
第三章 机器学习概述

机器学习是对能通过经验自动改进的计算机算法的研究。

— Mitchell [1997]

在介绍人工神经网络之前，我们先来了解下机器学习的基本概念。

通俗地讲，**机器学习**（Machine Learning, ML）就是如何让计算机从数据中进行学习，得到某种知识（或规律）。稍微正式一点的定义为：机器学习是指一类问题以及解决这类问题的方法，主要是研究如何从观测数据（样本）中寻找规律，并利用学习到的规律（模型）对未知或无法观测的数据进行预测。

机器学习问题在早期的工程领域也经常称为**模式识别**（Pattern Recognition, PR），但模式识别更偏向于具体的应用任务，比如光学字符识别、语音识别、人脸识别等。这些任务的特点是我们人类很容易做到，但我们不知道自己是如何做到的，因此也很难人工设计一个计算机程序来解决这些任务。一个可行的方法是设计一个算法可以让计算机自己从有标注的样本上学习其中的规律，并用来完成各种识别任务。随着机器学习技术的应用越来越广，现在机器学习的概念逐渐替代模式识别，成为这一类方法的统称。

以手写体数字识别为例，我们需要让计算机能自动识别手写的数字。比如像图3.1中的例子，将5识别为5，将6识别为6。手写数字识别是一个经典的机器学习任务，对人类来说很简单，但对计算机来说十分困难。设计一套识别算法几乎是不可能的任务。我们很难总结每个数字的手写体特征，或者区分不同数字的规则。在现实生活中，很多问题都类似于手写体数字识别这类问题，比如物体识别、语音识别等。我们不知道如何如何来让程序实现，即使可以通过



图 3.1: 手写体数字识别示例。图片来源: MNIST 数据集 [LeCun et al., 1998]

一些启发式规则来实现也是极其复杂的。因此,人们尝试采用另一种思路,即让计算机“看”大量的样本,并从中学习到一些经验,然后用这些经验来识别新的样本。这就是机器学习的思路。

还是以手写体数字识别为例,机器学习方法首先需要人工标注的大量手写体数字图像(即每个图像都人工标记了它是什么数字)作为训练样本集合,然后通过学习算法自动生成一套模型,并依靠它来识别新的手写体数字。这和人类学习过程也比较类似,我们教小孩子识别数字也是这样的过程。

3.1 机器学习定义

狭义地讲,机器学习是给定一些训练样本 $(x^{(i)}, y^{(i)}), 1 \leq i \leq N$ (其中, $x^{(i)}$ 是观测样本, $y^{(i)}$ 是需要预测的目标), 让计算机自动寻找一个决策函数 $f(\cdot)$ 来建立 $x^{(i)}$ 和 $y^{(i)}$ 之间的关系。这样对于一个新的输入样本 x , 我们可以通过决策函数来预测目标 y 。

$$y = f(\phi(x), \theta), \quad (3.1)$$

这里, θ 表示决策函数 $f(\cdot)$ 的参数, $\phi(x) \in \mathbb{R}^d$ 表示样本 x 对应的特征, 一般为向量形式。

在机器学习的应用中, 样本 x 的类型多种多样, 比较有代表性的类型为文本、音频、图像、视频等。在数字手写体识别中, 样本 x 为待识别的图像, 类

样本 x 不一定是数值型的, 因此需要通过 $\phi(x)$ 将 x 转换为数值型表示。

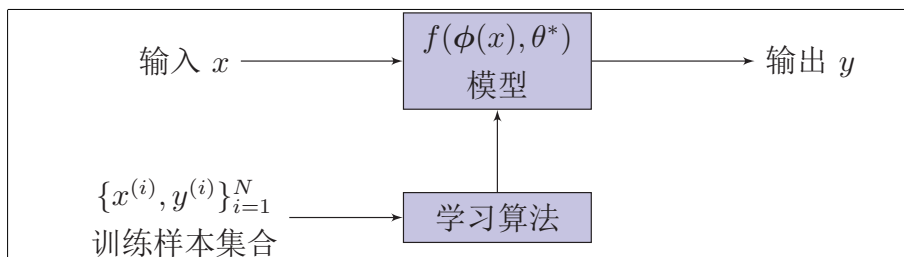


图 3.2: 机器学习系统示例

别 $y \in \{0, 1, \dots, 9\}$ 分别对应 10 个不同数字。为了识别 x 是什么数字，我们可以从图像中抽取一些特征，这些特征可以是直方图、宽高比、水平、垂直、对角线方向上的像素百分比、笔画数等。假设我们总共抽取了 d 个特征，这些特征可以表示为一个向量 $\phi(x) \in \mathbb{R}^d$ 。在情感分类中，样本 x 为自然语言文本，类别 $y \in \{+1, -1\}$ 分别表示正面或负面的评价。为了将样本 x 从文本形式转为为向量形式，我们可以使用**词袋模型** (Bag-of-Words, BoW) 模型。假设训练样本对应的词典 \mathcal{V} 中包含 v 个词，则每个文本可以表示为一个维度为 v 的向量 $\phi(x) \in \mathbb{R}^v$ ，向量中每一维对应词典中的一个词。如果向量中某一维对应的词在文本中出现，其值为 1，否则为 0。

机器学习系统的示例见图 3.2。对一个预测任务，输入为 x ，输出为 y 。我们需要提取特征 $\phi(x)$ ，再选择一个模型，即决策函数 $f(\phi(x), \theta)$ ，其中 θ 为模型参数，需要通过学习算法在一组训练样本上来得到一个最优的参数 θ^* 。通过训练样本来求解参数的过程也叫做机器学习的**训练过程**。有了模型 $f(\phi(x), \theta^*)$ ，我们就可以对任何输入 x 进行预测。

为了简单起见，我们直接用向量 \mathbf{x} 来表示样本 x 的特征向量，即 $\phi(x) = \mathbf{x}$ 。公式 3.1 也可以直接写为

$$y = f(\mathbf{x}, \theta). \quad (3.2)$$

3.1.1 机器学习类型

到目前为止已经非常多的机器学习模型。按照训练数据提供的信息以及反馈方式的不同，机器学习算法一般可以分为以下几类：

1. **监督学习** (Supervised Learning) 是利用一组已知输入 x 和目标标签 y 的

表 3.1: 三种机器学习类型的比较

| | 监督学习 | 无监督学习 | 增强学习 |
|------|-------|--------|------|
| 输入 | 有标记数据 | 无标记数据 | 决策过程 |
| 反馈方式 | 直接反馈 | 无反馈 | 奖励 |
| 目标 | 分类、预测 | 发现隐藏结构 | 动作 |

数据来学习模型的参数,使得模型预测的目标标签和真实标签尽可能的一致。

根据目标标签的类型不同,有监督学习又可以分为**回归**和**分类**两类。

- (a) **回归 (Regression) 问题**: 目标标签 y 是连续值 (实数或连续整数), $f(\mathbf{x}, \theta)$ 的输出也是连续值。对于所有已知或未知的 (x, y) ,使得 $f(\mathbf{x}, \theta)$ 和 y 尽可能地一致。
 - (b) **分类 (Classification) 问题**: 目标标签 y 是离散的类别 (符号)。在分类问题中,通过学习得到的决策函数 $f(\mathbf{x}, \theta)$ 也叫**分类器**。分类问题根据其类别数量可分为**两类分类 (Binary Classification)**和**多类分类 (Multi-class Classification)**。
2. **无监督学习 (Unsupervised Learning)** 是用来学习的数据不包含目标标签,需要学习算法自动学习到一些有价值的信息。典型的无监督学习问题有**聚类 (Clustering)**、**降维**、**密度估计**等。
 3. **增强学习 (Reinforcement Learning)** 也叫**强化学习**,强调如何基于环境做出一系列的**动作**,以取得最大化的累积收益。每做出一个动作,并不一定立刻得到收益。增强学习和有监督学习的不同在于增强学习不需要显式地以输入/输出对的方式给出训练样本,是一种在线的学习机制。

有监督的学习方法需要每个数据记录都有类标号,而无监督的学习方法则不考虑任何指导性信息。一般而言,一个监督学习模型需要大量的有标记数据集,而这些数据集是需要人工标注的。因此,也出现了很多**弱监督学习**和**半监督学习**的方法,希望从大规模的未标记数据中充分挖掘有用的信息,降低对标记数据数量的要求。

表3.1给出了三种机器学习类型的比较。

3.1.2 常见的基本概念

在机器学习中，经常会提及了一些基本概念，比如“数据”，“样本”，“特征”，“数据集”等。我们首先来解释下这些概念。

数据 在计算机科学中，**数据**是指所有能计算机程序处理的对象的总称，可以是数字、字母和符号等。在不同的任务中，表现形式不一样，比如图像、声音、文字、传感器数据等。

样本 **样本** (Sample)，也叫**实例** (Instance)，是按照一定的抽样规则或概率分布 P 从全部样本空间 \mathcal{X} 中取出的一部分数据，是实际观测得到的数据。不同形式数据的样本空间不同。比如，一张灰度图像（像素数为 n ）的样本空间为 $[0, 255]^n$ ，一个自然语言句子（长度为 L ）的样本空间为 \mathcal{V}^L ，其中 \mathcal{V} 为词表集合。

标签

正例和负例 对于两类分类问题，类别 $y \in \{+1, -1\}$ 可以直接用正负号表示。因此，常用**正例**和**负例**来分别表示属于不同类别的样本。

特征 每个样本通常表示为一组固定的属性或特征，这些属性通常是并不都以连续变量或离散变量的形式存在的。而机器学习中的很多算法的输入要求是数学上可计算的，因此需要从样本的属性中抽取出一些可以表征它们的数值型特征，一般为向量形式，也称为**特征向量**。对于一个样本 x ，假设我们提取了 d 个特征，这些特征可以构成一个 d 维的向量 $\phi(x) \in \mathbb{R}^d$ 。

如果样本 x 是一张大小为 $m \times n$ 的图片，其特征向量可以简单地表示为 mn 维的向量，每一维的值为图片中对应像素的灰度值。为了提高模型准确率，也会经常加入一个额外的特征，比如纹理特征、边缘特征等。如果样本 x 是一段文本，我们可以用词袋模型来得到其特征向量。

为了简单起见，和很多机器学习的文献一样，我们直接用特征向量来表示样本，并简写为 \mathbf{x} 。样本 $\mathbf{x} = [x_1; \dots; x_d]$ 为 d 维向量，其中每个 x_i 为样本 x 的第 i 个特征。

数据集 一组样本集合就称为**数据集**。一般情况下，假设数据集中的每个样本都是独立地从相同的数据分布中抽取的。我们用 $\mathcal{D} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$ 来表示由 N 个样本组成的数据集。

为了检验机器学习算法的好坏，一般将数据集分为两部分：训练集和测试集。**训练集**（Training Set）中的样本是用来训练模型的，也叫训练样本，而**测试集**（Test Set）中的样本是用来检验模型好坏的，也叫测试样本。

通过学习算法，在训练集得到一个模型，这个模型可以对测试集 \mathcal{D}_{Test} 上的每个样本 $(\mathbf{x}^{(i)}, y^{(i)})$ 预测一个类别标签 $\hat{y}^{(i)} = f(\mathbf{x}^{(i)}, \theta)$ 。模型的正确率为：

$$Acc = \frac{1}{|\mathcal{D}_{Test}|} \sum_{i=1}^{|\mathcal{D}_{Test}|} |\hat{y}^{(i)} = y^{(i)}|, \quad (3.3)$$

其中 $|\mathcal{D}_{Test}|$ 为测试集的大小。第3.6节中会介绍更多的评价方法。

除了训练集和测试集之外，有时也会使用一个**验证集**（Validation Set）（也叫**开发集**（Development Set））来进行模型选择。

在很多领域，数据集也经常称为**语料库**（Corpus）。

参数与超参数 机器学习可以归结为学习一个映射函数 $f: \mathbf{x} \rightarrow y$ ，将输入变量 \mathbf{x} 映射为输出变量 y 。一般我们可以假设映射函数为 $y = f(x, \theta)$ 。其中 θ 即为函数的参数。参数可以通过学习算法进行学习。

除了可学习的参数之外，还有一类参数是用来定义模型结构或训练策略的，这类参数叫做**超参数**（Hyper-Parameter）。超参数和可学习的参数不同，通常是按照人的经验设定，或者通过网格搜索对一组超参数组合进行不断试错调整。

常见的超参数有：聚类算法中的类别个数、梯度下降法的步长、正则项的系数、神经网络的层数、支持向量机中的核函数等。超参数的选取一般都是组合优化问题，很难通过优化算法来自动学习。因此，优化超参数是机器学习的一个经验性很强的技术。

特征学习 原始数据的特征有很多，但是并不是所有的特征都是有用的。并且，很多特征通常是冗余并且易变的。我们需要抽取有效的、稳定的特征。传统的特征提取是通过人工方式进行的，这需要大量的人工和专家知识。即使这样，人工设计的特征在很多任务上也不能满足需要。因此，如何自动地学习有效的特

在贝叶斯方法中，超参数可以理解为由参数的参数，即控制模型参数分布的参数。

征也成为机器学习中一个重要的研究内容，也就是**特征学习**，也叫**表示学习**。特征学习可以简化模型、缩短训练时间、提高模型泛华能力、避免过拟合等。

特征学习分成两种：特征选择和特征抽取。**特征选择**（Feature Selection）是选取原始特征集合的一个有效子集，使得基于这个特征子集训练出来的模型准确率最高。简单地说，特征选择就是保留有用特征，移除冗余或无关的特征。假设原始的特征个数为 d ，则共有 2^d 个候选子集。特征选择的目标是选择一个最优的候选子集。最暴力的做法是测试每个特征子集，看机器学习模型哪个子集上的准确率最高。但这种方式效率太低。常用的方法是采样贪心的策略，由空集合开始，每一轮添加该轮最优的特征；或者从原始特征集合开始，每次删除最无用的特征。

特征抽取（Feature Extraction）是构造一个新的特征空间，并将原始特征投影在新的空间中。以线性投影为例，原始特征向量 $\mathbf{x} \in \mathbb{R}^d$ ，经过线性投影后得到在新空间中的特征向量 \mathbf{x}' 。

$$\mathbf{x}' = P\mathbf{x}, \quad (3.4)$$

其中 $P \in \mathbb{R}^{k \times d}$ 为映射矩阵。

经典的特征抽取方法有主成分分析（Principle Components Analysis, PCA）和线性判别分析（Linear Discriminant Analysis, LDA）等。

特征选择和特征抽取又都可以分为有监督和和无监督学习。有监督学习目标是在抽取对模型最有用的特征，而无监督学习的目标是减少冗余信息。

经过特征选择或特征抽取后，特征的数量一般会减少，因此特征选择和特征抽取也经常称为维数约减（Dimension Reduction）或降维。

特征选择和特征抽取的优点是可以较少的特征来表示原始特征中的大部分相关信息，去掉噪声信息，并进而提高计算效率和减小维度灾难（Curse of Dimensionality）。对于很多没有正则化的模型，特征选择和特征抽取非常必要。

例子

正则化参考第3.3节。

3.2 机器学习模型

机器学习是从有限的观测数据中学习（或“猜测”）出具有一般性的规律，并可以将总结出来的规律推广应用到未观测样本上。

以监督学习为例，给定一组训练样本 $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$ ，我们假定存在一个输入空间 \mathcal{X} 和一个输出空间 \mathcal{Y} 。不同机器学习任务的主要区别在于输出空间不同。如果是两类分类问题， $\mathcal{Y} = \{+1, -1\}$ 。如果是 K 类分类问题， $\mathcal{Y} = \{1, 2, \dots, K\}$ 。如果是回归问题， $\mathcal{Y} = \mathbb{R}$ 。

训练集中的每个样本 $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$ 是从 \mathcal{X} 和 \mathcal{Y} 的联合空间中按照某个未知分布 $p(\mathbf{x}, y)$ 独立地随机产生的。这里要求样本分布 $p(\mathbf{x}, y)$ 未知但必须是固定的，不会随时间而变化。如果 P 本身可变的话，我们就无法通过这些数据进行学习。也就是说，训练集 \mathcal{D} 是由 N 个独立同分布（Identically and Independently Distributed, IID）的样本组成。

对于样本空间中的样本 $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$ ，假定存在一个未知的目标函数（Target Function） $g: \mathcal{X} \rightarrow \mathcal{Y}$ 使得

$$y = g(\mathbf{x}), \quad (3.5)$$

机器学习的目标是通过观测训练集 \mathcal{D} 中的样本，找到一个函数 $f(\mathbf{x})$ 使得

$$f(\mathbf{x}) \sim y, \quad \forall (\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}, \quad (3.6)$$

$f(\mathbf{x})$ 也叫作假设（Hypothesis）。一个好的函数 $f(\mathbf{x})$ 应该在所有 (\mathbf{x}, y) 的可能取值上都和目标函数一致。我们定义函数 f 在潜在的数据分布 $p(\mathbf{x}, y)$ 上的期望误差（Expected Error）为

$$\mathcal{R}(f) = \mathbb{E}_{(\mathbf{x}, y) \sim p(\mathbf{x}, y)}[\mathcal{L}(f(\mathbf{x}), y)], \quad (3.7)$$

其中， $\mathcal{L}(f(\mathbf{x}), y)$ 为损失函数，用来量化两个变量之间的差异； $p(\mathbf{x}, y)$ 未知的联合分布。

这样，一个好的函数 f 应当有一个比较小的期望误差。但是，由于我们不知道真实的数据分布和目标函数，实际上无法计算期望误差 $\mathcal{R}(f)$ 。我们可以计算的是经验误差（Empirical Error），即在训练集上的错误率。

$$\hat{\mathcal{R}}_{\mathcal{D}}(f) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}^{(i)}), y^{(i)}), \quad (3.8)$$

这里，“输入空间”默认为样本的特征空间。

这里，“目标函数”为需要去估计的未知函数。和优化中目标函数（Objective Function）意义不同。

“损失函数”的定义见第3.3节。

期望误差和经验误差之间的差距就叫做泛化误差 (Generalization Error), 或泛化错误。

$$G = \mathcal{R}(f) - \hat{\mathcal{R}}_{\mathcal{D}}(f) \quad (3.9)$$

泛化误差可以衡量一个机器学习模型是否可以很好地泛化到未知数据。

由于我们不知道真实目标函数 $g(\mathbf{x})$ 的具体形式, 只能根据经验来确定一个假设函数集合 \mathcal{H} , 称为假设空间 (Hypothesis Space), 然后通过观测其在训练集 \mathcal{D} 上的特性, 从中选择一个理想的假设 $f(\mathbf{x}) \in \mathcal{H}$ 。这里的“理想”是指泛化错误率比较低。

给定一个训练集 \mathcal{D} , 机器学习的目标是学习如何对未知的样本进行预测, 特别是在不在训练集中出现的样本。一个学习算法或学习器 \mathcal{A} 观察了一系列的训练样本之后, 从假设空间 \mathcal{H} 中得到最终的假设 $f(\mathbf{x})$ 。

因为我们不知道真实的目标函数 g , 所以不也知道 f 是否真的接近 g 。那么我们如何验证假设函数与真实目标函数之间的关系呢?

“泛化误差”在有些文献中也指“期望误差”。

泛化错误是指在未知样本上的错误, 具体解释见下。

PAC 理论

3.3 风险函数与损失函数

给定一个训练集 \mathcal{D} 和一个模型 $f(\mathbf{x}, \theta)$, 机器学习算法需要学习一个最优的参数 θ 。参数 θ 的好还通过建立一些准则来衡量。在很多机器学习算法中, 一般是定义一个损失函数 $\mathcal{L}(y, f(x, \theta))$, 然后在所有的训练样本上来评价决策函数的风险。

损失函数

给定一个实例 (x, y) , 真实目标是 y , 机器学习模型的预测为 $f(x, \theta)$ 。如果预测错误时 ($f(x, \theta) \neq y$), 我们需要定义一个度量函数来定量地计算错误的程度。常见的损失函数有如下几类:

0-1 损失函数 0-1 损失函数 (0-1 loss function) 是

$$\mathcal{L}(y, f(x, \theta)) = \begin{cases} 0 & \text{if } y = f(x, \theta) \\ 1 & \text{if } y \neq f(x, \theta) \end{cases} \quad (3.10)$$

$$= I(y \neq f(x, \theta)), \quad (3.11)$$

这里 I 是指示函数。

平方损失函数 平方损失函数 (quadratic loss function) 是

$$\mathcal{L}(y, \hat{y}) = (y - f(x, \theta))^2 \quad (3.12)$$

交叉熵损失函数 对于分类问题, 预测目标 y 为离散的类别, 模型输出 $f(x, \theta)$ 为每个类的条件概率。

假设 $y \in \{1, \dots, C\}$, 模型预测的第 i 个类的条件概率 $P(y = i|x) = f_i(x, \theta)$, 则 $f(x, \theta)$ 满足

$$f_i(x, \theta) \in [0, 1], \quad \sum_{i=1}^C f_i(x, \theta) = 1 \quad (3.13)$$

$f_y(x, \theta)$ 可以看作真实类别 y 的似然函数。参数可以直接用最大似然估计来优化。考虑到计算问题, 我们经常使用最小化负对数似然, 也就是**负对数似然损失函数** (Negative Log Likelihood function)。

$$\mathcal{L}(y, f(x, \theta)) = -\log f_y(x, \theta). \quad (3.14)$$

如果我们用 one-hot 向量 \mathbf{y} 来表示目标类别 c , 其中只有 $y_c = 1$, 其余的向量元素都为 0。

负对数似然函数也可以写为:

$$\mathcal{L}(y, f(x, \theta)) = -\sum_{i=1}^C y_i \log f_i(x, \theta). \quad (3.15)$$

y_i 也可以看成是真实类别的分布, 这样公式 3.15 恰好是交叉熵的形式。因此, 负对数似然损失函数也常叫做**交叉熵损失函数** (Cross Entropy Loss function) 是负对数似然函数的一种改进。

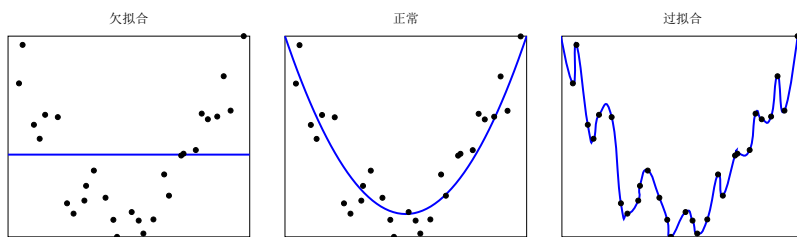


图 3.3: 欠拟合和过拟合示例

Hinge 损失函数 对于两类分类问题,假设 y 和 $f(x, \theta)$ 的取值为 $\{-1, +1\}$ 。Hinge 损失函数 (Hinge Loss Function) 的定义如下:

$$\mathcal{L}(y, f(x, \theta)) = \max(0, 1 - yf(x, \theta)) \quad (3.16)$$

$$= |1 - yf(x, \theta)|_+. \quad (3.17)$$

风险函数

$$\mathcal{R}_{emp}(\theta) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y^{(i)}, f(x^{(i)}, \theta)). \quad (3.18)$$

这里, 风险函数 $\mathcal{R}_{emp}(\theta)$ 是在已知的训练样本 (经验数据) 上计算得来的, 因此被称为**经验风险**。用对参数求经验风险来逐渐逼近理想的期望风险的最小值, 就是我们常说的**经验风险最小化原则** (Empirical Risk Minimization)。这样, 我们的目标就是变成了找到一个参数 θ^* 使得经验风险最小。

$$\theta^* = \arg \min_{\theta} \mathcal{R}_{emp}(\theta). \quad (3.19)$$

因为用来训练的样本往往是真实数据的一个很小的子集或者包含一定的噪声数据, 不能很好地反映全部数据的真实分布。经验风险最小化原则很容易导致模型在训练集上错误率很低, 但是在未知数据上错误率很高。这就是所谓的**过拟合**。过拟合问题往往是由于训练数据少和噪声等原因造成的。过拟合的标准定义为: 给定一个假设空间 \mathcal{H} , 一个假设 h 属于 \mathcal{H} , 如果存在其他的假设 h' 也属于 \mathcal{H} , 使得在训练集上 h 的损失比 h' 小, 但在整个样本空间上 h' 比 h 的损失小, 那么就假设 h 过度拟合训练数据 [Mitchell, 1997]。

为了解决过拟合问题，一般在经验风险最小化的原则上加上参数的正则化 (Regularization)，也叫结构风险最小化原则 (Structure Risk Minimization)。

$$\theta^* = \arg \min_{\theta} \mathcal{R}_{struct}(\theta) \quad (3.20)$$

$$= \arg \min_{\theta} \mathcal{R}_{emp}(\theta) + \lambda \|\theta\|_2^2 \quad (3.21)$$

$$= \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y^{(i)}, f(x^{(i)}, \theta)) + \lambda \|\theta\|_2^2. \quad (3.22)$$

这里， $\|\theta\|_2$ 是 L_2 范数的正则化项，用来减少参数空间，避免过拟合。 λ 用来控制正则化的强度。

正则化项也可以使用其它函数，比如 L_1 范数。 L_1 范数的引入通常会使得参数有一定稀疏性，因此在很多算法中也经常使用。在 Bayes 估计的角度来讲，正则化是假设了参数的先验分布，不完全依赖训练数据。

判别函数

经过特征抽取后，一个样本可以表示为 k 维特征空间中的一个点。为了对这个特征空间中的点进行区分，就需要寻找一些超平面来将这个特征空间分为一些互不重叠的子区域，使得不同类别的点分布在不同的子区域中，这些超平面就成为判别界面。

为了定义这些用来进行空间分割的超平面，就需要引入判别函数的概念。假设变量 $\mathbf{z} \in \mathbb{R}^m$ 为特征空间中的点，这个超平面由所有满足函数 $f(\mathbf{z}) = 0$ 的点组成。这里的 $f(\mathbf{z})$ 就称为判别函数。

有了判别函数，分类就变得很简单，就是看一个样本在特征空间中位于哪个区域，从而确定这个样本的类别。

3.4 参数学习算法

学习算法就是如何从训练集的样本中，自动学习决策函数的参数。不同机器学习算法的区别在于决策函数和学习算法的差异。相同的决策函数可以有不同的学习算法。比如，线性分类器有感知器，logistic 回归以及支持向量机等。感知器使用在线被动更新方法来学习参数，logistic 回归使用梯度下降法来学习参

数，而支持向量机使用SMO（Sequential minimal optimization）优化算法来学习参数。

通过一个学习算法进行自动学习参数的过程也叫作**训练过程**。

在选择合适的风险函数后，我们寻找一个参数 θ^* ，使得风险函数最小化。

$$\theta^* = \arg \min_{\theta} \mathcal{R}(\theta) \quad (3.23)$$

$$= \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y^{(i)}, f(x^{(i)}, \theta)). \quad (3.24)$$

这样，我们就将机器学习问题转化成为一个最优化问题。

批量随机梯度下降法

在最优化问题中，最常用的方法就是梯度下降法。

梯度下降是求得所有样本上的风险函数最小值，也叫做**批量梯度下降法**。

$$\theta_{t+1} = \theta_t - \alpha \frac{\partial \mathcal{R}(\theta)}{\partial \theta_t} \quad (3.25)$$

$$= \theta_t - \alpha \frac{1}{N} \sum_{i=1}^N \frac{\partial \mathcal{R}(\theta_t; x^{(i)}, y^{(i)})}{\partial \theta}. \quad (3.26)$$

在机器学习，搜索步长 α 中也叫作**学习率**（Learning Rate）。

随机梯度下降法

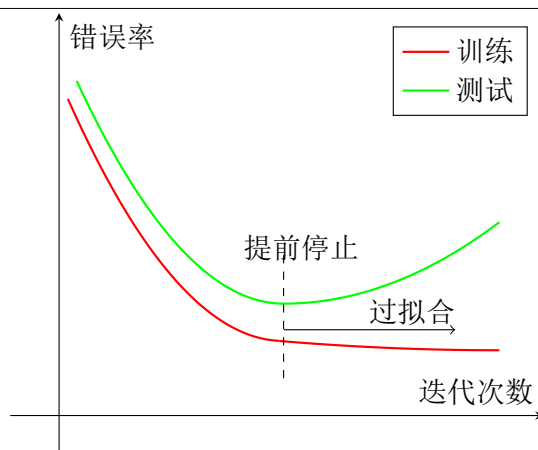
若样本个数 N 很大，输入 \mathbf{x} 的维数也很大时，那么批量梯度下降法每次迭代要处理所有的样本，效率会较低。为此，有一种改进的方法即**随机梯度下降法**。

随机梯度下降法（Stochastic Gradient Descent, SGD）也叫**增量梯度下降**，每个样本都进行更新

$$\theta_{t+1} = \theta_t - \alpha \frac{\partial \mathcal{L}(\theta_t; x^{(t)}, y^{(t)})}{\partial \theta}, \quad (3.27)$$

$x^{(t)}, y^{(t)}$ 是第 t 次迭代选取的样本。

梯度下降法见第??页第??节。



批量梯度下降和随机梯度下降之间的区别在于每次迭代的风险是对所有样本汇总的风险还是单个样本的风险。随机梯度下降因为实现简单，收敛速度也非常快，因此使用非常广泛。

还有一种折中的方法就是小批量（Mini-Batch）随机梯度下降法，每次迭代时，只采用一小部分的训练样本，兼顾了批量梯度下降和随机梯度下降的优点。

提前停止

在梯度下降训练的过程中，由于过拟合的原因，在训练样本上收敛的参数，并不一定在测试集上最优。因此，我们使用一个**验证集**（Validation Dataset）（也叫**开发集**（Development Dataset））来测试每一次迭代的参数在验证集上是否最优。如果在验证集上的错误率不再下降，就停止迭代。这种策略叫**提前停止**（Early-Stop）。如果没有验证集，可以在训练集上进行**交叉验证**。

学习率设置

在梯度下降中，学习率的取值非常关键，如果过大就不会收敛，如果过小则收敛速度太慢。一般步长可以由线性搜索算法来确定。在机器学习中，经常使用自适应调整学习率的方法。

动量法 动量法 (Momentum Method) [Rumelhart et al., 1988] 是对当前迭代的更新中加入上一次迭代的更新。我们记 $\nabla\theta_t = \theta_t - \theta_{t-1}$ 。在第 t 迭代时，

$$\theta_t = \theta_{t-1} + (\rho\nabla\theta_{t-1} - \lambda g_t), \quad (3.28)$$

其中， ρ 为动量因子，通常设为 0.9。这样，在迭代初期，使用前一次的梯度进行加速。在迭代后期的收敛值附近，因为两次更新方向基本相反，增加稳定性。

从某种角度来说，当前梯度叠加上部分的上次梯度，一定程度上可以近似看作二阶梯度。

优化、二阶梯度

AdaGrad AdaGrad (Adaptive Gradient) 算法 [Duchi et al., 2011] 是借鉴 L2 正则化的思想。在第 t 迭代时，

$$\theta_t = \theta_{t-1} - \frac{\rho}{\sqrt{\sum_{\tau=1}^t g_\tau^2}} g_t, \quad (3.29)$$

其中， ρ 是初始的学习率， $g_\tau \in \mathbb{R}^{|\theta|}$ 是第 τ 次迭代时的梯度。

随着迭代次数的增加，梯度逐渐缩小。

AdaM

AdaDelta AdaDelta 算法 [Zeiler, 2012] 用指数衰减的移动平均来累积历史的梯度信息。第 t 次迭代的梯度的期望 $E(g^2)_t$ 为：

$$E(g^2)_t = \rho E(g^2)_{t-1} + (1 - \rho) g_t^2, \quad (3.30)$$

其中， ρ 是衰减常数。

本次迭代的更新为

$$\nabla\theta_t = -\frac{\sqrt{E(\nabla\theta^2)_{t-1} + \epsilon}}{\sqrt{E(g^2)_t + \epsilon}} g_t \quad (3.31)$$

其中， $E(\nabla\theta^2)_t$ 为前一次迭代时 $\nabla\theta^2$ 的移动平均， ϵ 为常数。

累计更新本次迭代 $\nabla\theta^2$ 的移动平均

$$E(\nabla\theta^2)_t = \rho E(\nabla\theta^2)_{t-1} + (1 - \rho) \nabla\theta_t^2. \quad (3.32)$$

最后更新参数

$$\theta_t = \theta_{t-1} + \nabla \theta_t. \quad (3.33)$$

3.5 常见机器学习算法

本节介绍一些常见的机器学习算法。

3.5.1 线性回归

如果输入 \mathbf{x} 是列向量，目标 y 是连续值（实数或连续整数），预测函数 $f(\mathbf{x})$ 的输出也是连续值。这种机器学习问题是回归问题。

如果我们定义 $f(\mathbf{x})$ 是线性函数，

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b, \quad (3.34)$$

这就是线性回归问题（Linear Regression）。

为了简单起见，我们将公式3.34写为

$$f(\mathbf{x}) = \hat{\mathbf{w}}^T \hat{\mathbf{x}}, \quad (3.35)$$

其中 $\hat{\mathbf{w}}$ 和 $\hat{\mathbf{x}}$ 分别称为增广权重向量和增广特征向量。

$$\hat{\mathbf{x}} = \mathbf{x} \oplus 1 \triangleq \begin{bmatrix} x_1 \\ \vdots \\ x_k \\ 1 \end{bmatrix} = \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}, \quad (3.36)$$

$$\hat{\mathbf{w}} = \mathbf{w} \oplus b \triangleq \begin{bmatrix} w_1 \\ \vdots \\ w_k \\ b \end{bmatrix} = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}, \quad (3.37)$$

这里， \oplus 定义为两个向量的拼接操作。

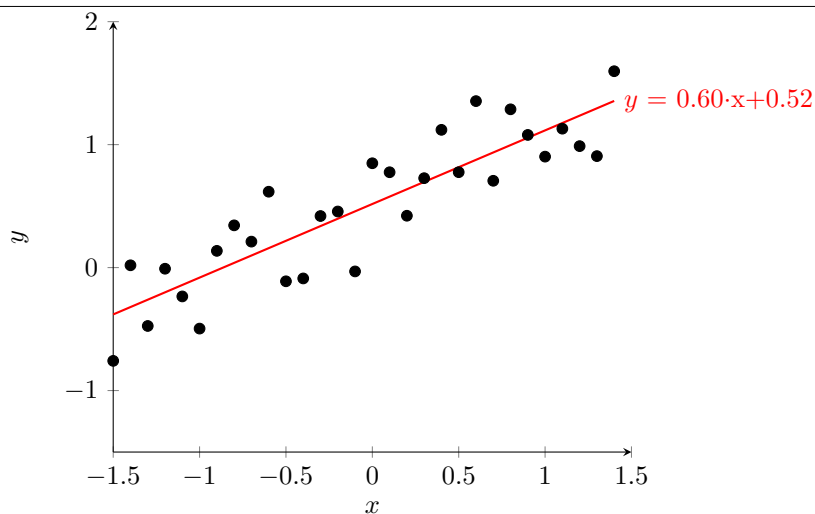


图 3.5: 线性回归示例

不失一般性，在本章后面的描述中我们采用简化的表示方法，直接用 \mathbf{w} 和 \mathbf{x} 来表示增广权重向量和增广特征向量。

线性回归的损失函数通常定义为平方损失函数。给定 N 个样本 $(\mathbf{x}^{(i)}, y^{(i)})$, $1 \leq i \leq N$ ，模型的经验风险为 平方损失参见第3.3节

$$\mathcal{R}(\mathbf{y}, f(X, \mathbf{w})) = \sum_{i=1}^N \mathcal{L}(y^{(i)}, f(\mathbf{x}^{(i)}, \mathbf{w})) \quad (3.38)$$

$$= \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2 \quad (3.39)$$

$$= \|\mathbf{X}^T \mathbf{w} - \mathbf{y}\|^2, \quad (3.40)$$

其中， $\mathbf{y} \in \mathbb{R}^N$ 是一个由目标值 $y^{(1)}, \dots, y^{(N)}$ 组成的列向量， $\mathbf{X} \in \mathbb{R}^{(k+1) \times N}$ 是所有输入 $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$ 组成的矩阵

$$\mathbf{X} = \begin{pmatrix} x_1^{(1)} & x_1^{(2)} & \cdots & x_1^{(N)} \\ \vdots & \vdots & \ddots & \vdots \\ x_k^{(1)} & x_k^{(2)} & \cdots & x_k^{(N)} \\ 1 & 1 & \cdots & 1 \end{pmatrix}. \quad (3.41)$$

要最小化 $\mathcal{R}(\mathbf{y}, f(X, \mathbf{w}))$ ，我们要计算 $\mathcal{R}(\mathbf{y}, f(X, \mathbf{w}))$ 对 \mathbf{w} 的导数

$$\frac{\partial \mathcal{R}(Y, f(X, \mathbf{w}))}{\partial \mathbf{w}} = \frac{\partial \|X^T \mathbf{w} - \mathbf{y}\|^2}{\partial \mathbf{w}} \quad (3.42)$$

$$= X(X^T \mathbf{w} - \mathbf{y}). \quad (3.43)$$

让 $\frac{\partial \mathcal{R}(\mathbf{y}, f(X, \mathbf{w}))}{\partial \mathbf{w}} = 0$ ，则可以得到

$$\mathbf{w} = (X X^T)^{-1} X \mathbf{y} \quad (3.44)$$

$$= \left(\sum_{i=1}^N \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T \right)^{-1} \left(\sum_{i=1}^N \mathbf{x}^{(i)} y^{(i)} \right). \quad (3.45)$$

这里要求 $X X^T$ 是满秩的，存在逆矩阵，也就是要求 \mathbf{x} 的每一维之间是线性不相关的。这样的参数求解方法也叫**最小二乘法估计**。

图3.5给出了用最小二乘法估计方法来进行参数学习的示例。

如果 $X X^T$ 不可求逆矩阵，说明在训练数据上，输入的不同特征之间是线性相关的。我们可以预处理数据，先使用主成分分析等特征提取方法来消除不同特征之间的相关性，然后再使用最小二乘估计方法来求解。

另外一种方法是通过用梯度下降法来求解。初始化 $\mathbf{w}_0 = 0$ ，迭代公式为：

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha \frac{\partial \mathcal{R}(Y, f(X, \mathbf{w}))}{\partial \mathbf{w}} \quad (3.46)$$

$$= \mathbf{w}_t + \alpha X(X^T \mathbf{w} - \mathbf{y}), \quad (3.47)$$

其中 α 是学习率。

3.5.2 Logistic 回归

Logistic 回归 (Logistic Regression) 是一种两类分类问题的参数学习算法。和回归问题不同，分类问题中的目标标签 y 是离散的类别（符号）。因此，分类问题中的决策函数 $f(\mathbf{x}, \theta)$ 需要输出离散值，这一般可以通过一个连续判别函数加上一个阈值函数来实现。

我们先来看两类分类问题，其预测的目标标签 y 只有两种取值。一般可设置为 $y \in \{0, 1\}$ 或者 $y \in \{-1, +1\}$ 。

对于分类问题，使用线性回归算法来求解是不合适的。图3.6a使用线性回归算法来解决一维的两类分类问题示例。

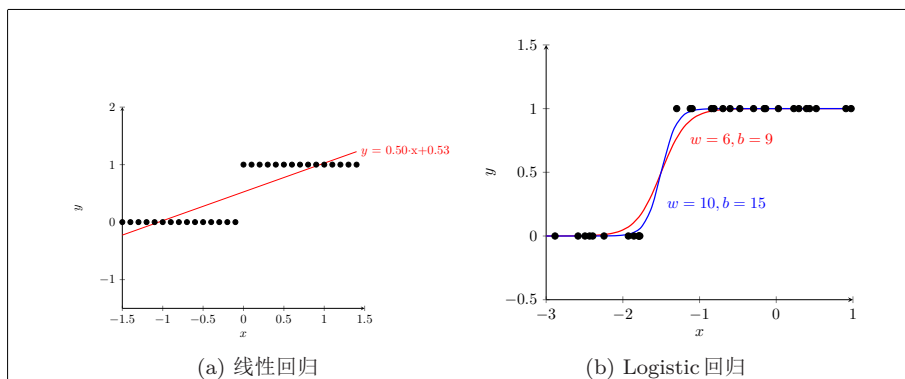


图 3.6: 一维两类问题示例

线性分类器是机器学习中应用最广泛的一类分类器，使用线性函数来作为判别函数。

$$\hat{y} = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} > 0 \\ 0 & \text{if } \mathbf{w}^T \mathbf{x} \leq 0 \end{cases} \quad (3.48)$$

$$= I(\mathbf{w}^T \mathbf{x} > 0), \quad (3.49)$$

其中， $I()$ 是指示函数。

公式3.71定义了一个两类分类问题的线性判别函数。在高维的特征空间中，所有满足 $\mathbf{f}(z) = 0$ 的点组成用一个超平面，这个超平面将特征空间一分为二，划分成两个区域。这两个区域分别对应两个类别。

图3.7中给了一个二维数据的判别函数以及对应的判别界面。在二维空间中，分类界面为一个直线。在三维空间中，分类界面为一个平面。在高维空间中，分类界面为一个超平面。对于线性函数来说，权重向量在线性空间中垂直于分类界面的向量。

线性分类函数的参数 \mathbf{w} 有很多种学习方式，比如第??章中介绍的感知器。这里使用另一种常用的学习算法：**Logistic** 回归。

我们定义目标类别 $y = 1$ 的后验概率为：

$$P(y = 1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) \quad (3.50)$$

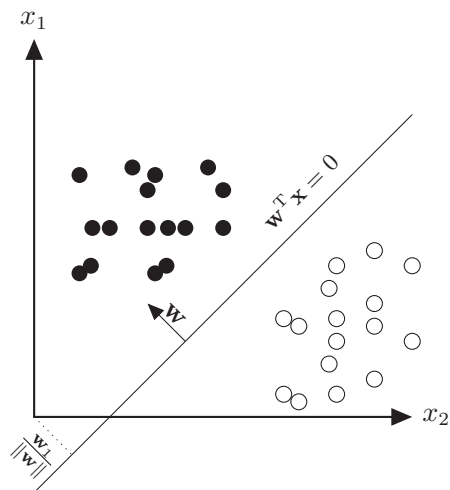


图 3.7: 两类分类线性判别函数

$$\sigma(z) \triangleq \frac{1}{1 + \exp(-z)}, \quad (3.51)$$

其中, $\sigma(\cdot)$ 为 logistic 函数, \mathbf{x} 和 \mathbf{w} 为增广的输入向量和权重向量。

$y = 0$ 的后验概率为

$$P(y = 0|\mathbf{x}) = 1 - P(y = 1|\mathbf{x}) \quad (3.52)$$

$$= \frac{\exp(-\mathbf{w}^T \mathbf{x})}{1 + \exp(-\mathbf{w}^T \mathbf{x})}. \quad (3.53)$$

根据公式3.51,

$$\mathbf{w}^T \mathbf{x} = \log \frac{P(y = 1|\mathbf{x})}{1 - P(y = 1|\mathbf{x})} \quad (3.54)$$

$$= \log \frac{P(y = 1|\mathbf{x})}{P(y = 0|\mathbf{x})}. \quad (3.55)$$

基于梯度下降法的训练算法

给定 N 个样本 $(x^{(i)}, y^{(i)})$, $1 \leq i \leq N$, 我们使用交叉熵损失函数, 模型在训练集的风险函数为:

$$\mathcal{R}(\mathbf{w}) = - \sum_{i=1}^N \left(y^{(i)} \log \left(\sigma(\mathbf{w}^T \mathbf{x}^{(i)}) \right) + (1 - y^{(i)}) \log \left(1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) \right) \right) \quad (3.56)$$

$$= - \sum_{i=1}^N \left(y^{(i)} \log \left(\frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x}^{(i)})} \right) + (1 - y^{(i)}) \log \left(1 - \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x}^{(i)})} \right) \right) \quad (3.57)$$

$$= - \sum_{i=1}^N \left(y^{(i)} \log \left(\frac{\exp(\mathbf{w}^T \mathbf{x}^{(i)})}{1 + \exp(\mathbf{w}^T \mathbf{x}^{(i)})} \right) + (1 - y^{(i)}) \log \left(\frac{1}{1 + \exp(\mathbf{w}^T \mathbf{x}^{(i)})} \right) \right) \quad (3.58)$$

$$= - \sum_{i=1}^N \left(\mathbf{w}^T \mathbf{x}^{(i)} y^{(i)} - \log \left(1 + \exp(\mathbf{w}^T \mathbf{x}^{(i)}) \right) \right) \quad (3.59)$$

$$(3.60)$$

采样梯度下降法， $\mathcal{R}(\mathbf{w})$ 关于 \mathbf{w} 的梯度为：

$$\frac{\partial \mathcal{R}(\mathbf{w})}{\partial \mathbf{w}} = - \sum_{i=1}^N \left(\mathbf{x}^{(i)} y^{(i)} - \mathbf{x}^{(i)} \cdot \exp(\mathbf{w}^T \mathbf{x}^{(i)}) \cdot \frac{1}{1 + \exp(\mathbf{w}^T \mathbf{x}^{(i)})} \right) \quad (3.61)$$

$$= - \sum_{i=1}^N \left(\mathbf{x}^{(i)} y^{(i)} - \mathbf{x}^{(i)} \cdot \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x}^{(i)})} \right) \quad (3.62)$$

$$= - \sum_{i=1}^N \left(\mathbf{x}^{(i)} \cdot \left(y^{(i)} - \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) \right) \right) \quad (3.63)$$

$$= \sum_{i=1}^N \left(\mathbf{x}^{(i)} \cdot \left(\sigma(\mathbf{w}^T \mathbf{x}^{(i)}) - y^{(i)} \right) \right) \quad (3.64)$$

我们可以初始化 $\mathbf{w}_0 = 0$ ，然后用梯度下降法进行更新

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \lambda \frac{\partial \mathcal{R}(\mathbf{w})}{\partial \mathbf{w}}, \quad (3.65)$$

其中 λ 是学习率。

3.5.3 Softmax 回归

Softmax 回归是 logistic 回归的多类推广。

多类分类问题

对于多类分类问题（假设类别数为 $C(C > 2)$ ），有很多种利用判别函数的方法。比较常用的是把多类分类问题转换为两类分类问题，在得到两类分类结果后，还需要通过投票方法进一步确定多类的分类结果。一般有两种多类转两类的转换方式：

1. 把多类分类问题转换为 C 个两类分类问题，构建 C 个一对多的分类器。每个两类分类问题都是把某一类和其他类用一个超平面分开。
2. 把多类分类问题转换为 $C(C-1)/2$ 个两类分类问题，构建 $C(C-1)/2$ 个两两分类器。每个两类分类问题都是把 C 类中某两类用一个超平面分开。

当对一个样本进行分类时，首先用多个两个分类器进行分类，然后进行投票，选择一个得分最高的类别。

但是上面两种转换方法都存在一个缺陷：空间中的存在一些区域，这些区域中点的类别是不能区分确定的。因为如果用多个两个分类器对这些点进行分，然后进行投票时，会发现有两个或更多个类别的得分是一样的。

为了避免上述缺陷，可以使用一个更加有效的决策规则，直接建立多类线性分类器。假设 $y = \{1, \dots, C\}$ 共 C 个类别，首先定义 C 个判别函数：

$$f_c(\mathbf{x}) = \mathbf{w}_c^T \mathbf{x}, \quad c = 1, \dots, C, \quad (3.66)$$

这里 \mathbf{w}_c 为类 c 的权重向量。

这样，对于空间中的一个点 \mathbf{x} ，如果存在类别 c ，对于所有的其他类别 $\tilde{c}(\mathbf{w}_c^T \mathbf{x} \neq c)$ 都满足 $f_c(\mathbf{x}) > f_{\tilde{c}}(\mathbf{x})$ ，那么 \mathbf{x} 属于类别 c 。相应的分类函数可以表示为：

$$\hat{y} = \arg \max_{c=1}^C \mathbf{w}_c^T \mathbf{x} \quad (3.67)$$

当 $C = 2$ 时，

$$\begin{aligned} \hat{y} &= \arg \max_{y \in \{0,1\}} f_y(\mathbf{x}) \\ &= I(f_1(\mathbf{x}) - f_0(\mathbf{x}) > 0) \\ &= I(\mathbf{w}_1^T \mathbf{x} - \mathbf{w}_0^T \mathbf{x} > 0) \\ &= I((\mathbf{w}_1 - \mathbf{w}_0)^T \mathbf{x} > 0) \end{aligned} \quad (3.68)$$

这里, $I()$ 是指示函数。对比公式3.71中的两类分类判别函数, 可以发现 $\mathbf{w} = \mathbf{w}_1 - \mathbf{w}_0$ 。

多类线性分类函数的参数 \mathbf{w} 也有很多种学习方式。这里我们介绍一种常用的学习算法: **softmax 回归**。在 softmax 回归中, 机器学习模型预测目标为每一个类别的后验概率。这就需要用到 softmax 函数。

利用 softmax 函数, 我们定义目标类别 $y = c$ 的后验概率为:

$$P(y = c|\mathbf{x}) = \text{softmax}(\mathbf{w}_c^T \mathbf{x}) \quad (3.69)$$

$$= \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_{i=1}^C \exp(\mathbf{w}_i^T \mathbf{x})}. \quad (3.70)$$

对于样本 (\mathbf{x}, y) , 输出目标 $y = \{1, \dots, C\}$, 我们用 C 维的 one-hot 向量 \mathbf{y} 来表示输出目标。对于类别 c ,

$$\mathbf{y} = [I(1 = c), I(2 = c), \dots, I(C = c)]^T, \quad (3.71)$$

这里, $I()$ 是指示函数。

同时, 我们将公式3.69重新定义一下, 直接输出 k 维向量。

$$\begin{aligned} \hat{\mathbf{y}} &= \text{softmax}(W^T \mathbf{x}) \\ &= \frac{\exp(W^T \mathbf{x})}{\mathbf{1}^T \exp(W^T \mathbf{x})} \\ &= \frac{\exp(\mathbf{z})}{\mathbf{1}^T \exp(\mathbf{z})}, \end{aligned} \quad (3.72)$$

其中, $W = [\mathbf{w}_1, \dots, \mathbf{w}_C]$ 是 C 个类对应权重向量组成的矩阵。 $\hat{\mathbf{y}}$ 的第 c 维的值是第 c 类的预测后验概率。其中, $\mathbf{z} = W^T \mathbf{x}$, 为 softmax 函数的输入向量。

基于梯度下降法的训练算法

给定 N 给样本 $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$, $1 \leq i \leq N$, 我们使用交叉熵损失函数, 模型在训练集的风险函数为:

$$\begin{aligned} \mathcal{R}(W) &= - \sum_{i=1}^N \sum_{c=1}^C \mathbf{y}_c^{(i)} \log \hat{\mathbf{y}}_c^{(i)} \\ &= - \sum_{i=1}^N (\mathbf{y}^{(i)})^T \log \hat{\mathbf{y}}^{(i)} \end{aligned}$$

$$\begin{aligned}
&= - \sum_{i=1}^N (\mathbf{y}^{(i)})^T \log (\text{softmax}(\mathbf{z}^{(i)})) \\
&= - \sum_{i=1}^N (\mathbf{y}^{(i)})^T \log (\text{softmax}(W^T \mathbf{x}^{(i)})). \quad (3.73)
\end{aligned}$$

采样梯度下降法，我们要计算 $\mathcal{R}(W)$ 关于 W 的梯度。首先，我们列下要用到的公式。

(1) softmax 函数的导数为

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \frac{\partial \text{softmax}(\mathbf{x})}{\partial \mathbf{x}} \quad (3.74)$$

$$= \mathbf{diag}(\text{softmax}(\mathbf{x})) - \text{softmax}(\mathbf{x}) \text{softmax}(\mathbf{x})^T \quad (3.75)$$

$$= \mathbf{diag}(\mathbf{y}) - \mathbf{y}\mathbf{y}^T. \quad (3.76)$$

(2) $\mathbf{z} = W^T \mathbf{x}$ ，则

$$\frac{\partial \mathbf{z}}{\partial \mathbf{w}_c} = M(\mathbf{x}, c), \quad (3.77)$$

$M(\mathbf{x}, c)$ 为第 c 列为 \mathbf{x} ，其余为 0 的矩阵。

(3) $\mathbf{x}^T \mathbf{diag}(\mathbf{x})^{-1} = \mathbf{1}_C^T$

(4) 因为 \mathbf{y} 为 onehot 向量，所以 $\mathbf{1}_K^T \mathbf{y} = 1$

采样梯度下降法， $\mathcal{R}(W)$ 关于 \mathbf{w}_c 的梯度为：

$$\begin{aligned}
\frac{\partial \mathcal{R}(W)}{\partial \mathbf{w}_c} &= - \sum_{i=1}^N \frac{\partial ((\mathbf{y}^{(i)})^T \log(\hat{\mathbf{y}}^{(i)}))}{\partial \mathbf{w}_c} \\
&= - \sum_{i=1}^N \frac{\partial \mathbf{z}^{(i)}}{\partial \mathbf{w}_c} \frac{\partial \hat{\mathbf{y}}^{(i)}}{\partial \mathbf{z}^{(i)}} \frac{\partial \log \hat{\mathbf{y}}^{(i)}}{\partial \hat{\mathbf{y}}^{(i)}} \mathbf{y}^{(i)} \\
&= - \sum_{i=1}^N \frac{\partial \mathbf{z}^{(i)}}{\partial \mathbf{w}_c} \frac{\partial \text{softmax}(\mathbf{z}^{(i)})}{\partial \mathbf{z}^{(i)}} \frac{\partial \log \hat{\mathbf{y}}^{(i)}}{\partial \hat{\mathbf{y}}^{(i)}} \mathbf{y}^{(i)} \\
&= - \sum_{i=1}^N M(\mathbf{x}^{(i)}, c) \left(\mathbf{diag}(\hat{\mathbf{y}}^{(i)}) - \hat{\mathbf{y}}^{(i)} (\hat{\mathbf{y}}^{(i)})^T \right) \left(\mathbf{diag}(\hat{\mathbf{y}}^{(i)}) \right)^{-1} \mathbf{y}^{(i)} \\
&= - \sum_{i=1}^N M(\mathbf{x}^{(i)}, c) \left(I - \hat{\mathbf{y}}^{(i)} \mathbf{1}_C^T \right) \mathbf{y}^{(i)}
\end{aligned}$$

$$\begin{aligned}
&= -\sum_{i=1}^N M(\mathbf{x}^{(i)}, c) \left((\mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)}) \mathbf{1}_C^T \mathbf{y}^{(i)} \right) \\
&= -\sum_{i=1}^N M(\mathbf{x}^{(i)}, c) \left(\mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)} \right) \\
&= -\sum_{i=1}^N \mathbf{x}^{(i)} \left(\mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)} \right)_c
\end{aligned} \tag{3.78}$$

如果采用 $y \in \{1, \dots, C\}$ 的离散表示形式，梯度公式也可以写为：

$$\frac{\partial \mathcal{R}(W)}{\partial \mathbf{w}_c} = -\frac{1}{N} \sum_{i=1}^C \mathbf{x}^{(i)} \left(I(y^{(i)} = c) - p(y^{(i)} = c | \mathbf{x}^{(i)}; W) \right), \tag{3.79}$$

其中， $I()$ 为指示函数。

根据公式3.5.3，我们可以得到

$$\frac{\partial \mathcal{R}(W)}{\partial W} = -\sum_{i=1}^N \mathbf{x}^{(i)} \left(\mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)} \right)^T \tag{3.80}$$

$$= \sum_{i=1}^N \mathbf{x}^{(i)} \left(\hat{\mathbf{y}}^{(i)} - \mathbf{y}^{(i)} \right)^T \tag{3.81}$$

我们可以初始化 $W_0 = 0$ ，然后用梯度下降法进行更新

$$W_{t+1} = W_t + \alpha \sum_{i=1}^N \mathbf{x}^{(i)} \left(\hat{\mathbf{y}}^{(i)} - \mathbf{y}^{(i)} \right)^T, \tag{3.82}$$

其中 α 是学习率。

3.5.4 决策树

决策树（Decision Tree）是一种简单但是广泛使用的分类器。通过训练数据构建决策树，可以高效的对未知的数据进行分类。决策树有两大优点：1）决策树模型可读性好，具有描述性，有助于人工分析；2）效率高，决策树只需要一次构建，反复使用，每一次预测的最大计算次数不超过决策树的深度。

决策树的典型算法有 ID3, C4.5, CART 等。相对于其它算法，决策树易于理解和实现，人们在通过解释后都有能力去理解决策树所表达的意义。决策树

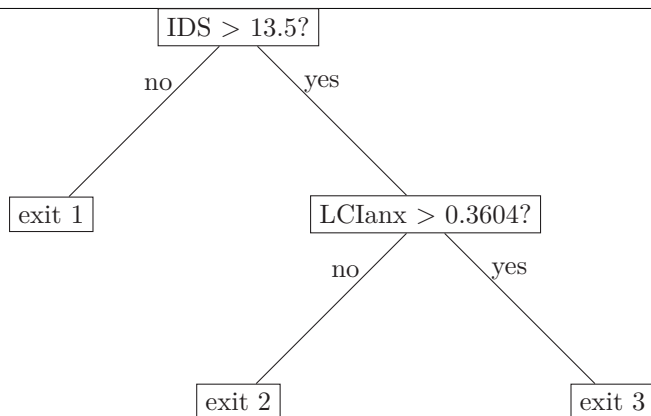


图 3.8: 决策树示例

可以同时处理不同类型的属性,并且在相对短的时间内能够对大型数据源做出可行且效果良好的结果。

3.5.5 朴素贝叶斯分类器

贝叶斯分类器根据在给定样本情况下不同类别的后验概率进行决策分类。

$$\hat{y} = \arg \max_{c=1}^C p(c|\mathbf{x}), \quad (3.83)$$

其中, $p(c|\mathbf{x})$ 为第 c 类的后验概率。

根据贝叶斯公式,上式可写为:

$$\hat{y} = \arg \max_{c=1}^C p(c|\mathbf{x}) \quad (3.84)$$

$$= \arg \max_{c=1}^C \frac{p(\mathbf{x}|c)p(c)}{\sum_c p(\mathbf{x}|c)p(c)} \quad (3.85)$$

$$\propto \arg \max_{c=1}^C p(\mathbf{x}|c)p(c). \quad (3.86)$$

这里需要估计样本 \mathbf{x} 的多维变量的联合条件概率 $p(\mathbf{x}|c)$ 。当 \mathbf{x} 为高维向量时,这是比较困难的工作。因此在实际应用中,经常假设样本每一维之间是彼此独立的。 $p(\mathbf{x}|c)$ 近似为多维向量 \mathbf{x} 中每一维变量的条件概率 $p(\mathbf{x}_i|c)$ 的乘积。

$$p(\mathbf{x}|c) = \prod_{i=1}^m p(\mathbf{x}_i|c). \quad (3.87)$$

其中， m 是样本 \mathbf{x} 的维数， \mathbf{x}_i 是 \mathbf{x} 第 i 维的值。

这样，公式3.86就可以写为

$$\hat{y} \propto \arg \max_{c=1}^C \prod_{i=1}^m p(\mathbf{x}_i|c). \quad (3.88)$$

公式3.88就称为朴素贝叶斯 (Naïve Bayes, NB) 分类器。其主要假设就是样本的每一维特征之间是彼此独立的。但这种情况在实际情况中一般不成立，因此其分类准确率可能会下降。但在许多场合，朴素贝叶斯分类算法的假设依然取得很好的性能，并且十分简单，可以与很多复杂的分类算法相媲美，是机器学习中最为常用的算法之一。

3.5.6 k 最近邻算法

k 最近邻算法 (k-Nearest Neighbor, kNN) 是根据待分样本的最相似的几个样本的类别来决定其所属类别，是最简单有效的机器学习算法之一。

对于目标函数形式不作过多的假设的算法称为非参数机器学习算法。通过不做假设，算法可以自由地从训练数据中学习任意形式的函数。

当你拥有许多数据而先验知识很少时，非参数学习通常很有用，此时你不需要关注于参数的选取。

对于一个样本，我们可以在特征空间中找到和它最相似 (即特征空间中最近邻) 的 k 个样本，如果 k 个样本中的大多数属于某一个类别，则该样本也属于这个类别。当 $k = 1$ 时，也称为**最近邻 (NN) 算法**。

kNN 方法相当于非参数密度估计方法，在决策时只与极少量的相邻样本有关。由于 kNN 方法主要靠周围有限的邻近的样本，因此对于类域的交叉或重叠较多的非线性可分数据来说，kNN 方法较其他方法更为适合，可以很好的克服了线性不可分问题的缺陷。

在 kNN 算法中，所选择的邻居都是已经有类别标记的样本，这也意味着 kNN 算法不需要训练阶段。因此，kNN 算法也很适用于分类标准随时会产生变化的需求。只要删除旧的标记样本，添加新的标记样本，就改变了分类准则。

但是，kNN 算法也存在一些不足之处。

1. 在判断一个样本的类别时，需要把它与所有已知类别的样本都比较一遍，这样计算开销是相当大的。虽然可以通过对样本空间建立索引来提高找到最近邻的效率，但是效率依然低于其他分类器。
2. 当样本不平衡时，如一个类的样本容量很大，而其他类样本容量很小时，有可能导致当输入一个新样本时，该样本的k个邻居中大容量类的样本占多数，导致分类错误。
3. kNN 算法很容易受到噪声特征的影响，在进行 kNN 分类之前需要进行特征选择。

3.5.7 支持向量机

支持向量机（Support Vector Machine, SVM）是一个经典的监督学习算法，它在解决很多任务中表现很强的优势。

支持向量机引入**边际距离**（Margin）的概念，给定训练样本 $\{(x_i, y_i)\}_{i=1}^n$ ，我们定义

$$\gamma_i(\mathbf{w}) = \mathbf{w}^T \mathbf{x} - \arg \max_{\hat{y} \neq y} \theta^T \phi(x, \hat{y}). \quad (3.89)$$

支持向量机的目标是寻找一个权重向量 θ^* 使得，

$$\theta^* = \arg \max_{\theta: \|\theta\|=1} \min_{i=1}^n \gamma_i(\theta). \quad (3.90)$$

公式3.90可以写为：

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{2} \|\theta\|^2 \\ \text{s.t.} \quad & \gamma_i(w) \geq 1, (\forall i) \end{aligned} \quad (3.91)$$

公式3.91的约束条件比较强，为了能够容忍部分不满足约束的样本，可以引入松弛变量 ξ ：

$$\begin{aligned} \min_{\theta, \xi} \quad & \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & \gamma_i(w) \geq 1 - \xi_i, \xi_i \geq 0, (\forall i) \end{aligned} \quad (3.92)$$

支持向量机的求解可以通过二次优化方法得到全局最优解，这使它有着其他统计学习技术难以比拟的优越性。同时，还使用核函数将原始的样本空间向高维空间进行变换，能够解决原始样本线性不可分的问题。支持向量机的优点在于通用性较好，且分类精度高、分类速度快、分类速度与训练样本个数无关，是最常用的分类器之一。

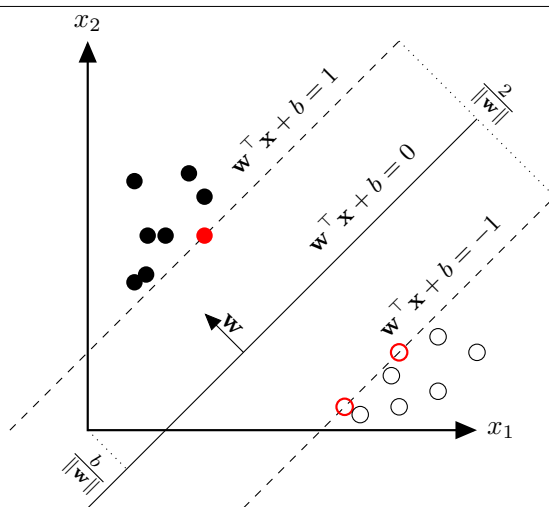


图 3.9: 支持向量机

3.6 评价方法

为了衡量一个分类算法好坏，需要给定一个测试集，用分类器对测试集中的每一个样本进行分类，并根据分类结果计算评价分数。常见的评价标准有正确率、准确率、召回率和F值等。

给定测试集 $T = (x_1, y_1), \dots, (x_N, y_N)$ ，对于所有的 $y_i \in \{\omega_1, \dots, \omega_C\}$ 。假设分类结果为 $Y = \hat{y}_1, \dots, \hat{y}_N$ 。

则**正确率** (Accuracy, Correct Rate) 为：

$$Acc = \frac{\sum_{i=1}^N |y_i = \hat{y}_i|}{N} \quad (3.93)$$

其中， $|\cdot|$ 为指示函数

和正确率相对应的就是**错误率** (Error Rate)。

$$Err = \frac{\sum_{i=1}^N |y_i \neq \hat{y}_i|}{N} \quad (3.94)$$

正确率是平均的整体性能。

在很多情况下，我们需要对每个类都进行性能估计，这就需要计算准确率和召回率。正确率和召回率是广泛用于信息检索和统计学分类领域的两个度量值，在机器学习的评价中也被大量使用。

准确率 (Precision, P)，也叫查准率，精确率或精度，是识别出的个体总数中正确识别的个体总数的比例。对于类 c 来说，

$$P_c = \frac{\sum_{\substack{i=1 \\ \hat{y}_i=c}}^N |y_i = \hat{y}_i|}{\sum_{\substack{i=1 \\ \hat{y}_i=c}}^N 1} \quad (3.95)$$

召回率 (Recall, R)，也叫查全率，是测试集中存在的个体总数中正确识别的个体总数的比例。

$$R_c = \frac{\sum_{\substack{i=1 \\ y_i=c}}^N |y_i = \hat{y}_i|}{\sum_{\substack{i=1 \\ y_i=c}}^N 1} \quad (3.96)$$

F1 值是根据正确率和召回率二者给出的一个综合的评价指标，具体定义如下：

$$F1_c = \frac{P_c * R_c * 2}{(P_c + R_c)} \quad (3.97)$$

为了计算分类算法在整个数据集上的总体准确率、召回率和 F1 值，经常使用两种平均方法，分别称为**宏平均** (macro average) 和**微平均** (micro average)。

宏平均是每一个类的性能指标的算术平均值，

$$R_{macro} = \sum_{i=1}^C R_c / C, \quad (3.98)$$

$$P_{macro} = \sum_{i=1}^C P_c / C, \quad (3.99)$$

$$F1_{macro} = \frac{P_{macro} * R_{macro} * 2}{(P_{macro} + R_{macro})}. \quad (3.100)$$

而微平均是每一个样本的性能指标的算术平均。对于单个样本而言，它的准确率和召回率是相同的（要么都是 1，要么都是 0）因此准确率和召回率的微平均是相同的，根据 F1 值公式，对于同一个数据集它的准确率、召回率和 F1 的微平均指标是相同的。

3.7 定理

在机器学习中，有一些非常有名的定理。这些定理对理解机器学习的内在特性非常有帮助。

3.7.1 没有免费午餐定理

没有免费午餐定理(No Free Lunch Theorem, NFL)是由 Wolpert 和 Macerday 在最优化理论中提出的。没有免费午餐定理证明：对于基于迭代的最优化算法，不存在某种算法对所有问题（有限的搜索空间内）都有效。如果一个算法对某些问题有效，那么它一定在另外一些问题上比纯随机搜索算法更差。也就是说，不能脱离具体问题来谈论算法的优劣，任何算法都有局限性。必须要“具体问题具体分析”。

没有免费午餐定理对于机器学习算法也同样适用。不存在一种机器学习算法适合于任何领域或任务。如果有人宣称自己的模型在所有问题上都好于其他模型，那么他肯定是在吹牛。

3.7.2 丑小鸭定理

丑小鸭定理 (Ugly Duckling) 1960年代，美籍日本学者渡边慧提出了丑小鸭定理：“丑小鸭与白天鹅之间的区别和两只白天鹅之间的区别一样大”。这个定理初看好像不符合常识，但是仔细思考后是非常有道理的。因为世界上不存在相似性的客观标准，一切相似性的标准都是主观的。如果以体型大小的角度来看，丑小鸭和白天鹅的区别大于两只白天鹅的区别；但是如果以基因的角度来看，丑小鸭与它父母的差别要小于他父母和其他白天鹅之间的差别。

渡边慧，模式识别研究的鼻祖之一。

这里的“丑小鸭”是指白天鹅的幼雏，而不是“丑陋的小鸭子”。

3.8 总结和深入阅读

本章简单地介绍了机器学习的理论知识，主要为后面讲解人工神经网络铺垫一些基础知识。如果需要快速全面地了解机器学习的基本概念可以阅读《Pattern Classification》[Duda et al., 2001]和《Pattern Recognition and Machine

Learning》[Bishop, 2006], 进一步深入了解可以阅读《The Elements of Statistical Learning》[Hastie et al., 2001] 以及《Learning in Graphical Models》[Jordan, 1998]。

参考文献

- C.M. Bishop. *Pattern recognition and machine learning*. Springer New York., 2006.
- 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.
- T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer, New York, 2001.
- M.I. Jordan. *Learning in Graphical Models*. Kluwer Academic Publishers, 1998.
- Yann LeCun, Corinna Cortes, and Christopher JC Burges. MNIST handwritten digit database. Online, 1998. URL <http://yann.lecun.com/exdb/mnist>.
- T.M. Mitchell. *Machine learning*. Burr Ridge, IL: McGraw Hill, 1997.
- David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *Cognitive modeling*, 5:3, 1988.
- Matthew D Zeiler. Adadelta: An adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.
-