

$$\left(\frac{\partial x}{\partial \mathbf{A}}\right)_{ij} = \frac{\partial x}{\partial A_{ij}}. \quad (\text{A.19})$$

对于函数 $f(\mathbf{x})$, 假定其对向量的元素可导, 则 $f(\mathbf{x})$ 关于 \mathbf{x} 的一阶导数是一个向量, 其第 i 个分量为

$$(\nabla f(\mathbf{x}))_i = \frac{\partial f(\mathbf{x})}{\partial x_i}, \quad (\text{A.20})$$

$f(\mathbf{x})$ 关于 \mathbf{x} 的二阶导数是称为海森矩阵(Hessian matrix)的一个方阵, 其第 i 行第 j 列上的元素为

$$(\nabla^2 f(\mathbf{x}))_{ij} = \frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j}. \quad (\text{A.21})$$

向量和矩阵的导数满足乘法法则(product rule)

\mathbf{a} 相对于 \mathbf{x} 为常向量.

$$\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}, \quad (\text{A.22})$$

$$\frac{\partial \mathbf{A} \mathbf{B}}{\partial \mathbf{x}} = \frac{\partial \mathbf{A}}{\partial \mathbf{x}} \mathbf{B} + \mathbf{A} \frac{\partial \mathbf{B}}{\partial \mathbf{x}}. \quad (\text{A.23})$$

由 $\mathbf{A}^{-1} \mathbf{A} = \mathbf{I}$ 和 式(A.23), 逆矩阵的导数可表示为

$$\frac{\partial \mathbf{A}^{-1}}{\partial \mathbf{x}} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial \mathbf{x}} \mathbf{A}^{-1}. \quad (\text{A.24})$$

若求导的标量是矩阵 \mathbf{A} 的元素, 则有

$$\frac{\partial \text{tr}(\mathbf{A} \mathbf{B})}{\partial A_{ij}} = B_{ji}, \quad (\text{A.25})$$

$$\frac{\partial \text{tr}(\mathbf{A} \mathbf{B})}{\partial \mathbf{A}} = \mathbf{B}^T. \quad (\text{A.26})$$

进而有

$$\frac{\partial \text{tr}(\mathbf{A}^T \mathbf{B})}{\partial \mathbf{A}} = \mathbf{B}, \quad (\text{A.27})$$

$$\frac{\partial \text{tr}(\mathbf{A})}{\partial \mathbf{A}} = \mathbf{I}, \quad (\text{A.28})$$

$$\frac{\partial \text{tr}(\mathbf{A} \mathbf{B} \mathbf{A}^T)}{\partial \mathbf{A}} = \mathbf{A}(\mathbf{B} + \mathbf{B}^T). \quad (\text{A.29})$$

由式(A.15)和(A.29)有

$$\frac{\partial \|\mathbf{A}\|_F^2}{\partial \mathbf{A}} = \frac{\partial \text{tr}(\mathbf{A}\mathbf{A}^T)}{\partial \mathbf{A}} = 2\mathbf{A} . \quad (\text{A.30})$$

链式法则(chain rule)是计算复杂导数时的重要工具. 简单地说, 若函数 f 是 g 和 h 的复合, 即 $f(x) = g(h(x))$, 则有

$$\frac{\partial f(x)}{\partial x} = \frac{\partial g(h(x))}{\partial h(x)} \cdot \frac{\partial h(x)}{\partial x} . \quad (\text{A.31})$$

例如在计算下式时, 将 $\mathbf{Ax} - \mathbf{b}$ 看作一个整体可简化计算:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{x}} (\mathbf{Ax} - \mathbf{b})^T \mathbf{W} (\mathbf{Ax} - \mathbf{b}) &= \frac{\partial (\mathbf{Ax} - \mathbf{b})}{\partial \mathbf{x}} \cdot 2\mathbf{W} (\mathbf{Ax} - \mathbf{b}) \\ &= 2\mathbf{AW} (\mathbf{Ax} - \mathbf{b}) . \end{aligned} \quad (\text{A.32})$$

A.3 奇异值分解

任意实矩阵 $\mathbf{A} \in \mathbb{R}^{m \times n}$ 都可分解为

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T , \quad (\text{A.33})$$

其中, $\mathbf{U} \in \mathbb{R}^{m \times m}$ 是满足 $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ 的 m 阶酉矩阵(unitary matrix); $\mathbf{V} \in \mathbb{R}^{n \times n}$ 是满足 $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ 的 n 阶酉矩阵; $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$ 是 $m \times n$ 的矩阵, 其中 $(\mathbf{\Sigma})_{ii} = \sigma_i$ 且其他位置的元素均为 0, σ_i 为非负实数且满足 $\sigma_1 \geq \sigma_2 \geq \dots \geq 0$.

常将奇异值按降序排列以确保 $\mathbf{\Sigma}$ 的唯一性.

当 \mathbf{A} 为对称正定矩阵时, 奇异值分解与特征值分解结果相同.

式(A.33)中的分解称为奇异值分解(Singular Value Decomposition, 简称 SVD), 其中 \mathbf{U} 的列向量 $\mathbf{u}_i \in \mathbb{R}^m$ 称为 \mathbf{A} 的左奇异向量(left-singular vector), \mathbf{V} 的列向量 $\mathbf{v}_i \in \mathbb{R}^n$ 称为 \mathbf{A} 的右奇异向量(right-singular vector), σ_i 称为奇异值(singular value). 矩阵 \mathbf{A} 的秩(rank)就等于非零奇异值的个数.

奇异值分解有广泛的用途, 例如对于低秩矩阵近似(low-rank matrix approximation)问题, 给定一个秩为 r 的矩阵 \mathbf{A} , 欲求其最优 k 秩近似矩阵 $\tilde{\mathbf{A}}$, $k \leq r$, 该问题可形式化为

$$\begin{aligned} \min_{\tilde{\mathbf{A}} \in \mathbb{R}^{m \times n}} \quad & \|\mathbf{A} - \tilde{\mathbf{A}}\|_F \\ \text{s.t.} \quad & \text{rank}(\tilde{\mathbf{A}}) = k . \end{aligned} \quad (\text{A.34})$$

奇异值分解提供了上述问题的解析解: 对矩阵 \mathbf{A} 进行奇异值分解后, 将矩阵 Σ 中的 $r - k$ 个最小的奇异值置零获得矩阵 Σ_k , 即仅保留最大的 k 个奇异值, 则

$$\mathbf{A}_k = \mathbf{U}_k \Sigma_k \mathbf{V}_k^T \quad (\text{A.35})$$

就是式(A.34)的最优解, 其中 \mathbf{U}_k 和 \mathbf{V}_k 分别是式(A.33)中的前 k 列组成的矩阵. 这个结果称为 Eckart-Young-Mirsky 定理.

B 优化

B.1 拉格朗日乘子法

拉格朗日乘子法(Lagrange multipliers)是一种寻找多元函数在一组约束下的极值的方法. 通过引入拉格朗日乘子, 可将有 d 个变量与 k 个约束条件的最优化问题转化为具有 $d + k$ 个变量的无约束优化问题求解.

先考虑一个等式约束的优化问题. 假定 \mathbf{x} 为 d 维向量, 欲寻找 \mathbf{x} 的某个取值 \mathbf{x}^* , 使目标函数 $f(\mathbf{x})$ 最小且同时满足 $g(\mathbf{x}) = 0$ 的约束. 从几何角度看, 该问题的目标是在由方程 $g(\mathbf{x}) = 0$ 确定的 $d - 1$ 维曲面上寻找能使目标函数 $f(\mathbf{x})$ 最小化的点. 此时不难得到如下结论:

函数等值线与约束曲面相切.

可通过反证法证明: 若梯度 $\nabla f(\mathbf{x}^*)$ 与约束曲面不正交, 则仍可在约束曲面上移动该点使函数值进一步下降.

- 对于约束表面上的任意点 \mathbf{x} , 该点的梯度 $\nabla g(\mathbf{x})$ 正交于约束曲面;
- 在最优点 \mathbf{x}^* , 目标函数在该点的梯度 $\nabla f(\mathbf{x}^*)$ 正交于约束曲面.

由此可知, 在最优点 \mathbf{x}^* , 如附图B.1所示, 梯度 $\nabla g(\mathbf{x})$ 和 $\nabla f(\mathbf{x})$ 的方向必相同或相反, 即存在 $\lambda \neq 0$ 使得

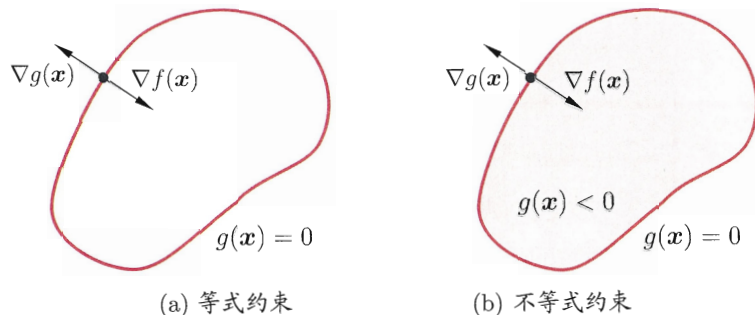
$$\nabla f(\mathbf{x}^*) + \lambda \nabla g(\mathbf{x}^*) = 0, \quad (\text{B.1})$$

对等式约束, λ 可能为正也可能为负.

λ 称为拉格朗日乘子. 定义拉格朗日函数

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x}), \quad (\text{B.2})$$

不难发现, 将其对 \mathbf{x} 的偏导数 $\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda)$ 置零即得式(B.1), 同时, 将其对 λ 的偏导数 $\nabla_{\lambda} L(\mathbf{x}, \lambda)$ 置零即得约束条件 $g(\mathbf{x}) = 0$. 于是, 原约束优化问题可转化为对拉格朗日函数 $L(\mathbf{x}, \lambda)$ 的无约束优化问题.



附图B. 1 拉格朗日乘子法的几何含义: 在 (a) 等式约束 $g(\mathbf{x}) = 0$ 或 (b) 不等式约束 $g(\mathbf{x}) \leq 0$ 下, 最小化目标函数 $f(\mathbf{x})$. 红色曲线表示 $g(\mathbf{x}) = 0$ 构成的曲面, 而其围成的阴影区域表示 $g(\mathbf{x}) < 0$.

现在考虑不等式约束 $g(\mathbf{x}) \leq 0$, 如附图B. 1 所示, 此时最优点 \mathbf{x}^* 或在 $g(\mathbf{x}) < 0$ 的区域中, 或在边界 $g(\mathbf{x}) = 0$ 上. 对于 $g(\mathbf{x}) < 0$ 的情形, 约束 $g(\mathbf{x}) \leq 0$ 不起作用, 可直接通过条件 $\nabla f(\mathbf{x}) = 0$ 来获得最优点; 这等价于将 λ 置零然后对 $\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda)$ 置零得到最优点. $g(\mathbf{x}) = 0$ 的情形类似于上面等式约束的分析, 但需注意的是, 此时 $\nabla f(\mathbf{x}^*)$ 的方向必与 $\nabla g(\mathbf{x}^*)$ 相反, 即存在常数 $\lambda > 0$ 使得 $\nabla f(\mathbf{x}^*) + \lambda \nabla g(\mathbf{x}^*) = 0$. 整合这两种情形, 必满足 $\lambda g(\mathbf{x}) = 0$. 因此, 在约束 $g(\mathbf{x}) \leq 0$ 下最小化 $f(\mathbf{x})$, 可转化为在如下约束下最小化式(B.2) 的拉格朗日函数:

$$\begin{cases} g(\mathbf{x}) \leq 0; \\ \lambda \geq 0; \\ \mu_j g_j(\mathbf{x}) = 0. \end{cases} \quad (\text{B.3})$$

式(B.3)称为 Karush-Kuhn-Tucker (简称KKT)条件.

上述做法可推广到多个约束. 考虑具有 m 个等式约束和 n 个不等式约束, 且可行域 $\mathbb{D} \subset \mathbb{R}^d$ 非空的优化问题

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & h_i(\mathbf{x}) = 0 \quad (i = 1, \dots, m), \\ & g_j(\mathbf{x}) \leq 0 \quad (j = 1, \dots, n). \end{aligned} \quad (\text{B.4})$$

引入拉格朗日乘子 $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \dots, \lambda_m)^T$ 和 $\boldsymbol{\mu} = (\mu_1, \mu_2, \dots, \mu_n)^T$, 相应的拉格

朗日函数为

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i h_i(\mathbf{x}) + \sum_{j=1}^n \mu_j g_j(\mathbf{x}), \quad (\text{B.5})$$

由不等式约束引入的 KKT 条件($j = 1, 2, \dots, n$)为

$$\begin{cases} g_j(\mathbf{x}) \leq 0; \\ \mu_j \geq 0; \\ \mu_j g_j(\mathbf{x}) = 0. \end{cases} \quad (\text{B.6})$$

一个优化问题可以从两个角度来考察, 即“主问题”(primal problem)和“对偶问题”(dual problem). 对主问题(B.4), 基于式(B.5), 其拉格朗日“对偶函数”(dual function) $\Gamma: \mathbb{R}^m \times \mathbb{R}^n \mapsto \mathbb{R}$ 定义为

在推导对偶问题时, 常通过将拉格朗日乘子 $L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ 对 \mathbf{x} 求导并令导数为 0, 来获得对偶函数的表达形式.

$$\begin{aligned} \Gamma(\boldsymbol{\lambda}, \boldsymbol{\mu}) &= \inf_{\mathbf{x} \in \mathbb{D}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \\ &= \inf_{\mathbf{x} \in \mathbb{D}} \left(f(\mathbf{x}) + \sum_{i=1}^m \lambda_i h_i(\mathbf{x}) + \sum_{j=1}^n \mu_j g_j(\mathbf{x}) \right). \end{aligned} \quad (\text{B.7})$$

$\boldsymbol{\mu} \geq 0$ 表示 $\boldsymbol{\mu}$ 的分量均为非负.

若 $\tilde{\mathbf{x}} \in \mathbb{D}$ 为主问题(B.4)可行域中的点, 则对任意 $\boldsymbol{\mu} \geq 0$ 和 $\boldsymbol{\lambda}$ 都有

$$\sum_{i=1}^m \lambda_i h_i(\mathbf{x}) + \sum_{j=1}^n \mu_j g_j(\mathbf{x}) \leq 0, \quad (\text{B.8})$$

进而有

$$\Gamma(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\mathbf{x} \in \mathbb{D}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq L(\tilde{\mathbf{x}}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq f(\tilde{\mathbf{x}}). \quad (\text{B.9})$$

若主问题(B.4)的最优值为 p^* , 则对任意 $\boldsymbol{\mu} \geq 0$ 和 $\boldsymbol{\lambda}$ 都有

$$\Gamma(\boldsymbol{\lambda}, \boldsymbol{\mu}) \leq p^*, \quad (\text{B.10})$$

即对偶函数给出了主问题最优值的下界. 显然, 这个下界取决于 $\boldsymbol{\mu}$ 和 $\boldsymbol{\lambda}$ 的值. 于是, 一个很自然的问题是: 基于对偶函数能获得的最好下界是什么? 这就引出了优化问题

$$\max_{\lambda, \mu} \Gamma(\lambda, \mu) \quad \text{s.t. } \mu \geq 0. \quad (\text{B.11})$$

式(B.11)就是主问题(B.4)的对偶问题, 其中 λ 和 μ 称为“对偶变量”(dual variable). 无论主问题(B.4)的凸性如何, 对偶问题(B.11)始终是凸优化问题.

这称为 Slater 条件.

考虑式(B.11)的最优值 d^* , 显然有 $d^* \leq p^*$, 这称为“弱对偶性”(weak duality)成立; 若 $d^* = p^*$, 则称为“强对偶性”(strong duality)成立, 此时由对偶问题能获得主问题的最优下界. 对于一般的优化问题, 强对偶性通常不成立. 但是, 若主问题为凸优化问题, 如式(B.4)中 $f(x)$ 和 $g_j(x)$ 均为凸函数, $h_i(x)$ 为仿射函数, 且其可行域中至少有一点使不等式约束严格成立, 则此时强对偶性成立. 值得注意的是, 在强对偶性成立时, 将拉格朗日函数分别对原变量和对偶变量求导, 再并令导数等于零, 即可得到原变量与对偶变量的数值关系. 于是, 对偶问题解决了, 主问题也就解决了.

B.2 二次规划

二次规划(Quadratic Programming, 简称 QP)是一类典型的优化问题, 包括凸二次优化和非凸二次优化. 在此类问题中, 目标函数是变量的二次函数, 而约束条件是变量的线性不等式.

非标准二次规划问题中可以包含等式约束. 注意到等式约束能用两个不等式约束来代替; 不等式约束可通过增加松弛变量的方式转化为等式约束.

假定变量个数为 d , 约束条件的个数为 m , 则标准的二次规划问题形如

$$\begin{aligned} \min_x \quad & \frac{1}{2} x^T Q x + c^T x \\ \text{s.t.} \quad & A x \leq b, \end{aligned} \quad (\text{B.12})$$

其中 x 为 d 维向量, $Q \in \mathbb{R}^{d \times d}$ 为实对称矩阵, $A \in \mathbb{R}^{m \times d}$ 为实矩阵, $b \in \mathbb{R}^m$ 和 $c \in \mathbb{R}^d$ 为实向量, $A x \leq b$ 的每一行对应一个约束.

若 Q 为半正定矩阵, 则式(B.12)目标函数是凸函数, 相应的二次规划是凸二次优化问题; 此时若约束条件 $A x \leq b$ 定义的可行域不为空, 且目标函数在此可行域有下界, 则该问题将有全局最小值. 若 Q 为正定矩阵, 则该问题有唯一的全局最小值. 若 Q 为非正定矩阵, 则式(B.12)是有多个平稳点和局部极小点的 NP 难问题.

常用的二次规划解法有椭球法(ellipsoid method)、内点法(interior point)、增广拉格朗日法(augmented Lagrangian)、梯度投影法(gradient projection)等. 若 Q 为正定矩阵, 则相应的二次规划问题可由椭球法在多项式时间内求解.

B.3 半正定规划

半正定规划(Semi-Definite Programming, 简称SDP)是一类凸优化问题, 其中的变量可组织成半正定对称矩阵形式, 且优化问题的目标函数和约束都是这些变量的线性函数.

给定 $d \times d$ 的对称矩阵 \mathbf{X} 、 \mathbf{C} ,

$$\mathbf{C} \cdot \mathbf{X} = \sum_{i=1}^d \sum_{j=1}^d C_{ij} X_{ij}, \quad (\text{B.13})$$

若 \mathbf{A}_i ($i = 1, 2, \dots, m$) 也是 $d \times d$ 的对称矩阵, b_i ($i = 1, 2, \dots, m$) 为 m 个实数, 则半正定规划问题形如

$$\min_{\mathbf{X}} \quad \mathbf{C} \cdot \mathbf{X} \quad (\text{B.14})$$

$$\text{s.t.} \quad \mathbf{A}_i \cdot \mathbf{X} = b_i, \quad i = 1, 2, \dots, m$$

$\mathbf{X} \succeq 0$ 表示 \mathbf{X} 半正定.

$$\mathbf{X} \succeq 0.$$

半正定规划与线性规划都拥有线性的目标函数和约束, 但半正定规划中的约束 $\mathbf{X} \succeq 0$ 是一个非线性、非光滑约束条件. 在优化理论中, 半正定规划具有一定的一般性, 能将几种标准的优化问题(如线性规划、二次规划)统一起来.

常见的用于求解线性规划的内点法经过少许改造即可求解半正定规划问题, 但半正定规划的计算复杂度较高, 难以直接用于大规模问题.

B.4 梯度下降法

一阶方法仅使用目标函数的一阶导数, 不利用其高阶导数.

梯度下降法(gradient descent)是一种常用的一阶(first-order)优化方法, 是求解无约束优化问题最简单、最经典的方法之一.

考虑无约束优化问题 $\min_{\mathbf{x}} f(\mathbf{x})$, 其中 $f(\mathbf{x})$ 为连续可微函数. 若能构造一个序列 $\mathbf{x}^0, \mathbf{x}^1, \mathbf{x}^2, \dots$ 满足

$$f(\mathbf{x}^{t+1}) < f(\mathbf{x}^t), \quad t = 0, 1, 2, \dots \quad (\text{B.15})$$

则不断执行该过程即可收敛到局部极小点. 欲满足式(B.15), 根据泰勒展式有

$$f(\mathbf{x} + \Delta \mathbf{x}) \simeq f(\mathbf{x}) + \Delta \mathbf{x}^T \nabla f(\mathbf{x}), \quad (\text{B.16})$$

于是, 欲满足 $f(\mathbf{x} + \Delta\mathbf{x}) < f(\mathbf{x})$, 可选择

$$\Delta\mathbf{x} = -\gamma\nabla f(\mathbf{x}), \quad (\text{B.17})$$

每步的步长 γ_k 可不同.

其中步长 γ 是一个小常数. 这就是梯度下降法.

L -Lipschitz条件是指对于任意 \mathbf{x} , 存在常数 L 使得 $\|\nabla f(\mathbf{x})\| \leq L$ 成立.

若目标函数 $f(\mathbf{x})$ 满足一些条件, 则通过选取合适的步长, 就能确保通过梯度下降收敛到局部极小点. 例如若 $f(\mathbf{x})$ 满足 L -Lipschitz 条件, 则将步长设置为 $1/(2L)$ 即可确保收敛到局部极小点. 当目标函数为凸函数时, 局部极小点就对应着函数的全局最小点, 此时梯度下降法可确保收敛到全局最优解.

当目标函数 $f(\mathbf{x})$ 二阶连续可微时, 可将式(B.16)替换为更精确的二阶泰勒展开式, 这样就得到了牛顿法(Newton's method). 牛顿法是典型的二阶方法, 其迭代轮数远小于梯度下降法. 但牛顿法使用了二阶导数 $\nabla^2 f(\mathbf{x})$, 其每轮迭代中涉及到海森矩阵(A.21)的求逆, 计算复杂度相当高, 尤其在高维问题中几乎不可行. 若能以较低的计算代价寻找海森矩阵的近似逆矩阵, 则可显著降低计算开销, 这就是拟牛顿法(quasi-Newton method).

B.5 坐标下降法

求解极大值问题时亦称“坐标上升法”(coordinate ascent).

坐标下降法(coordinate descent)是一种非梯度优化方法, 它在每步迭代中沿一个坐标方向进行搜索, 通过循环使用不同的坐标方向来达到目标函数的局部极小值.

不妨假设目标是求解函数 $f(\mathbf{x})$ 的极小值, 其中 $\mathbf{x} = (x_1, x_2, \dots, x_d)^T \in \mathbb{R}^d$ 是一个 d 维向量. 从初始点 \mathbf{x}^0 开始, 坐标下降法通过迭代地构造序列 $\mathbf{x}^0, \mathbf{x}^1, \mathbf{x}^2, \dots$ 来求解该问题, \mathbf{x}^{t+1} 的第 i 个分量 x_i^{t+1} 构造为

$$x_i^{t+1} = \arg \min_{y \in \mathbb{R}} f(x_1^{t+1}, \dots, x_{i-1}^{t+1}, y, x_{i+1}^t, \dots, x_d^t). \quad (\text{B.18})$$

通过执行此操作, 显然有

$$f(\mathbf{x}^0) \geq f(\mathbf{x}^1) \geq f(\mathbf{x}^2) \geq \dots \quad (\text{B.19})$$

与梯度下降法类似, 通过迭代执行该过程, 序列 $\mathbf{x}^0, \mathbf{x}^1, \mathbf{x}^2, \dots$ 能收敛到所期望的局部极小点或驻点(stationary point).

坐标下降法不需计算目标函数的梯度, 在每步迭代中仅需求解一维搜索问题, 对于某些复杂问题计算较为简便. 但若目标函数不光滑, 则坐标下降法有可能陷入非驻点(non-stationary point).

C 概率分布

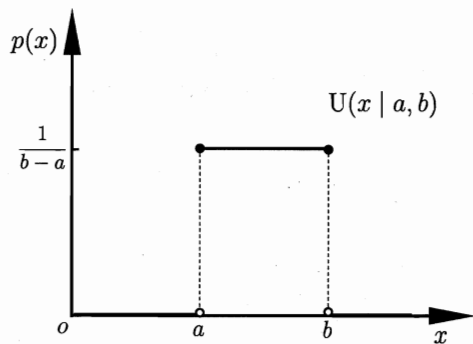
C.1 常见概率分布

本节简要介绍几种常见概率分布. 对于每种分布, 我们将给出概率密度函数以及期望 $\mathbb{E}[\cdot]$ 、方差 $\text{var}[\cdot]$ 和协方差 $\text{cov}[\cdot, \cdot]$ 等几个主要的统计量.

C.1.1 均匀分布

这里仅介绍连续均匀分布.

均匀分布(uniform distribution)是关于定义在区间 $[a, b]$ ($a < b$) 上连续变量的简单概率分布, 其概率密度函数如附图C.1 所示.



附图C.1 均匀分布的概率密度函数

$$p(x | a, b) = U(x | a, b) = \frac{1}{b-a}; \quad (\text{C.1})$$

$$\mathbb{E}[x] = \frac{a+b}{2}; \quad (\text{C.2})$$

$$\text{var}[x] = \frac{(b-a)^2}{12}. \quad (\text{C.3})$$

不难发现, 若变量 x 服从均匀分布 $U(x | 0, 1)$ 且 $a < b$, 则 $a + (b-a)x$ 服从均匀分布 $U(x | a, b)$.

C.1.2 伯努利分布

以瑞士数学家雅各布·伯努利 (Jacob Bernoulli, 1654–1705) 的名字命名.

伯努利分布(Bernoulli distribution)是关于布尔变量 $x \in \{0, 1\}$ 的概率分布, 其连续参数 $\mu \in [0, 1]$ 表示变量 $x = 1$ 的概率.

$$P(x | \mu) = \text{Bern}(x | \mu) = \mu^x (1 - \mu)^{1-x}; \quad (\text{C.4})$$

$$\mathbb{E}[x] = \mu ; \quad (\text{C.5})$$

$$\text{var}[x] = \mu(1 - \mu) . \quad (\text{C.6})$$

C.1.3 二项分布

二项分布(binomial distribution)用以描述 N 次独立的伯努利实验中有 m 次成功(即 $x = 1$)的概率, 其中每次伯努利实验成功的概率为 $\mu \in [0, 1]$.

$$P(m | N, \mu) = \text{Bin}(m | N, \mu) = \binom{N}{m} \mu^m (1 - \mu)^{N-m} ; \quad (\text{C.7})$$

$$\mathbb{E}[x] = N\mu ; \quad (\text{C.8})$$

$$\text{var}[x] = N\mu(1 - \mu) . \quad (\text{C.9})$$

对于参数 μ , 二项分布的共轭先验分布是贝塔分布. 共轭分布参见 C.2.

当 $N = 1$ 时, 二项分布退化为伯努利分布.

C.1.4 多项分布

若将伯努利分布由单变量扩展为 d 维向量 \mathbf{x} , 其中 $x_i \in \{0, 1\}$ 且 $\sum_{i=1}^d x_i = 1$, 并假设 x_i 取 1 的概率为 $\mu_i \in [0, 1]$, $\sum_{i=1}^d \mu_i = 1$, 则将得到离散概率分布

$$P(\mathbf{x} | \boldsymbol{\mu}) = \prod_{i=1}^d \mu_i^{x_i} ; \quad (\text{C.10})$$

$$\mathbb{E}[x_i] = \mu_i ; \quad (\text{C.11})$$

$$\text{var}[x_i] = \mu_i(1 - \mu_i) ; \quad (\text{C.12})$$

$$\text{cov}[x_j, x_i] = \mathbb{I}[j = i] \mu_i . \quad (\text{C.13})$$

在此基础上扩展二项分布则得到多项分布(multinomial distribution), 它描述了在 N 次独立实验中有 m_i 次 $x_i = 1$ 的概率.

对于参数 $\boldsymbol{\mu}$, 多项分布的共轭先验分布是狄利克雷分布. 共轭分布参见 C.2.

$$\begin{aligned} P(m_1, m_2, \dots, m_d | N, \boldsymbol{\mu}) &= \text{Mult}(m_1, m_2, \dots, m_d | N, \boldsymbol{\mu}) \\ &= \frac{N!}{m_1! m_2! \dots m_d!} \prod_{i=1}^d \mu_i^{m_i} ; \end{aligned} \quad (\text{C.14})$$

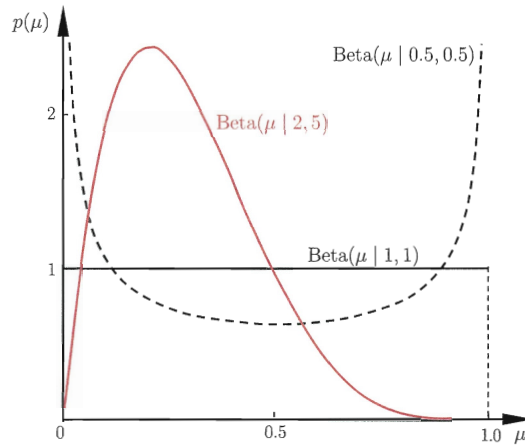
$$\mathbb{E}[m_i] = N\mu_i ; \quad (\text{C.15})$$

$$\text{var}[m_i] = N\mu_i(1 - \mu_i) ; \quad (\text{C.16})$$

$$\text{cov}[m_j, m_i] = -N\mu_j\mu_i . \quad (\text{C.17})$$

C.1.5 贝塔分布

贝塔分布(Beta distribution)是关于连续变量 $\mu \in [0, 1]$ 的概率分布, 它由两个参数 $a > 0$ 和 $b > 0$ 确定, 其概率密度函数如附图C.2 所示.



附图C. 2 贝塔分布的概率密度函数

$$\begin{aligned} p(\mu | a, b) &= \text{Beta}(\mu | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} \\ &= \frac{1}{B(a, b)} \mu^{a-1} (1-\mu)^{b-1} ; \end{aligned} \quad (\text{C.18})$$

$$\mathbb{E}[\mu] = \frac{a}{a+b} ; \quad (\text{C.19})$$

$$\text{var}[\mu] = \frac{ab}{(a+b)^2(a+b+1)} , \quad (\text{C.20})$$

其中 $\Gamma(a)$ 为 Gamma 函数

$$\Gamma(a) = \int_0^{+\infty} t^{a-1} e^{-t} dt , \quad (\text{C.21})$$

$B(a, b)$ 为 Beta 函数

$$B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} . \quad (\text{C.22})$$

当 $a = b = 1$ 时, 贝塔分布退化为均匀分布.

C.1.6 狄利克雷分布

以德国数学家狄利克雷 (1805—1859) 的名字命名.

狄利克雷分布 (Dirichlet distribution) 是关于一组 d 个连续变量 $\mu_i \in [0, 1]$ 的概率分布, $\sum_{i=1}^d \mu_i = 1$. 令 $\boldsymbol{\mu} = (\mu_1; \mu_2; \dots; \mu_d)$, 参数 $\boldsymbol{\alpha} = (\alpha_1; \alpha_2; \dots; \alpha_d)$, $\alpha_i > 0$, $\hat{\alpha} = \sum_{i=1}^d \alpha_i$.

$$p(\boldsymbol{\mu} | \boldsymbol{\alpha}) = \text{Dir}(\boldsymbol{\mu} | \boldsymbol{\alpha}) = \frac{\Gamma(\hat{\alpha})}{\Gamma(\alpha_1) \dots \Gamma(\alpha_d)} \prod_{i=1}^d \mu_i^{\alpha_i-1}; \quad (\text{C.23})$$

$$\mathbb{E}[\mu_i] = \frac{\alpha_i}{\hat{\alpha}}; \quad (\text{C.24})$$

$$\text{var}[\mu_i] = \frac{\alpha_i(\hat{\alpha} - \alpha_i)}{\hat{\alpha}^2(\hat{\alpha} + 1)}; \quad (\text{C.25})$$

$$\text{cov}[\mu_j, \mu_i] = \frac{\alpha_j \alpha_i}{\hat{\alpha}^2(\hat{\alpha} + 1)}. \quad (\text{C.26})$$

当 $d = 2$ 时, 狄利克雷分布退化为贝塔分布.

C.1.7 高斯分布

高斯分布 (Gaussian distribution) 亦称正态分布 (normal distribution), 是应用最为广泛的连续概率分布.

对于单变量 $x \in (-\infty, \infty)$, 高斯分布的参数为均值 $\mu \in (-\infty, \infty)$ 和方差 $\sigma^2 > 0$. 附图 C.3 给出了在几组不同参数下高斯分布的概率密度函数.

σ 为标准差.

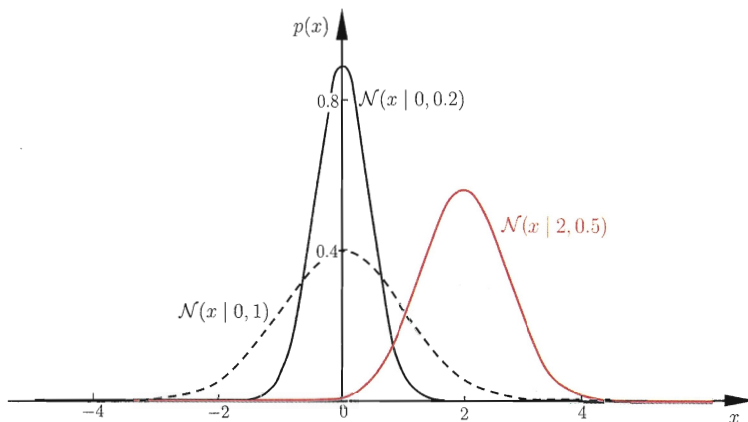
$$p(x | \mu, \sigma^2) = \mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\}; \quad (\text{C.27})$$

$$\mathbb{E}[x] = \mu; \quad (\text{C.28})$$

$$\text{var}[x] = \sigma^2. \quad (\text{C.29})$$

对于 d 维向量 \mathbf{x} , 多元高斯分布的参数为 d 维均值向量 $\boldsymbol{\mu}$ 和 $d \times d$ 的对称正定协方差矩阵 $\boldsymbol{\Sigma}$.

$$\begin{aligned} p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ &= \frac{1}{\sqrt{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\}; \end{aligned} \quad (\text{C.30})$$



附图C.3 高斯分布的概率密度函数

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}; \quad (\text{C.31})$$

$$\text{cov}[\mathbf{x}] = \boldsymbol{\Sigma}. \quad (\text{C.32})$$

C.2 共轭分布

假设变量 x 服从分布 $P(x | \Theta)$, 其中 Θ 为参数, $X = \{x_1, x_2, \dots, x_m\}$ 为变量 x 的观测样本, 假设参数 Θ 服从先验分布 $\Pi(\Theta)$. 若由先验分布 $\Pi(\Theta)$ 和抽样分布 $P(X | \Theta)$ 决定的后验分布 $F(\Theta | X)$ 与 $\Pi(\Theta)$ 是同种类型的分布, 则称先验分布 $\Pi(\Theta)$ 为分布 $P(x | \Theta)$ 或 $P(X | \Theta)$ 的共轭分布(conjugate distribution).

例如, 假设 $x \sim \text{Bern}(x | \mu)$, $X = \{x_1, x_2, \dots, x_m\}$ 为观测样本, \bar{x} 为观测样本的均值, $\mu \sim \text{Beta}(\mu | a, b)$, 其中 a, b 为已知参数, 则 μ 的后验分布

$$\begin{aligned} F(\mu | X) &\propto \text{Beta}(\mu | a, b) P(X | \mu) \\ &= \frac{\mu^{a-1} (1-\mu)^{b-1}}{B(a, b)} \mu^{m\bar{x}} (1-\mu)^{m-m\bar{x}} \\ &= \frac{1}{B(a+m\bar{x}, b+m-m\bar{x})} \mu^{a+m\bar{x}-1} (1-\mu)^{b+m-m\bar{x}-1} \\ &= \text{Beta}(\mu | a', b'), \end{aligned} \quad (\text{C.33})$$

亦为贝塔分布, 其中 $a' = a + m\bar{x}$, $b' = b + m - m\bar{x}$, 这意味着贝塔分布与伯努利分布共轭. 类似可知, 多项分布的共轭分布是狄利克雷分布, 而高斯分布的共轭分布仍是高斯分布.

这里仅考虑高斯分布方差已知、均值服从先验的情形.

先验分布反映了某种先验信息, 后验分布既反映了先验分布提供的信息、又反映了样本提供的信息. 当先验分布与抽样分布共轭时, 后验分布与先验分布属于同种类型, 这意味着先验信息与样本提供的信息具有某种同一性. 于是, 若使用后验分布作为进一步抽样的先验分布, 则新的后验分布仍将属于同种类型. 因此, 共轭分布在不少情形下会使问题得以简化. 例如在式(C.33)的例子中, 对服从伯努利分布的事件 X 使用贝塔先验分布, 则贝塔分布的参数值 a 和 b 可视为对伯努利分布的真实情况(事件发生和不发生)的预估. 随着“证据”(样本)的不断到来, 贝塔分布的参数值从 a, b 变化为 $a + m\bar{x}, b + m - m\bar{x}$, 且 $a/(a+b)$ 将随着 m 的增大趋近于伯努利分布的真实参数值 \bar{x} . 显然, 使用共轭先验之后, 只需调整 a 和 b 这两个预估值即可方便地进行模型更新.

C.3 KL散度

KL散度(Kullback-Leibler divergence), 亦称相对熵(relative entropy)或信息散度(information divergence), 可用于度量两个概率分布之间的差异. 给定两个概率分布 P 和 Q , 二者之间的KL散度定义为

$$\text{KL}(P\|Q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx, \quad (\text{C.34})$$

这里假设两个分布均为连续型概率分布; 对于离散型概率分布, 只需将定义中的积分替换为对所有离散值遍历求和.

其中 $p(x)$ 和 $q(x)$ 分别为 P 和 Q 的概率密度函数.

KL散度满足非负性, 即

$$\text{KL}(P\|Q) \geq 0, \quad (\text{C.35})$$

当且仅当 $P = Q$ 时 $\text{KL}(P\|Q) = 0$. 但是, KL散度不满足对称性, 即

$$\text{KL}(P\|Q) \neq \text{KL}(Q\|P), \quad (\text{C.36})$$

度量应满足四个基本性质, 参见9.3节.

因此, KL散度不是一个度量(metric).

若将KL散度的定义(C.34)展开, 可得

$$\begin{aligned} \text{KL}(P\|Q) &= \int_{-\infty}^{\infty} p(x) \log p(x) dx - \int_{-\infty}^{\infty} p(x) \log q(x) dx \\ &= -H(P) + H(P, Q), \end{aligned} \quad (\text{C.37})$$

其中 $H(P)$ 为熵(entropy), $H(P, Q)$ 为 P 和 Q 的交叉熵(cross entropy). 在信

息论中, 熵 $H(P)$ 表示对来自 P 的随机变量进行编码所需的最小字节数, 而交叉熵 $H(P, Q)$ 则表示使用基于 Q 的编码对来自 P 的变量进行编码所需的字节数. 因此, KL散度可认为是使用基于 Q 的编码对来自 P 的变量进行编码所需的“额外”字节数; 显然, 额外字节数必然非负, 当且仅当 $P = Q$ 时额外字节数为零.

后 记

写作本书的主因,是2016年准备在南京大学开设“机器学习”课.十五年前笔者曾主张开设此课,但那时国内对机器学习闻之不多,不少人听到这个名字的第一反应是“学习什么机器?”学校估计学生兴趣不大,于是笔者开设了“数据挖掘”这门名字听上去就觉得很有用的课.被评为省优秀研究生课程后,又给本科生单开了一门“数据挖掘导论”.这两门课很受欢迎,选修学生很多,包括不少外来蹭听生.虽然课上有一多半其实在讲机器学习,但笔者仍一直希望专开一门机器学习课,因笔者以为机器学习迟早会变成计算机学科的基础内容.

图灵奖得主 E. W. Dijkstra 曾说“计算机科学并不仅是关于计算机,就像天文学并不仅是关于望远镜”.正如天文学早期的研究关注如何制造望远镜,计算机科学早期研究是在关注如何令计算机运转.到了今天,建造强大的天文望远镜虽仍重要,但天文学更要紧的是“用”望远镜来开展研究.类似地,计算机科学发展至今,也该到了从关注“造”计算机转入更关注“用”计算机来认识和改造世界的阶段,其中最重要的无疑是用计算机对数据进行分析,因为这是计算的主要目的,而这就离不开机器学习.十多年前在国内某次重要论坛上笔者刚抛出此观点就被专家迎头指斥,但今日来看,甚至很多计算机学科外人士都已对机器学习的重大价值津津乐道,现在才开设机器学习基础课似乎已有点嫌晚了.

1995年在南大图书馆偶然翻看了《机器学习:通往人工智能的途径》,这算是笔者接触机器学习的开始.那时机器学习在国内问津者寥,甚至连科研人员申请基金项目也无合适代码方向可报.周边无专家可求教,又因国内科研经费匮乏而几无国际交流,加之学校尚无互联网和电子文献库,能看到的最新文献仅是两年前出版且页数不全的某 IEEE 汇刊……可谓举步维艰,经历的困惑和陷阱不可胜数.笔者切身体会到,入门阶段接触的书籍是何等重要,对自学者尤甚.一本好书能让人少走许多弯路,材料不佳则后续要花费数倍精力方能纠偏.中文书当然要国人自己来写.虽已不需靠“写书出名”,且深知写教科书极耗时间精力,但踌躇后笔者仍决定动手写这本书,唯望为初学者略尽绵薄之力.

有人说“一千个人眼中就有一千个哈姆雷特”,一个学科何尝不是如此.之所以不欲使用市面上流行的教科书(主要是英文的),除了觉得对大多数中国学生来说中文教科书更便于学习,另一个原因则是希望从笔者自己的视角来展现机器学习.

2013年中开始规划提纲,由此进入了焦躁的两年.该写哪些内容、先写什么后写什么、从哪个角度写、写到什么程度,总有千丝万缕需考虑.及至写作进行,更是战战兢兢,深恐不慎误人子弟.写书难,写教科书更难.两年下来,甘苦自知.子曰:“取乎其上,得乎其中;取乎其中,得乎其下”,且以顶级的态度,出一本勉强入得方家法眼之书.

本书贯穿以西瓜为例,一则因为瓜果中笔者尤喜西瓜,二则因为西瓜在笔者所生活的区域有个有趣的蕴义。朋友小聚、请客吃饭,菜已全而主未知,或饌未齐而人待走,都挺尴尬。于是聪明人发明了“潜规则”:席终上西瓜。无论整盘抑或小菜,宾主见瓜至,则心领神会准备起身,皆大欢喜。久而久之,无论菜肴价格贵贱、场所雅鄙,宴必有西瓜。若将宴席比作(未来)应用系统,菜肴比作所涉技术,则机器学习好似那必有的西瓜,它可能不是最“高大上”的,但却是离不了的、没用上总觉得不甘心的。

本书写作过程从材料搜集,到习题设计,再到阅读校勘,都得到了笔者的很多学生、同事和学术界朋友的支持和帮助,在此谨列出他们的姓名以致谢意(姓氏拼音序):陈松灿,戴望州,高阳,高尉,黄圣君,黎铭,李楠,李武军,李宇峰,钱超,王魏,王威廉,吴建鑫,徐淼,俞扬,詹德川,张利军,张敏灵,朱军。书稿在 LAMDA 组学生 2015 年暑期讨论班上试讲,高斌斌、郭翔宇、李绍园、钱鸿、沈芷玉、叶翰嘉、张腾等同学又帮助发现了许多笔误。特别感谢李楠把笔者简陋的手绘图转变为精致的插图,俞扬帮助调整排版格式和索引,刘冲把笔者对封面设计的想法具体表现出来。

中国计算机学会终身成就奖得主、中国科学院院士陆汝钤先生是我国人工智能事业的开拓者之一,他在 1988 年和 1996 年出版的《人工智能》(上、下册)曾给予笔者很多启发。承蒙陆老师厚爱,在百忙中为本书作序,不胜惶恐之至。陆老师在序言中提出的问题很值得读者在本书之后的进阶学习与研究中深思。

感谢清华大学出版社薛慧老师为本书出版所做的努力。十二年前笔者入选国家杰出青年科学基金时薛老师即邀著书,笔者以年纪尚轻、学力未逮婉辞。十年前“机器学习及其应用”研讨会(MLA)从陆汝钤院士肇始的复旦大学智能信息处理重点实验室移师南京,参会人数从复旦最初的 20 人,发展到 2010 年 400 余人,此后在清华、复旦、西电达 800 余人,今年再回南大竟至 1300 余人,场面热烈。MLA 倡导“学术至上、其余从简”,不搞繁文缛节,参会免费。但即便如此,仍有很多感兴趣的师生因旅费不菲而难以参加。于是笔者提议每两年以《机器学习及其应用》为题出版一本报告选集以飨读者。这个主意得到了薛老师、陆老师以及和笔者一起长期组织 MLA、去年因病去世的王珏老师的大力支持。此类专业性学术文集销量不大,出版社多半要贴钱。笔者曾跟薛老师说,自著的第一本中文书必交由薛老师在清华出版,或可稍为出版社找补。转眼《机器学习及其应用》系列已出到第六本,薛老师或以为十年前是玩笑话,某日告之书快完稿时她蓦然惊喜。

最后要感谢笔者的家人,本书几乎耗尽了两年来笔者所有的节假日和空闲时间。写作时垂髫小子常跑来案边,不是问“爸爸去哪儿?”而是看几眼然后问“爸爸你又写了几页?”为了给他满意的答复,笔者埋头努力。

周志华

2015 年 11 月于南京渐宽斋

索引

- 0/1损失函数, 130, 147
- 5×2 交叉验证, 41
- ϵ -贪心, 374
- AdaBoost, 173
- ART网络, 108
- Bagging, 178
- Bellman等式, 380
- Boltzmann分布, 111
- Boltzmann机, 111
- Boosting, 173, 190
- BP算法, 101
- BP网络, 101
- C4.5决策树, 78, 83
- CART决策树, 79
- ECOC, 64
- Elman网络, 111
- EM算法, 162, 208, 295, 335
- $F1$, 32
- Fisher判别分析, 60
- Friedman检验, 42
- Frobenius 范数, 400
- Hoeffding不等式, 192, 268
- hinge损失, 130
- ID3决策树, 75
- ILP, 357, 364
- Jensen不等式, 268
- K -摇臂赌博机, 373
- KKT条件, 124, 132, 135
- KL散度, 335, 414
- Kohonen网络, 109
- k 折交叉验证, 26
- k 近邻, 225
- k 均值算法, 202, 218
- L_1 正则化, 253
- L_2 正则化, 253
- LASSO, 252, 261
- Lipschitz条件, 253
- LVQ, 204, 218
- M-P神经元模型, 97
- McDiarmid不等式, 268
- MCMC, 331
- McNemar检验, 41
- MDP, 371
- Mercer定理, 137, 139
- MH算法, 333
- MvM, 63
- Nemenyi后续检验, 43
- OvO, 63
- OvR, 63
- P-R曲线, 31
- PAC辨识, 269
- PAC可学习, 269
- PAC学习算法, 270
- PCA, 229
- Q-学习, 387, 393
- Rademacher复杂度, 279
- RBF网络, 108
- ReLU, 114
- RIPPER, 353
- RKHS, 128
- ROC曲线, 33, 46
- S3VM, 298
- Sarsa 算法, 387, 390
- Sigmoid函数, 58, 98, 102
- Softmax, 375
- SOM网络, 109
- Stacking, 184
- SVM, 123
- TD学习, 386, 393
- Tikhonov正则化, 252
- VC维, 273, 274
- V型结构, 158
- WEKA, 16

- 奥卡姆剃刀, 7, 17
- 版本空间, 5
- 半监督聚类, 240, 307
- 半监督学习, 293, 294
- 半监督支持向量机, 298
- 半朴素贝叶斯分类器, 154
- 半正定规划, 407
- 包裹式特征选择, 250
- 包外估计, 27, 179
- 贝塔分布, 411
- 贝叶斯定理, 148
- 贝叶斯分类器, 164
- 贝叶斯风险, 147
- 贝叶斯决策论, 147
- 贝叶斯模型平均, 185
- 贝叶斯网, 156, 319, 339
- 贝叶斯学习, 164
- 贝叶斯最优分类器, 147
- 本真低维空间, 232
- 本真距离, 234
- 必连约束, 239, 307
- 边际独立性, 158
- 边际分布, 328
- 边际化, 158, 328
- 边际似然, 163
- 编码矩阵, 65
- 变分推断, 334
- 变量消去, 328
- 标记, 2
- 标记传播, 302
- 标记空间, 3
- 表格值函数, 388
- 表示定理, 137
- 表示学习, 114
- 伯努利分布, 409
- 不可分, 269, 272
- 不可知PAC可学习, 273
- 不一致, 269
- 参数估计, 54
- 参数调节, 28
- 参数空间, 106
- 策略, 372
- 策略迭代, 381
- 测地线距离, 234
- 测试, 3
- 测试样本, 3
- 层次聚类, 214
- 查全率, 30
- 查询, 293
- 查准率, 30
- 差异性度量, 187
- 超父, 155
- 成对马尔可夫性, 325
- 冲突消解, 348
- 重采样, 177
- 重赋权, 177
- 簇, 3, 197
- 错误率, 23, 29
- 打散, 273
- 带序规则, 348
- 代价, 35, 47
- 代价矩阵, 35
- 代价敏感, 36, 67
- 代价曲线, 36
- 单隐层网络, 101
- 道德图, 158
- 等度量映射, 234
- 低密度分隔, 298
- 低维嵌入, 226
- 低秩矩阵近似, 402
- 狄利克雷分布, 412
- 递归神经网络, 111
- 典型相关分析, 240
- 独立同分布, 3, 267
- 独依赖估计, 154
- 度量学习, 237
- 端正图, 158
- 对比散度, 112
- 对分, 273
- 对率函数, 58
- 对率回归, 58, 132, 325
- 对率损失, 130
- 对偶函数, 405
- 对偶问题, 123, 405
- 对数几率函数, 58, 98
- 对数几率回归, 57
- 对数似然, 59, 149
- 对数线性回归, 56
- 多变量决策树, 88, 92
- 多标记学习, 68
- 多层前馈神经网络, 100
- 多分类, 3
- 多分类器系统, 171
- 多分类学习, 63
- 多核学习, 140
- 多视图学习, 240, 304
- 多维缩放, 227

多项分布, 410
多样性, 172
多样性度量, 187
多元线性回归, 55

二次规划, 406
二分类, 3
二项分布, 410
二项检验, 38

发散, 113
罚函数法, 133
反向传播算法, 101
泛化, 3, 121, 350
泛化误差, 23, 267
非参数化方法, 340
非度量距离, 201
非线性降维, 232
非线性可分, 99
分层采样, 25
分而治之, 74
分类, 3
分歧, 185, 304
风险, 147
符号主义, 10, 363

概率近似正确, 268
概率模型, 206, 319
概率图模型, 156, 319
概念类, 268
概念学习, 4, 17
感知机, 98
高斯分布, 412
高斯核, 128
高斯混合, 206, 296
割平面法, 139
个体学习器, 171
功能神经元, 99
共轭分布, 413
关系学习, 363
广义 δ 规则, 115
广义瑞利商, 61
广义线性模型, 57
规范化, 36, 183
规则, 347
规则学习, 347
归结商, 362
归纳, 359
归纳逻辑程序设计, 357, 364
归纳偏好, 6

归纳学习, 4, 11
归一化, 36
过采样, 67
过滤式特征选择, 249
过拟合, 23, 104, 191, 352
过配, 23

行列式, 399
豪斯多夫距离, 220
核范数, 260
核方法, 137
核函数, 126
核化, 137, 232
核化线性降维, 232
核技巧, 127
核矩阵, 128, 138, 233
核线性判别分析, 137
核主成分分析, 232
合一, 361
宏 $F1$, 32
宏查全率, 32
宏查准率, 32
后剪枝, 79, 352
划分超平面, 121, 298
划分选择, 75, 92
话题模型, 337
回归, 3
汇合, 114
混合属性, 201
混合专家, 191, 313
混淆矩阵, 30

基尼指数, 79
基学习器, 171
基学习算法, 171
基于分歧的方法, 304
基于能量的模型, 111
机械学习, 11
迹, 399
迹范数, 260
激活函数, 98
吉布斯采样, 161, 334
极大似然法, 59, 149, 297
极大似然估计, 149, 328
集成修剪, 191
集成学习, 171, 311
急切学习, 225
级联相关, 110
挤压函数, 98
几率, 58

- 计算学习理论, 267
- 加权距离, 201
- 加权平均, 182, 225
- 加权投票, 183, 225
- 加性模型, 173
- 假设, 2, 269
- 假设检验, 37
- 假设空间, 268
- 监督学习, 3
- 间隔, 122
- 简单平均, 182
- 剪枝, 79, 352
- 奖赏, 371
- 降维, 227
- 交叉验证成对 t 检验, 40
- 交叉验证法, 26
- 交叉熵, 415
- 街区距离, 200
- 阶跃函数, 57, 98
- 结构风险, 133
- 近端梯度下降, 253, 259
- 近邻成分分析, 238
- 近似动态规划, 393
- 近似推断, 161, 328, 331
- 精度, 23, 29
- 精确推断, 161, 328
- 经验风险, 133
- 经验风险最小化, 278
- 经验误差, 23, 267
- 径向基函数, 108
- 竞争型学习, 108
- 纠错输出码, 64
- 局部极小, 106
- 局部马尔可夫性, 324
- 局部线性嵌入, 235
- 矩阵补全, 259
- 聚类, 3, 197
- 聚类集成, 219
- 聚类假设, 294
- 距离度量, 199
- 距离度量学习, 201, 237
- 卷积神经网络, 113
- 决策树, 73, 363
- 决策树桩, 82
- 绝对多数投票, 182
- 均方误差, 29, 54
- 均匀分布, 409
- 均匀稳定性, 285
- 可分, 269, 270
- 可解释性, 115, 191
- 可塑性-稳定性窘境, 109
- 拉格朗日乘子法, 403
- 拉普拉斯修正, 153
- 拉斯维加斯方法, 251
- 懒惰学习, 154, 225, 240
- 累积误差逆传播, 105
- 类比学习, 11
- 类别不平衡, 66, 299
- 类间散度矩阵, 61, 138
- 类内散度矩阵, 61, 138
- 离散化, 83
- 离散属性, 200
- 联系函数, 57
- 连接权, 101, 104
- 连接主义, 10
- 连续属性, 200
- 链式法则, 103, 402
- 列联表, 41, 187
- 列名属性, 200
- 岭回归, 252
- 留出法, 25
- 流形假设, 240, 294
- 流形学习, 234
- 流形正则化, 240
- 逻辑文字, 347
- 码书学习, 255
- 马尔可夫决策过程, 371
- 马尔可夫链, 161, 320
- 马尔可夫随机场, 322
- 马尔可夫毯, 325
- 马尔可夫网, 319
- 曼哈顿距离, 200
- 没有免费的午餐定理, 9
- 蒙特卡罗方法, 251, 340, 384
- 密采样, 226
- 密度聚类, 211
- 免模型学习, 382
- 闵可夫斯基距离, 200, 220
- 命题规则, 348
- 模仿学习, 390
- 模拟退火, 107
- 模型选择, 24
- 默认规则, 348
- 逆归结, 359
- 逆强化学习, 391
- 欧氏距离, 200

- 盘式记法, 334
- 判别式模型, 148, 325
- 判定树, 73
- 偏差-方差分解, 44, 177
- 偏好, 6
- 平方损失, 54
- 平衡点, 31
- 平均场, 337
- 平均法, 225
- 平稳分布, 161
- 朴素贝叶斯分类器, 150

- 奇异值分解, 231, 402
- 恰PAC可学习, 270
- 迁移学习, 17
- 嵌入式特征选择, 252
- 欠采样, 67
- 欠拟合, 23
- 欠配, 23
- 强化学习, 371
- 切比雪夫距离, 200
- 亲和矩阵, 301
- 权共享, 113
- 全局马尔可夫性, 323
- 全局散度矩阵, 62
- 全局最小, 106
- 缺省规则, 348
- 缺失值, 85

- 人工神经网络, 97
- 人工智能, 10
- 冗余特征, 247
- 软间隔, 129
- 软间隔支持向量机, 131
- 弱学习器, 171

- 熵, 415
- 上采样, 67
- 深度学习, 113
- 神经网络, 97
- 神经元, 97
- 生成式模型, 148, 295, 325
- 胜者通吃, 108
- 时间复杂度, 269
- 时序差分学习, 386, 393
- 示教学习, 11
- 示例, 2
- 势函数, 322
- 视图, 304
- 收敛, 99

- 受限 Boltzmann 机, 112
- 属性, 2
- 属性空间, 2
- 属性子集, 189
- 数据集, 2
- 数据挖掘, 14
- 数据预处理, 247
- 数值属性, 200
- 似然, 148
- 似然率, 352
- 松弛变量, 130
- 随机森林, 179
- 随机子空间, 189

- 探索-利用窘境, 374
- 特化, 350
- 特征, 2, 247
- 特征工程, 114
- 特征向量, 2
- 特征选择, 247
- 特征学习, 114
- 梯度下降, 102, 254, 389, 407
- 替代函数, 58
- 替代损失, 130
- 条件独立性假设, 150, 305
- 条件风险, 147
- 条件随机场, 325
- 同父, 158
- 统计关系学习, 364
- 统计学习, 12, 139
- 投票法, 172, 225
- 图半监督学习, 300
- 推断, 319

- 微 $F1$, 32
- 微查全率, 32
- 微查准率, 32
- 维数约简, 227
- 维数灾难, 227, 247
- 未标记样本, 293
- 稳定基学习器, 189
- 稳定性, 284
- 无导师学习, 3
- 无关特征, 247
- 无监督学习, 3, 197
- 无监督逐层训练, 113
- 无序属性, 200
- 勿连约束, 239, 307
- 误差, 23
- 误差-分歧分解, 185

- 误差逆传播, 101
- 稀疏编码, 255
- 稀疏表示, 67, 255
- 稀疏性, 67
- 下采样, 67
- 先验, 148
- 线性超平面, 99
- 线性核, 128
- 线性回归, 53, 252
- 线性降维, 229
- 线性可分, 99, 126
- 线性模型, 53
- 线性判别分析, 60, 139
- 相对多数投票, 183
- 相对熵, 414
- 相关特征, 247
- 相似度度量, 201
- 协同过滤, 259
- 协同训练, 304
- 斜决策树, 90
- 信念传播, 330, 340
- 信念网, 156
- 信息散度, 414
- 信息增益, 75, 248
- 信息熵, 75
- 序贯覆盖, 349
- 选择性集成, 191
- 学习, 2
- 学习率, 99
- 学习器, 2
- 学习向量量化, 204
- 训练, 2
- 训练集, 2
- 训练误差, 23
- 训练样本, 2
- 压缩感知, 257
- 哑结点, 99
- 演绎, 359
- 验证集, 28, 105
- 样本, 2
- 样本复杂度, 270
- 样本空间, 2
- 样例, 2
- 一阶规则, 348
- 一致性, 140
- 遗传算法, 107
- 异常检测, 219
- 因子, 322
- 隐变量, 162, 319
- 隐狄利克雷分配模型, 337
- 隐马尔可夫模型, 319
- 硬间隔, 129
- 优先级规则, 348
- 有标记样本, 293
- 有导师学习, 3
- 有模型学习, 377
- 有限假设空间, 270
- 有向分离, 158
- 有效性指标, 197
- 有序属性, 200
- 预剪枝, 79, 352
- 阈值, 97, 104
- 阈值逻辑单元, 98
- 阈值移动, 67
- 元规则, 348
- 原型聚类, 202
- 原子命题, 348
- 再励学习, 371
- 再平衡, 67
- 再生核希尔伯特空间, 128
- 再缩放, 67
- 在线学习, 109, 241, 393
- 早停, 105
- 增长函数, 273
- 增量学习, 92, 109
- 增益率, 77
- 召回率, 30
- 真相, 2
- 振荡, 99
- 正态分布, 412
- 正则化, 56, 105, 133
- 证据, 148
- 支持向量, 122
- 支持向量回归, 133
- 支持向量机, 123
- 支持向量展式, 127
- 直推学习, 295
- 值迭代, 382
- 值函数近似, 388
- 指数损失, 130, 173
- 置换, 361
- 置信度, 38
- 主成分分析, 229
- 主动学习, 293
- 状态-动作值函数, 377
- 状态值函数, 377
- 准确率, 30

子集评价, 248
子集搜索, 248
子空间, 189, 227
自适应谐振理论, 108
自助采样法, 178
自助法, 27
自组织映射, 109
字典学习, 255

总体代价, 36
最近邻分类器, 225
最小二乘法, 54, 72
最小描述长度, 159
最小一般泛化, 358
最一般合一置换, 361
坐标下降, 163, 408

MACHINE LEARNING

机器学习

周志华 著

清华大学出版社
北 京

内 容 简 介

机器学习是计算机科学的重要分支领域。本书作为该领域的入门教材,在内容上尽可能涵盖机器学习基础知识的各方面。全书共16章,大致分为3个部分:第1部分(第1~3章)介绍机器学习的基础知识;第2部分(第4~10章)讨论一些经典而常用的机器学习方法(决策树、神经网络、支持向量机、贝叶斯分类器、集成学习、聚类、降维与度量学习);第3部分(第11~16章)为进阶知识,内容涉及特征选择与稀疏学习、计算学习理论、半监督学习、概率图模型、规则学习以及强化学习等。每章都附有习题并介绍了相关阅读材料,以便有兴趣的读者进一步钻研探索。

本书可作为高等院校计算机、自动化及相关专业的本科生或研究生教材,也可供对机器学习感兴趣的研究人员和工程技术人员阅读参考。

本书封面贴有清华大学出版社防伪标签,无标签者不得销售。

版权所有,侵权必究。侵权举报电话:010-62782989 13701121933

图书在版编目(CIP)数据

机器学习/周志华著.--北京:清华大学出版社,2016

ISBN 978-7-302-42328-7

I. ①机… II. ①周… III. ①机器学习 IV. ①TP181

中国版本图书馆 CIP 数据核字(2015)第 287090 号

责任编辑:薛 慧

封面设计:何凤霞

责任校对:刘玉霞

责任印制:宋 林

出版发行:清华大学出版社

网 址: <http://www.tup.com.cn>, <http://www.wqbook.com>

地 址:北京清华大学学研大厦 A 座

邮 编:100084

社 总 机:010-62770175

邮 购:010-62786544

投稿与读者服务:010-62776969, c-service@tup.tsinghua.edu.cn

质量反馈:010-62772015, zhiliang@tup.tsinghua.edu.cn

印 装 者:北京亿浓世纪彩色印刷有限公司

经 销:全国新华书店

开 本:210mm×235mm 印 张:27.75 字 数:626 千字

版 次:2016 年 1 月第 1 版 印 次:2016 年 1 月第 1 次印刷

印 数:1~5000

定 价:88.00 元

产品编号:064027-01

前言

这是一本面向中文读者的机器学习教科书,为了使尽可能多的读者通过本书对机器学习有所了解,作者试图尽可能少地使用数学知识.然而,少量的概率、统计、代数、优化、逻辑知识似乎不可避免.因此,本书更适合大学三年级以上的理工科本科生和研究生,以及具有类似背景的对机器学习感兴趣的人士.为方便读者,本书附录给出了一些相关数学基础知识简介.

全书共 16 章,大体上可分为 3 个部分:第 1 部分包括第 1~3 章,介绍机器学习基础知识;第 2 部分包括第 4~10 章,介绍一些经典而常用的机器学习方法;第 3 部分包括第 11~16 章,介绍一些进阶知识.前 3 章之外的后续各章均相对独立,读者可根据自己的兴趣和时间情况选择使用.根据课时情况,一个学期的本科生课程可考虑讲授前 9 章或前 10 章;研究生课程则不妨使用全书.

书中除第 1 章外,每章都给出了十道习题.有的习题是帮助读者巩固本章学习,有的是为了引导读者扩展相关知识.一学期的一般课程可使用这些习题,再辅以两到三个针对具体数据集的大作业.带星号的习题则有相当难度,有些并无现成答案,谨供富有进取心的读者启发思考.

本书在内容上尽可能涵盖机器学习基础知识的各方面,但作为机器学习入门读物且因授课时间的考虑,很多重要、前沿的材料未能覆盖,即便覆盖到的部分也仅是管中窥豹,更多的内容留待读者在进阶课程中学习.为便于有兴趣的读者进一步钻研探索,本书每章均介绍了一些阅读材料,谨供读者参考.

笔者以为,对学科相关的重要人物和事件有一定了解,将会增进读者对该学科的认识.本书在每章最后都写了一个与该章内容相关的小故事,希望有助于读者增广见闻,并且在紧张的学习过程中稍微放松调剂一下.

书中不可避免地涉及大量外国人名,若全部译为中文,则读者在日后进一步阅读文献时或许会对不少人名产生陌生感,不利于进一步学习.因此,本书仅对一般读者耳熟能详的名字如“图灵”等加以直接使用,对故事中的一些主要人物给出了译名,其他则保持外文名.

机器学习发展极迅速,目前已成为一个广袤的学科,罕有人士能对其众多分支领域均有精深理解.笔者自认才疏学浅,仅略知皮毛,更兼时间和精力所限,书中错谬之处在所难免,若蒙读者诸君不吝告知,将不胜感激.

周志华

2015 年 6 月

序 言

在人工智能界有一种说法,认为机器学习是人工智能领域中最能够体现智能的一个分支.从历史来看,机器学习似乎也是人工智能中发展最快的分支之一.在二十世纪八十年代的时候,符号学习可能还是机器学习的主流,而自从二十世纪九十年代以来,就一直是统计机器学习的天下了.不知道是否可以这样认为:从主流为符号机器学习发展到主流为统计机器学习,反映了机器学习从纯粹的理论研究和模型研究发展到以解决现实生活中实际问题为目的的应用研究,这是科学研究的一种进步.有关机器学习的专著国内出版的不是很多.前两年有李航教授的《统计学习方法》出版,以简要的方式介绍了一批重要和常用的机器学习方法.此次周志华教授的鸿篇巨著《机器学习》则全面而详细地介绍了机器学习的各个分支,既可作为教材,又可作为自学用书和科研参考书.

翻阅书稿的过程引起了一些自己的思考,平时由于和机器学习界的朋友接触多了,经常获得一些道听途说的信息以及专家们对机器学习现状及其发展前途的评论.在此过程中,难免会产生一些自己的疑问.我借此机会把它写下来放在这里,算是一种“外行求教机器学习”.

问题一:在人工智能发展早期,机器学习的技术内涵几乎全部是符号学习.可是从二十世纪九十年代开始,统计机器学习犹如一匹黑马横空出世,迅速压倒并取代了符号学习的地位.人们可能会问:在满目的统计学习期刊和会议文章面前,符号学习是否被彻底忽略了?它还能成为机器学习的研究对象吗?它是否将继续在统计学习的阴影里生活并苟延残喘?对这个问题有三种可能的答案:一是告诉符号学习:“你就是该退出历史舞台,认命吧!”二是告诉统计学习:“你的一言堂应该关门了!”单纯的统计学习已经走到了尽头,再想往前走就要把统计学习和符号学习结合起来.三是事物发展总会有“三十年河东,三十年河西”的现象,符号学习还有“翻身”的日子.第一种观点我没有听人明说过,但是我想恐怕有可能已经被许多人默认了.第二种观点我曾听王珏教授多次说过.他并不认为统计学习会衰退,而只是认为机器学习已经到了一个转折点,从今往后,统计学习应该和知识的利用相结合,这是一种“螺旋式上升,进入更高级的形式”,否则,统计学习可能会停留于现状而止步不前.王珏教授还认为:进入转折点的标志就是 Koller 等的《概率图模型》一书的出版.至于第三种观点,恰好我收到老朋友,美国人工智能资深学者、俄亥俄大学 Chandrasekaran 教授的来信,他正好谈起符号智能被统计智能“打压”的现象,并且正好表达了河东河西的观点.我请求他允许我把这段话引进正在撰写的序言中,他爽快地同意了,仅仅修改了几处私人通信的口吻.全文如下:“最近几年,人工智能在很大程度上集中于统计学和大数据.我同意由于计算能力的大幅提高,这些技术曾经取得过某些令人印象深刻的成果.但是我们完全有理由相信,虽然这些技术还会继续改进、提高,总有一天这个领域(指 AI)会对它们说再见,并转向更加基本的认知科学研究.尽管钟摆的摆回去还需要一段时间,我

相信定有必要把统计技术和对认知结构的深刻理解结合起来。”看来, Chandrasekaran 教授也并不认为若干年以后 AI 真会回到河西, 他的意见和王珏教授的意见基本一致, 但不仅限于机器学习, 而是涉及整个人工智能领域. 只是王珏教授强调知识, 而 Chandrasekaran 教授强调更加基本的“认知”.

问题二: 王珏教授认为统计机器学习不会“一路顺风”的判据是: 统计机器学习算法都是基于样本数据独立同分布的假设. 但是自然界现象千变万化, 王珏教授认为“哪有那么多独立同分布?”这就引来了下一个问题: “独立同分布”条件对于机器学习来讲真是必需的吗? 独立同分布的不存在一定是一个不可逾越的障碍吗? 无独立同分布条件下的机器学习也许只是一个难题, 而不是不可解问题. 我有一个“胡思乱想”, 认为前些时候出现的“迁移学习”也许会对这个问题的解决带来一线曙光. 尽管现在的迁移学习还要求迁移双方具备“独立同分布”条件, 但是不同分布之间的迁移学习, 同分布和异分布之间的迁移学习也许迟早会出现?

问题三: 近年来出现了一些新的动向, 例如“深度学习”、“无终止学习”等等, 社会上给予了特别关注, 尤其是深度学习. 但它们真的代表了机器学习的新的方向吗? 包括本书作者周志华教授在内的一些学者认为: 深度学习掀起的热潮也许大过它本身真正的贡献, 在理论和技术上并没有太多的创新, 只不过是硬件技术的革命, 计算机的速度大大提高了, 使得人们有可能采用原来复杂度很高的算法, 从而得到比过去更精细的结果. 当然这对于推动机器学习应用于实践有很大意义. 但我们不禁要斗胆问一句: 深度学习是否又要取代统计学习了? 事实上, 确有专家已经感受到来自深度学习的压力, 指出统计学习正在被深度学习所打压, 正如我们早就看到的符号学习被统计学习所打压. 不过我觉得这种打压还远没有强大到像统计学习打压符号学习的程度. 这—是因为深度学习的“理论创新”还不明显; 二是因为目前的深度学习主要适合于神经网络, 在各种机器学习方法百花盛开的今天, 它的应用范围还有限, 还不能直接说是连接主义方法的回归; 三是因为统计学习仍然在机器学习中被有效地普遍采用, “得道多助”, 想抛弃它不容易.

问题四: 机器学习研究出现以来, 我们看到的主要是从符号方法到统计方法的演变, 用到的数学主要是概率统计. 但是, 数学之大, 就像大海. 难道只有统计方法适合于在机器学习方面应用吗? 当然, 我们也看到了一些其他数学分支在机器学习上的应用的好例子, 例如微分几何在流形学习上的应用, 微分方程在归纳学习上的应用. 但如果和统计方法相比, 它们都只能算是配角. 还有的数学分支如代数可能应用得更广, 但在机器学习中代数一般是作为基础工具来使用, 例如矩阵理论和特征值理论. 又如微分方程求解最终往往归结为代数问题求解. 它们可以算是幕后英雄: “出头露面的是概率和统计, 埋头苦干的是代数和逻辑”. 是否可以想象以数学方法为主角, 以统计方法为配角的机器学习理论呢? 在这方面, 流形学习已经“有点意思”了, 而彭实戈院士的倒排随机微分方程理论之预测金融走势, 也许是用高深数学推动新的机器学习模式的更好例子. 但是从宏观的角度看, 数学理论的介入程度还远远不够. 这里指的主要是深刻的、现代的数学理论, 我们期待着有更多数学家的参与, 开辟机器学习的新模式、新理论、新方向.

问题五: 上一个问题的延续: 符号机器学习时代主要以离散方法处理问题, 统计机器学习时代主要以连续方法处理问题. 这两种方法之间应该没有一条鸿沟. 流形学习中李群、李代数方法的引入给我们以很好的启示. 从微分流形到李群, 再从李群到李代数, 就是一个沟通连续和离散的过程. 然而, 现有的方法在数学上并不完美. 浏览流形学习的文献可知, 许多论文直接把任意数据集看成微分流形, 从而就认定测地线的存在并讨论起降维来了. 这样的例子也许不是个别的, 足可说明数学家介入机器学习研究之必要.

问题六: 大数据时代的出现, 有没有给机器学习带来本质性的影响? 理论上讲, 似乎“大数据”给统计机器学习提供了更多的机遇, 因为海量的数据更加需要统计、抽样的方法. 业界人士估计, 大数据的出现将使人工智能的作用更加突出. 有人把大数据处理分成三个阶段: 收集、分析和预测. 收集和分析的工作相对来说已经做得相当好了, 现在关注的焦点是要有科学的预测, 机器学习技术在这里不可或缺. 这一点大概毋庸置疑. 然而, 同样是使用统计、抽样方法, 同样是收集、分析和预测, 大数据时代使用这类方法和以前使用这类方法有什么本质的不同吗? 量变到质变是辩证法的一个普遍规律. 那么, 从前大数据时代到大数据时代, 数理统计方法有没有发生本质的变化? 反映到它们在机器学习上的应用有无本质变化? 大数据时代正在呼唤什么样的机器学习方法的产生? 哪些机器学习方法又是由于大数据研究的驱动而产生的呢?

以上这些话也许说得远了, 我们还是回到本书上来. 本书的作者周志华教授在机器学习的许多领域都有出色的贡献, 是中国机器学习研究的领军人物之一, 在国际学术界有着很高的声誉. 他在机器学习的一些重要领域, 例如集成学习、半监督学习、多示例和多标记学习等方面都做出了在国际上有重要影响的工作, 其中一些可以认为是中国学者在国际上的代表性贡献. 除了自身的学术研究以外, 他在推动中国的机器学习发展方面也做了许多工作. 例如他和不久前刚过世的王珏教授从 2002 年开始, 组织了系列化的“机器学习及其应用”研讨会. 初在复旦, 后移至南大举行, 越办越兴旺, 从单一的专家报告发展到专家报告、学生论坛和张贴论文三种方式同时举行, 参会者从数十人发展到数百人, 活动搞得有声有色, 如火如荼. 最近更是把研讨会推向全国高校轮流举行. 他和王珏教授紧密合作, 南北呼应, 人称“南周北王”. 王珏教授的离去使我们深感悲伤. 令我们欣慰的是国内不但有周志华教授这样的机器学习领军人物, 而且比周教授更年轻的许多机器学习青年才俊也成长起来了. 中国的机器学习大有希望.

陆汝钊

中国科学院数学与系统科学研究院

2015 年 8 月于北京

主要符号表

x	标量
\boldsymbol{x}	向量
\mathbf{x}	变量集
\mathbf{A}	矩阵
\mathbf{I}	单位阵
\mathcal{X}	样本空间或状态空间
\mathcal{D}	概率分布
D	数据样本 (数据集)
\mathcal{H}	假设空间
H	假设集
\mathcal{L}	学习算法
(\cdot, \cdot, \cdot)	行向量
$(; ; \cdot)$	列向量
$(\cdot)^T$	向量或矩阵转置
$\{\cdots\}$	集合
$ \{\cdots\} $	集合 $\{\cdots\}$ 中元素个数
$\ \cdot\ _p$	L_p 范数, p 缺省时为 L_2 范数
$P(\cdot), P(\cdot \cdot)$	概率质量函数, 条件概率质量函数
$p(\cdot), p(\cdot \cdot)$	概率密度函数, 条件概率密度函数
$\mathbb{E}_{\cdot \sim \mathcal{D}}[f(\cdot)]$	函数 $f(\cdot)$ 对 \cdot 在分布 \mathcal{D} 下的数学期望; 意义明确时将省略 \mathcal{D} 和(或) \cdot
$\sup(\cdot)$	上确界
$\mathbb{I}(\cdot)$	指示函数, 在 \cdot 为真和假时分别取值为 1, 0
$\text{sign}(\cdot)$	符号函数, 在 $\cdot < 0, = 0, > 0$ 时分别取值为 -1, 0, 1

目 录

第 1 章 绪论	1
1.1 引言	1
1.2 基本术语	2
1.3 假设空间	4
1.4 归纳偏好	6
1.5 发展历程	10
1.6 应用现状	13
1.7 阅读材料	16
习题	19
参考文献	20
休息一会儿	22
第 2 章 模型评估与选择	23
2.1 经验误差与过拟合	23
2.2 评估方法	24
2.3 性能度量	28
2.4 比较检验	37
2.5 偏差与方差	44
2.6 阅读材料	46
习题	48
参考文献	49
休息一会儿	51
第 3 章 线性模型	53
3.1 基本形式	53
3.2 线性回归	53
3.3 对数几率回归	57
3.4 线性判别分析	60
3.5 多分类学习	63

3.6 类别不平衡问题	66
3.7 阅读材料	67
习题	69
参考文献	70
休息一会儿	72
第4章 决策树	73
4.1 基本流程	73
4.2 划分选择	75
4.3 剪枝处理	79
4.4 连续与缺失值	83
4.5 多变量决策树	88
4.6 阅读材料	92
习题	93
参考文献	94
休息一会儿	95
第5章 神经网络	97
5.1 神经元模型	97
5.2 感知机与多层网络	98
5.3 误差逆传播算法	101
5.4 全局最小与局部极小	106
5.5 其他常见神经网络	108
5.6 深度学习	113
5.7 阅读材料	115
习题	116
参考文献	117
休息一会儿	120
第6章 支持向量机	121
6.1 间隔与支持向量	121
6.2 对偶问题	123
6.3 核函数	126
6.4 软间隔与正则化	129
6.5 支持向量回归	133

6.6 核方法	137
6.7 阅读材料	139
习题	141
参考文献	142
休息一会儿	145
第7章 贝叶斯分类器	147
7.1 贝叶斯决策论	147
7.2 极大似然估计	149
7.3 朴素贝叶斯分类器	150
7.4 半朴素贝叶斯分类器	154
7.5 贝叶斯网	156
7.6 EM算法	162
7.7 阅读材料	164
习题	166
参考文献	167
休息一会儿	169
第8章 集成学习	171
8.1 个体与集成	171
8.2 Boosting	173
8.3 Bagging与随机森林	178
8.4 结合策略	181
8.5 多样性	185
8.6 阅读材料	190
习题	192
参考文献	193
休息一会儿	196
第9章 聚类	197
9.1 聚类任务	197
9.2 性能度量	197
9.3 距离计算	199
9.4 原型聚类	202
9.5 密度聚类	211

9.6 层次聚类	214
9.7 阅读材料	217
习题	220
参考文献	221
休息一会儿	224
第 10 章 降维与度量学习	225
10.1 k 近邻学习	225
10.2 低维嵌入	226
10.3 主成分分析	229
10.4 核化线性降维	232
10.5 流形学习	234
10.6 度量学习	237
10.7 阅读材料	240
习题	242
参考文献	243
休息一会儿	246
第 11 章 特征选择与稀疏学习	247
11.1 子集搜索与评价	247
11.2 过滤式选择	249
11.3 包裹式选择	250
11.4 嵌入式选择与 L_1 正则化	252
11.5 稀疏表示与字典学习	254
11.6 压缩感知	257
11.7 阅读材料	260
习题	262
参考文献	263
休息一会儿	266
第 12 章 计算学习理论	267
12.1 基础知识	267
12.2 PAC学习	268
12.3 有限假设空间	270
12.4 VC维	273

12.5 Rademacher复杂度	279
12.6 稳定性	284
12.7 阅读材料	287
习题	289
参考文献	290
休息一会儿	292
第 13 章 半监督学习	293
13.1 未标记样本	293
13.2 生成式方法	295
13.3 半监督SVM	298
13.4 图半监督学习	300
13.5 基于分歧的方法	304
13.6 半监督聚类	307
13.7 阅读材料	311
习题	313
参考文献	314
休息一会儿	317
第 14 章 概率图模型	319
14.1 隐马尔可夫模型	319
14.2 马尔可夫随机场	322
14.3 条件随机场	325
14.4 学习与推断	328
14.5 近似推断	331
14.6 话题模型	337
14.7 阅读材料	339
习题	341
参考文献	342
休息一会儿	345
第 15 章 规则学习	347
15.1 基本概念	347
15.2 序贯覆盖	349
15.3 剪枝优化	352

15.4 一阶规则学习	354
15.5 归纳逻辑程序设计	357
15.6 阅读材料	363
习题	365
参考文献	366
休息一会儿	369
第 16 章 强化学习	371
16.1 任务与奖赏	371
16.2 K -摇臂赌博机	373
16.3 有模型学习	377
16.4 免模型学习	382
16.5 值函数近似	388
16.6 模仿学习	390
16.7 阅读材料	393
习题	394
参考文献	395
休息一会儿	397
附录	399
A 矩阵	399
B 优化	403
C 概率分布	409
后记	417
索引	419