



人工智能：机器学习 III

饶洋辉

计算机学院,

中山大学

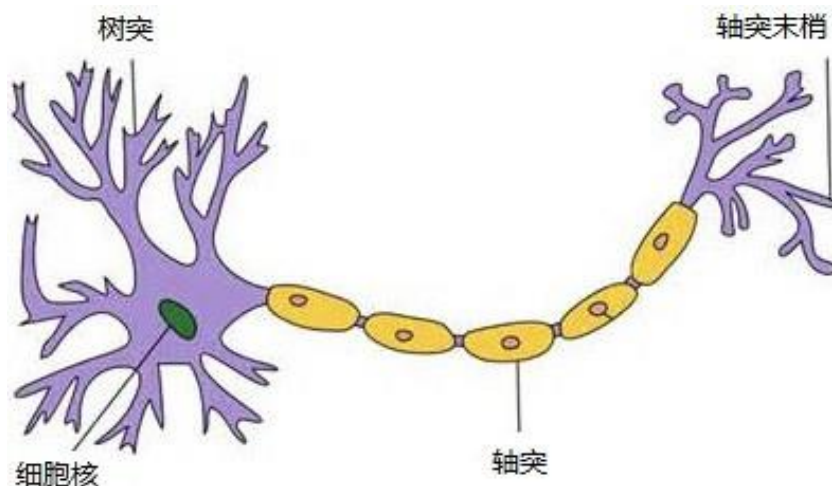
raoyangh@mail.sysu.edu.cn

<http://cse.sysu.edu.cn/node/2471>

课件来源：中山大学陈川副教授等

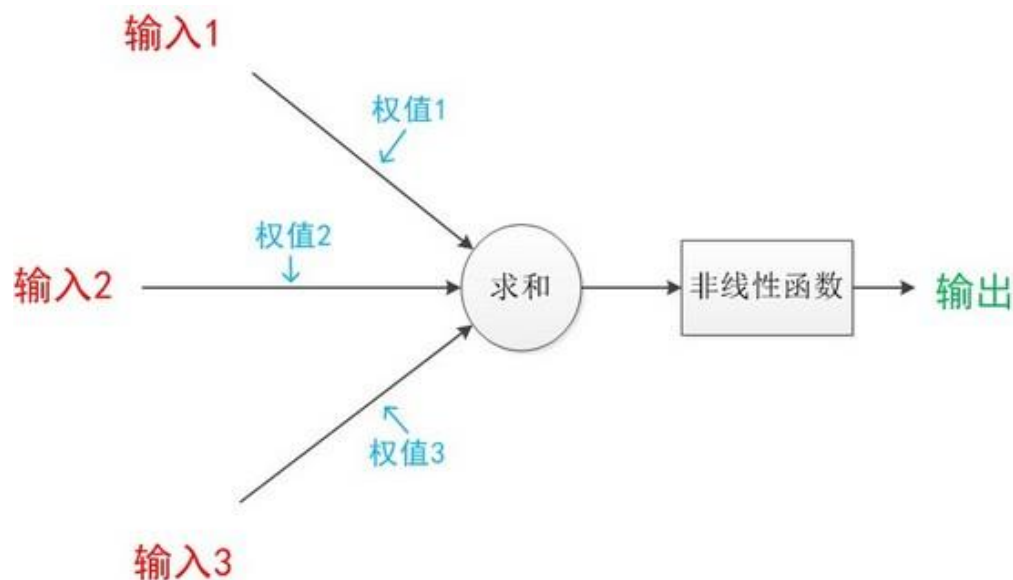
神经网络模型

- （人工）神经网络模型受到生物神经系统的启发。
- 人脑主要由称为神经元的神经细胞组成。1904年，生物学家已经知晓神经元的组成结构。
- 一个神经元通常具有多个树突，主要用来接受传入信息；而轴突只有一条，轴突尾端有许多轴突末梢可以给其他多个神经元传递信息。轴突末梢跟其他神经元的树突产生连接，从而传递信号。这个连接的位置在生物学上叫做“突触”。



神经网络模型

- 类比生物结构：输入可以类比为神经元的树突，而输出可以类比为神经元的轴突，计算则可以类比为细胞核。神经元模型是一个包含输入，输出与计算功能的模型。
- 下图是一个典型的神经元模型，它包含3个输入，1个输出，以及2个计算功能。
- 中间的箭头线被称为“连接”，它们附有一个“权值”。



神经网络模型

- 神经网络模型由大量简单的交互节点组成(人工神经元)
- 知识由这些节点之间的连接强度来表示
- 知识是通过学习过程调整连接而获得的
- 所有的神经元同时并独立地处理他们的输入

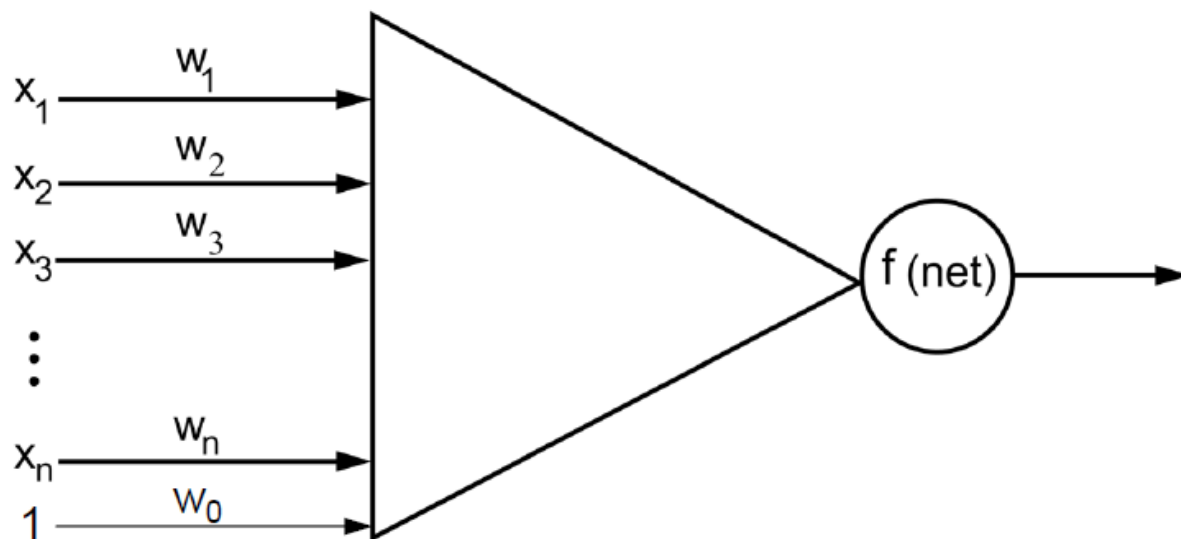
人工神经元

- 神经网络中的计算单位是人工神经元
- 人工神经元的组成部分为：
 - 输入信号 x_i . 这些信号代表来自环境的数据或其他神经元的激活(activation)
 - 一组实值(real-value)权重 w_i . 这些权重的值代表连接强度
 - 一个激活层 $\sum_i w_i x_i$. 神经元的激活层由加权输入的总和确定
 - 一个阈值函数 f . 该函数通过确定激活是低于还是高于阈值来计算最终输出

人工神经元

- 给定一个激活值 $net = \sum_i w_i x_i$, 该神经元的输出为:

$$f(net) = \begin{cases} +1 & \text{if } \sum_i w_i x_i \geq 0 \\ -1 & \text{if } \sum_i w_i x_i < 0 \end{cases}$$



示例

- 可以使用人工神经元来计算逻辑“与”功能
 - 神经元有三个输入
 - x_1 和 x_2 表示原始输入
 - 第三个是具有固定值+1的偏置输入
 - 输入数据和偏差分别具有+1，+1和-1.5的权重
- 那么逻辑或功能呢？

人工神经元

- 感知机学习算法 (PLA) 可用于调整人工神经元的权重
- 调整权重，直到神经元的输出与训练样例的真实输出一致
- 使用以下规则

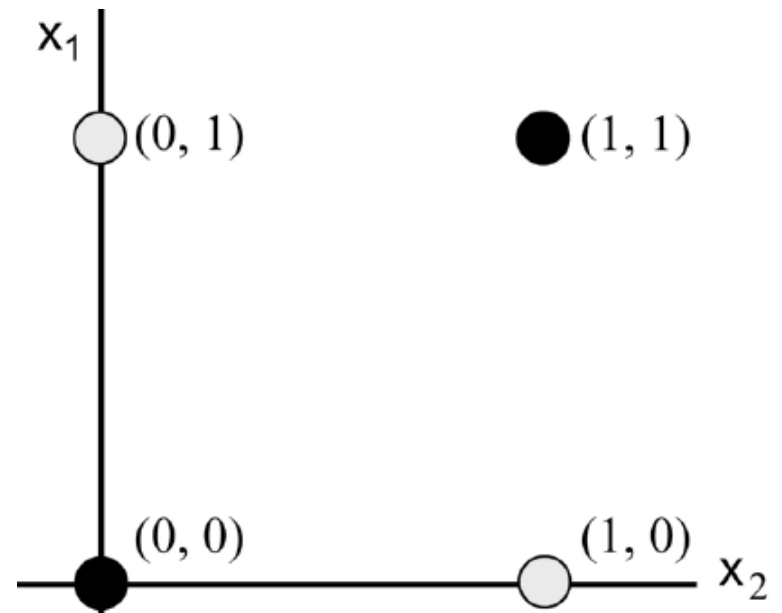
$$\mathbf{w}_{(t+1)} \leftarrow \mathbf{w}_{(t)} + y_{n(t)} \mathbf{x}_{n(t)}$$

人工神经元

- 感知机学习算法不能解决模式不可线性分离的问题
- 一个简单的例子是异或问题
- 需要多层网络来解决这类问题

示例

x_1	x_2	Output
1	1	-1
1	0	1
0	1	1
0	0	-1



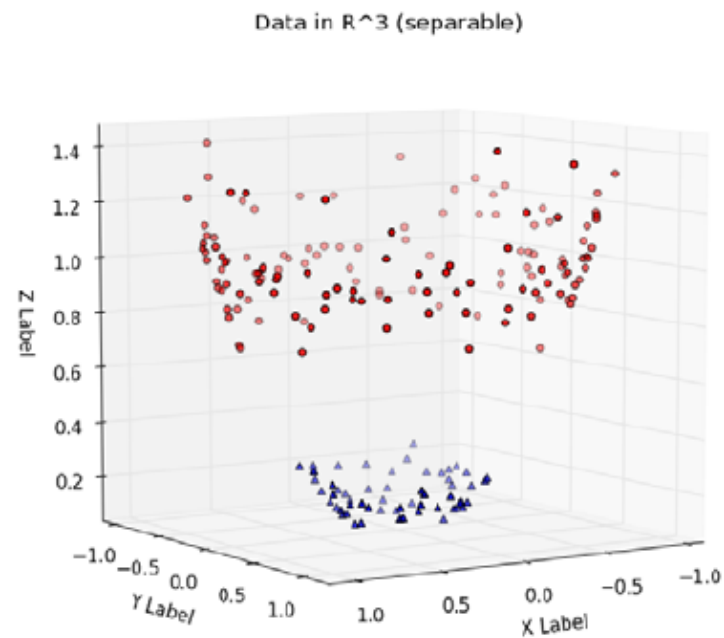
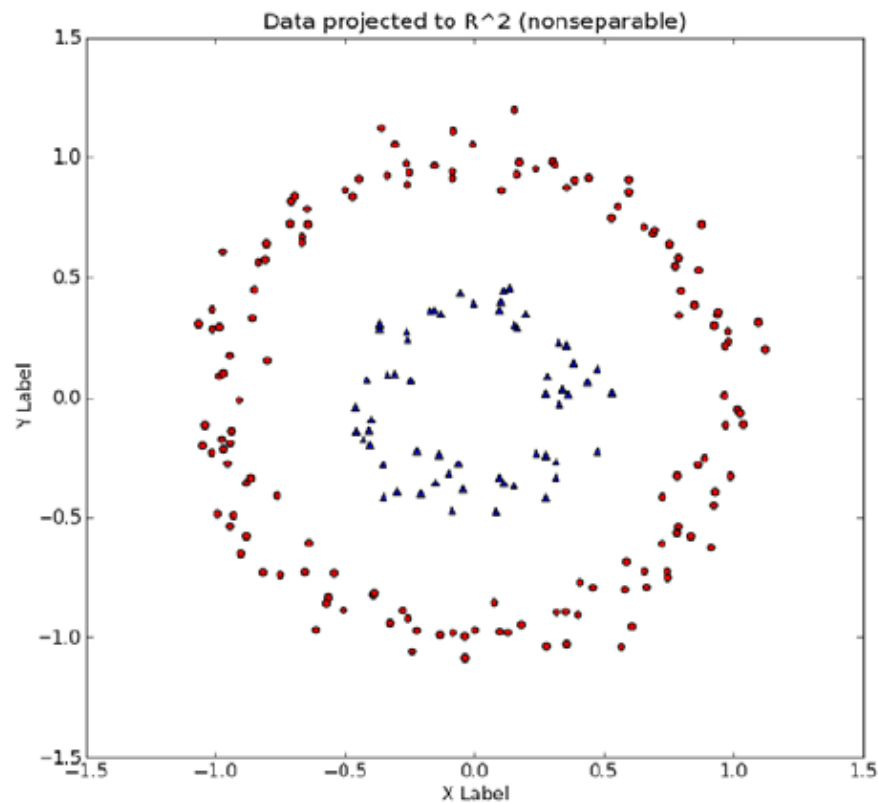


Figure 5: (Left) A dataset in \mathbb{R}^2 , not linearly separable. (Right) The same dataset transformed by the transformation: $[x_1, x_2] = [x_1, x_2, x_1^2 + x_2^2]$.

两层神经网络（多层感知机）

- 单层神经网络无法解决异或问题，但是当增加一个计算层之后，两层神经网络不仅可以解决异或问题，而且具有非常好的非线性分类效果。不过两层神经网络的计算是一个问题，没有一个较好的解法。
- 1986年，Rumelhart和Hinton等人提出了反向传播（Backpropagation, BP）算法，解决了两层神经网络所需要的复杂计算量问题，从而带动了业界使用两层神经网络研究的热潮。
- 这时候的Hinton还很年轻，30年以后，正是他重新定义了神经网络，带来了神经网络复苏。



David Rumelhart



Geoffery Hinton

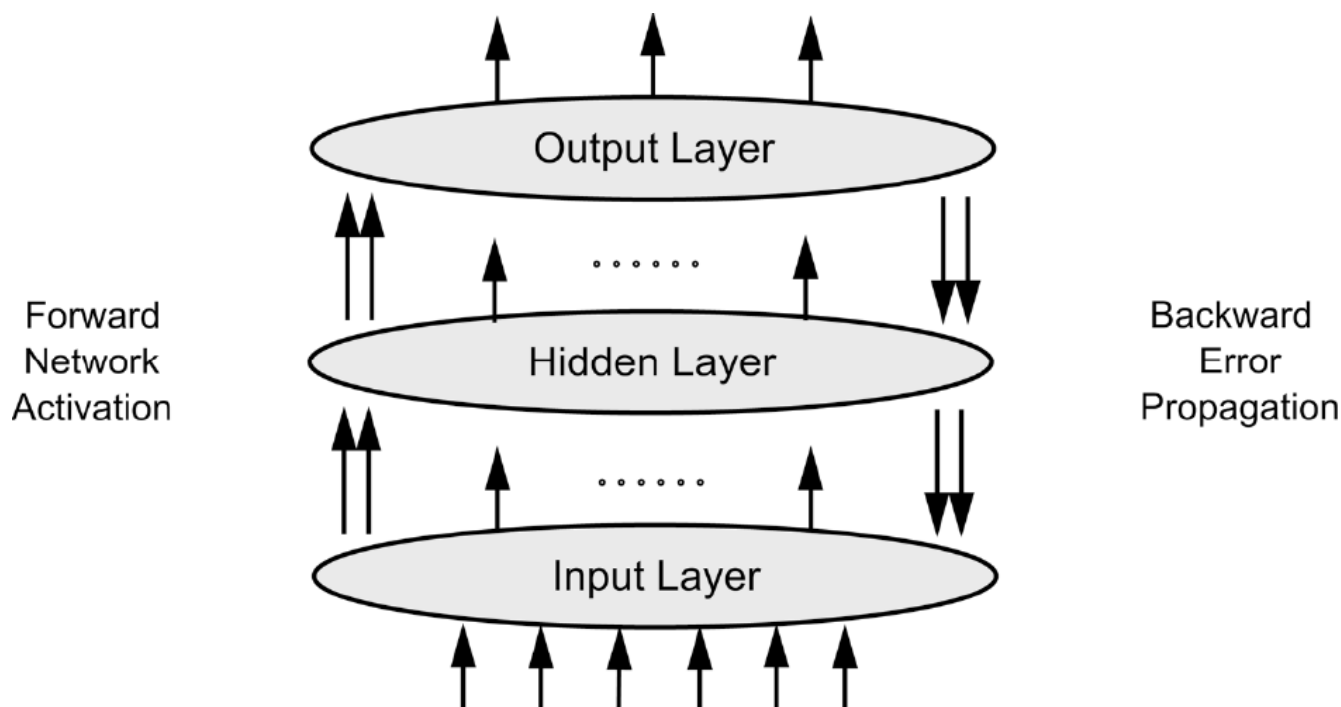
www.cybercontrols.org

神经网络模型

- 输入和输出节点之间的附加层称为隐藏层
- 嵌入在这些层中的节点被称为隐藏节点
- 一个前馈神经网络 (feedforward neural networks), 其中一层中的节点仅连接到下一层中的节点
- 反向传播学习算法是专门为多层神经网络设计的

神经网络模型

- 训练算法的每次迭代中有两个阶段
 - 前向传递阶段(forward phase)
 - 反向传递阶段(backward phase)



神经网络模型

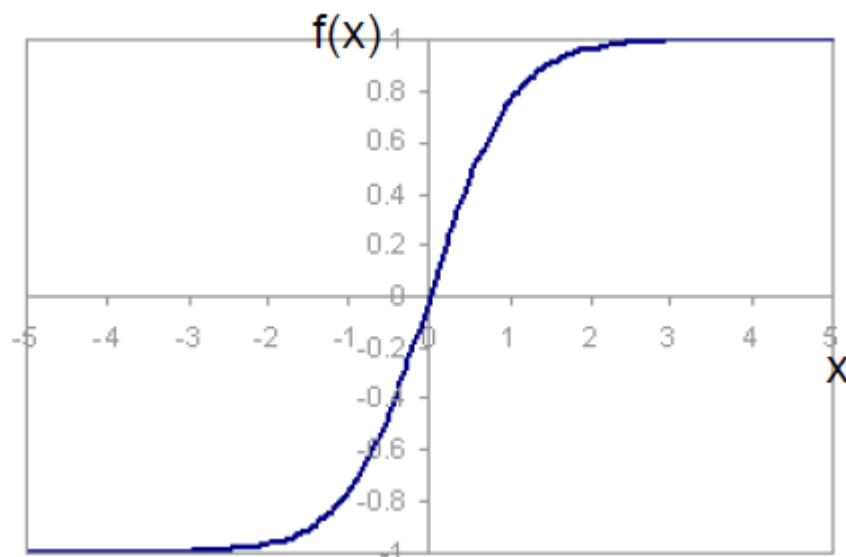
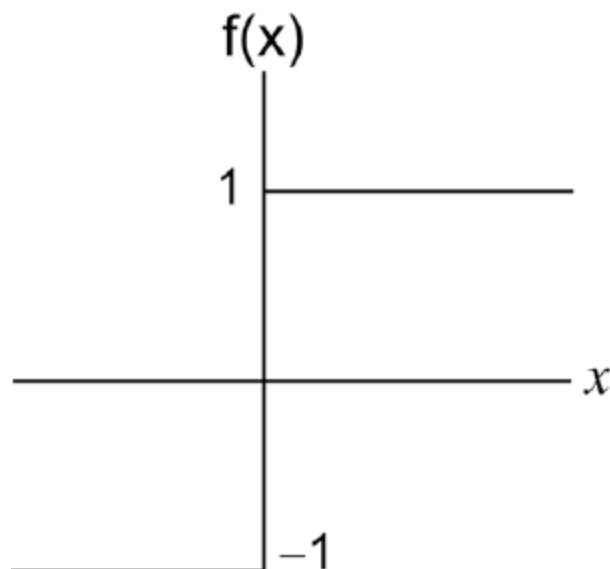
- 在前向传递阶段，从前一次迭代获得的权重被用来计算每个神经元的输出值
- 在计算第 $(l+1)$ 层的输出之前计算第 (l) 层处的神经元的输出

神经网络模型

- 在反向传递阶段，权重更新方程应用于相反的方向
- 换句话说，第 $(l+1)$ 层的权重在更新第 (l) 层的权重之前被更新
- 这种学习算法允许我们使用第 $(l+1)$ 层神经元的误差来估计第 (l) 层神经元的误差

神经网络模型

- 对于这种类型的网络，不采用PLA的阈值函数，而是使用另一个激活函数(activation function).



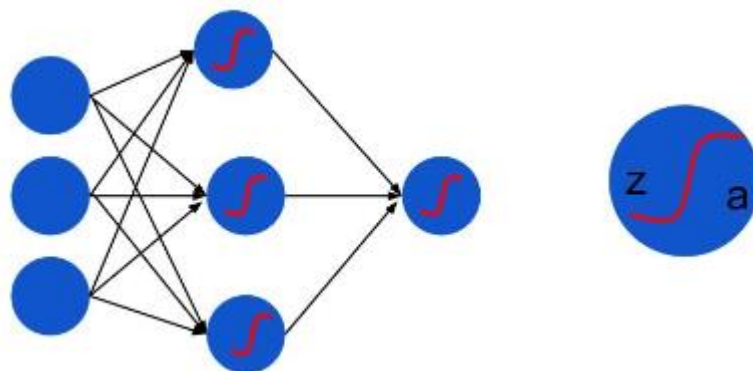
神经网络模型

- 常见的激活函数是双曲正切函数(tanh)

$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^{2x} - 1}{e^{2x} + 1}$$

- 这个函数的一个重要属性是可求导的

$$f'(x) = 1 - f(x)^2$$



神经网络模型

- 另一个常见的激活函数是sigmoid函数

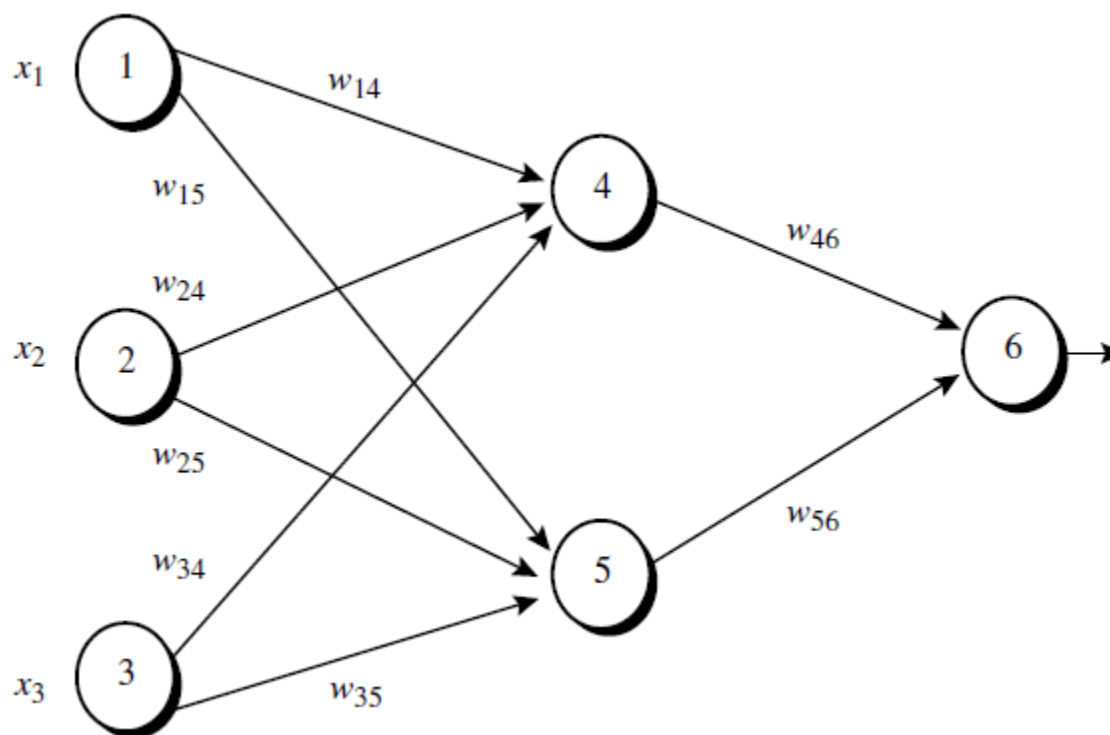
$$f(x) = \frac{1}{1 + e^{-x}}$$

- 这个函数也具备可求导的属性

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

示例

- 一个训练元组 $X=(1,0,1)$, 标签为1



J. Han, M. Kamber, J. Pei著; 范明等译. 数据挖掘: 概念与技术 (第3版). 机械工业出版社, 2012.

示例

Initial input, weight, and bias values.

x_1	x_2	x_3	w_{14}	w_{15}	w_{24}	w_{25}	w_{34}	w_{35}	w_{46}	w_{56}	θ_4	θ_5	θ_6
1	0	1	0.2	-0.3	0.4	0.1	-0.5	0.2	-0.3	-0.2	-0.4	0.2	0.1

The net input and output calculations.

Unit j	Net input, I_j	Output, O_j
4	$0.2 + 0 - 0.5 - 0.4 = -0.7$	$1/(1 + e^{0.7}) = 0.332$
5	$-0.3 + 0 + 0.2 + 0.2 = 0.1$	$1/(1 + e^{-0.1}) = 0.525$
6	$(-0.3)(0.332) - (0.2)(0.525) + 0.1 = -0.105$	$1/(1 + e^{0.105}) = 0.474$

示例

Calculation of the error at each node.

<i>Unit j</i>	<i>Err_j</i>
6	$(0.474)(1 - 0.474)(1 - 0.474) = 0.1311$
5	$(0.525)(1 - 0.525)(0.1311)(-0.2) = -0.0065$
4	$(0.332)(1 - 0.332)(0.1311)(-0.3) = -0.0087$

Calculations for weight and bias updating.

<i>Weight or bias</i>	<i>New value</i>
w_{46}	$-0.3 + (0.9)(0.1311)(0.332) = -0.261$
w_{56}	$-0.2 + (0.9)(0.1311)(0.525) = -0.138$
w_{14}	$0.2 + (0.9)(-0.0087)(1) = 0.192$
w_{15}	$-0.3 + (0.9)(-0.0065)(1) = -0.306$
w_{24}	$0.4 + (0.9)(-0.0087)(0) = 0.4$
w_{25}	$0.1 + (0.9)(-0.0065)(0) = 0.1$
w_{34}	$-0.5 + (0.9)(-0.0087)(1) = -0.508$
w_{35}	$0.2 + (0.9)(-0.0065)(1) = 0.194$
θ_6	$0.1 + (0.9)(0.1311) = 0.218$
θ_5	$0.2 + (0.9)(-0.0065) = 0.194$
θ_4	$-0.4 + (0.9)(-0.0087) = -0.408$

神经网络模型

- 给定隐藏层或输出层的单元 j , 单位 j 的净输入 I_j 为:

$$I_j = \sum_i w_{ij} O_i + \theta_j$$

w_{ij} 是从上一层单元 i 到单元 j 的连接权重; O_i 是上一层单元 i 的输出;
 θ_j 单元的偏置

- 给定单元 j 的输入 I_j , 则单元 j 的输出 O_j 的公式为:

$$O_j = \frac{1}{1 + e^{-I_j}}$$

- 对于输出层中的单元 k , 误差 Err_k 由下式计算:

$$Err_k = O_k (1 - O_k) (T_k - O_k)$$

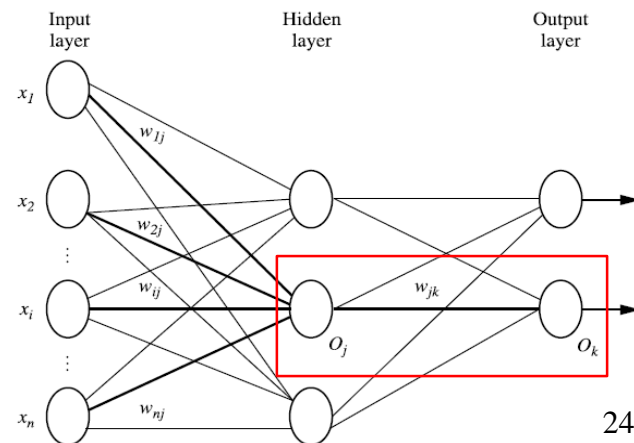
- 隐藏层单元 j 的误差为:

$$Err_j = O_j (1 - O_j) \sum_k Err_k w_{jk}$$

- 权重的更新公式为

$$w_{jk} = w_{jk} + \eta Err_k O_j$$

$$\theta_k = \theta_k + \eta Err_k$$

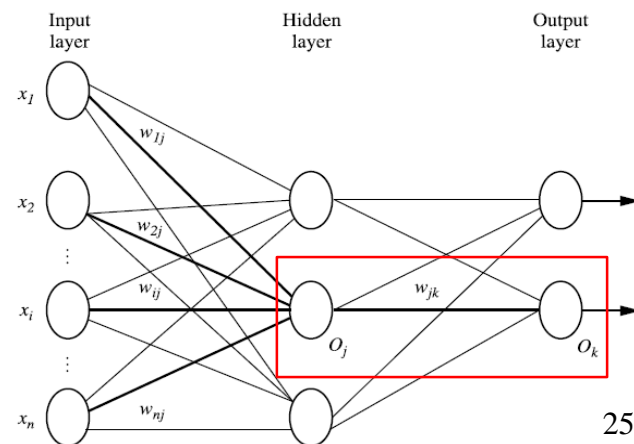


神经网络模型

- 最小化节点 O_k 的误差
- 我们将其定义为 $E = \frac{1}{2}e^2 = \frac{1}{2}(T - O)^2$
- 为了调整权重 w_{jk} 我们首先计算在 w_{jk} 上 E 的偏导数

$$\begin{aligned}\frac{\partial E}{\partial w_{jk}} &= \frac{\partial E}{\partial e} \times \frac{\partial e}{\partial O_k} \times \frac{\partial O_k}{\partial w_{jk}} \\ &= -(e) \times (O_k(1 - O_k)) \times (O_j) \\ &= -(T_k - O_k)O_k(1 - O_k)O_j\end{aligned}$$

- 然后我们使用梯度下降法



神经网络模型

- 反向传播学习算法是基于错误曲面的思想
- 表面代表了作为网络权重函数的数据集上的累积误差
- 曲面表示数据集上的累积误差每个可能的网络权重配置由表面上的点表示为网络权重的函数

神经网络模型

- 学习算法的目标是确定一组最小化错误的权重
- 学习算法应该被设计成能够找到表面上最快速减少错误的方向
- 这可以通过在每个表面点处的梯度矢量的相反方向上移动来实现（即，通过采用梯度下降学习方法）

神经网络模型

- 优点
 - 高鲁棒性
 - 非常适合于连续值输入和输出
 - 在广泛的现实数据上取得成功

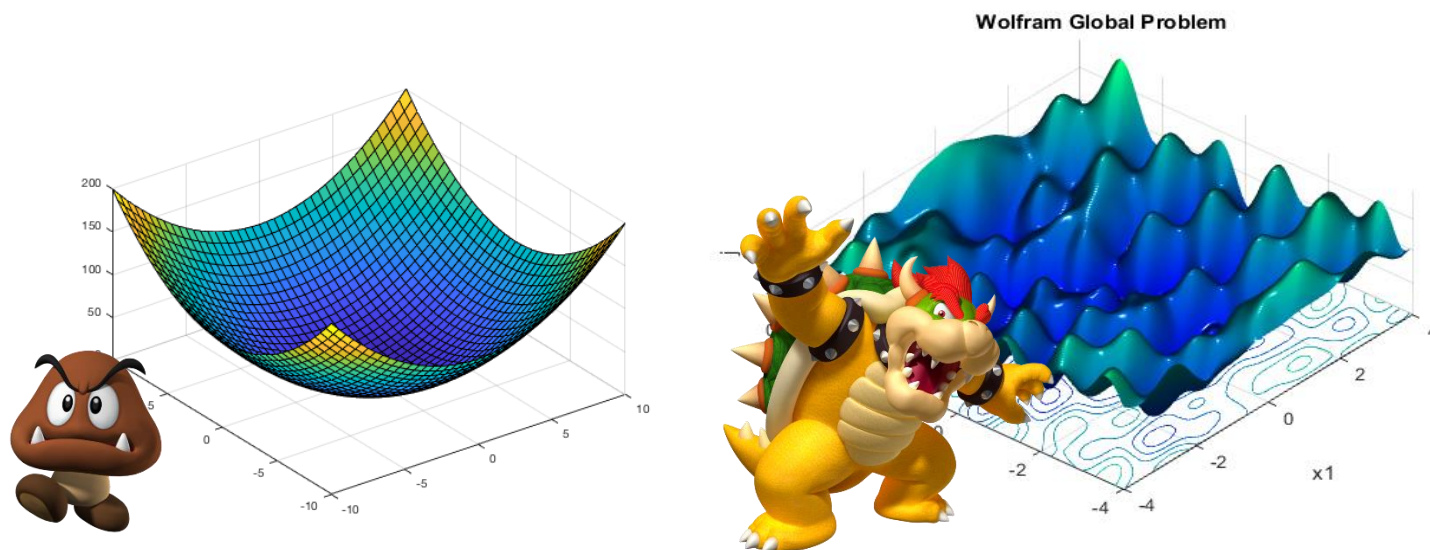
神经网络模型

- 缺点

- 训练时间长：尽管使用了BP算法，一次多层神经网络的训练仍然耗时太久
- 局部最优解和梯度消失问题，这使得神经网络的优化较为困难
- 通常需要经验最佳确定的多个参数，例如网络拓扑或“结构”
- 可解释性差：很难解释网络中学习到的权重和“隐藏层神经元”背后的象征意义

非凸模型优化问题

- 非凸模型优化的一个问题就是局部最优解问题，这使得神经网络的优化较为困难。



- 优化效率慢也限制了感知器规模。

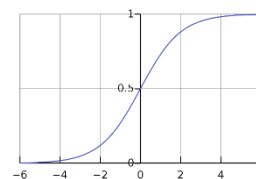
梯度消失问题

- 在深度网络中，不同的层的学习速度差异很大。使用sigmoid作为神经元的输入输出函数，在BP反向传播梯度时，每传递一层梯度衰减为原来的0.25。层数一多，梯度指数衰减后低层基本上接受不到有效的训练信号。

$$\frac{\partial E_{\text{total}}}{\partial w_1} = \frac{\partial E_{\text{total}}}{\partial \text{out}_{h1}} * \frac{\partial \text{out}_{h1}}{\partial \text{net}_{h1}} * \frac{\partial \text{net}_{h1}}{\partial w_1}$$

$$\frac{\partial E_{\text{total}}}{\partial \text{out}_{h1}} = \frac{\partial E_{o1}}{\partial \text{out}_{h1}} + \frac{\partial E_{o2}}{\partial \text{out}_{h1}}$$

$$\frac{\partial E_{o1}}{\partial \text{out}_{h1}} = \frac{\partial E_{o1}}{\partial \text{out}_{o1}} * \frac{\partial \text{out}_{o1}}{\partial \text{net}_{o1}} * \frac{\partial \text{net}_{o1}}{\partial \text{out}_{h1}}$$



the sigmoid function as $\sigma(x) = \frac{1}{1 + e^{-x}}$.

$$0 \leq \frac{d}{dx} \sigma(x) = \sigma(x)(1 - \sigma(x)) \leq 0.25$$

- 梯度消失，使得多层网络等价于浅层网络。

冰河期

- 90 年代中期，由 Vapnik 等人发明的 SVM（Support Vector Machine，支持向量机）算法诞生，很快就在若干个方面体现出了优势：无需调参；高效；全局最优解。基于以上种种理由，SVM 迅速打败了神经网络算法成为主流。



Vladimir Vapnik



- 神经网络的研究再次陷入了冰河期。当时，只要你的论文中包含神经网络相关的字眼，很可能被拒收，研究界那时对神经网络的不待见可想而知。

深度学习

- 在被人摒弃的10年中，有几个学者仍然在坚持研究，包括加拿大多伦多大学的Geoffery Hinton教授
- 2006年，Hinton在《Science》和相关期刊上发表了论文，首次提出了“深度信念网络”的概念
- 与传统的训练方式不同，“深度信念网络”有一个“预训练”（pre-training）的过程，这可以方便的让神经网络中的权值找到一个接近最优解的值，之后再使用“微调”（fine-tuning）技术来对整个网络进行优化训练。这两个技术的运用大幅度减少了训练多层神经网络的时间。他给多层神经网络相关的学习方法赋予了一个新名词--“深度学习”（深度学习是多个算法的总称）

深度学习

- 2006年Hinton提出的逐层预训练方法
- pre-training + fine-tuning...
- 使用Rectified Linear Unit (Relu)激活函数
- sigmoid等函数涉及指数运算，计算量大，反向传播求误差梯度时，求导涉及除法，计算量相对大，而采用Relu激活函数，整个过程的计算量节省很多
- sigmoid函数反向传播时，容易出现梯度消失
- Relu会使一部分神经元的输出为0，这样就造成了网络的稀疏性，并且减少了参数的相互依存关系，缓解了过拟合问题的发生
- Dropout: 模型训练时随机让网络某些隐含层节点的权重不工作，防止过拟合

深度学习

- 参数越多的模型复杂度越高
 - 优点：“容量” (capacity) 越大，能完成更复杂的任务
 - 缺点：训练效率低，数据量少时易陷入过拟合
- 缺点解决方案
 - 云计算和大数据时代下，计算能力大幅提升，可缓解训练低效性
 - 训练数据大幅增加，降低过拟合风险
- 以深度学习为代表的复杂模型开始受人关注

深度学习

- 很快，深度学习在语音识别领域暂露头角。2012年，Hinton与他的学生在ImageNet竞赛中，用多层的卷积神经网络成功地对包含一千类别的一百万张图片进行了训练，取得了分类错误率15%的好成绩，这个成绩比第二名高了近11个百分点，充分证明了多层神经网络识别效果的优越性。
- 此后，关于深度神经网络的研究与应用不断涌现。



Geoffery Hinton
1986, BP算法



Geoffery Hinton
图灵奖获得者