

## 第五章 前馈神经网络

在人工智能领域，**人工神经网络**（Artificial Neural Network，ANN）是指一系列受生物学和神经学启发的数学模型。这些模型主要是通过对人脑的神经网络进行抽象，构建人工神经元，并按照一定拓扑结构来建立人工神经元之间的连接，来模拟生物神经网络。在人工智能领域，人工神经网络也常常简称为神经网络（Neural Network，NN）或神经模型。

到目前为止，研究者已经提出了多种多样的神经网络模型，比如：感知器、前馈网络、卷积网络、循环网络、自组织映射、Hopfield网络、Boltzmann机等。这些神经网络模型的差异主要在于神经元的激活规则、网络连接的拓扑结构以及参数的学习规则等。

### 5.1 人工神经网络

人工神经网络主要由大量的神经元以及它们之间的有向连接构成。因此考虑三方面：

- **神经元的激活规则**：主要是指神经元输入到输出之间的映射关系，一般为非线性函数。
- **网络的拓扑结构**：不同神经元之间的连接关系。神经元间的连接也称为权重，即需要学习的参数。
- **学习算法**：通过训练数据来学习神经网络的参数。

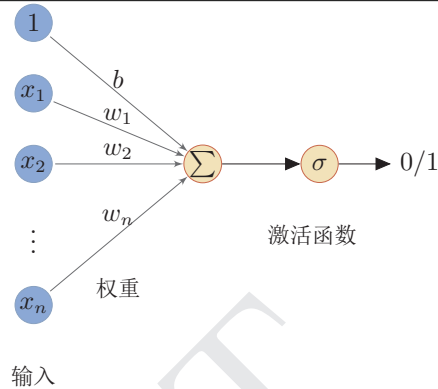


图 5.1: 人工神经元模型

### 5.1.1 神经元

人工神经元 (Artificial Neuron) 是构成人工神经网络的基本单元。人工神经元和感知器非常类似，也是模拟生物神经元特性，接受一组输入信号并产生输出。生物神经元有一个阈值，当神经元所获得的输入信号的积累效果超过阈值时，它就处于兴奋状态；否则，应该处于抑制状态。

净输入也叫净活性值 (net activation)。

我们用  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  来表示人工神经元的一组输入，用净输入  $z$  表示一个神经元所获得的输入信号  $x$  的加权和，

$$z = \sum_{i=1}^n w_i x_i + b \quad (5.1)$$

$$= \mathbf{w}^T \mathbf{x} + b, \quad (5.2)$$

其中， $\mathbf{w} = (w_1, w_2, \dots, w_n)$  是  $n$  维的权重向量， $b$  是偏置。

净输入  $z$  在经过一个非线性函数后，人工神经元输出它的活性值 (Activation)  $a$ ，

$$a = f(z), \quad (5.3)$$

这里的非线性函数  $f$  也称为激活函数 (Activation Function)，有 sigmoid 型函数、非线性斜面函数等。

人工神经元的结构如图 5.1 所示。如果我们设激活函数  $f$  为 0 或 1 的阶跃函数，人工神经元就是感知器。

### 数学小知识 | 饱和

对于函数  $f(x)$ ，若  $x \rightarrow -\infty$  时，其导数  $f'(x) \rightarrow 0$ ，则称其为**左饱和**。若  $x \rightarrow +\infty$  时，其导数  $f'(x) \rightarrow 0$ ，则称其为**右饱和**。当同时满足左、右饱和时，就称其为**饱和**。

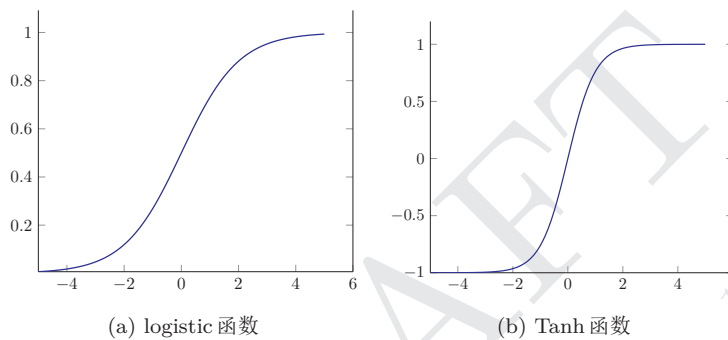


图 5.2: Sigmoid 型激活函数

### 5.1.2 激活函数

为了增强网络的表达能力，我们需要使用连续非线性**激活函数**（Activation Function）。因为连续非线性激活函数可导，所以可以用最优化的方法来求解。

下面介绍几个在神经网络中常用的激活函数。

#### Logistic 函数

Logistic 函数是一种 sigmoid 型函数。其定义为  $\sigma(x)$

$$\sigma(x) = \frac{1}{1 + e^{-x}}. \quad (5.4)$$

图5.3a给出了 logistic 函数的形状。Logistic 函数可以看成是一个“挤压”函数，把一个实数域的输入“挤压”到  $(0, 1)$ 。当输入值在 0 附近时，sigmoid 型函数近似为线性函数；当输入值靠近两端时，对输入进行抑制。输入越小，越接

Sigmoid 型函数是指一类 S 型曲线函数，常用的 sigmoid 型函数有 logistic 函数和 tanh 函数。

近于0；输入越大，越接近于1。这样的特点也和生物神经元类似，对一些输入会产生兴奋（输出为1），对另一些输入产生抑制（输出为0）。和感知器使用的阶跃激活函数相比，logistic函数是连续可导的，其数学性质更好。

**Tanh 函数** Tanh函数也是一种sigmoid型函数。其定义为

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (5.5)$$

Tanh函数可以看作是放大并平移的logistic函数：

$$\tanh(x) = 2\sigma(2x) - 1. \quad (5.6)$$

图5.3b给出了tanh函数的形状，其值域是 $(-1, 1)$ 。

### Hard-Logistic 和 Hard-Tanh 函数

Logistic函数和tanh函数都是sigmoid型函数，具有饱和性，但是计算开销较大。因为这两个函数都是在中间（0附近）近似线性，两端饱和。因此，这两个函数可以通过分段函数来近似。

以logistic函数 $\sigma(x)$ 为例，其导数为 $\sigma'(x) = \sigma(x)(1 - \sigma(x))$ 。logistic函数在0附近的一阶泰勒展开（Taylor expansion）为

$$g_l(x) \approx \sigma(0) + \sigma'(0) \quad (5.7)$$

$$= 0.25x + 0.5. \quad (5.8)$$

$$\text{hard-logistic}(x) = \begin{cases} 1 & g_l(x) \geq 1 \\ g_l & 0 < g_l(x) < 1 \\ 0 & g_l(x) \leq 0 \end{cases} \quad (5.9)$$

$$= \max(\min(g_l(x), 1), 0) \quad (5.10)$$

$$= \max(\min(0.25x + 0.5, 1), 0). \quad (5.11)$$

同样，tanh函数在0附近的一阶泰勒展开为

$$g_t(x) \approx \tanh(0) + \tanh'(0) \quad (5.12)$$

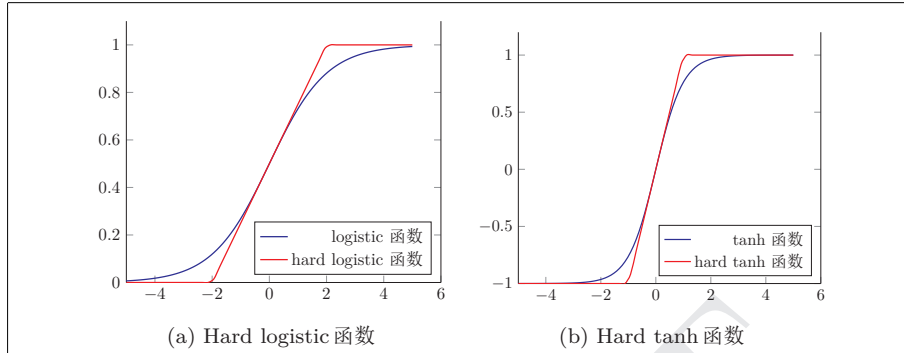


图 5.3: Sigmoid 型激活函数近似

$$= x. \quad (5.13)$$

这样,  $\tanh$  函数也可以用分段函数  $\text{hard-tanh}(x)$  来近似。

$$\text{hard-tanh}(x) = \max(\min(g_t(x), 1), -1) \quad (5.14)$$

$$= \max(\min(x, 1), -1). \quad (5.15)$$

图5.3给出了  $\text{hard-logistic}$  和  $\text{hard-tanh}$  函数两种函数的形状。

### 修正线性单元

**修正线性单元** (rectified linear unit, ReLU) [Nair and Hinton, 2010], 也叫 **rectifier 函数** [Glorot et al., 2011], 是深层神经网络中经常使用的激活函数。ReLU 实际上是一个斜坡 (ramp) 函数, 定义为

$$\text{rectifier}(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases} \quad (5.16)$$

$$= \max(0, x). \quad (5.17)$$

采用 ReLU 的神经网络只需要进行加、乘和比较的操作, 计算上也更加高效。此外, rectifier 函数被认为有生物上的解释性。神经科学家发现神经元具有**单侧抑制**、**宽兴奋边界**、**稀疏激活性**等特性。

需要参见第82页的  
第5.3节内容

**稀疏激活性** 生物神经元只对少数输入信号选择性响应，大量信号被屏蔽了。因为随机初始化的原因，sigmoid 系函数同时近乎有 50% 的神经元被激活，这不符合神经科学的发现。而 rectifier 函数却具备很好的稀疏性。

Rectifier 函数为左饱和函数，在  $x > 0$  时导数为 1，在  $x < 0$  时导数为 0。这样在训练时，如果学习率设置过大，在一次更新参数后，一个采用 ReLU 的神经元在所有的训练数据上都不能被激活。那么，这个神经元在以后的训练过程中永远不能被激活，这个神经元的梯度就永远都会是 0。

在实际使用中，为了避免上述情况，ReLU 有几种被广泛使用的变种。

**带泄露的 ReLU** 带泄露的 ReLU (Leaky ReLU) 在输入  $x < 0$  时，保持一个很小的梯度  $\lambda$ 。这样当神经元非激活时也能有一个非零的梯度可以更新参数，避免永远不能被激活 [Maas et al., 2013]。带泄露的 ReLU 的定义如下：

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \lambda x & \text{if } x \leq 0 \end{cases} \quad (5.18)$$

$$= \max(0, x) + \lambda \min(0, x), \quad (5.19)$$

这里， $\lambda \in (0, 1)$  是一个很小的常数，比如 0.01。

**带参数的 ReLU** 带参数的 ReLU (Parametric ReLU, PReLU) 引入一个可学习的参数，不同神经元可以有不同的参数 [He et al., 2015]。对于第  $i$  个神经元，其 PReLU 的定义为

$$\text{PReLU}_i(x) = \begin{cases} x & \text{if } x > 0 \\ a_i x & \text{if } x \leq 0 \end{cases} \quad (5.20)$$

$$= \max(0, x) + a_i \min(0, x), \quad (5.21)$$

其中， $a_i$  为  $x \leq 0$  时函数的斜率。因此，PReLU 是非饱和函数。如果  $a_i = 0$ ，那么 PReLU 就退化为 ReLU。如果  $a_i$  为一个很小的常数，则 PReLU 可以看作带泄露的 ReLU。

这里，不同神经元对应的参数都可以是不同的，也可以一组神经元共享一个参数。

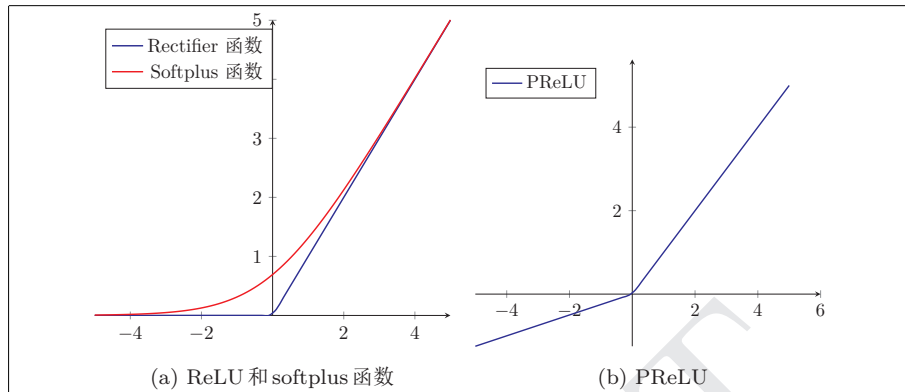


图 5.4: ReLU、PReLU 和 softplus 函数

### Softplus 函数

Softplus 函数 [Dugas et al., 2001] 可以看作是 rectifier 函数的平滑版本，其定义为

$$\text{softplus}(x) = \log(1 + e^x) \quad (5.22)$$

softplus 函数其导数刚好是 logistic 函数。softplus 虽然也有具有单侧抑制、宽兴奋边界的特性，却没有稀疏激活性。

图5.4给出了 ReLU、PReLU 以及 softplus 函数的示例。

**Maxout 单元** Maxout 单元 [Goodfellow et al., 2013] 也是一种分段线性函数。Sigmoid 型函数、ReLU 等激活函数的输入是神经元的净输入  $z$ ，是一个标量。而 maxout 单元的输入不是经过加权叠加后的净输入，而是神经元的全部原始输入，是一个向量  $\mathbf{x} = [x_1, x_2, \dots, x_n]^\top$ 。

每个 maxout 单元有  $k$  个权重向量  $\mathbf{w}_i \in \mathbb{R}^n$  和偏置  $b_i$  ( $i \in [1, k]$ )。对于输入  $\mathbf{x}$ ，可以得到  $k$  个净输入  $z$

$$z_i = \mathbf{w}_i \mathbf{x} + b_i \quad (5.23)$$

$$= \sum_{j=1}^n w_{i,j} x_j + b_i, \quad (5.24)$$

其中， $\mathbf{w}_i = [w_{i,1}, \dots, w_{i,n}]^\top$  为第  $i$  个权重向量。

激活函数	函数	导数
Logistic 函数	$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
Tanh 函数	$f(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
ReLU	$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
SoftPlus 函数	$f(x) = \ln(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$

表 5.1: 常见激活函数及其导数

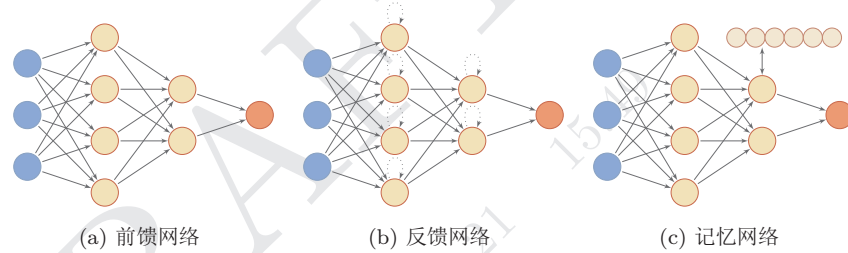


图 5.5: 三种不同的网络模型

Maxout 单元的非线性函数定义为

$$\text{maxout}(\mathbf{x}) = \max_{i \in [1, k]} (z_i). \quad (5.25)$$

Maxout 单元不单是净输入到输出之间的非线性映射，而是整体学习输入到输出之间的非线性映射关系。Maxout 激活函数可以看作任意凸函数的分段线性近似，并且在有限的点上是不可微的。采用 maxout 单元的神经网络也就做 maxout 网络。

### 5.1.3 网络结构

一个生物神经细胞的功能比较简单，而人工神经元只是生物神经细胞的理想化和简单实现，功能更加简单。要想模拟人脑的能力，单一的神经元是远远不够的，需要通过很多神经元一起协作来完成复杂的功能。这样通过一定的连



接方式或信息传递方式进行协作的神经元可以看作是一个网络，叫做**人工神经网络**，简称**神经网络**。

到目前为止，研究者已经发明了各种各样的神经网络结构。目前常用的神经网络结构有以下三种：

**前馈网络** 网络中各个神经元按接受信息的先后分为不同的组。每一组可以看作一个神经层。每一层中的神经元接受前一层神经元的输出，并输出到下一层神经元。整个网络中的信息是朝一个方向传播，没有反向的信息传播。前馈网络可以用一个有向无环路图表示。前馈网络可以看作一个**函数**，通过简单非线性函数的多次复合，实现输入空间到输出空间的复杂映射。这种网络结构简单，易于实现。前馈网络包括全连接前馈网络和卷积神经网络等。

**反馈网络** 网络中神经元不但可以接收其它神经元的信号，也可以接收自己的反馈信号。和前馈网络相比，反馈网络在不同的时刻具有不同的状态，具有记忆功能，因此反馈网络可以看作一个**程序**，也具有更强的计算能力。反馈神经网络可用一个完备的无向图来表示。

**记忆网络** 记忆网络在前馈网络或反馈网络的基础上，引入一组记忆单元，用来保存中间状态。同时，根据一定的取址、读写机制，来增强网络能力。和反馈网络相比，反馈网络具有更强的记忆功能。

图5.5给出了前馈网络、反馈网络和记忆网络的网络结构示例。

## 5.2 前馈神经网络

给定一组神经元，我们可以以神经元为节点来构建一个网络。不同的神经网络模型有着不同网络连接的拓扑结构。一种比较直接的拓扑结构是前馈网络。**前馈神经网络**（Feedforward Neural Network）是最早发明的简单人工神经网络。

在前馈神经网络中，各神经元分别属于不同的层。每一层的神经元可以接收前一层神经元的信号，并产生信号输出到下一层。第一层叫**输入层**，最后一层叫**输出层**，其它中间层叫做**隐藏层**。整个网络中无反馈，信号从输入层向输出层单向传播，可用一个有向无环图表示。

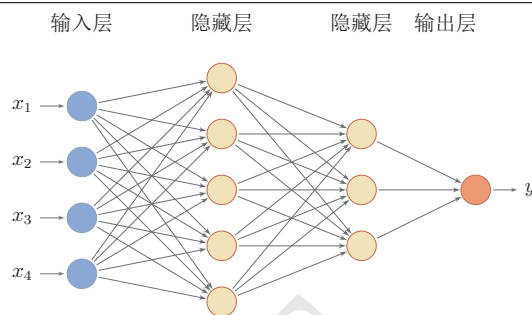


图 5.6: 多层神经网络

前馈神经网络也经常称为多层感知器 (Multilayer Perceptron, MLP)。但多层感知器的叫法并不是十分合理, 因为前馈神经网络其实是由多层的 logistic 回归模型 (连续的非线性函数) 组成, 而不是由多层的感知器 (不连续的非线性函数) 组成 [Bishop, 2006]。

图5.6给出了前馈神经网络的示例。

### 5.2.1 前馈计算

给定一个前馈神经网络, 我们用下面的记号来描述这样网络。

- $L$ : 表示神经网络的层数;
- $n^l$ : 表示第  $l$  层神经元的个数;
- $f_l(\cdot)$ : 表示  $l$  层神经元的激活函数;
- $W^{(l)} \in \mathbb{R}^{n^l \times n^{l-1}}$ : 表示  $l-1$  层到第  $l$  层的权重矩阵;
- $\mathbf{b}^{(l)} \in \mathbb{R}^{n^l}$ : 表示  $l-1$  层到第  $l$  层的偏置;
- $\mathbf{z}^{(l)} \in \mathbb{R}^{n^l}$ : 表示  $l$  层神经元的净输入 (净活性值);
- $\mathbf{a}^{(l)} \in \mathbb{R}^{n^l}$ : 表示  $l$  层神经元的输出 (活性值)。

前馈神经网络通过下面公式进行信息传播。

$$\mathbf{z}^{(l)} = W^{(l)} \cdot \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)} \quad (5.26)$$

$$\mathbf{a}^{(l)} = f_l(\mathbf{z}^{(l)}) \quad (5.27)$$

公式5.26和5.27也可以合并写为：

$$\mathbf{z}^{(l)} = W^{(l)} \cdot f_{l-1}(\mathbf{z}^{(l-1)}) + \mathbf{b}^{(l)} \quad (5.28)$$

或者

$$\mathbf{a}^{(l)} = f_l(W^{(l)} \cdot \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)}). \quad (5.29)$$

这样，前馈神经网络可以通过逐层的信息传递，得到网络最后的输出  $\mathbf{a}^L$ 。整个网络可以看作一个复合函数  $f(\mathbf{x}; W, \mathbf{b})$ ，将输入  $x$  作为第1层的输入  $\mathbf{a}^{(0)}$ ，将第  $L$  层的输出  $\mathbf{a}^{(L)}$  作为整个函数的输出。

$$\mathbf{x} = \mathbf{a}^{(0)} \rightarrow \mathbf{z}^{(1)} \rightarrow \mathbf{a}^{(1)} \rightarrow \mathbf{z}^{(2)} \rightarrow \dots \rightarrow \mathbf{a}^{(L-1)} \rightarrow \mathbf{z}^{(L)} \rightarrow \mathbf{a}^{(L)} = f(\mathbf{x}; W, \mathbf{b}). \quad (5.30)$$

### 5.2.2 应用到机器学习

在机器学习中，输入样本的特征对分类器的影响很大。以监督学习为例，好的特征可以极大提高分类器的性能。因此，要取得好的分类效果，需要样本的原始特征向量  $\mathbf{x}$  转换到更有效的特征向量  $\phi(\mathbf{x})$ ，这个过程叫做特征抽取或特征转换。

多层前馈神经网络可以看作是一个非线性复合函数  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ ，将输入  $\mathbf{x} \in \mathbb{R}^n$  映射到输出  $\mathbf{r} \in \mathbb{R}^m$ 。因此，多层前馈神经网络也可以看成是一种特征转换方法，其输出  $\mathbf{r}$  作为分类器的输入进行分类。

给定一个训练样本  $(\mathbf{x}, \mathbf{y})$ ，先利用多层前馈神经网络将  $\mathbf{x}$  映射到  $\mathbf{r}$ ，然后再将  $\mathbf{r}$  输入到分类器  $g(\cdot)$ 。

$$\hat{\mathbf{y}} = g(\mathbf{r}, \theta) \quad (5.31)$$

$$= g(f(\mathbf{x}; W, \mathbf{b}), \theta), \quad (5.32)$$

这里， $g(\cdot)$  为线性或非线性的分类器， $\theta$  为分类器  $g(\cdot)$  的参数， $\hat{\mathbf{y}}$  为分类器的输出。

特别地，如果分类器  $g(\cdot)$  为 logistic 回归或 softmax 回归，那么  $g(\cdot)$  也可以看成是网络的最后一层，即神经网络直接输出不同类别的后验概率。

Logistic 回归 参考第 3.2.2 节。

Softmax 回 归 参 考第 3.2.3 节。

对于两类分类问题  $y \in \{0, 1\}$ , logistic 回归分类器可以看成神经网络的最后一层。也就是说, 网络的最后一层只用一个神经元, 并且其激活函数为 logistic 函数。网络的输出可以直接可以作为两个类别的后验概率。

$$P(y = 1|\mathbf{x}) = f(\mathbf{x}; W, \mathbf{b}) \quad (5.33)$$

其中,  $f(\mathbf{x}; W, \mathbf{b}) \in \mathbb{R}$ 。

对于多类分类问题  $y \in \{1, \dots, C\}$ , 如果使用 softmax 回归分类器, 相当于网络最后一层设置  $C$  个神经元, 其输出经过 softmax 函数进行归一化后可以作为每个类的后验概率。

$$\hat{\mathbf{y}} = \text{softmax}(f(\mathbf{x}; W, \mathbf{b})), \quad (5.34)$$

其中,  $f(\mathbf{x}; W, \mathbf{b}) \in \mathbb{R}^C$  为神经网络输出,  $\hat{\mathbf{y}} \in \mathbb{R}^C$  为预测的不同类别后验概率组成的向量。

如果采用交叉熵损失函数, 对于样本  $(\mathbf{x}, \mathbf{y})$ , 其损失函数为

$$\mathcal{L}(y, f(\mathbf{x}; W, \mathbf{b})) = -\mathbf{y}^\top \log \hat{\mathbf{y}} \quad (5.35)$$

这样, 给定一组训练样本, 网络的参数就可以通过梯度下降法来进行学习。梯度下降法需要计算损失函数对参数的导数, 如果通过链式法则逐一进行求导效率比较低。因此在神经网络的训练中经常使用反向传播算法来计算梯度。

## 5.3 反向传播算法

给定一组样本  $(\mathbf{x}^{(i)}, y^{(i)}), 1 \leq i \leq N$ , 用前馈神经网络的输出为  $f(\mathbf{x}|\mathbf{w}, \mathbf{b})$ , 目标函数为:

$$\mathcal{R}(W, \mathbf{b}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y^{(i)}, f(\mathbf{x}^{(i)}; W, \mathbf{b})) + \frac{1}{2} \lambda \|W\|_F^2, \quad (5.36)$$

$$= \frac{1}{N} \sum_{i=1}^N \mathcal{L}(W, \mathbf{b}; \mathbf{x}^{(i)}, y^{(i)}) + \frac{1}{2} \lambda \|W\|_F^2, \quad (5.37)$$

这里,  $W$  和  $\mathbf{b}$  包含了每一层的权重矩阵和偏置向量;  $\|W\|_F^2$  是正则化项, 用来防止过拟合;  $\lambda$  是为正数的超参。  $\lambda$  越大,  $W$  越接近于 0。这里的  $\|W\|_F^2$  一般使

注意这里的正则化项只包含权重参数  $W$ , 而不

用 Frobenius 范数：

$$\|W\|_F^2 = \sum_{l=1}^L \sum_{i=1}^{n^l} \sum_{j=1}^{n^{l-1}} (W_{ij}^{(l)})^2. \quad (5.38)$$

我们的目标是最小化  $\mathcal{R}(W, \mathbf{b}; \mathbf{x}, y)$ 。如果采用梯度下降方法，我们可以用如下方法更新参数：

$$W^{(l)} = W^{(l)} - \alpha \frac{\partial \mathcal{R}(W, \mathbf{b})}{\partial W^{(l)}}, \quad (5.39)$$

$$= W^{(l)} - \alpha \left( \frac{1}{N} \sum_{i=1}^N \left( \frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}^{(i)}, y^{(i)})}{\partial W^{(l)}} \right) - \lambda W^{(l)} \right), \quad (5.40)$$

$$\mathbf{b}^{(l)} = \mathbf{b}^{(l)} - \alpha \frac{\partial \mathcal{R}(W, \mathbf{b})}{\partial \mathbf{b}^{(l)}}, \quad (5.41)$$

$$= \mathbf{b}^{(l)} - \alpha \left( \frac{1}{N} \sum_{i=1}^N \frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}^{(i)}, y^{(i)})}{\partial \mathbf{b}^{(l)}} \right), \quad (5.42)$$

$$(5.43)$$

这里  $\alpha$  是参数的更新率。

我们首先来看下  $\frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}, y)}{\partial W^{(l)}}$  怎么计算。

根据链式法则， $\frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}, y)}{\partial W_{ij}^{(l)}}$  可以写为

$$\frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}, y)}{\partial W_{ij}^{(l)}} = \left( \frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}, y)}{\partial \mathbf{z}^{(l)}} \right)^\top \frac{\partial \mathbf{z}^{(l)}}{\partial W_{ij}^{(l)}}. \quad (5.44)$$

公式5.44的第一项是为目标函数关于第  $l$  层的神经元  $\mathbf{z}^{(l)}$  的偏导数，我们称为误差项  $\delta^{(l)}$ ：

$$\delta^{(l)} = \frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}, y)}{\partial \mathbf{z}^{(l)}} \in \mathbb{R}^{n^{(l)}}. \quad (5.45)$$

误差项  $\delta^{(l)}$  来表示第  $l$  层的神经元对最终误差的影响，也反映了最终的输出对第  $l$  层的神经元对最终误差的敏感程度。

公式5.44中的第二项是  $l$  层的神经元  $\mathbf{z}^{(l)}$  关于参数  $W_{ij}^{(l)}$  的偏导数。

我们分别来计算这两个偏导数。

链式法则参见第12页的公式2.26

计算误差项  $\delta^{(l)}$  我们先来看下第  $l$  层的误差项  $\delta^{(l)} = \frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}, y)}{\partial \mathbf{z}^{(l)}}$  怎么计算。

$$\delta^{(l)} \triangleq \frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}, y)}{\partial \mathbf{z}^{(l)}} \quad (5.46)$$

$$= \frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{z}^{(l)}} \cdot \frac{\partial \mathbf{z}^{(l+1)}}{\partial \mathbf{a}^{(l)}} \cdot \frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}, y)}{\partial \mathbf{z}^{(l+1)}} \quad (5.47)$$

$$= \text{diag}(f'_l(\mathbf{z}^{(l)})) \cdot (W^{(l+1)})^\top \cdot \delta^{(l+1)} \quad (5.48)$$

$$= f'_l(\mathbf{z}^{(l)}) \odot ((W^{(l+1)})^\top \delta^{(l+1)}), \quad (5.49)$$

其中  $\odot$  是向量的点积运算符，表示每个元素相乘。

上述推导中，关键是公式5.47到公式5.48的推导。公式5.47中有三项，第三项根据定义为  $\delta^{(l+1)}$ 。

第二项因为  $\mathbf{z}^{(l+1)} = W^{(l+1)} \cdot \mathbf{a}^{(l)} + \mathbf{b}^{(l)}$ ，所以

$$\frac{\partial \mathbf{z}^{(l+1)}}{\partial \mathbf{a}^{(l)}} = (W^{(l+1)})^\top. \quad (5.50)$$

第一项因为  $\mathbf{a}^{(l)} = f_l(\mathbf{z}^{(l)})$ ，而  $f_l(\cdot)$  为按位计算的函数。因此

$$\frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{z}^{(l)}} = \frac{\partial f_l(\mathbf{z}^{(l)})}{\partial \mathbf{z}^{(l)}} \quad (5.51)$$

$$= \text{diag}(f'_l(\mathbf{z}^{(l)})). \quad (5.52)$$

从公式5.49可以看出，第  $l$  层的误差项可以通过第  $l+1$  层的误差项计算得到。这就是误差的**反向传播**（Backpropagation, BP）。反向传播算法的含义是：第  $l$  层的一个神经元的误差项（或敏感性）是所有与该神经元相连的第  $l+1$  层的神经元的误差项的权重和。然后，再乘上该神经元激活函数的梯度。

计算  $\frac{\partial \mathbf{z}^{(l)}}{\partial W_{ij}^{(l)}}$  我们先来计算公式5.44中的第二项  $\frac{\partial \mathbf{z}^{(l)}}{\partial W_{ij}^{(l)}}$ 。

因为  $\mathbf{z}^{(l)} = W^{(l)} \cdot \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)}$ ，所以

$$\frac{\partial \mathbf{z}^{(l)}}{\partial W_{ij}^{(l)}} = \frac{\partial (W^{(l)} \cdot \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)})}{\partial W_{ij}^{(l)}} = \begin{bmatrix} 0 \\ \vdots \\ a_j^{(l-1)} \\ \vdots \\ 0 \end{bmatrix}. \quad (5.53)$$

在得到上面两个偏导数之后，公式5.44可以写为

$$\frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}, y)}{\partial W_{ij}^{(l)}} = (\delta^{(l)})^\top \mathbf{a}^{(l-1)} \quad (5.54)$$

$$= \delta_i^{(l)} a_j^{(l-1)}. \quad (5.55)$$

进一步， $\mathcal{L}(W, \mathbf{b}; \mathbf{x}, y)$  关于第  $l$  层权重  $W^{(l)}$  的梯度为

$$\frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}, y)}{\partial W^{(l)}} = \delta^{(l)} (\mathbf{a}^{(l-1)})^\top. \quad (5.56)$$

同理可得， $\mathcal{L}(W, \mathbf{b}; \mathbf{x}, y)$  关于第  $l$  层偏置  $\mathbf{b}^{(l)}$  的梯度为

$$\frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}, y)}{\partial \mathbf{b}^{(l)}} = \delta^{(l)}. \quad (5.57)$$

在计算出每一层的误差项之后，我们就可以得到每一层参数的梯度。因此，前馈神经网络的训练过程可以分为以下三步：

1. 前馈计算每一层的状态和激活值，直到最后一层；
2. 反向传播计算每一层的误差项；
3. 计算每一层参数的偏导数，并更新参数。

← 第  $i$  行

具体的训练过程如算法5.1所示，也叫反向传播算法。

#### 算法 5.1: 反向传播算法

输入: 训练集:  $(\mathbf{x}^{(i)}, y^{(i)}), i = 1, \dots, N$ , 最大迭代次数:  $T$   
 输出:  $W, \mathbf{b}$

```

1 初始化  $W, \mathbf{b}$ ;
2 for  $t = 1 \dots T$  do
3   for  $i = 1 \dots N$  do
4     (1) 前馈计算每一层的状态和激活值, 直到最后一层;
5     (2) 用公式5.49反向传播计算每一层的误差  $\delta^{(l)}$ ;
6     (3) 用公式5.56和5.57每一层参数的导数;
7        $\frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}^{(i)}, y^{(i)})}{\partial W^{(l)}} = \delta^{(l)} (\mathbf{a}^{(l-1)})^\top$ ;
8        $\frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}^{(i)}, y^{(i)})}{\partial \mathbf{b}^{(l)}} = \delta^{(l)}$ ;
9     (4) 更新参数;
10       $W^{(l)} = W^{(l)} - \alpha \sum_{i=1}^N \left( \frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}^{(i)}, y^{(i)})}{\partial W^{(l)}} \right) - \alpha \lambda W^{(l)}$ ;
11       $\mathbf{b}^{(l)} = \mathbf{b}^{(l)} - \alpha \sum_{i=1}^N \left( \frac{\partial \mathcal{L}(W, \mathbf{b}; \mathbf{x}^{(i)}, y^{(i)})}{\partial \mathbf{b}^{(l)}} \right)$ ;
12   end
13 end
```

## 5.4 自动梯度计算

神经网络的参数主要通过梯度下降来进行优化的。当确定了风险函数以及网络结构后,我们就可以手动用链式法则来计算风险函数对每个参数的梯度,并用代码进行实现。但是手动求导并转换为计算机程序的过程非常琐碎并容易出错,导致实现神经网络变得十分低效。目前,几乎所有的主流深度学习框架都包含了自动梯度计算的功能,即我们可以只考虑网络结构并用代码实现,其梯度可以自动进行计算,无需人工干预。这样开发的效率就大大提高了。

自动求导的方法可以分为以下三类:

**数值微分** 数值微分 (Numerical Differentiation) 是用数值方法来计算函数  $f(x)$  的导数。函数  $f(x)$  的点  $x$  的导数定义为

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}. \quad (5.58)$$



要计算函数  $f(x)$  在点  $x$  的导数，可以对  $x$  加上一个很少的非零的扰动  $\Delta x$ ，通过上述定义来直接计算函数  $f(x)$  的梯度。数值微分方法非常容易实现，但是找到一个合适的扰动  $\Delta x$  却十分困难。如果  $\Delta x$  过小，会引起数值计算问题，比如舍入误差；如果  $\Delta x$  过大，会增加截断误差，使得导数计算不准确。因此，数值微分的实用性比较差。在实际应用，经常使用下面公式来计算梯度，可以减少截断误差。

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x - \Delta x)}{2\Delta x}. \quad (5.59)$$

**符号微分** 符号微分 (Symbolic Differentiation) 是一种基于符号计算的自动求导方法。

**符号计算** (Symbolic Computation)，也叫**代数计算**，是指用计算机来处理带有变量的数学表达式。这里的变量看作是符号 (Symbols)，一般不需要代入具体的值。符号计算的输入和输出都是数学表达式，一般包括对数学表达式的化简、因式分解、微分、积分、解代数方程、求解常微分方程等运算。

比如数学表达式的化简：

$$\text{输入: } 3x - x + 2x + 1 \quad (5.60)$$

$$\text{输出: } 4x + 1. \quad (5.61)$$

符号计算一般来讲是对输入的表达式，通过不断使用一些事先定义的规则进行转换。当转换结果不能再继续使用变换规则时，便停止计算。

Theano [Bergstra et al., 2010] 和 Tensorflow [Abadi et al., 2016] 都采用了符号微分的方法进行自动求解梯度。符号微分可以在编译时就计算梯度的数学表示，并进一步利用符号计算方法进行优化。此外，符号计算的一个优点是符号计算和平台无关，可以在 CPU 或 GPU 上运行。

符号微分也有一些不足之处。一是编译时间较长，特别是对于循环，需要很长时间进行编译。二是为了进行符号微分，一般需要设计一种专门的语言来表示数学表达式，并且要对变量（符号）进行预先声明。三是程序很难调试。

**自动微分** 自动微分 (Automatic Differentiation) 是一种可以对一个（程序）函数进行计算导数的方法。符号微分的处理对象是数学表达式，而自动微分的处理对象是一个函数或一段程序。而自动微分可以直接在原始程序代码进行微分。自

**舍入误差** (Round-off Error) 是指数值计算中由于数字舍入造成的近似值和精确值之间的差异，比如用浮点数来表示实数。

**截断误差** (Truncation Error) 是数学模型的理论与数值计算问题的精确解之间的误差。

和符号计算相对应的概念是**数值计算**，即将数值代入数学表示中进行计算。

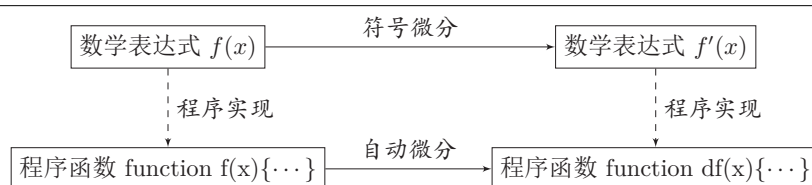


图 5.7: 符号微分与自动微分对比

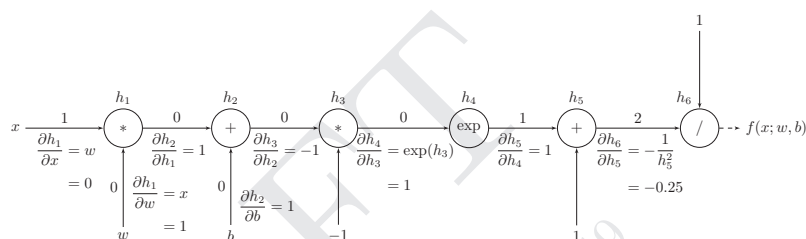


图 5.8: 复合函数  $f(x; w, b) = \sigma(wx + b)$  的计算图

自动微分的基本原理是所有的数值计算可以分解为一些基本操作，包含  $+$ ,  $-$ ,  $*$ ,  $/$  和一些初等函数  $\exp$ ,  $\log$ ,  $\sin$ ,  $\cos$  等。图5.7给出了符号微分与自动微分的关联。

自动微分也是利用链式法则来自动计算一个复合函数的梯度。以神经网络中常见的复合函数  $f(x; w, b) = \sigma(wx + b)$  为例，

$$f(x; w, b) = \sigma(wx + b) \quad (5.62)$$

$$= \frac{1}{\exp(-(wx + b)) + 1} \quad (5.63)$$

为了简单起见，我们设复合函数  $f(x; w, b)$  的输入  $x$  为标量，参数为权重  $w$  和偏置  $b$ 。我们可以将其分解为一系列的基本操作，并构成一个计算图（Computational Graph），其中每个节点为一个基本操作，如图5.8所示。每一步基本操作的导数都可以通过简单的规则来实现。

这样函数  $f(x; w, b)$  关于参数  $w$  和  $b$  的导数可以通过计算图上的节点  $f(x; w, b)$  与参数  $w$  和  $b$  之间路径上所有的导数连乘来得到，即

$$\frac{\partial f(x; w, b)}{\partial w} = \frac{\partial f(x; w, b)}{\partial h_6} \frac{\partial h_6}{\partial h_5} \frac{\partial h_5}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w}, \quad (5.64)$$

$$\frac{\partial f(x; w, b)}{\partial b} = \frac{\partial f(x; w, b)}{\partial h_6} \frac{\partial h_6}{\partial h_5} \frac{\partial h_5}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial b}. \quad (5.65)$$

计算图是数学运算的图形化表示。计算图中节点表示一个变量或基本操作。

以  $\frac{\partial f(x;w,b)}{\partial w}$  为例，当  $x = 1, w = 0, b = 0$  时，可以得到

$$\frac{\partial f(x;w,b)}{\partial w} \Big|_{x=1, w=0, b=0} = \frac{\partial f(x;w,b)}{\partial h_6} \frac{\partial h_6}{\partial h_5} \frac{\partial h_5}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w} \quad (5.66)$$

$$= 1 * -0.25 * 1 * 1 * -1 * 1 * 1 \quad (5.67)$$

$$= 0.25. \quad (5.68)$$

如果函数和参数之间有多条路径，可以将这多条路径上的导数再进相加，得到最终的梯度。

按照计算导数的顺序，自动微分可以分为两种模式：前向模式和反向模式。**前向模式**是按计算图中计算方向的相同方向来递归的计算梯度。以  $\frac{\partial f(x;w,b)}{\partial w}$  为例，当  $x = 1, w = 0, b = 0$  时，前向模式的累积计算顺序如下：

$$\frac{\partial h_1}{\partial w} = x = 1 \quad (5.69)$$

$$\frac{\partial h_2}{\partial w} = \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w} = 1 * 1 = 1 \quad (5.70)$$

$$\frac{\partial h_3}{\partial w} = \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial w} = -1 * 1 \quad (5.71)$$

$$\vdots \quad \vdots \quad (5.72)$$

$$\frac{\partial h_6}{\partial w} = \frac{\partial h_6}{\partial h_5} \frac{\partial h_5}{\partial w} = -0.25 * -1 = 0.25 \quad (5.73)$$

$$\frac{\partial f(x;w,b)}{\partial w} = \frac{\partial f(x;w,b)}{\partial h_6} \frac{\partial h_6}{\partial w} = 1 * 0.25 = 0.25 \quad (5.74)$$

**反向模式**是按计算图中计算方向的相反方向来递归的计算梯度。以  $\frac{\partial f(x;w,b)}{\partial w}$  为例，当  $x = 1, w = 0, b = 0$  时，反向模式的累积计算顺序如下：

$$\frac{\partial f(x;w,b)}{\partial h_6} = 1 \quad (5.75)$$

$$\frac{\partial f(x;w,b)}{\partial h_5} = \frac{\partial f(x;w,b)}{\partial h_6} \frac{\partial h_6}{\partial h_5} = 1 * -0.25 \quad (5.76)$$

$$\frac{\partial f(x;w,b)}{\partial h_4} = \frac{\partial f(x;w,b)}{\partial h_5} \frac{\partial h_5}{\partial h_4} = -0.25 * 1 = -0.25 \quad (5.77)$$

$$\vdots \quad \vdots \quad (5.78)$$

$$\frac{\partial f(x;w,b)}{\partial w} = \frac{\partial f(x;w,b)}{\partial h_1} \frac{\partial h_1}{\partial w} = 0.25 * 1 = 0.25 \quad (5.79)$$

前向模式和反向模式可以看作是应用链式法则的两种梯度累积损失。从上面反向模式的计算顺序可以看出，反向模式和反向传播的计算梯度的方式相同。

对于一般的函数形式  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ ，前向模式需要对每一个输入变量都进行一遍遍历，共需要  $n$  遍。而反向模式需要对每一个输出都进行一个遍历，共需要  $m$  遍。当  $n > m$  时，反向模式更高效。在前馈神经网络的参数学习中，风险函数为  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ，输出为标量，因此采用反向模式为最有效的计算方式，只需要一遍计算。

在深度学习框架里，Python 语言 autograd<sup>1</sup> 和 Lua 语言的 Torch-autograd<sup>2</sup> 采用了自动微分方法，可以对各自的表达式、程序进行自动求导。比如，Autograd 可以直接对 Python 和 Numpy 代码进行求导。

符号微分和自动微分都利用计算图和链式法则来自动求解导数。符号微分在编译阶段先构造一个复合函数的计算图，通过符号计算得到导数的表达式，还可以对导数表达式进行优化，在程序运行阶段才代入变量的具体数值进行计算导数。而自动微分则无需事先编译，在程序运行阶段边计算边记录计算图，计算图上的局部梯度都直接代入数值进行计算，然后用前向或反向模式来计算最终的梯度。

## 5.5 梯度消失问题

在神经网络中误差反向传播的迭代公式为

$$\delta^{(l)} = f'_l(\mathbf{z}^{(l)}) \odot ((W^{(l+1)})^\top \delta^{(l+1)}), \quad (5.80)$$

其中需要用到激活函数  $f_l$  的导数。

误差从输出层反向传播时，在每一层都要乘以该层的激活函数的导数。

当我们使用 sigmoid 型函数：logistic 函数  $\sigma(x)$  或 tanh 函数时，其导数为

$$\sigma'(x) = \sigma(x)(1 - \sigma(x)) \in [0, 0.25] \quad (5.81)$$

$$\tanh'(x) = 1 - (\tanh(x))^2 \in [0, 1]. \quad (5.82)$$

<sup>1</sup> <https://github.com/HIPS/autograd>

<sup>2</sup> <https://github.com/twitter/torch-autograd>

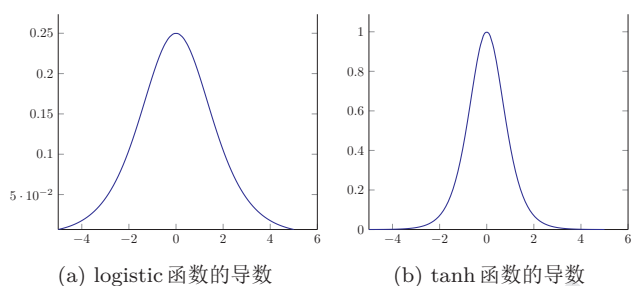


图 5.9: 激活函数的导数

我们可以看到，sigmoid 型函数导数的值域都小于 1。并且由于 sigmoid 型函数的饱和性，饱和区的导数更是接近于 0。这样，误差经过每一层传递都会不断衰减。当网络层数很深时，梯度就会不停的衰减，甚至消失，使得整个网络很难训练。这就是所谓的**梯度消失问题**（Vanishing Gradient Problem），也叫**梯度弥散问题**。梯度消失问题在过去的二三十年里一直没有有效地解决，是阻碍神经网络发展的重要原因之一。

在深层神经网络中，减轻梯度消失问题的一个方法就是使用导数比较大的激活函数，比如 ReLU（rectifier 函数）。这样误差可以很好地传播，训练速度得到了很大的提高。

## 5.6 训练方法

**随机梯度下降** 对于大规模数据，很难使用批量梯度下降。随机梯度下降（Stochastic Gradient Descent, SGD）

随机梯度下降：每次迭代（随机）选择一个样本计算梯度，并更新参数。

$$\mathbf{a}_{t+1} = \mathbf{a}_t - \lambda \frac{\partial \mathcal{R}(\theta_t; x^{(t)}, y^{(t)})}{\partial \theta}, \quad (5.83)$$

**Mini-batch 随机梯度下降** Mini-batch 随机梯度下降：批量梯度下降和随机梯度下降的折中。

每次迭代选取  $m (1 \leq m \leq N)$  个样本进行参数更新。

$$\mathbf{a}_{t+1} = \mathbf{a}_t - \lambda \frac{1}{m} \sum_{i \in I_t} \frac{\partial \mathcal{R}(\theta_t; x^{(i)}, y^{(i)})}{\partial \theta}, \quad (5.84)$$

## 5.7 经验

因为深度神经网络的表达能力很强，所以很容易产生过拟合。另外，大量的参数会导致训练比较慢。在训练深度神经网络时，同时也需要掌握一定的技巧。目前，人们在大量的实践中总结了一些经验技巧，可以从以下几个方面来提高学习效率并得到一个好的网络模型：1) 数据增强；2) 数据预处理；3) 网络参数初始化；4) 正则化；5) 超参数优化等。

### 5.7.1 数据增强

深层神经网络一般都需要大量的训练数据才能获得比较理想的结果。

### 5.7.2 数据预处理

一般而言，原始的训练数据中，每一维特征的来源以及度量单位不同，会造成这些特征值的分布范围往往差异很大。当我们计算不同样本之间的欧式距离时，取值范围大的特征会起到主导作用。这样，对于基于相似度比较的机器学习方法（比如最近邻分类器），必须先对样本进行预处理，将各个维度的特征归一化到同一个取值区间，并且消除不同特征之间的相关性，才能获得比较理想的结果。虽然神经网络可以通过参数的调整来适应不同特征的取值范围，但是会导致训练效率比较低。

假设一个只有一层的网络  $y = \tanh(w_1x_1 + w_2x_2 + b)$ ，其中  $x_1 \in [0, 10]$ ,  $x_2 \in [0, 1]$ 。之前我们提到  $\tanh()$  的导数在区间  $[-2, 2]$  上是敏感的，其余的导数接近于 0。因此，如果  $w_1x_1 + w_2x_2 + b$  过大或过小，都会导致梯度过小，难以训练。为了提高训练效率，我们需要使  $w_1x_1 + w_2x_2 + b$  在  $[-2, 2]$  区间，我们需要将  $w_1$  设得小一点，比如在  $[-0.1, 0.1]$  之间。可以想象，如果数据维数很多时，我们很难这样精心去选择每一个参数。因此，如果每一个特征的取值范围都在相似的区间，比如  $[0, 1]$  或者  $[-1, 1]$ ，我们就不太需要区别对待每一个参数，减少人工干预。

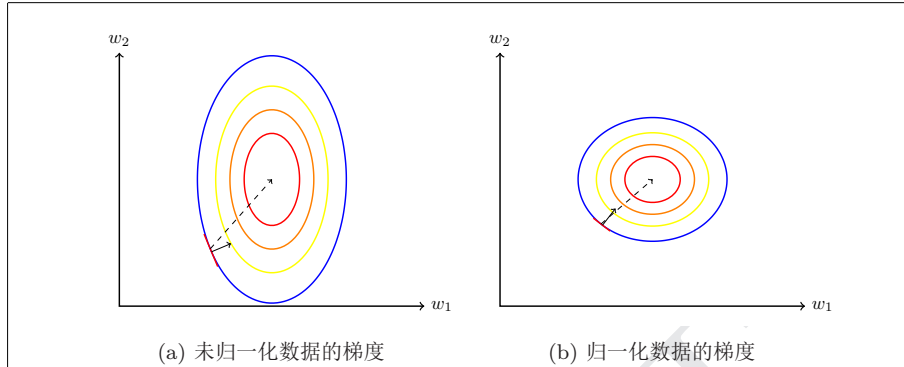


图 5.10: 数据归一化对梯度的影响。

除了参数初始化之外，不同特征取值范围差异比较大时还会梯度下降法的搜索效率。图5.10给出了数据归一化对梯度的影响。其中，图5.10a为未归一化数据的等高线图。取值范围不同会造成在大多数位置上的梯度方向并不是最优的搜索方向。当使用梯度下降法寻求最优解时，会导致需要很多次迭代才能收敛。如果我们把数据归一化为取值范围相同，如图5.10b所示，大部分位置的梯度方向近似于最优搜索方向。这样，在梯度下降求解时，每一步梯度的方向都基本指向最小值，训练效率会大大提高。

归一化的方法有很多种，比如之前我们介绍的sigmoid型函数等都可以将不同取值范围的特征挤压到一个比较受限的区间。这里，我们介绍几种在神经网络中经常使用的归一化方法。

**标准归一化** 标准归一化也叫z-score归一化，来源于统计上的标准分数。将每一个维特征都处理为符合标准正态分布（均值为0，标准差为1）。假设有  $N$  个样本  $\{\mathbf{x}^{(i)}\}, i = 1, \dots, N$ ，对于每一维特征  $x$ ，我们先计算它的均值和标准差：

$$\mu = \frac{1}{N} \sum_{i=1}^N x^{(i)}, \quad (5.85)$$

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \mu)^2. \quad (5.86)$$

然后，将特征  $x^{(i)}$  减去均值，并除以标准差，得到新的特征值  $\hat{x}^{(i)}$ 。

$$\hat{x}^{(i)} = \frac{x^{(i)} - \mu}{\sigma}, \quad (5.87)$$

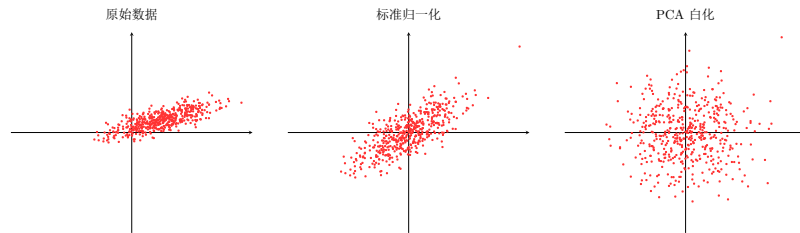


图 5.11: 数据归一化示例

这里， $\sigma$  不能为 0。如果标准差为 0，说明这一维特征没有任务区分性，可以直接删掉。

在标准归一化之后，每一维特征都服从标准正态分布。

**缩放归一化** 另外一种非常简单的归一化是通过缩放将特征取值范围归一到  $[0, 1]$  或  $[-1, 1]$  之间：

$$\hat{x}^{(i)} = \frac{x^{(i)} - \min(x)}{\max(x) - \min(x)}, \quad (5.88)$$

其中， $\min(x)$  和  $\max(x)$  分别为这一维特征在所有样本上的最小值和最大值。

**PCA** 使用 PCA (Principal Component Analysis) 方法可以去除掉各个成分之间的相关性。

### 5.7.3 批量归一化

在传统机器学习中，一个常见的问题的**协变量偏移** (Covariate Shift)。协变量是一个统计学概念，是可能影响预测结果的变量。在机器学习中，协变量可以看作是输入变量。一般的机器学习算法都要求输入变量在训练集和测试集上的分布是相似的。如果不满足这个要求，这些学习算法在测试集的表现会比较差。

在多层的神经网络中，中间某一层的输入是其前面网络的输出。因为前面网络的参数在每次迭代时也会被更新，而一旦这些参数变化会造成这一个中间层的输入也发生变化，其分布往往会和参数更新之前差异比较大。换句话说，从



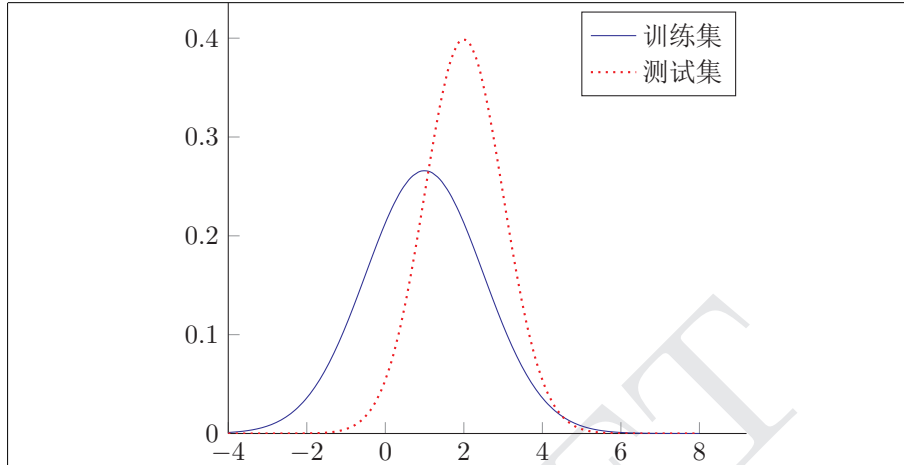


图 5.12: 协变量偏移。

这一个中间层开始，之后的网络参数白学了，需要重新学习。这个中间层的深度很大时，这种现象就越明显。这种现象叫做**内部协变量偏移**（Internal Covariate Shift）。

为了解决这个问题，通过对每一层的输入进行归一化使其分布保存稳定。这种方法称为**批量归一化方法**（Batch Normalization）[Ioffe and Szegedy, 2015]。

归一化方法可以采用第5.7.2节中介绍的几种归一化方法。因为每一层都要进行归一化，所以要求归一化方法的速度要快。这样，PCA白化的归一化方法就不太合适。为了提高归一化效率，这里使用标准归一化，对每一维特征都归一到标准正态分布。相当于每一层都进行一次数据预处理，从而加速收敛速度。

因为标准归一化会使得输入的取值集中的0附近，如果使用sigmoid型激活函数时，这个取值区间刚好是接近线性变换的区间，减弱的神经网络的非线性性质。因此，为了使得归一化不对网络的表示能力造成负面影响，我们可以通过一个附加的缩放和平移变换改变取值区间。从最保守的角度考虑，可以通过来标准归一化的逆变换来使得归一化后的变量可以被还原为原来的值。

$$y_k = \gamma_k \hat{x}_k + \beta_k, \quad (5.89)$$

这里  $\gamma_k$  和  $\beta_k$  分别代表缩放和平移的参数。当  $\gamma_k = \sqrt{\sigma[x_k]}$ ,  $\beta_k = \mu[x_k]$  时， $y_k$  即为原始的  $x_k$ 。

#### 算法 5.2: 批量归一化

输入: 一次 mini-batch 中的所有样本集合:

$$\mathcal{B} = \{\mathbf{x}^{(i)}\}, i = 1, \dots, m;$$

参数:  $\gamma, \beta$ ;

输出:  $\{y^{(i)} = \mathbf{BN}_{\gamma, \beta}(\mathbf{x}^{(i)})\}$

1 for  $k = 1 \dots K$  do

2

$$\mu_{\mathcal{B}, k} = \frac{1}{m} \sum_{i=1}^m x_k^{(i)}, \quad // \text{ mini-batch 均值}$$

$$\sigma_{\mathcal{B}, k}^2 = \frac{1}{m} \sum_{i=1}^m (x_k^{(i)} - \mu_{\mathcal{B}, k})^2. \quad // \text{ mini-batch 方差}$$

$$\hat{x}_k^{(i)} = \frac{x_k^{(i)} - \mu_{\mathcal{B}, k}}{\sqrt{\sigma_{\mathcal{B}, k}^2 + \epsilon}}, \forall i \quad // \text{ 归一化}$$

$$3 \quad y_k^{(i)} = \gamma \hat{x}_k^{(i)} + \beta \equiv \mathbf{BN}_{\gamma, \beta}(\mathbf{x}^{(i)}), \forall i \quad // \text{ 缩放和平移}$$

4 end

当训练完成时, 用整个数据集上的均值  $\mu_k$  和方差  $\sigma_k^2$  来分别代替  $\mu_{\mathcal{B}, k}$  和方差  $\sigma_{\mathcal{B}, k}^2$ 。

通过每一层的归一化, 从而减少前面网络参数更新对后面网络输入带来的内部协变量偏移问题, 提高训练效率。

#### 5.7.4 参数初始化

神经网络的训练过程中的参数学习是基于梯度下降法进行优化的。梯度下降法需要在开始训练时给每一个参数赋一个初始值。这个初始值的选取十分关键。在感知器和 logistic 回归的训练中, 我们一般将参数全部初始化为 0。但是这在神经网络的训练中会存在一些问题。因为如果参数都为 0, 在第一遍前向计算时, 所有的隐层神经元的激活值都相同。这样会导致深层神经元没有区分性。这种现象也称为对称权重现象。

为了打破这个平衡, 比较好的方式是对每个参数都随机初始化, 这样使得

不同神经元之间的区分性更好。

但是一个问题是如何选取随机初始化的区间呢？如果参数太小，会导致神经元的输入过小。经过多层之后信号就慢慢消失了。参数过小还会使得sigmoid型激活函数丢失非线性的能力。以logistic函数为例，在0附近基本上是近似线性的。这样多层神经网络的优势也就不存在了。如果参数取得太大，会导致输入状态过大。对于sigmoid型激活函数来说，激活值变得饱和，从而导致梯度接近于0。

因此，如果要高质量地训练一个网络，给参数选取一个合适的初始化区间是非常重要的。经常使用的初始化方法有以下几种：

**Gaussian 初始化方法** Gaussian 初始化方法是最简单的初始化方法，参数从一个固定均值（比如0）和固定方差（比如0.01）的Gaussian分布进行随机初始化。

**Xavier 初始化方法** Glorot and Bengio [2010] 提出一个初始化方法，参数可以在区间  $[-r, r]$  内采用均匀分布进行初始化。

对于logistic函数，第  $l-1$  到  $l$  层的权重，

$$r = \sqrt{\frac{6}{n^{l-1} + n^l}}, \quad (5.90)$$

这里  $n^l$  是第  $l$  层神经元个数， $n^{l-1}$  是第  $l-1$  层神经元个数。

对于tanh函数，

$$r = 4\sqrt{\frac{6}{n^{l-1} + n^l}}. \quad (5.91)$$

假设第  $l$  层的一个隐藏层神经元  $z^l$ ，其接受前一层的  $n^{l-1}$  个神经元的输出  $a_i^{(l-1)}$ ， $i \in [1, n^{(l-1)}]$ ，

$$z^l = \sum_{i=1}^n w_i^l a_i^{(l-1)} \quad (5.92)$$

为了避免初始化参数使得激活值变得饱和，我们需要尽量使得  $z^l$  处于激活函数的线性区间，也就是其绝对值比较小的值。这时，该神经元的激活值为  $a^l \equiv f(z^l) \approx z^l$ 。

Xavier 初始化方法中，Xavier 是发明者 Glorot 的名字。

假设  $w_i^l$  和  $a_i^{(l-1)}$  都是相互独立，并且均值都为0，则  $a$  的均值为

$$E(a^l) = E\left(\sum_{i=1}^n w_i^l a_i^{(l-1)}\right) = \sum_{i=1}^d E(\mathbf{w}_i) E(a_i^{(l-1)}) = 0. \quad (5.93)$$

$a^l$  的方差为

$$\text{Var}(a^l) = \text{Var}\left(\sum_{i=1}^{n^{(l-1)}} w_i^l a_i^{(l-1)}\right) \quad (5.94)$$

$$= \sum_{i=1}^{n^{(l-1)}} \text{Var}(w_i^l) \text{Var}(a_i^{(l-1)}) \quad (5.95)$$

$$= n^{(l-1)} \text{Var}(w_i^l) \text{Var}(a_i^{(l-1)}). \quad (5.96)$$

也就是说,输入信号的方差在经过该神经元后被放大或缩小了  $n^{(l-1)} \text{Var}(w_i^l)$  倍。为了使得在经过多层网络后，信号不被过分放大或过分减弱，我们尽可能保持每个神经元的输入和输出的方差一致。这样  $n^{(l-1)} \text{Var}(w_i^l)$  设为1比较合理，即

$$\text{Var}(w_i^l) = \frac{1}{n^{(l-1)}}. \quad (5.97)$$

同理，为了使得在反向传播中，误差信号也不被放大或缩小，需要将  $w_i^l$  的方差保持为

$$\text{Var}(w_i^l) = \frac{1}{n^{(l)}}. \quad (5.98)$$

作为折中，同时考虑信号在前向和反向传播中都不被放大或缩小，可以设置

$$\text{Var}(w_i^l) = \frac{2}{n^{(l-1)} + n^{(l)}}. \quad (5.99)$$

假设随机变量  $x$  在区间  $[a, b]$  内均匀分布，则其方差为：

$$\text{Var}(x) = \frac{(b-a)^2}{12}. \quad (5.100)$$

因此，若让  $w_i^l \in [-r, r]$ ，则  $r$  的取值为

$$r = \sqrt{\frac{6}{n^{l-1} + n^l}}. \quad (5.101)$$

### 5.7.5 正则化

深度神经网络很容易产生过拟合现象，因为增加的抽象层使得模型能够对训练数据中较为罕见的依赖关系进行建模。对此，权重递减（ $\ell_2$  正规化）或者稀疏（ $\ell_1$  -正规化）等方法可以利用在训练过程中以减小过拟合现象 [Bengio et al., 2013]。

### 5.7.6 Dropout

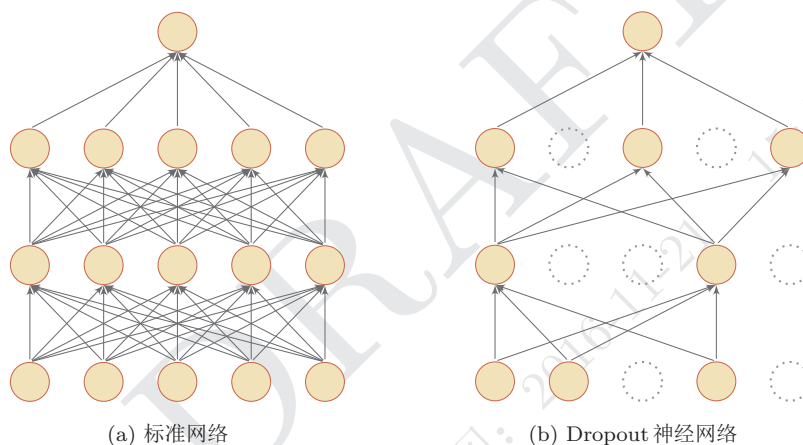


图 5.13: Dropout 神经网络模型

训练一个大规模的神经网络经常容易过拟合。过拟合在很多机器学习中都会出现。一般解决过拟合的方法有正则化、早期终止、集成学习以及使用验证集等。Srivastava et al. [2014] 提出了适用于神经网络的避免过拟合的方法，叫 **dropout** 方法（丢弃法），即在训练中随机丢弃一部分神经元（同时丢弃其对应的连接边）来避免过拟合。图5.13给出了 dropout 网络的示例。

每做一次 dropout，相当于从原始的网络中采样得到一个更瘦的网络。如果一个神经网络有  $n$  个神经元，那么总共可以采样出  $2^n$  个子网络。每次迭代都相当于训练一个不同的子网络，这些子网络都共享原始网络的参数。那么，最终的神经网络可以近似看作是集成了指数级个不同网络的组合模型。每次选择丢弃的神经元是随机的。最简单的方法是设置一个固定的概率  $p$ 。对每一个神经元都一

个概率  $p$  来判定要不要保留。 $p$  可以通过验证集来选取一个最优的值。或者,  $p$  也可以设为 0.5, 这对大部分的网络和任务有比较有效。

一般情况下 dropout 是针对神经元进行随机丢弃, 但是也可以扩展到对每条边进行随机丢弃, 或每一层进行随机丢弃。

当  $p = 0.5$  时, 在训练时有一半的神经元被丢弃, 只剩余一半的神经元是可以激活的。而在测试时, 所有的神经元都是可以激活的。因此每个神经元训练时的净输入值平均比测试时小一半左右。这会造成训练和测试时网络的输出不一致。为了缓解这个问题, 在测试时需要将每一个神经元的输出都折半, 也相当于把不同的神经网络做了平均。

一般来讲, 对于隐藏层的神经元, 其 dropout 率等于 0.5 时效果最好, 因为此时通过 dropout 方法, 随机生成的网络结构最具多样性。对于输入层的神经元, 其 dropout 率通常设为更接近 1 的数, 使得输入变化不会太大。当对输入层神经元进行 dropout 时, 其可以被给数据增加噪声, 以此来提高网络的鲁棒性。

### 5.7.7 超参数优化

在构建前馈神经网络时, 有下面的超参数需要设置。

- 网络层数
- 每层的神经元数量
- 激活函数的类型
- 学习率 (以及动态调整算法)
- 正则化系数
- 每次小批量梯度下降的

超参数的优化问题是一个组合优化问题, 没有通用的优化方法。对于超参数的设置, 一般用**网格搜索** (Grid Search) 或者人工搜索的方法来进行。假设总共有  $K$  个超参数, 第  $k$  个超参数的可以取  $m_k$  个值。如果参数是连续的, 可以将参数离散化, 选择几个“经验”值。比如学习率  $\alpha$ , 我们可以设置

$$\alpha \in \{0.01, 0.1, 0.5, 1.0\}.$$

这样，这些超参数可以有  $m_1 \times m_2 \times \cdots \times m_K$  个取值组合。所谓网格搜索就是根据这些超参数的不同组合分别训练一个模型，然后评价这些模型在检验数据集上的性能，选取一组性能最好的组合。

如果每个超参数对模型性能的影响有很大差异，有些超参数对模型性能的影响有限，而有些则非常大。在这种情况下，网格搜索可能会有些浪费。采用随机搜索会比网格搜索更加有效，而且在实践中也更容易实现 [Bergstra and Bengio, 2012]。

## 5.8 总结和深入阅读

多层前馈神经网络作为一种机器学习方法在很多模式识别和机器学习的教材中都有介绍，比如《Pattern Recognition and Machine Learning》[Bishop, 2006]，《Pattern Classification》[Duda et al., 2001] 等。虽然多层前馈网络在2000年以前就已被广泛使用，但是基本上都是三层网络（即只有一个隐藏层），神经元的激活函数基本上都是 sigmoid 型函数，并且使用的损失函数大多数是平方损失。

## 参考文献

- Martin Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*, 2016.
- James A Anderson and Edward Rosenfeld. *Talking nets: An oral history of neural networks*. MiT Press, 2000.
- Yoshua Bengio. Learning deep architectures for AI. *Foundations and trends® in Machine Learning*, 2(1):1–127, 2009.
- Yoshua Bengio, Patrice Simard, and Paolo Frasconi. Learning long-term dependencies with gradient descent is difficult. *Neural Networks, IEEE Transactions on*, 5(2):157–166, 1994.
- Yoshua Bengio, Jean-Sébastien Senécal, et al. Quick training of probabilistic neural nets by importance sampling. In *AISTATS Conference*, 2003.
- Yoshua Bengio, Nicolas Boulanger-Lewandowski, and Razvan Pascanu. Advances in optimizing recurrent networks. In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 8624–8628. IEEE, 2013.
- James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13(Feb):281–305, 2012.
- James Bergstra, Olivier Breuleux, Frédéric Bastien, Pascal Lamblin, Razvan Pascanu, Guillaume Desjardins, Joseph Turian, David Warde-Farley, and Yoshua Bengio. Theano: a cpu and gpu math expression compiler. In *Proceedings of the Python for Scientific Computing Conference (SciPy)*, 2010.
- C.M. Bishop. *Pattern recognition and machine learning*. Springer New York., 2006.



- Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- P.F. Brown, P.V. Desouza, R.L. Mercer, V.J.D. Pietra, and J.C. Lai. Class-based n-gram models of natural language. *Computational linguistics*, 18(4): 467–479, 1992.
- Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- Michael Collins. Discriminative training methods for hidden markov models: Theory and experiments with perceptron algorithms. In *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing*, 2002.
- Hal Daumé III. A course in machine learning. <http://ciml.info/>. [Online].
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *The Journal of Machine Learning Research*, 12:2121–2159, 2011.
- R. Duda, P. Hart, and D. Stork. *Pattern Classification*. Wiley, New York, 2nd edition, 2001. ISBN 0471056693.
- Charles Dugas, Yoshua Bengio, François Bélisle, Claude Nadeau, and René Garcia. Incorporating second-order functional knowledge for better option pricing. *Advances in Neural Information Processing Systems*, pages 472–478, 2001.
- Jeffrey L Elman. Finding structure in time. *Cognitive science*, 14(2):179–211, 1990.
- Yoav Freund and Robert E Schapire. Large margin classification using the perceptron algorithm. *Machine learning*, 37(3):277–296, 1999.
- Kunihiko Fukushima. Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological cybernetics*, 36(4):193–202, 1980.
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *International conference on artificial intelligence and statistics*, pages 249–256, 2010.

- Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In *International Conference on Artificial Intelligence and Statistics*, pages 315–323, 2011.
- Ian Goodfellow, Aaron Courville, and Yoshua Bengio. Deep learning. Book in preparation for MIT Press, 2015. URL <http://goodfeli.github.io/dlbook/>.
- Ian J Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron C Courville, and Yoshua Bengio. Maxout networks. In *ICML*, volume 28, pages 1319–1327, 2013.
- T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer, New York, 2001.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1026–1034, 2015.
- Geoffrey Hinton, Li Deng, Dong Yu, George E Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N Sainath, et al. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *Signal Processing Magazine, IEEE*, 29(6):82–97, 2012.
- Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- Sepp Hochreiter, Yoshua Bengio, Paolo Frasconi, and Jürgen Schmidhuber. Gradient flow in recurrent nets: the difficulty of learning long-term dependencies, 2001.
- John J Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.
- David H Hubel and Torsten N Wiesel. Receptive fields of single neurones in the cat’s striate cortex. *The Journal of physiology*, 148(3):574–591, 1959.
- David H Hubel and Torsten N Wiesel. Receptive fields, binocular interaction and functional architecture in the cat’s visual cortex. *The Journal of physiology*, 160(1):106–154, 1962.

- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- M.I. Jordan. *Learning in Graphical Models*. Kluwer Academic Publishers, 1998.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acoustic models. In *Proc. ICML*, volume 30, 2013.
- Ryan McDonald, Keith Hall, and Gideon Mann. Distributed training strategies for the structured perceptron. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 456–464. Association for Computational Linguistics, 2010.
- Marvin Minsky and Papert Seymour. Perceptrons. 1969.
- Marvin L Minsky and Seymour A Papert. *Perceptrons - Expanded Edition: An Introduction to Computational Geometry*. MIT press Boston, MA:, 1987.
- T.M. Mitchell. *Machine learning*. Burr Ridge, IL: McGraw Hill, 1997.
- Andriy Mnih and Yee Whye Teh. A fast and simple algorithm for training neural probabilistic language models. *arXiv preprint arXiv:1206.6426*, 2012.
- Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 807–814, 2010.
- Jorge Nocedal and Stephen Wright. *Numerical optimization*. Springer Science & Business Media, 2006.
- Albert BJ Novikoff. On convergence proofs for perceptrons. Technical report, DTIC Document, 1963.
- F. Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386–408, 1958.

- David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *Cognitive modeling*, 5:3, 1988.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- Ilya Sutskever, Oriol Vinyals, and Quoc VV Le. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, pages 3104–3112, 2014.
- Paul Werbos. Beyond regression: New tools for prediction and analysis in the behavioral sciences. 1974.
- Paul J Werbos. Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*, 78(10):1550–1560, 1990.
- DE Rumelhart GE Hinton RJ Williams and GE Hinton. Learning representations by back-propagating errors. *Nature*, pages 323–533, 1986.
- Matthew D Zeiler. Adadelta: An adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.