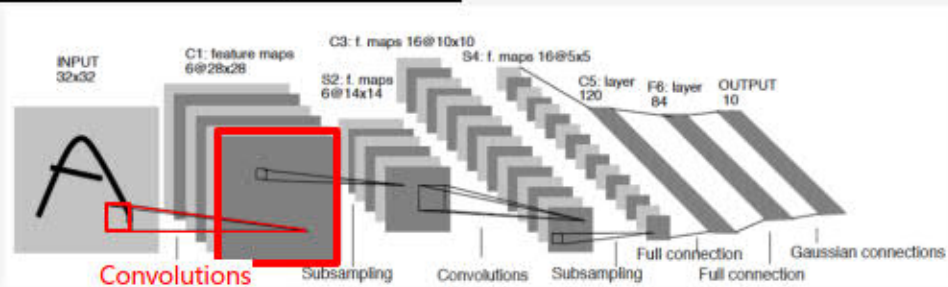
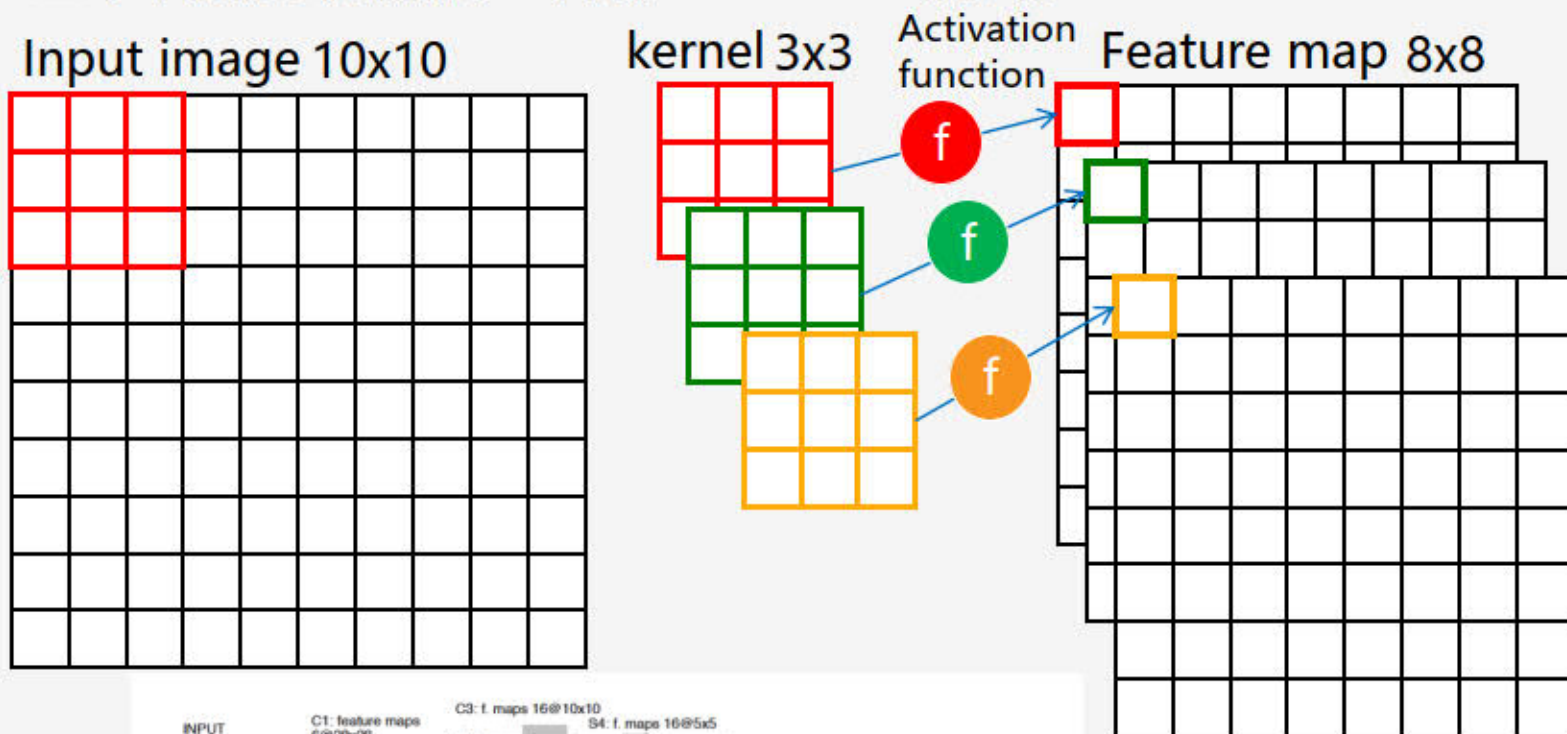


Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. "Gradient-based learning applied to document recognition." Proceedings of the IEEE, 86(11) 1998

5.4 卷积神经网络

卷积层

- 范例：具有 3×3 过滤器的 10×10 输入图像产生 8×8 输出图像
- 三个不同的过滤器(权重不同)产生3张 8×8 的图像



5.4 卷积神经网络

下采样池化层

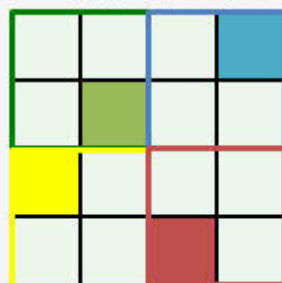
池化

- 将输入图像划分为一组非重叠矩形，并且对于每个这样的子区域，输出最大值或平均值

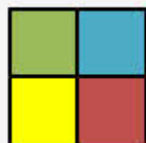
最大池化

- 优点
 - 减少计算
 - 是一种获取给定兴趣区域响应最强烈的节点的方法
- 缺点
 - 可能会导致准确的空间信息丢失

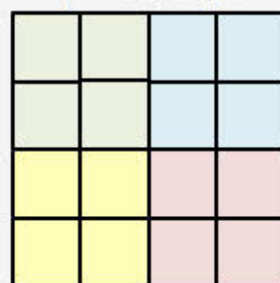
最大池化



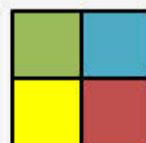
2x2过滤器
中最大值



平均池化



2x2过滤器
中平均值

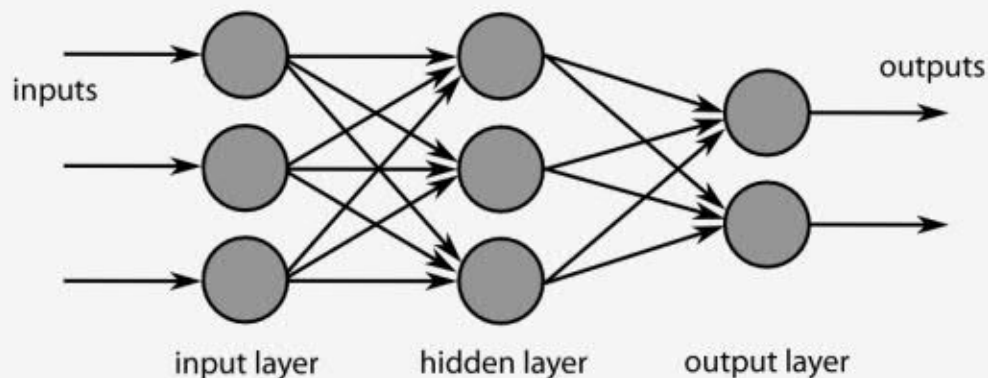
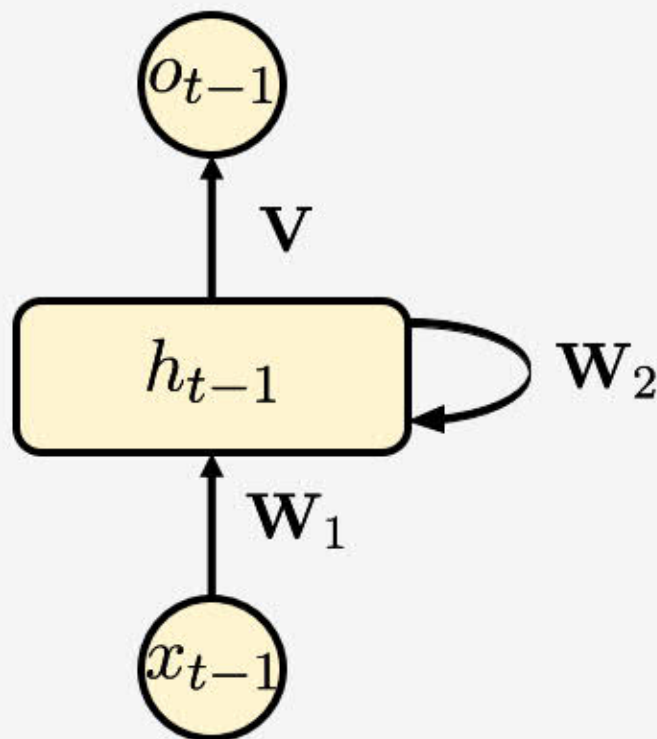


5.5 循环神经网络

循环神经网络 (Recurrent Neural Network)

□ 从前馈神经网络到循环神经网络

- 前馈神经网络一个有向无环图，无自连边
- 循环神经网络有一条自连的边，将隐层信息随着时间步往前传



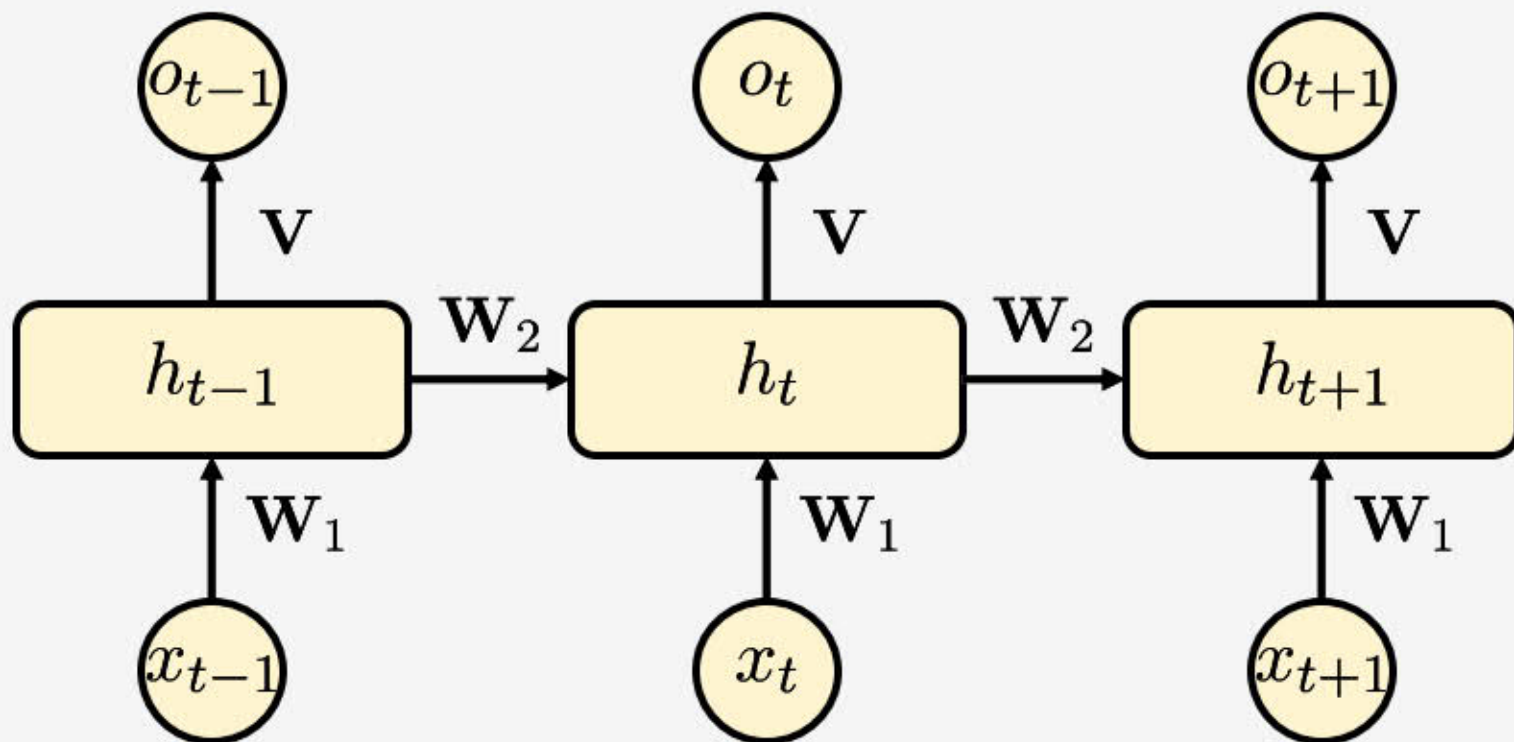
前馈神经网络

5.5 循环神经网络

循环神经网络 (Recurrent Neural Network)

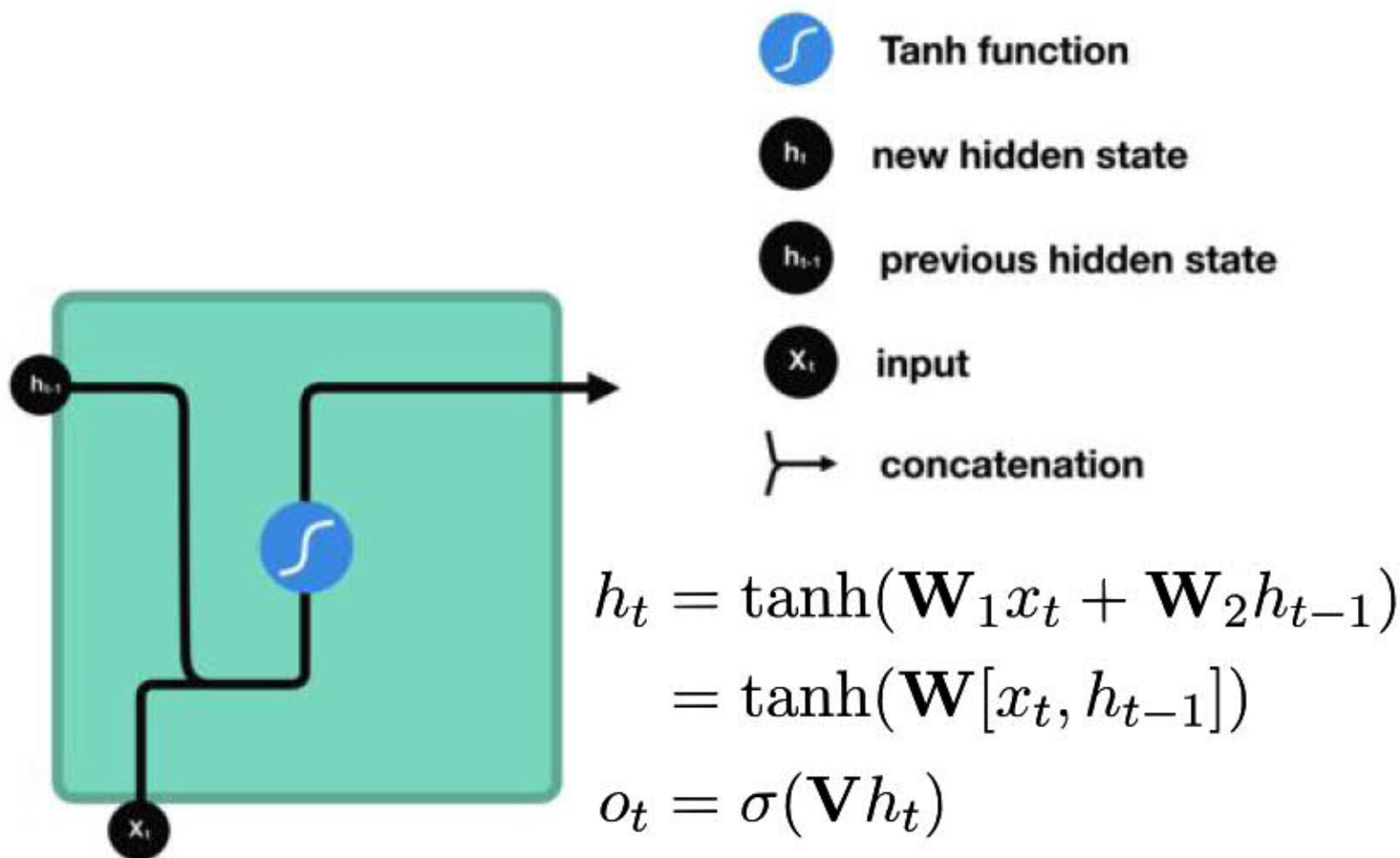
□ 从前馈神经网络到循环神经网络

- 前馈神经网络一个有向无环图，无自连边
- 循环神经网络有一条自连的边，将隐层信息随着时间步往前传



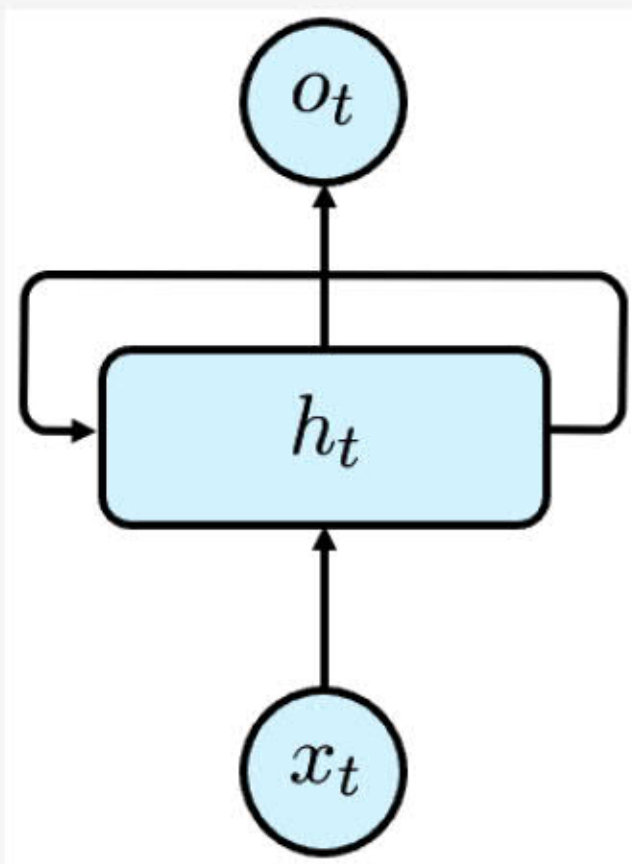
5.5 循环神经网络

RNN工作原理图示



5.5 循环神经网络

RNN的工作方式



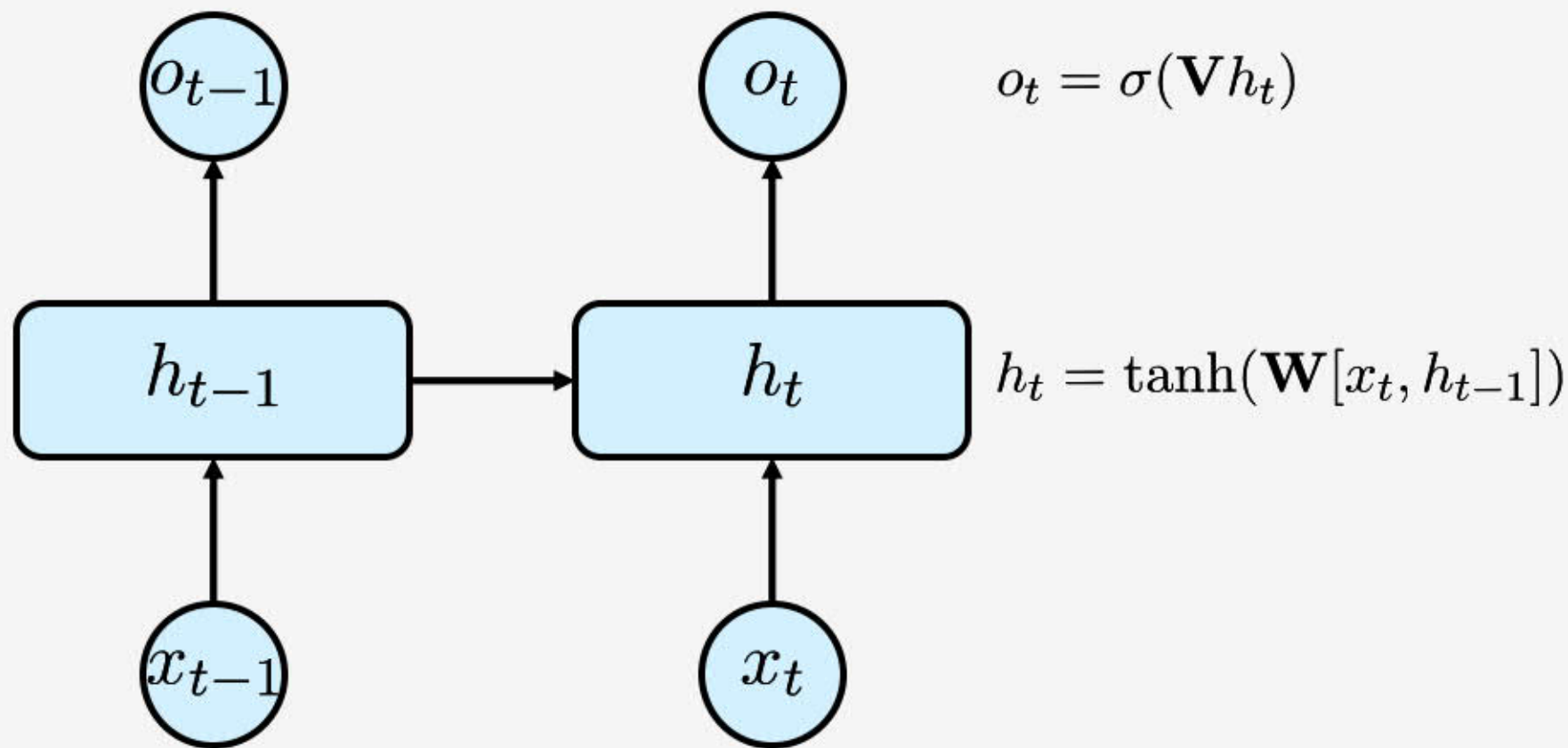
$$o_t = \sigma(\mathbf{V}h_t)$$

$$\begin{aligned} h_t &= \tanh(\mathbf{W}_1x_t + \mathbf{W}_2h_{t-1}) \\ &= \tanh(\mathbf{W}[x_t, h_{t-1}]) \end{aligned}$$

把 x_t, h_{t-1} 接起来

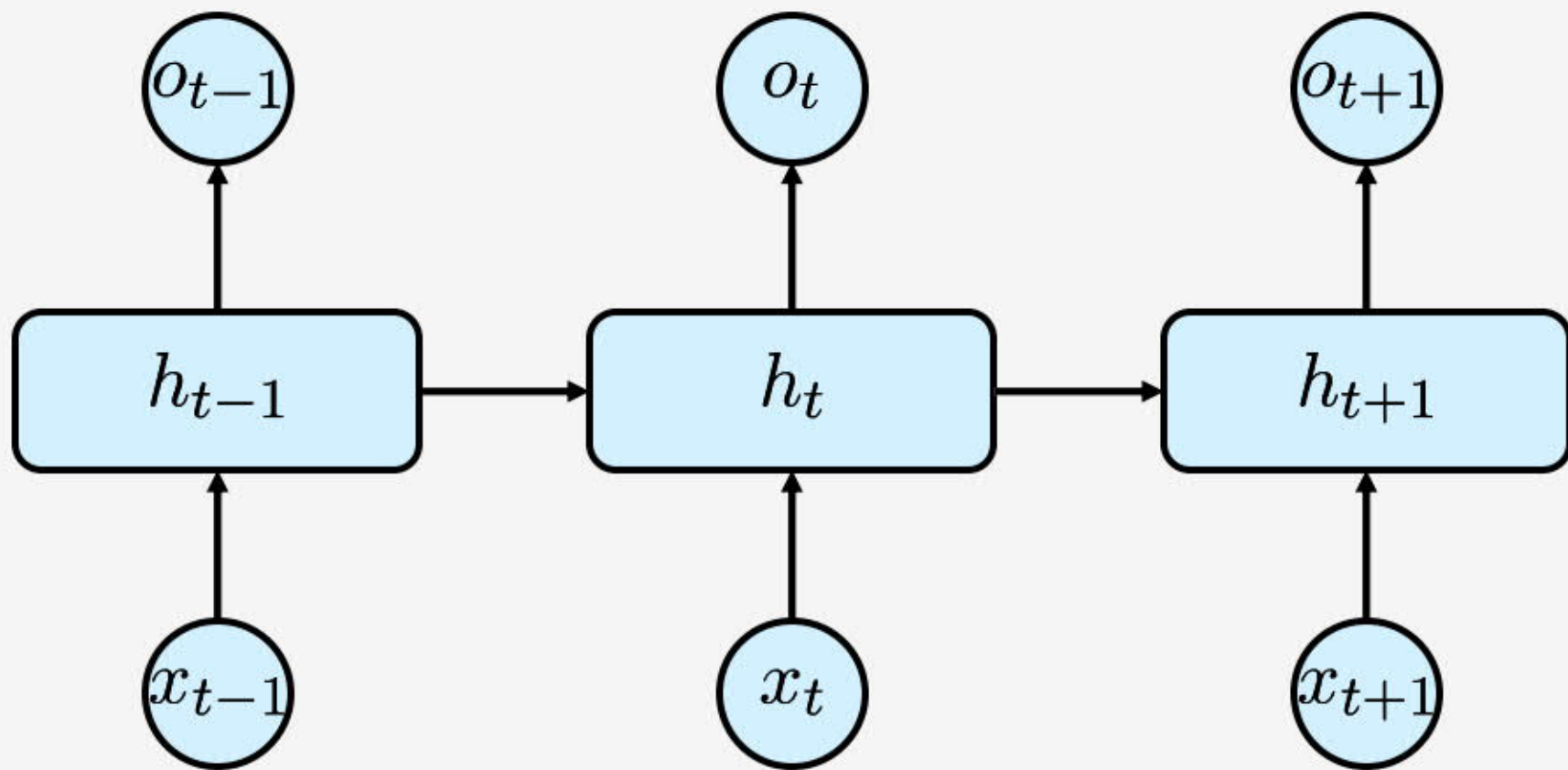
5.5 循环神经网络

RNN的工作方式



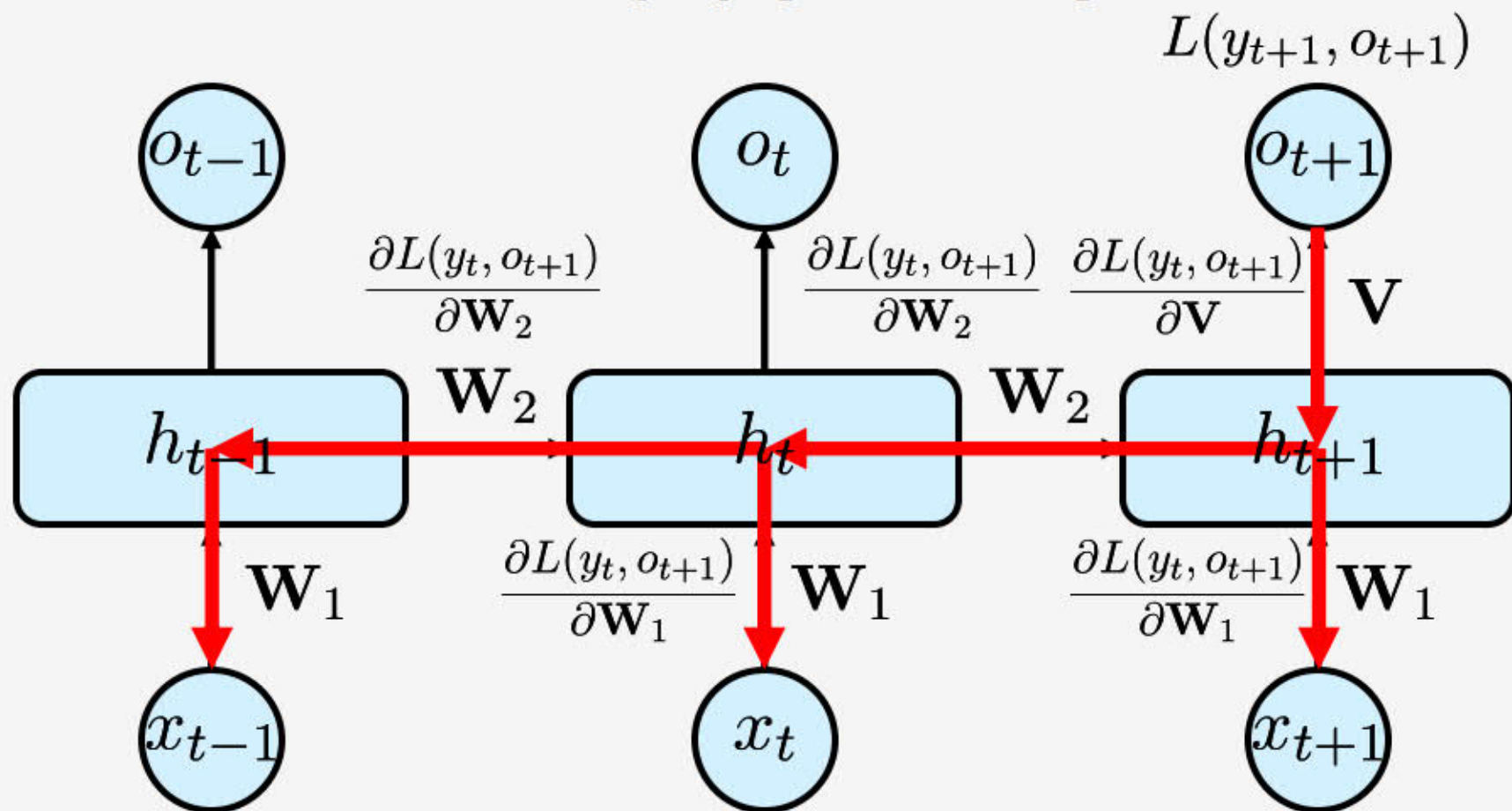
5.5 循环神经网络

RNN的工作方式



5.5 循环神经网络

RNN的反向传播学习 (Backpropagation through time)



在上图例子中, BPTT让 $t+1$ 时间步的损失梯度更新了:

- V 一次, W_2 两次, W_1 三次

神经网络总结

- 通用近似：双层神经网络可以逼近任何函数
- 到目前为止，反向传播是多层神经网络最重要的训练方案
- 深度学习，即用大数据训练的深度神经网络，效果非常好
 - MLP (DNN) 、CNN、RNN、LSTM
 - 感兴趣的同学可以关注Attention和Transformer
- 结合其他机器学习模型的神经网络取得了进一步成功