

第五章 矩阵计算问题

在众多科学与工程学科，如物理、化学工程、统计学、经济学、生物学、信号处理、自动控制、系统理论、医学和军事工程等中，许多问题都可用数学建模成矩阵方程 $A\mathbf{x} = \mathbf{b}$ 。根据数据向量 $\mathbf{b} \in \mathbb{R}^{m \times 1}$ 和数据矩阵 $A \in \mathbb{R}^{m \times n}$ 的不同, 矩阵方程主要有以下三种类型:

(1) **适定方程组**: 方程的个数与未知量的个数相等即 $m = n$, 并且 A 满秩可逆, 此时 \mathbf{x} 有唯一的解。

(2) **超定方程组**: 当上述 $m > n$ 时, 并且数据矩阵 A 和数据向量 \mathbf{b} 均已知, 其中之一或者二者可能存在误差或者干扰。

(3) **欠定方程组**: 当上述 $m < n$ 时, 数据矩阵 A 和数据向量 \mathbf{b} 均已知, 但未知向量 \mathbf{x} 可能要求为稀疏向量。

我们这里引进线性方程并给出它的标准形式 $A\mathbf{x} = \mathbf{y}$, 其中 $\mathbf{x} \in \mathbb{R}^n$ 是未知变量, $A \in \mathbb{R}^{m \times n}$ 是参数矩阵, $\mathbf{y} \in \mathbb{R}^m$ 是已知向量。线性方程构成了数值线性代数的基础, 它们的解法在许多优化方法的关键。事实上, 解线性方程组问题 $A\mathbf{x} = \mathbf{y}$ 可以被看成**优化问题**, 即关于 \mathbf{x} , 最小化 $\|A\mathbf{x} - \mathbf{y}\|^2$ 。我们描述线性方程组解得集合并且当线性方程组正确解不存在的情况下, 讨论求解线性方程组近似解的方法。随后引出最小二乘问题以及它的变体、解的数值敏感性及其解决方法, 它们与矩阵分解的关系 (例如 QR 分解和 SVD) 也将被介绍。

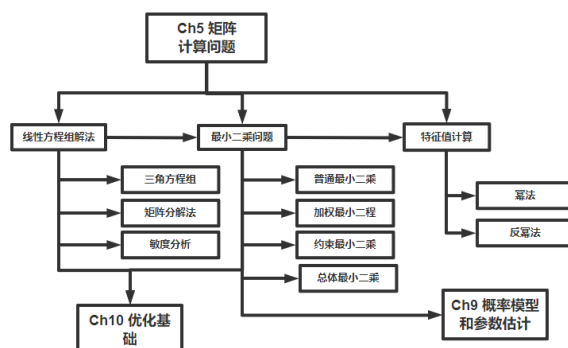


图 5.1: 本章导图

5.1 线性方程组的直接解法

5.1.1 线性方程组问题

在工程问题中，线性方程组描述了变量之间最基本的关系。线性方程在各个科学分支中无处不在，例如弹性力学、电阻网络、曲线拟合等。线性方程构成了线性代数的核心并时常作为优化问题的约束条件。由于许多优化算法的迭代过程非常依赖线性方程组的解，所以它也是许多优化算法的基础。下面我们以一个例子来说明，线性方程组如何解决上面的问题。

例 5.1.1. (三点测距问题) 三角测量是一种确定点位置的方法，给定距离到已知控制点 (锚点)，三边测量可以应用于许多不同的领域，如地理测绘、地震学、导航 (例如 GPS 系统) 等。在图 5.2 中，三个测距点 $a_1, a_2, a_3 \in \mathbb{R}^2$ 的坐标是已知的，并且从点 $\mathbf{x} = (x_1, x_2)^T$ 到测距点的距离为 d_1, d_2, d_3 ， \mathbf{x} 的未知坐标与距离测量有关，可以由下面非线性方程组描述

$$\|\mathbf{x} - \mathbf{a}_1\|_2^2 = d_1^2, \quad \|\mathbf{x} - \mathbf{a}_2\|_2^2 = d_2^2, \quad \|\mathbf{x} - \mathbf{a}_3\|_2^2 = d_3^2 \quad (5.1)$$

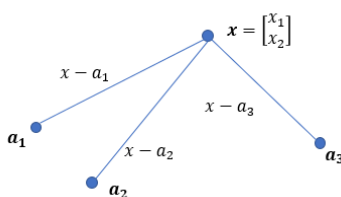


图 5.2: 三点测量位置图例。点 \mathbf{x} 处，我们测量距三个测距点 $\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3$ 的距离，以便确定 \mathbf{x} 的坐标。

通过第一个方程减去另外两个方程，我们获得了两个 \mathbf{x} 的线性方程组。

$$\begin{aligned} 2(\mathbf{a}_2 - \mathbf{a}_1)^T \mathbf{x} &= d_1^2 - d_2^2 + \|\mathbf{a}_2\|_2^2 - \|\mathbf{a}_1\|_2^2 \\ 2(\mathbf{a}_3 - \mathbf{a}_1)^T \mathbf{x} &= d_1^2 - d_3^2 + \|\mathbf{a}_3\|_2^2 - \|\mathbf{a}_1\|_2^2 \end{aligned}$$

也就是说，原始非线性方程组(5.1)的每个解也可以看作线性方程组的解。使用方程组标准形式 $A\mathbf{x} = \mathbf{y}$ (标准形式的定义在下一小节给出) 可以描述为：

$$A = \begin{bmatrix} 2(\mathbf{a}_2 - \mathbf{a}_1)^T \\ 2(\mathbf{a}_3 - \mathbf{a}_1)^T \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} d_1^2 - d_2^2 + \|\mathbf{a}_2\|_2^2 - \|\mathbf{a}_1\|_2^2 \\ d_1^2 - d_3^2 + \|\mathbf{a}_3\|_2^2 - \|\mathbf{a}_1\|_2^2 \end{bmatrix} \quad (5.2)$$

上述问题的解，将在后面详细讨论。

我们先回顾线性方程组中的一般概念。

[illegible]
$$\mathbf{A} = (a_{ij})^{m \times n} \square \mathbf{x} = (x_1, x_2, \dots, x_n)^\top \square \mathbf{b} = (b_1, b_2, \dots, b_m)^\top, \quad (5.4)$$
$$Ax = b. \quad (5.5)$$
$$Ax = 0. \quad (5.6)$$

定义 5.1.2. 给定方程组

$$\bar{A}x = \bar{b}. \quad (5.7)$$

定理 5.1.1. 对方程组(5.5)的系数矩阵 A 及右端作相同的行初等变换, 所得到的新方程组与原方程组同解。

定理 5.1.2. 设齐次线性方程组(5.6)的系数矩阵 A 的秩为 r , 此时

- (1) 方程组(5.6)有非零解的必要充分条件是 $r < n$ 。
- (2) 若 $r < n$, 则方程组(5.6)一定有基础解系。基础解系不是唯一的, 但任两个基础解系必等价, 且每一个基础解系所含解向量的个数都等于 $n - r$ 。
- (3) 若 $r < n$, 设 $\eta_1, \eta_2, \dots, \eta_{n-r}$ 是方程组(5.6)的一个基础解系, 则它的一般解为

$$\eta = \lambda_1 \eta_1 + \lambda_2 \eta_2 + \cdots + \lambda_{n-r} \eta_{n-r}, \quad (5.8)$$

其中 $\lambda_i (i = 1, 2, \dots, n-r)$ 是数域 \mathbb{K} 中的任意常数.

定理 5.1.3. 方程组(5.5)有解的必要充分条件是: $\text{rank}(\mathbf{A}) = \text{rank}(\tilde{\mathbf{A}})$. 矩阵 $\tilde{\mathbf{A}}$ 为它的增广矩阵.

定理 5.1.4. 设 $\text{rank}(\mathbf{A}) = \text{rank}(\tilde{\mathbf{A}}) = r$, γ_0 是非齐次方程组(5.5)的一个解向量 (常称为 特解), $\eta_1, \eta_2, \dots, \eta_{n-r}$ 是其导出组(5.6)的一个基础解系, 则方程组(5.5)的解向量均可表为:

$$\gamma = \gamma_0 + \eta = \gamma_0 + \lambda_1 \eta_1 + \lambda_2 \eta_2 + \dots + \lambda_{n-r} \eta_{n-r},$$

其中 $\lambda_i (i = 1, 2, \dots, n-r)$ 是数域 \mathbb{K} 中的任意常数 (这种形式的解向量常称为一般解).

对增广矩阵 $\tilde{\mathbf{A}}$ 进行初等行变换将其化为阶梯型矩阵, 写出相应的阶梯型方程组.

(1) 若 $r = n$, 则阶梯型方程组形如:

$$\begin{cases} c_{11}x_1 + c_{12}x_2 + \dots + c_{1n}x_n = d_1, \\ c_{22}x_2 + \dots + c_{2n}x_n = d_2, \\ \dots\dots\dots \\ c_{nn}x_n = d_n. \end{cases} \quad (5.9)$$

其中 $c_{ii} \neq 0 (i = 1, 2, \dots, n)$. 依次由第 n 个, 第 $n-1$ 个, \dots , 第一个方程组可解得 x_n, x_{n-1}, \dots, x_1 , 由此即得方程组(5.3)的唯一解 x_1, x_2, \dots, x_n .

(2) 若 $r < n$, 则阶梯型方程组可表为:

$$\begin{cases} c_{11}x_1 + c_{12}x_2 + \dots + c_{1r}x_r = d_1 - c_{1,r+1}x_{r+1} - \dots - c_{1n}x_n, \\ c_{22}x_2 + \dots + c_{2r}x_r = d_2 - c_{2,r+1}x_{r+1} - \dots - c_{2n}x_n, \\ \dots\dots\dots \\ c_{rr}x_r = d_r - c_{r,r+1}x_{r+1} - \dots - c_{rn}x_n. \end{cases} \quad (5.10)$$

其中 $c_{ii} \neq 0 (i = 1, 2, \dots, r)$. 此时方程组(5.3)有无穷多组解. 若令 $x_{r+1} = x_{r+2} = \dots = x_n = 0$, 则可由方程组(5.10)求得一个特解 $\gamma_0 = (\delta_1, \delta_2, \dots, \delta_r, 0, 0, \dots, 0)$, 再由其导出组的一个基础解系 $\eta_1, \eta_2, \dots, \eta_{n-r}$, 可得方程组(5.3)的一般解.

线性方程组的解集被定义为:

$$S \doteq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{Ax} = \mathbf{y}\} \quad (5.11)$$

用 $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n \in \mathbb{R}^n$ 表示矩阵 \mathbf{A} 的列, 即 $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n]$. \mathbf{Ax} 仅仅表示矩阵 \mathbf{A} 的列与向量 \mathbf{x} 中各个元素的加权和:

$$\mathbf{Ax} = x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \dots + x_n\mathbf{a}_n \quad (5.12)$$

通过定义, 我们能够看出, 无论 \mathbf{x} 的值是什么, \mathbf{Ax} 生成了由矩阵 \mathbf{A} 的列张成的子空间. 向量 $\mathbf{Ax} \in \text{Range}(\mathbf{A})$. 若 $\mathbf{y} \notin \text{Range}(\mathbf{A})$, 则线性方程组没有解. 因此解集 S 为空. 等价地, 线性方程组有解当且仅当 $\mathbf{y} \in \text{Range}(\mathbf{A})$.

回顾定义, 设 $\mathbf{A} \in \mathbb{R}^{m \times n}$. \mathbf{A} 的值域定义为

$$\text{Range}(\mathbf{A}) = \{\mathbf{y} \in \mathbb{R}^m : \mathbf{y} = \mathbf{Ax}, \mathbf{x} \in \mathbb{R}^n\} \quad (5.13)$$

易证 $\text{Range}(\mathbf{A}) = \text{span}\{\mathbf{a}_1, \dots, \mathbf{a}_n\}$, 其中 \mathbf{a}_i 为 \mathbf{A} 的列向量. \mathbf{A} 的零空间定义为:

$$\text{Null}(\mathbf{A}) = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{0}\} \quad (5.14)$$

它的维数记为 $\text{nullity}(\mathbf{A})$

一个子空间 $\mathbb{S} \subset \mathbb{R}^n$ 的正交补定义为:

$$\mathbb{S}^\perp = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{y}^\text{T}\mathbf{x} = 0, \forall \mathbf{x} \in \mathbb{S}\} \quad (5.15)$$

从矩阵值域的角度, 给出定理5.1.3的证明:

定理 5.1.5. 方程组的解(5.5)存在的充分必要条件是 $\text{rank}(\mathbf{A}) = \text{rank}([\mathbf{A}, \mathbf{b}])$.

证明. 必要性: 设存在 \mathbf{x} 使 $\mathbf{A}\mathbf{x} = \mathbf{b}$, 则 \mathbf{b} 是 \mathbf{A} 的列向量的线性组合, 即 $\mathbf{b} \in \mathcal{R}(\mathbf{A})$. 这说明 $\mathcal{R}([\mathbf{A}, \mathbf{b}]) = \mathcal{R}(\mathbf{A})$, 所以有 $\text{rank}(\mathbf{A}) = \text{rank}([\mathbf{A}, \mathbf{b}])$.

充分性: 若 $\text{rank}(\mathbf{A}) = \text{rank}([\mathbf{A}, \mathbf{b}])$ 成立, 则 $\mathbf{b} \in \mathcal{R}(\mathbf{A})$, 即 \mathbf{b} 可表示为 $\mathbf{b} = \sum_{i=1}^n x_i \mathbf{a}_i$, 这里 $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]$, 所以令 $\mathbf{x} = (x_1, \dots, x_n)^\text{T}$, 即有 $\mathbf{A}\mathbf{x} = \mathbf{b}$. \square

定理 5.1.6. 假定方程组(5.5)的解存在, 并且假定 \mathbf{x} 是其任一给定的解, 则(5.5)全部解的集合是

$$\mathbf{x} + \text{Null}(\mathbf{A}) \quad (5.16)$$

证明. 如果 \mathbf{y} 满足(5.5), 则 $\mathbf{A}(\mathbf{y} - \mathbf{x}) = \mathbf{0}$, 即 $(\mathbf{y} - \mathbf{x}) \in \text{Null}(\mathbf{A})$, 于是有 $\mathbf{y} = \mathbf{x} + (\mathbf{y} - \mathbf{x}) \in \mathbf{x} + \text{Null}(\mathbf{A})$. 反之, 如果 $\mathbf{y} \in \mathbf{x} + \text{Null}(\mathbf{A})$, 则存在 $\mathbf{z} \in \text{Null}(\mathbf{A})$, 使 $\mathbf{y} = \mathbf{x} + \mathbf{z}$, 从而有 $\mathbf{A}\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{A}\mathbf{z} = \mathbf{A}\mathbf{x} = \mathbf{b}$. \square

定理5.1.6告诉我们, 只要知道了方程组(5.5)的一个解, 便可以用它及 $\text{Null}(\mathbf{A})$ 中向量的和得到(5.5)的全部解. 由此可知, 方程组(5.5)的解唯一, 只有当 $\text{Null}(\mathbf{A})$ 中仅有零向量才行.

推论 5.1.1. 方程组(5.5)的解唯一的充分必要条件是 $\text{nullity}(\mathbf{A}) = 0$.

对于不同类型的方程组, 解的数量情况有各自的特点。

超定方程组

线性方程组 $\mathbf{A}\mathbf{x} = \mathbf{y}$ 中线性方程的个数大于未知变量的个数时, 我们说 $\mathbf{A}\mathbf{x} = \mathbf{y}$ 是超定的. 或者说矩阵 $\mathbf{A}^{m \times n}$ 的行数大于列数: $m > n$. 假设 \mathbf{A} 是一个列满秩矩阵, 也就是说 $\text{rank}(\mathbf{A}) = n$, 则我们可以得出 $\text{Null}(\mathbf{A}) = \mathbf{0}$.

因此线性方程组的解要么没有解, 要么有唯一解. 在超定方程组中, $\mathbf{y} \notin \text{Range}(\mathbf{A})$ 是很常见的, 因此引入近似解的概念, 近似解使得 $\mathbf{A}\mathbf{x}$ 与 \mathbf{y} 在合适的度量下距离最小.

欠定方程组

线性方程组 $A\mathbf{x} = \mathbf{y}$ 中未知变量的个数小于方程组的个数时, 我们说线性方程组是欠定的. 或者说 $A^{m \times n}$ 的列数大于行数: $m < n$. 假设 A 是一个行满秩矩阵, 也就是说 $\text{rank}(A) = m$, $\text{Range}(A) = \mathbb{R}^m$. 则根据定理 5.1.6:

$$\text{rank}(A) + \dim(\text{Null}(A)) = n \quad (5.17)$$

因此 $\dim(\text{Null}(A)) = n - m > 0$. 此时线性方程组有解且有无限多个解, 并且解集的维度是 $n - m$. 在所有可能的解中, 我们总是对具有最小范数的解很感兴趣.(后面详细讨论).

适定方程组

线性方程组 $A\mathbf{x} = \mathbf{y}$ 中线性方程的个数等于未知变量的个数时, 我们说 $A\mathbf{x} = \mathbf{y}$ 是适定方程组. 或者说 $A^{m \times n}$ 的列数等于行数: $m = n$. 如果系数矩阵是满秩的即 A 可逆, A^{-1} 唯一且有 $A^{-1}A = I$. 在这种情况下, 线性方程组的解是唯一的:

$$\mathbf{x} = A^{-1}\mathbf{y} \quad (5.18)$$

注意, 实际中我们几乎不会通过先求 A^{-1} 再乘以向量 \mathbf{y} 的方式求解 \mathbf{x} . 而是通过数值方法 (比如之前学过的 LU 分解, Cholesky 分解) 来计算线性方程组非奇异方程组的解.

5.1.2 三角形线性方程组

在前面, 我们考虑了理论上如何对一个一般的线性方程组求解。

接下来, 我们将考虑如何使用计算机对一个线性方程组求解, 尤其是一个规模巨大的线性方程组。

首先来考虑一个稍微简单些的情况, 三角形线性方程组。

定义 5.1.4. 如果一个矩阵 A 主对角线以上所有元素为 0, 则称其为下三角矩阵。

如果一个矩阵 A 主对角线以下的所以元素为 0, 则称其为上三角矩阵。

定义 5.1.5. 如果一个线性方程组 $A\mathbf{x} = \mathbf{b}$ 的系数矩阵 A 是上三角形矩阵或者下三角矩阵, 我们则称其为上三角形线性方程组或者下三角形线性方程组。

针对上三角形线性方程组和下三角形线性方程组, 我们可以分别用两种特别的方法解出方程组的解。

接下来, 我们分别介绍这两种方法。

我们利用前代法计算下三角形线性方程组。

注意, 我们要求系数矩阵主对角线上元素均非 0。从而保证方程组有且仅有一个解。

$$\begin{pmatrix} a_{11} & & & & \\ a_{21} & a_{22} & & & \\ a_{31} & a_{32} & a_{33} & & \\ \vdots & \vdots & \vdots & \ddots & \\ a_{n1} & a_{n2} & a_{n3} & \dots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_n \end{pmatrix}$$

其中 $a_{11}, a_{22}, \dots, a_{nn}$ 非 0.

在前代法的第 (k) 个循环中, 我们将会遇到下面这样一个形式

$$\begin{pmatrix} a_{11} & & & & & & \\ 0 & a_{22} & & & & & \\ 0 & 0 & a_{33} & & & & \\ \vdots & \vdots & \vdots & \ddots & & & \\ 0 & 0 & 0 & \dots & a_{kk} & & \\ 0 & 0 & 0 & \dots & a_{k+1k} & a_{k+1k+1} & \\ \vdots & \vdots & \vdots & & \vdots & \vdots & \ddots \\ 0 & 0 & 0 & \dots & a_{nk} & a_{nk+1} & \dots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2^{(1)} \\ b_3^{(2)} \\ \vdots \\ b_k^{(k-1)} \\ b_{k+1}^{(k-1)} \\ \vdots \\ b_n^{(k-1)} \end{pmatrix}$$

此时我们将第 k 列从第 $k+1$ 行到第 n 行化为 0, 同时更新 b 。

$$\begin{pmatrix} a_{11} & & & & & & \\ 0 & a_{22} & & & & & \\ 0 & 0 & a_{33} & & & & \\ \vdots & \vdots & \vdots & \ddots & & & \\ 0 & 0 & 0 & \dots & a_{kk} & & \\ 0 & 0 & 0 & \dots & 0 & a_{k+1k+1} & \\ \vdots & \vdots & \vdots & & \vdots & \vdots & \ddots \\ 0 & 0 & 0 & \dots & 0 & a_{nk+1} & \dots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2^{(1)} \\ b_3^{(2)} \\ \vdots \\ b_k^{(k-1)} \\ b_{k+1}^{(k)} \\ \vdots \\ b_n^{(k)} \end{pmatrix}$$

所以前代法, 就从前 (x_1) 往后 (x_n) 来依次求解。

回代法则恰好相反, 他是从后往前一次求解。

回代法是用于上三角形的线性方程组求解。

同样我们要求其系数矩阵对角线上元素非 0.

$$\begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ & a_{22} & \cdots & a_{2n} \\ & & \ddots & \vdots \\ & & & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}$$

与前代法类似，在回代法第 $(n - k + 1)$ 个循环内。

$$\begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1k-1} & a_{1k} & 0 & \cdots & 0 \\ & a_{22} & \cdots & a_{2k-1} & a_{2k} & 0 & \cdots & 0 \\ & & \ddots & & \vdots & \vdots & \vdots & \\ & & & a_{k-1k-1} & a_{k-1k} & 0 & \cdots & 0 \\ & & & & a_{kk} & 0 & \cdots & 0 \\ & & & & & a_{k+1k+1} & \cdots & 0 \\ & & & & & & \ddots & \vdots \\ & & & & & & & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1^{(n-k)} \\ b_2^{(n-k)} \\ \vdots \\ b_{k-1}^{(n-k)} \\ b_k^{(n-k)} \\ b_{k+1}^{(n-k-1)} \\ \vdots \\ b_n \end{pmatrix}$$

此时我们将第 k 列从第 1 行到第 $k - 1$ 行化为 0，同时更新 b 。

$$\begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1k-1} & 0 & 0 & \cdots & 0 \\ & a_{22} & \cdots & a_{2k-1} & 0 & 0 & \cdots & 0 \\ & & \ddots & & \vdots & \vdots & \vdots & \\ & & & a_{k-1k-1} & 0 & 0 & \cdots & 0 \\ & & & & a_{kk} & 0 & \cdots & 0 \\ & & & & & a_{k+1k+1} & \cdots & 0 \\ & & & & & & \ddots & \vdots \\ & & & & & & & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1^{(n-k+1)} \\ b_2^{(n-k+1)} \\ \vdots \\ b_{k-1}^{(n-k+1)} \\ b_k^{(n-k)} \\ b_{k+1}^{(n-k-1)} \\ \vdots \\ b_n \end{pmatrix}$$

5.1.3 矩阵分解解线性方程组

LU 分解解方程

如果线性方程组 $Ax = b$ 的系数矩阵 A 可以进行 LU 分解，则我们可以对 A 先进行 LU 分解。

令 $A = LU$ ，我们则可以得到方程组 $LUx = b$ 。令 $Ux = y$ ，则我们可以得到两个线性方程组。

$$Ly = b$$

$$Ux = y$$

然后利用前代法求出 \mathbf{y} , 再利用回代法求解出原线性方程组的解 \mathbf{x} 。

因为在计算机上求解方程, 我们还需要考虑资源问题, 为了节约资源, 下面给出一种紧凑的求解方式。

给定矩阵 \mathbf{A} 和向量 \mathbf{b} , 我们先对 \mathbf{A} 进行 LU 分解。并且使用 \mathbf{A} 的上三角部分存储上三角矩阵, 用下三角部分存储下三角矩阵。

比如矩阵

$$\begin{pmatrix} 3 & 2 & -1 \\ 6 & 6 & -2 \\ -3 & 2 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & 2 & 1 \end{pmatrix} \begin{pmatrix} 3 & 2 & -1 \\ 0 & 2 & 0 \\ 0 & 0 & -1 \end{pmatrix}$$

就可以使用

$$\begin{pmatrix} 3 & 2 & -1 \\ 2 & 2 & 0 \\ -1 & 2 & -1 \end{pmatrix}$$

来存储。

然后我们来重新给出 LU 分解的算法流程。

$$u_{1i} = a_{1i}, \quad i = 1, \dots, n;$$

$$l_{i1} = a_{i1}/u_{11}, i = 2, \dots, n;$$

对 $k = 2, \dots, n$, 计算

$$u_{ki} = a_{ki} - \sum_{r=1}^{k-1} l_{kr} u_{ri}, \quad i = k, \dots, n;$$

$$l_{ik} = (a_{ik} - \sum_{r=1}^{k-1} l_{kr} u_{ri})/u_{kk}, \quad k = 2, \dots, n. \text{ 最后再使用前代法和回代法求出}$$

最终的解。

例 5.1.2. 求解 $\begin{pmatrix} 3 & 2 & -1 \\ 6 & 6 & -2 \\ -3 & 2 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ -2 \\ -5 \end{pmatrix}$

我们先对 $\begin{pmatrix} 3 & 2 & -1 \\ 6 & 6 & -2 \\ -3 & 2 & 0 \end{pmatrix}$ LU 分解。

$$\rightarrow \begin{pmatrix} 3 & 2 & -1 \\ 2 & 6 & -2 \\ -1 & 2 & 0 \end{pmatrix} \rightarrow \begin{pmatrix} 3 & 2 & -1 \\ 2 & 2 & 0 \\ -1 & 4 & -1 \end{pmatrix} \rightarrow \begin{pmatrix} 3 & 2 & -1 \\ 2 & 2 & 0 \\ -1 & 2 & -1 \end{pmatrix} \rightarrow \begin{pmatrix} 3 & 2 & -1 \\ 2 & 2 & 0 \\ -1 & 2 & -1 \end{pmatrix}$$

然后再进行前代法，注意此处，我们只关心 \mathbf{y} ，不关注 \mathbf{L} 如何变换，故 \mathbf{L} 并不需要实际上去化为 $\mathbf{0}$ 。可得

$$\hat{\mathbf{y}} = \begin{pmatrix} 0 \\ -2 \\ -1 \end{pmatrix}$$

最后进行回代法，便得解 $(1, -1, 1)^T$ 。

其他类型的分解解法

对于正定矩阵还可以使用 Cholesky 分解来进行求解。对于方程组

$$\mathbf{Ax} = \mathbf{b}$$

其中 \mathbf{A} 为正定矩阵。

我们可将其分解得到 $\mathbf{A} = \mathbf{L}^T \mathbf{L}$ 。然后再按照上述前代法和回代法的方式得到方程的解。

对于非方阵的情况，我们还可以使用 QR 分解来求解方程组。

对于方程组 $\mathbf{Ax} = \mathbf{b}$,

我们先对 \mathbf{A} 进行 QR 分解，得到 $\mathbf{QRx} = \mathbf{b}$ 即可得一上三角形线性方程组的同解方程组 $\mathbf{Rx} = \mathbf{Q}^T \mathbf{b}$ 。再利用回代法求解。

5.1.4 敏度分析与其他方法

在本节中，我们分析了数据中小扰动对非奇异方阵线性方程解的影响。

输入的扰动敏感性

令 \mathbf{x} 为线性方程 $\mathbf{Ax} = \mathbf{y}$ 的解，其中 \mathbf{A} 为非奇异方阵，且 $\mathbf{y} \neq \mathbf{0}$ 。假设我们通过向它添加一个小的扰动项 $\delta\mathbf{y}$ 来略微改变 \mathbf{y} ，并将 $\mathbf{x} + \delta\mathbf{x}$ 称为扰动方程组的解：

$$\mathbf{A}(\mathbf{x} + \delta\mathbf{x}) = \mathbf{y} + \delta\mathbf{y}$$

我们的关键问题是：如果 $\delta\mathbf{y}$ 变小， $\delta\mathbf{x}$ 将会不会变小？我们从上面的公式看出，并且从 $\mathbf{Ax} = \mathbf{y}$ 的事实看，扰动 $\delta\mathbf{x}$ 本身就是线性方程组的解。

$$\mathbf{A}\delta\mathbf{x} = \delta\mathbf{y}$$

并且，由于认为 \mathbf{A} 是可逆的，我们可以写成

$$\delta\mathbf{x} = \mathbf{A}^{-1}\delta\mathbf{y}$$

采用该方程两边的欧几里德范数得出

$$\|\delta\mathbf{x}\|_2 = \|\mathbf{A}^{-1}\delta\mathbf{y}\|_2 \leq \|\mathbf{A}^{-1}\|_2 \|\delta\mathbf{y}\|_2$$

其中 $\|A^{-1}\|_2$ 是 A^{-1} 的谱 (最大奇异值) 范数。类似地, 从 $A\mathbf{x} = \mathbf{y}$ 得出 $\|\mathbf{y}\|_2 = \|A\mathbf{x}\|_2 \leq \|A\|_2 \|\mathbf{x}\|_2$, 因此

$$\|\mathbf{x}\|_2^{-1} \leq \frac{\|A\|_2}{\|\mathbf{y}\|_2}$$

将上面两个公式相乘, 我们得到

$$\frac{\|\delta\mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \|A^{-1}\|_2 \|A\|_2 \frac{\|\delta\mathbf{y}\|_2}{\|\mathbf{y}\|_2}$$

这个结果是我们正在寻找的, 因为它将“输入项” \mathbf{y} 的相对变化与“输出” \mathbf{x} 的相对变化联系起来。数量

$$\kappa(A) = \|A^{-1}\|_2 \|A\|_2, 1 \leq \kappa(A) \leq \infty$$

是矩阵 A 的条件数, 见方程 (5.5)。大的 $\kappa(A)$ 意味着 \mathbf{y} 上的扰动在 \mathbf{x} 上被极大地放大, 即方程对输入数据的变化非常敏感。如果 A 是奇异的, 那么 $\kappa = \infty$ 。非常大的 $\kappa(A)$ 表明 A 接近奇异; 我们说在这种情况下 A 是病态的。我们在以下引理中总结了我们的发现。

引理 5.1.1. (对于输入的敏感性) 令 A 为非奇异方阵, $\mathbf{x}, \delta\mathbf{x}$ 满足

$$A\mathbf{x} = \mathbf{y}$$

$$A(\mathbf{x} + \delta\mathbf{x}) = \mathbf{y} + \delta\mathbf{y}$$

然后它认为

$$\frac{\|\delta\mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \kappa(A) \frac{\|\delta\mathbf{y}\|_2}{\|\mathbf{y}\|_2}$$

其中 $\kappa(A) = \|A^{-1}\|_2 \|A\|_2$ 是矩阵 A 的条件数

系数矩阵中的扰动敏感性

接下来我们考虑 A 矩阵的扰动对 \mathbf{x} 的影响。令 $A\mathbf{x} = \mathbf{y}$ 并且令 δA 为一个扰动, 满足下面等式

$$(A + \delta A)(\mathbf{x} + \delta\mathbf{x}) = \mathbf{y}, \quad \text{对于一些 } \delta\mathbf{x}$$

那么有

$$A\delta\mathbf{x} = -\delta A(\mathbf{x} + \delta\mathbf{x})$$

因此 $\delta\mathbf{x} = -A^{-1}\delta A(\mathbf{x} + \delta\mathbf{x})$ 。则

$$\|\delta\mathbf{x}\|_2 = \|A^{-1}\delta A(\mathbf{x} + \delta\mathbf{x})\|_2 \leq \|A^{-1}\|_2 \|\delta A\|_2 \|\mathbf{x} + \delta\mathbf{x}\|_2$$

并且

$$\frac{\|\delta\mathbf{x}\|_2}{\|\mathbf{x} + \delta\mathbf{x}\|_2} \leq \|A^{-1}\|_2 \|A\|_2 \frac{\|\delta A\|_2}{\|A\|_2}$$

我们再次看到只有在条件数不是太大时, 小扰动 $\frac{\|\delta A\|_2}{\|A\|_2} \ll 1$ 对 \mathbf{x} 的相对影响才很小。就是说, 它离 1 不太远, $\kappa(A) \simeq 1$ 。这个会在下一个引理中总结。

引理 5.1.2. (系数矩阵中的扰动敏感性) 令 A 为非奇异方阵, $x, \delta A, \delta x$ 满足

$$\begin{aligned} Ax &= y \\ (A + \delta A)(x + \delta x) &= y \end{aligned}$$

然后它认为

$$\frac{\|\delta x\|_2}{\|x + \delta x\|_2} \leq \kappa(A) \frac{\|\delta A\|_2}{\|A\|_2}$$

对 A, y 联合扰动的敏感性

我们最后考虑了 A 和 y 的同时扰动对 x 的影响。令 $Ax = y$, 并且令 $\delta A, \delta y$ 为扰动, 满足下面等式

$$(A + \delta A)(x + \delta x) = y + \delta y, \text{ 对于一些 } \delta x$$

然后, $A\delta x = \delta y - \delta A(x + \delta x)$, 因此 $\delta x = A^{-1}\delta y - A^{-1}\delta A(x + \delta x)$ 。则

$$\begin{aligned} \|\delta x\|_2 &= \|A^{-1}\delta y - A^{-1}\delta A(x + \delta x)\|_2 \\ &\leq \|A^{-1}\delta y\|_2 + \|A^{-1}\delta A(x + \delta x)\|_2 \\ &\leq \|A^{-1}\|_2 \|\delta y\|_2 + \|A^{-1}\|_2 \|\delta A\|_2 \|x + \delta x\|_2 \end{aligned}$$

接着, 上式除以 $\|x + \delta x\|_2$,

$$\frac{\|\delta x\|_2}{\|x + \delta x\|_2} \leq \|A^{-1}\|_2 \frac{\|\delta y\|_2}{\|y\|_2} \frac{\|y\|_2}{\|x + \delta x\|_2} + \kappa(A) \frac{\|\delta A\|_2}{\|A\|_2}$$

但是 $\|y\|_2 = \|Ax\|_2 \leq \|A\|_2 \|x\|_2$, 因此

$$\frac{\|\delta x\|_2}{\|x + \delta x\|_2} \leq \kappa(A) \frac{\|\delta y\|_2}{\|y\|_2} \frac{\|x\|_2}{\|x + \delta x\|_2} + \kappa(A) \frac{\|\delta A\|_2}{\|A\|_2}$$

下一步, 我们根据 $\|x\|_2 = \|x + \delta x - \delta x\|_2 \leq \|x + \delta x\|_2 + \|\delta x\|_2$ 去写

$$\frac{\|\delta x\|_2}{\|x + \delta x\|_2} \leq \kappa(A) \frac{\|\delta y\|_2}{\|y\|_2} \frac{\|x\|_2}{\|x + \delta x\|_2} + \kappa(A) \frac{\|\delta A\|_2}{\|A\|_2}$$

从中我们得到

$$\frac{\|\delta x\|_2}{\|x + \delta x\|_2} \leq \kappa(A) \frac{\|\delta y\|_2}{\|y\|_2} \left(1 + \frac{\|\delta x\|_2}{\|x + \delta x\|_2} \right) + \kappa(A) \frac{\|\delta A\|_2}{\|A\|_2}$$

因此

$$\frac{\|\delta x\|_2}{\|x + \delta x\|_2} \leq \frac{\kappa(A)}{1 - \kappa(A) \frac{\|\delta y\|_2}{\|y\|_2}} \left(\frac{\|\delta y\|_2}{\|y\|_2} + \frac{\|\delta A\|_2}{\|A\|_2} \right)$$

扰动的“放大因子”是受 $\frac{\kappa(A)}{1 - \kappa(A) \frac{\|\delta y\|_2}{\|y\|_2}}$ 的约束。因此, 该界限小于某些给定的 γ , 如果

$$\kappa(A) \leq \frac{\gamma}{1 + \gamma \frac{\|\delta y\|_2}{\|y\|_2}}$$

因此, 我们看到关节扰动的影响仍然由 A 的条件数控制, 如下所述

引理 5.1.3. (对 A, y 扰动的敏感性) 令 A 为非奇异方阵, 令 $\gamma > 1$ 已知, 并且令 $x, \delta y, \delta A, \delta x$ 满足下面等式

$$Ax = y$$

$$(A + \delta A)(x + \delta x) = y + \delta y$$

然后

$$\kappa(A) \leq \frac{\gamma}{1 + \gamma \frac{\|\delta y\|_2}{\|y\|_2}}$$

这意味着

$$\frac{\|\delta x\|_2}{\|x + \delta x\|_2} \leq \gamma \left(\frac{\|\delta y\|_2}{\|y\|_2} + \frac{\|\delta A\|_2}{\|A\|_2} \right)$$

5.2 最小二乘问题

本节的重点是如何通过解决最小二乘问题来解决超定方程组和欠定方程组。在本节开始我们先看一个在解决最小二乘问题中常用的定理:

定理 5.2.1. 下面 2 条性质成立:

1. $A \in \mathbb{R}^{m \times n}$ 是一个列满秩矩阵 (即 $\text{rank}(A) = n$) 当且仅当 $A^T A$ 是可逆的.
2. $A \in \mathbb{R}^{m \times n}$ 是一个行满秩矩阵 (即 $\text{rank}(A) = m$) 当且仅当 AA^T 是可逆的.

证明. 对于第一点. 如果 $A^T A$ 不是可逆的, 则存在 $x \neq 0$ 使得 $A^T Ax = 0$. $x^T A^T Ax = 0$, 因此 $Ax = 0$. 所以 A 不是一个列满秩矩阵. 反之, 如果 $A^T A$ 是可逆的, 对于每个 $x \neq 0, A^T Ax \neq 0$, 也能推出对于每一个非零的 $x, Ax \neq 0$. 第二点的证明过程与第一点的证明过程相似. \square

5.2.1 最小二乘问题与线性回归

最小二乘问题多产生于线性回归或者数据拟合问题。例如, 假定给出 m 个点 t_1, \dots, t_m 和这 m 个点上的试验或观测数据 y_1, \dots, y_m , 其中 $t_i \in \mathbb{R}^n, i = 1, \dots, m$ 。再假定给出在 t_i 上取值的 n 个已知函数 $a_1(t), \dots, a_n(t)$ 。考虑 a_i 的线性组合

$$f(x; t) = x_1 a_1(t) + x_2 a_2(t) + \dots + x_n a_n(t)$$

我们希望在 t_1, \dots, t_m 点上 $f(x; t)$ 能最佳地逼近这些数据 y_1, \dots, y_m 。为此, 若定义残量

$$r_i(x) = y_i - \sum_{j=1}^n x_j a_j(t_i), \quad i = 1, \dots, m$$

则问题成为: 估计参数 x_1, \dots, x_n , 使残量 r_1, \dots, r_m 尽可能小。可用矩阵向量形式表示为

$$r(x) = b - Ax$$

其中

$$\mathbf{A} = \begin{bmatrix} a_1(t_1) & \cdots & a_n(t_1) \\ \vdots & & \vdots \\ a_1(t_m) & \cdots & a_n(t_m) \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

$$\mathbf{x} = (x_1, \cdots, x_n)^T, \quad \mathbf{r}(\mathbf{x}) = (r_1(\mathbf{x}), \cdots, r_m(\mathbf{x}))$$

当 $m = n$ 时, 我们可以要求 $\mathbf{r}(\mathbf{x}) = 0$, 则估计 x 的问题就可用上一节关于线性方程组中讨论的方法解决。当 $m > n$ 时, 一般不可能使所有残量为 0, 但我们可要求残差向量在某种范数意义下达最小, 最小二乘问题就是求 \mathbf{x} 使得残差向量 $\mathbf{r}(\mathbf{x})$ 在 2 范数意义下最小。下面给出最小二乘问题的定义。

定义 5.2.1. 给定矩阵 $\mathbf{A} \in \mathbb{R}^{m \times n}$ 和向量 $\mathbf{b} \in \mathbb{R}^m$, 确定 $x \in \mathbb{R}^n$ 使得

$$\|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2 = \|\mathbf{r}(\mathbf{x})\|_2 = \min_{\mathbf{y} \in \mathbb{R}^n} \|\mathbf{r}(\mathbf{y})\|_2 = \min_{\mathbf{y} \in \mathbb{R}^n} \|\mathbf{A}\mathbf{y} - \mathbf{b}\|_2 \quad (5.19)$$

该问题称为最小二乘问题, 简记为 $LS(Least Squares)$ 问题, 其中 $\mathbf{r}(\mathbf{x})$ 称为残差向量。

如果残差向量 \mathbf{r} 线性依赖于 \mathbf{x} , 则称其为线性最小二乘问题; 如果 \mathbf{r} 非线性的依赖于 \mathbf{x} , 则称其为非线性最小二乘问题。我们主要讨论线性最小二乘问题, 简称最小二乘问题。对于残差向量选择不同的范数, 便得到不同的问题, 我们主要讨论残差向量选择 2 范数的情况。

最小二乘问题的解 \mathbf{x} 又称为线性方程组(5.20)的最小二乘解。

$$\mathbf{A}\mathbf{x} = \mathbf{b}, \quad \mathbf{A} \in \mathbb{R}^{m \times n} \quad (5.20)$$

即残差向量 $\mathbf{r}(\mathbf{x})$ 的 2 范数最小的意义下满足方程组(5.20)。

根据 m 与 n 以及矩阵 \mathbf{A} 的秩 $r(\mathbf{A})$ 的不同, 最小二乘问题可分为下面几种情况:

(1) $m = n$

$$(1a) r(\mathbf{A}) = m = n$$

$$(1b) r(\mathbf{A}) < m = n$$

(2) $m > n$

$$(2a) r(\mathbf{A}) = n < m$$

$$(2b) r(\mathbf{A}) < n < m$$

(3) $m < n$

$$(3a) r(\mathbf{A}) = m < n$$

$$(3b) r(\mathbf{A}) < m < n$$

其中 $m > n$ 对应的方程组称为超定方程组, $m < n$ 对应的方程组称为欠定方程组。不同的情况通常需要不同的方法去处理, 我们主要讨论 (2a) 情况的最小二乘问题的性质与求解。

定理 5.2.2. 线性最小二乘问题的解总是存在的, 而且其解唯一的充分必要条件是 $\text{null}(\mathbf{A}) = 0$ 。

证明. 因为 $\mathbb{R}^m = \mathcal{R}(\mathbf{A}) \oplus \mathcal{R}(\mathbf{A})^\perp$, 所以向量 \mathbf{b} 可以唯一地表示 $\mathbf{b} = \mathbf{b}_1 + \mathbf{b}_2$, 其中 $\mathbf{b}_1 \in \mathcal{R}(\mathbf{A}), \mathbf{b}_2 \in \mathcal{R}(\mathbf{A})^\perp$. 于是对任意 $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{b}_1 - \mathbf{A}\mathbf{x} \in \mathcal{R}(\mathbf{A})$ 且与 \mathbf{b}_2 正交, 从而

$$\begin{aligned}\|\mathbf{r}(\mathbf{x})\|_2^2 &= \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2 = \|(\mathbf{b}_1 - \mathbf{A}\mathbf{x}) + \mathbf{b}_2\|_2^2 \\ &= \|\mathbf{b}_1 - \mathbf{A}\mathbf{x}\|_2^2 + \|\mathbf{b}_2\|_2^2\end{aligned}$$

由此可得, $\|\mathbf{r}(\mathbf{x})\|_2^2$ 达到极小当且仅当 $\|\mathbf{b}_1 - \mathbf{A}\mathbf{x}\|_2^2$ 达到极小, 而 $\mathbf{b}_1 \in \mathcal{R}(\mathbf{A})$ 又蕴含着 $\|\mathbf{b}_1 - \mathbf{A}\mathbf{x}\|_2^2$ 达到极小地充要条件 $\mathbf{A}\mathbf{x} = \mathbf{b}_1$, 综上, 由 $\mathbf{b}_1 \in \mathcal{R}(\mathbf{A})$ 和推论 5.1.1 推出定理的结论成立。□

记最小二乘问题的解集为 \mathbb{S}_{LS} , 即

$$\mathbb{S}_{LS} = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} \text{ 是 LS 问题的解}\}$$

则由定理 5.2.2 知, \mathbb{S}_{LS} 总是非空的, 而且它仅有一个元素的充要条件是 \mathbf{A} 的列线性无关。

5.2.2 欠定问题的最小范数解

我们要考虑下一个情况: 当矩阵 \mathbf{A} 的列数比行数多: $m < n$. 假设矩阵 \mathbf{A} 是行满秩, 我们有 $\dim\{\text{Null}(\mathbf{A})\} = n - m > 0$, 因此得出 $\mathbf{y} = \mathbf{A}\mathbf{x}$ 有无数个解并且解的集合是 $\mathbb{S}_x = \{\mathbf{x} : \mathbf{x} = \tilde{\mathbf{x}} + \mathbf{z}, \mathbf{z} \in \text{Null}(\mathbf{A})\}$, 其中 $\tilde{\mathbf{x}}$ 是任意满足 $\mathbf{A}\tilde{\mathbf{x}} = \mathbf{y}$ 的向量。

(我们有兴趣从这个解的集合 \mathbb{S}_x 中挑选出一个满足最小欧几里得范数的解 \mathbf{x}^*)

我们想解决的问题:

$$\min_{\mathbf{x}: \mathbf{A}\mathbf{x}=\mathbf{y}} \|\mathbf{x}\|_2$$

这个式子等价于 $\min_{\mathbf{x} \in \mathbb{S}_x} \|\mathbf{x}\|_2$. (唯一的) 解 \mathbf{x}^* 必须与 $\text{Null}(\mathbf{A})$ 相互垂直, 等价地, $\mathbf{x}^* \in \text{Range}(\mathbf{A}^T)$ 这个意味着存在 ζ , $\mathbf{x}^* = \mathbf{A}^T \zeta$. 因为 \mathbf{x}^* 是方程组的解, 必须满足 $\mathbf{A}\mathbf{x}^* = \mathbf{y}$, 就等于说 $\mathbf{A}\mathbf{A}^T \zeta = \mathbf{y}$.

因为矩阵 \mathbf{A} 是行满秩, $\mathbf{A}\mathbf{A}^T$ 是可逆的并且有唯一的 ζ 是方程组的解。 $\zeta = (\mathbf{A}\mathbf{A}^T)^{-1}\mathbf{y}$.

这样我们得到了唯一的最小范数解:

$$\mathbf{x}^* = \mathbf{A}^T (\mathbf{A}\mathbf{A}^T)^{-1} \mathbf{y}$$

5.2.3 最小二乘问题的求解方法

正则化法

定理 5.2.3. $\mathbf{x} \in \mathbb{S}_{LS}$ 的充要条件是 \mathbf{x} 为正规方程组(5.21)的解

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b} \quad (5.21)$$

证明. 设 $\mathbf{x} \in \mathbb{S}_{LS}$, 由定理5.2.2证明知 $\mathbf{Ax} = \mathbf{b}_1$, 其中 $\mathbf{b}_1 \in \mathbb{R}(\mathbf{A})$, 而且 $\mathbf{r}(\mathbf{x}) = \mathbf{b} - \mathbf{Ax} = \mathbf{b} - \mathbf{b}_1 = \mathbf{b}_2 \in \mathbb{R}(\mathbf{A})^\perp$. 所以 $\mathbf{A}^T \mathbf{r}(\mathbf{x}) = \mathbf{A}^T \mathbf{b}_2 = \mathbf{0}$. 将 $\mathbf{r}(\mathbf{x}) = \mathbf{b} - \mathbf{Ax}$ 代入 $\mathbf{A}^T \mathbf{r}(\mathbf{x}) = \mathbf{0}$ 即得定理5.21。反之, 设 $\mathbf{x} \in \mathbb{R}^n$ 满足 $\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}$, 则对任意的 $\mathbf{y} \in \mathbb{R}^n$ 有

$$\begin{aligned}\|\mathbf{b} - \mathbf{A}(\mathbf{x} + \mathbf{y})\|_2^2 &= \|\mathbf{b} - \mathbf{Ax}\|_2^2 - 2\mathbf{y}^T \mathbf{A}^T (\mathbf{b} - \mathbf{Ax}) + \|\mathbf{Ay}\|_2^2 \\ &= \|\mathbf{b} - \mathbf{Ax}\|_2^2 + \|\mathbf{Ay}\|_2^2 \\ &\geq \|\mathbf{b} - \mathbf{Ax}\|_2^2\end{aligned}$$

由此即得 $\mathbf{x} \in \mathbb{S}_{LS}$ 。 \square

定理 5.2.4. 设 $\mathbf{A} \in \mathbb{R}^{m \times n} (m \geq n)$, 当 \mathbf{A} 的列满秩, $\mathbf{A}^T \mathbf{A}$ 对称正定可逆, 最小二乘问题的解是唯一的, 表达式为 $\mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$ 。其中 $(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ 又称为 \mathbf{A} 的广义逆, 用符号 \mathbf{A}^\dagger 表示。

证明. 因为 \mathbf{A} 列满秩, 所以 $\mathbf{A}^T \mathbf{A}$ 可逆, 正规方程 $\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}$ 两边同时乘 $(\mathbf{A}^T \mathbf{A})^{-1}$ 得到 $\mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$ 。 \square

矩阵的广义逆

不是所有的矩阵都有逆矩阵, 只有满秩的方阵才有逆矩阵, 即 $\text{rank}(\mathbf{A}) = n$, 那么对于不是方阵, 不满秩矩阵的情况, 数学中也引申出类似于逆矩阵的定义:

定义 5.2.2. 设 $\mathbf{A} = (a_{ij})_{m \times n}$ 是数域 \mathbb{R} 上的矩阵, 如果有唯一一个矩阵 \mathbf{X} , 满足

$$\begin{cases} \mathbf{AXA} = \mathbf{A}, \\ \mathbf{XAX} = \mathbf{X}, \\ \mathbf{AX} = (\mathbf{AX})^T, \\ \mathbf{XA} = (\mathbf{XA})^T \end{cases}$$

则矩阵 \mathbf{X} 称为矩阵 \mathbf{A} 的广义逆, 记作 \mathbf{A}^+ 。

例 5.2.1. 设矩阵 \mathbf{A}

$$\mathbf{A} = \begin{pmatrix} 1 & 2 & 3 \\ 0 & 1 & 2 \end{pmatrix}$$

则矩阵 \mathbf{A} 的广义逆为

$$\mathbf{A}^+ = \frac{1}{6} \begin{pmatrix} 5 & -8 \\ 2 & -2 \\ -1 & 4 \end{pmatrix}$$

广义逆扩充了逆矩阵的范围, 将逆矩阵扩大到不是方阵, 不满秩的情况。特别地, 如果矩阵 \mathbf{A} 本身有逆矩阵 \mathbf{A}^{-1} , 则 \mathbf{A}^{-1} 显然满足式 (2.1), 即 $\mathbf{A}^+ = \mathbf{A}^{-1}$ 。因此, 矩阵 \mathbf{A} 的广义逆矩阵 \mathbf{A}^+ 可视为 \mathbf{A} 的逆 \mathbf{A}^{-1} 的推广。

当 A 列满秩时, 我们可以使用 Cholesky 分解法来求解正规方程组。基本步骤如下:

- (1) 计算 $C = A^T A$, $d = A^T b$;
- (2) 用平方根法计算 C 的 Cholesky 分解: $C = LL^T$;
- (3) 求解三角方程组 $Ly = d$ 和 $L^T x = y$ 。

Cholesky 分解法求解的运算量 (乘法次数) 大约为 $mn^2 + \frac{1}{3}n^3 + O(n^2)$, 其中 mn^2 对应计算 $A^T A$, $\frac{1}{3}n^3$ 对应 n 矩阵的 Cholesky 分解, $O(n^2)$ 对应求解三角方程组。

QR 分解法

设 $A \in \mathbb{R}^{m \times n} (m \geq n)$ 列满秩, $b \in \mathbb{R}^m$, 将矩阵 A 分解成列正交矩阵 Q 和上三角矩阵 R 的乘积, 即

$$A = QR$$

然后将列正交矩阵 Q 扩充为正交矩阵 $\begin{bmatrix} Q & \hat{Q} \end{bmatrix} \in \mathbb{R}^{m \times m}$, 得到

$$A = \begin{bmatrix} Q & \hat{Q} \end{bmatrix} \begin{bmatrix} R \\ O \end{bmatrix}$$

根据 2 范数的正交不变性得

$$\begin{aligned} \|Ax - b\|_2^2 &= \left\| \begin{bmatrix} Q & \hat{Q} \end{bmatrix}^T (Ax - b) \right\|_2^2 \\ &= \left\| \begin{bmatrix} Q & \hat{Q} \end{bmatrix}^T (QRx - b) \right\|_2^2 \\ &= \left\| \begin{bmatrix} Rx - Q^T b \\ -\hat{Q}^T b \end{bmatrix} \right\|_2^2 \\ &= \|Rx - Q^T b\|_2^2 + \|\hat{Q}^T b\|_2^2 \\ &\geq \|\hat{Q}^T b\|_2^2 \end{aligned}$$

当且仅当 $Rx = Q^T b$ 时等号成立, 又因为 R 可逆, 所以最小二乘问题的解为

$$x = R^{-1} Q^T b.$$

QR 分解法求解最小二乘问题的基本步骤如下:

- (1) 计算 A 的 QR 分解: $A = QR$;
- (2) 计算 $c = Q^T b$;
- (3) 求解上三角方程组 $Rx = c$ 。

在矩阵分解部分, 我们介绍过 Gram-Schmidt 正交化、Householder 变换、Givens 变换三种方法进行 QR 分解。

在计算机中一般使用基于 Householder 变换的 QR 分解, 该算法有良好地数值性态, 结果通常要比正则化方法精确。但是运算量也比较大, 大约为 $2mn^2 - \frac{2}{3}n^3$ 。

我们也可以使用 Givens 变换来实现 QR 分解, 所需的运算量大约是 Householder 方法地两倍, 但是如果 A 有较多的零元素, 则灵活地使用 Givens 变换会使运算量大为减少。

奇异值分解法

设 $A \in \mathbb{R}^{m \times n} (m \geq n)$ 列满秩, $A = U \begin{bmatrix} \Sigma \\ O \end{bmatrix} V^T$ 是 A 的奇异值分解, 令 U_n 为 U 的前 n 列组成的矩阵, 即 $U = [U_n, \tilde{U}]$, 其中 U 是正交矩阵, 根据 2 范数的正交不变性得

$$\begin{aligned} \|Ax - b\|_2^2 &= \left\| U \begin{bmatrix} \Sigma \\ O \end{bmatrix} V^T x - b \right\|_2^2 \\ &= \left\| \begin{bmatrix} \Sigma \\ O \end{bmatrix} V^T x - [U_n, \tilde{U}]^T b \right\|_2^2 \\ &= \left\| \begin{bmatrix} \Sigma V^T x - U_n^T b \\ -\tilde{U}^T b \end{bmatrix} \right\|_2^2 \\ &= \|\Sigma V^T x - U_n^T b\|_2^2 + \|\tilde{U}^T b\|_2^2 \\ &\geq \|\tilde{U}^T b\|_2^2 \end{aligned}$$

等号当且仅当 $\Sigma V^T x - U_n^T b = 0$ 时成立, 即

$$x = (\Sigma V^T)^{-1} U_n^T b = V \Sigma^{-1} U_n^T b$$

5.2.4 最小二乘问题的变体

对最小二乘问题做一些修改, 会得到其他形式的最小二乘问题。

加权最小二乘

在普通的最小二乘法中, 我们想要最小化误差向量各项的平方和:

$$\|Ax - y\|_2^2 = \sum_{i=1}^m r_i^2, \quad r_i = a_i^T x - y_i$$

其中 $a_i^T, i = 1, \dots, m$ 是 A 的各列。但是, 在某些情形下, 方程的残差项并不是同样重要的, 相比其他方程, 有可能满足某一个方程更重要, 这样, 我们需要在残差项赋予权重:

$$f_0(x) = \sum_{i=1}^m w_i^2 r_i^2,$$

其中 $w_i \geq 0$ 是给定的权重。

这样最小化目标函数重写为：

$$f_0(\mathbf{x}) = \|\mathbf{W}(\mathbf{A}\mathbf{x} - \mathbf{y})\|_2^2 = \|\mathbf{A}_w\mathbf{x} - \mathbf{y}_w\|_2^2$$

其中

$$\mathbf{W} = \text{diag}(w_1, \dots, w_m), \mathbf{A}_w \doteq \mathbf{W}\mathbf{A}, \mathbf{y}_w \doteq \mathbf{W}\mathbf{y}$$

加权最小二乘仍然是普通最小二乘的形式，其权重最小二乘解为：

$$\begin{aligned} \hat{\mathbf{x}}_{\text{WLS}} &= (\mathbf{A}_w^T \mathbf{A}_w)^{-1} \mathbf{A}_w^T \mathbf{y}_w \\ &= (\mathbf{A}^T \mathbf{W}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{W}^T \mathbf{W} \mathbf{y} \end{aligned}$$

约束最小二乘

考虑带有约束的最小二乘问题

$$\begin{aligned} \min_{\mathbf{x}} \quad & \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 \\ \text{s.t.} \quad & \mathbf{B}\mathbf{x} = \mathbf{f} \end{aligned}$$

其中 $\mathbf{B}\mathbf{x} = \mathbf{f}$ 是约束条件。求解需要凸优化知识，在这里只列出解。如果 $\mathbf{A}^T \mathbf{A}$ 非奇异，且 \mathbf{B} 行满秩，则

$$\mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} (\mathbf{A}^T \mathbf{b} - \mathbf{B}^T \boldsymbol{\lambda})$$

其中 $\boldsymbol{\lambda} = [\mathbf{B}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{B}^T]^{-1} [\mathbf{B}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b} - \mathbf{f}]$ 。

总体最小二乘

考虑得到的数据矩阵和数据向量 \mathbf{A}, \mathbf{b} 都有误差，设实际观测的数据矩阵和数据向量

$$\mathbf{A} = \mathbf{A}_0 + \mathbf{E}, \quad \mathbf{b} = \mathbf{b}_0 + \mathbf{e}$$

其中 \mathbf{E} 和 \mathbf{e} 分别表示误差数据矩阵和误差数据向量。总体最小二乘的基本思想是：不仅用校正向量 $\Delta \mathbf{b}$ 去干扰数据向量 \mathbf{b} ，同时用校正矩阵 $\Delta \mathbf{A}$ 去干扰数据矩阵 \mathbf{A} ，以便对 \mathbf{A} 和 \mathbf{b} 二者内存在的误差或噪声进行联合补偿

$$\mathbf{b} + \Delta \mathbf{b} = \mathbf{b}_0 + \mathbf{e} + \Delta \mathbf{b} \rightarrow \mathbf{b}_0$$

$$\mathbf{A} + \Delta \mathbf{A} = \mathbf{A}_0 + \mathbf{E} + \Delta \mathbf{A} \rightarrow \mathbf{A}_0$$

以抑制观测误差或噪声对矩阵方程求解的影响，从而实现从有误差的矩阵方程到精确矩阵方程的求解的转换

$$(\mathbf{A} + \Delta \mathbf{A})\mathbf{x} = \mathbf{b} + \Delta \mathbf{b} \implies \mathbf{A}_0 \mathbf{x} = \mathbf{b}_0 \quad (5.22)$$

自然地，我们希望矫正数据矩阵和校正数据向量都尽量小。因此，总体最小二乘问题可以用约束优化问题叙述为：

$$\begin{aligned} \text{TLS:} \quad & \min_{\Delta \mathbf{A}, \Delta \mathbf{b}, \mathbf{x}} \|\Delta \mathbf{A}, \Delta \mathbf{b}\|_F^2 = \|\Delta \mathbf{A}\|_F^2 + \|\Delta \mathbf{b}\|_2^2 \\ \text{s.t.} \quad & (\mathbf{A} + \Delta \mathbf{A})\mathbf{x} = \mathbf{b} + \Delta \mathbf{b} \end{aligned}$$

约束条件有时也表示为 $(\mathbf{b} + \Delta\mathbf{b}) \in \text{Range}(\mathbf{A} + \Delta\mathbf{A})$

由(5.22), 校正过方程的解满足:

$$([\mathbf{A}, \mathbf{b}] + [\Delta\mathbf{A}, \Delta\mathbf{b}]) \begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix} = \mathbf{0} \quad (5.23)$$

如果 $([\mathbf{A}, \mathbf{b}] + [\Delta\mathbf{A}, \Delta\mathbf{b}])$ 是列满秩的矩阵, 记 $\hat{\mathbf{x}} = \begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix}$, 则以 $\hat{\mathbf{x}}$ 为未知量的方程:

$$([\mathbf{A}, \mathbf{b}] + [\Delta\mathbf{A}, \Delta\mathbf{b}])\hat{\mathbf{x}} = \mathbf{0} \quad (5.24)$$

只有零解, 与 $\hat{\mathbf{x}}$ 的最后一个分量为 -1 矛盾。因此, $([\mathbf{A}, \mathbf{b}] + [\Delta\mathbf{A}, \Delta\mathbf{b}])$ 是一个列亏损矩阵。问题转化为求一个最接近 $[\mathbf{A}, \mathbf{b}]$ 的列亏损矩阵。设 $[\mathbf{A}, \mathbf{b}]$ 的奇异值分解为

$$[\mathbf{A}, \mathbf{b}] = \sum_{i=1}^{n+1} \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

其中 σ_i 为 $[\mathbf{A}, \mathbf{b}]$ 的第 i 个奇异值, $\mathbf{u}_i, \mathbf{v}_i$ 分别为对应的左右奇异向量。

也就是 $[\Delta\mathbf{A}, \Delta\mathbf{b}] = \sigma_{n+1} \mathbf{u}_{n+1} \mathbf{v}_{n+1}^T$ 。

设 $\mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{x} \in \mathbb{R}^n, \mathbf{b} \in \mathbb{R}^m$, 其解为:

$$\hat{\mathbf{x}}_{\text{TLS}} = (\mathbf{A}^T \mathbf{A} - \sigma_{n+1}^2 \mathbf{I})^{-1} \mathbf{A}^T \mathbf{b}$$

其中 σ_{n+1} 为 $[-\mathbf{b}, \mathbf{A}]$ 的第 $n+1$ 个奇异值, $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_{n+1}$ 。

5.3 特征值计算

在数据科学中, 我们一般只讨论实矩阵的特征值问题。

应注意, 实矩阵的特征值和特征向量不一定是实数和实向量, 但实特征值一定对应于实特征向量 (方程 (2.14) 的解), 而一般的复特征值对应的特征向量一定不是实向量。此外, 由于特征方程为实系数方程, 若一个特征值不是实数, 则其复共轭也一定是它的特征值。

对于一个实对称矩阵来说, 它的 n 个特征值均为实数, 并且存在 n 个正交的实特征向量。

5.3.1 圆盘定理

在很多情况下, 我们并不需要确切地知道矩阵的每一个特征值的大小, 而是要估计出这个矩阵各个特征值大概的范围。

定理 5.3.1. 圆盘定理 设 $\mathbf{A} = (a_{kj}) \in \mathbb{C}^{n \times n}$, 则:

(I) \mathbf{A} 的每一个特征值必属于 \mathbf{A} 的格什戈林圆盘之中, 即对任一特征值 λ 必定存在 $k (1 \leq k \leq n)$, 使得

$$|\lambda - a_{kk}| \leq \sum_{j=1, j \neq k}^n |a_{kj}| \quad (5.25)$$

用集合的关系来说明, 这意味着 $\lambda(A) \subseteq \cup_{k=1}^n D_k$, 其中 $D_k = \{z | |z - a_{kk}| \leq \sum_{j=1, j \neq k}^n |a_{kj}|\}$

(2) 若 A 的格什戈林圆盘中有 m 个圆盘组成一连通并集 S , 且 S 与余下的 $n - m$ 个圆盘分离, 则 S 内恰好包含 A 的 m 个特征值 (重特征值按重数计)。

下面对定理5.3.1的结论 (1) 进行证明, 结论 (2) 的证明超出了本书的范围。

证明. 设 λ 为 A 的任一特征值, 则有 $Ax = \lambda x$. x 为非零常量. 设 x 中第 k 个分量最大, 即

$$|x_k| = \max_{1 \leq j \leq n} |x_j| > 0,$$

考虑线性方程中第 k 个方程

$$\sum_{j=1}^n a_{kj} x_j = \lambda x_k,$$

将其中与 x_k 有关的项移到等号左边, 其余移到右边, 再两边取模得

$$|\lambda - a_{kk}| |x_k| = \left| \sum_{\substack{j=1 \\ j \neq k}}^n a_{kj} x_j \right| \leq \sum_{\substack{j=1 \\ j \neq k}}^n |a_{kj}| |x_j| \leq |x_k| \sum_{\substack{j=1 \\ j \neq k}}^n |a_{kj}| \quad (5.26)$$

最后一个不等式的推导利用了“ x 中第 k 个分量最大”的假设. 将不等式(5.26)除以 $|x_k|$, 即得到式(5.25), 因此证明了定理5.3.1的结论 (1). 上述证明过程还说明, 若某个特征向量的第 k 个分量的模最大, 则相应的特征值必定属于第 k 个圆盘中. \square

还可以按照矩阵的每一列元素定义 n 个圆盘, 对于它们定理5.25仍然成立. 下面的定理是圆盘定理的重要推论, 其证明留给感兴趣的读者.

定理 5.3.2. 设 $A \in \mathbb{R}^{n \times n}$, 且 A 的对角元均大于 0, 则

(1) 若 A 严格对角占优, 则 A 的特征值的实部都大于 0.

(2) 若 A 为对角占优的对称矩阵, 则 A 一定是对称半正定矩阵, 若同时 A 非奇异, 则 A 为对称正定矩阵.

例 5.3.1. 试估计矩阵

$$\begin{pmatrix} 4 & 1 & 0 \\ 1 & 0 & -1 \\ 1 & 1 & -4 \end{pmatrix}$$

的特征值范围。

解. 直接应用圆盘定理, 该矩阵的三个圆盘如下:

$$D_1: |\lambda - 4| \leq 1, \quad D_2: |\lambda| \leq 2, \quad D_3: |\lambda + 4| \leq 2.$$

D_1 与其他圆盘分离, 则它仅含一个特征值, 且必定为实数 (若为虚数则其共轭也是特征值, 这与 D_1 仅含一个特征值矛盾)。所以对矩阵特征值的范围的估计是

$$3 \leq \lambda_1 \leq 5, \quad \lambda_2, \lambda_3 \in D_2 \cup D_3.$$

再对矩阵 A^T 应用圆盘定理, 则可以进一步优化上述结果。矩阵 A^T 对应的三个圆盘为

$$D'_1: |\lambda - 4| \leq 2, \quad D'_2: |\lambda| \leq 2, \quad D'_3: |\lambda + 4| \leq 1.$$

这说明 D'_3 中存在一个特征值, 且为实数, 它属于区间 $[-5, -3]$, 经过综合分析可知三个特征值均为实数, 它们的范围是

$$\lambda_1 \in [3, 5], \quad \lambda_2 \in [-2, 2], \quad \lambda_3 \in [-5, -3].$$

事实上, 使用 *MATLAB* 的 *eig* 命令可求出矩阵 A 的特征值为 4.2030, -0.4429, -3.7601.

在估计特征值范围的时候, 我们希望各个圆盘的半径越小越好。所以我们可以通过对矩阵 A 做相似变换, 例如取 X 为对角阵, 然后再应用圆盘定理估计特征值的范围.

例 5.3.2. (特征值范围的估计): 选取适当的矩阵 X , 应用定理 5.3.1 估计例 5.3.1 中矩阵的特征值范围.

解. 取

$$X^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0.9 \end{pmatrix}$$

则

$$A_1 = X^{-1}AX = \begin{pmatrix} 4 & 1 & 0 \\ 1 & 0 & -\frac{10}{9} \\ 0.9 & 0.9 & -4 \end{pmatrix}$$

的特征值与 A 的相同。对 A_1 应用圆盘定理, 得到三个分离的圆盘, 它们分别包含一个实特征值, 由此得到特征值的范围估计

$$\lambda_1 \in [3, 5], \lambda_2 \in \left[-\frac{19}{9}, \frac{19}{9}\right], \lambda_3 \in [-5.8, -2.2].$$

此外, 还可以进一步估计 $\rho(A)$ 的范围, 即 $3 \leq \rho(A) \leq 5.8$ 。

上述例子表明, 综合运用圆盘定理和矩阵特征值的性质, 可对特征值的范围进行一定的估计. 对具体例子, 可适当设置相似变换矩阵, 尽可能让圆盘相互分离, 从而提高估计的有效性。

5.3.2 幂法

幂法是通过求矩阵的特征向量来求出特征值的一种迭代法。它主要用来求按模最大的特征值和相应的特征向量。其优点是算法简单, 容易实现, 缺点是收敛速度慢, 其有效性依赖于矩阵特征值的分布情况。

适于使用幂法的常见情形是: A 的特征值可按模的大小排列为 $|\lambda_1| > |\lambda_2| \geq \cdots \geq |\lambda_n|$, 且其对应特征向量 $\xi_1, \xi_2, \dots, \xi_n$ 线性无关。此时, 任意非零向量 $x^{(0)}$ 均可用 $\xi_1, \xi_2, \dots, \xi_n$ 线性表示, 即

$$x^{(0)} = \alpha_1 \xi_1 + \alpha_2 \xi_2 + \cdots + \alpha_n \xi_n$$

且 $\alpha_1, \alpha_2, \dots, \alpha_n$ 不全为零。做向量序列 $x^{(k)} = A^k x^{(0)}$, 则

$$\begin{aligned} x^{(k)} &= A^k x^{(0)} = \alpha_1 A^k \xi_1 + \alpha_2 A^k \xi_2 + \cdots + \alpha_n A^k \xi_n \\ &= \alpha_1 \lambda_1^k \xi_1 + \alpha_2 \lambda_2^k \xi_2 + \cdots + \alpha_n \lambda_n^k \xi_n \\ &= \lambda_1^k [\alpha_1 \xi_1 + \alpha_2 (\frac{\lambda_2}{\lambda_1})^k \xi_2 + \cdots + \alpha_n (\frac{\lambda_n}{\lambda_1})^k \xi_n] \end{aligned}$$

由此可见, 若 $\alpha_1 \neq 0$, 则有

$$\lim_{k \rightarrow \infty} (\frac{\lambda_i}{\lambda_1})^k = 0, i = 2, \dots, n$$

故当 k 充分大的时候, 必有

$$x^{(k)} \approx \lambda_1^k \alpha_1 \xi_1$$

即 $x^{(k)}$ 可以近似看成 λ_1 对应的特征向量, 而 $x^{(k)}$ 与 $x^{(k-1)}$ 分量之比为

$$\frac{x^{(k)}}{x^{(k-1)}} \approx \frac{\lambda_1^k \alpha_1 \xi_1}{\lambda_1^{k-1} \alpha_1 \xi_1} = \lambda_1$$

于是利用向量序列 $\{x^{(k)}\}$ 即可求出按模最大的特征值 λ_1 , 又可以求出对应的特征向量 ξ_1 。

在实际计算中, 考虑到当 $|\lambda_1| > 1$ 时, $\lambda_1^k \rightarrow \infty$; $|\lambda_1| < 1$ 时, $\lambda_1^k \rightarrow 0$, 因而计算 $x^{(k)}$ 时可能会发生上溢或者下溢, 故每一步将 $x^{(k)}$ 归一化处理, 即将 $x^{(k)}$ 的各分量都除以模最大的分量, 使 $\|x^{(k)}\| = 1$, 于是求 A 按模最大的特征值 λ_1 和对应的特征向量 ξ_1 的算法, 可归纳为如下步骤。

- (1) 输入矩阵 A , 初始向量 $v^{(0)}$, 误差限 ϵ , 最大迭代次数 N 。记 m_0 是 $v^{(0)}$ 按模最大的分量, $x^{(0)} = v^{(0)}/m_0$ 。置 $k = 0$
- (2) 计算 $v^{(k+1)} = Ax^{(k)}$ 。记 m_{k+1} 是 $v^{(k+1)}$ 按模最大的分量, $x^{(k+1)} = v^{(k+1)}/m_{k+1}$
- (3) 若 $|m_{k+1} - m_k| < \epsilon$, 停止计算, 输出近似特征值 m_{k+1} 和近似特征向量 $x^{(k+1)}$ 否则转 (4)
- (4) 若 $k < N$ 置 $k = k + 1$ 转 (2) 否则输出计算失败信息, 停止计算

定理 5.3.3. 设矩阵 A 的特征值可按模的大小排列为 $|\lambda_1| > |\lambda_2| \geq \cdots \geq |\lambda_n|$, 且对应特征向量 $\xi_1, \xi_2, \dots, \xi_n$ 线性无关。序列 $\{x^{(k)}\}$ 有算法产生, 则有

$$\lim_{k \rightarrow \infty} x^{(k)} = \frac{\xi_1}{\max\{\xi_1\}} = \xi_1^0, \quad \lim_{k \rightarrow \infty} m_k = \lambda_1, \quad (5.27)$$

式中: ξ_1^0 为将 ξ_1 归一化后得到的向量; $\max\{\xi_1\}$ 为向量 ξ_1 模最大的分量。

证明. 由算法7.1.6的步2和步3知

$$\mathbf{x}^{(k)} = \frac{\mathbf{v}^{(k)}}{m_k} = \frac{\mathbf{A}\mathbf{x}^{(k-1)}}{m_k} = \frac{\mathbf{A}^2\mathbf{x}^{(k-2)}}{m_k m_{k-1}} = \cdots = \frac{\mathbf{A}^k \mathbf{x}^{(0)}}{m_k m_{k-1} \cdots m_1}.$$

由于 $\mathbf{x}^{(k)}$ 的最大分量为1, 即 $\max\{\mathbf{x}^{(k)}\} = 1$, 故

$$m_k m_{k-1} \cdots m_1 = \max\{\mathbf{A}^k \mathbf{x}^{(0)}\}$$

从而

$$\begin{aligned} \mathbf{x}^{(k)} &= \frac{\mathbf{A}^k \mathbf{x}^{(0)}}{\max\{\mathbf{A}^k \mathbf{x}^{(0)}\}} = \frac{\lambda_1^k [\alpha_1 \boldsymbol{\xi}_1 + \sum_{i=2}^n \alpha_i (\frac{\lambda_i}{\lambda_1})^k \boldsymbol{\xi}_i]}{\max\{\lambda_1^k [\alpha_1 \boldsymbol{\xi}_1 + \sum_{i=2}^n \alpha_i (\frac{\lambda_i}{\lambda_1})^k \boldsymbol{\xi}_i]\}} \\ &= \frac{\alpha_1 \boldsymbol{\xi}_1 + \sum_{i=2}^n \alpha_i (\frac{\lambda_i}{\lambda_1})^k \boldsymbol{\xi}_i}{\max\{\alpha_1 \boldsymbol{\xi}_1 + \sum_{i=2}^n \alpha_i (\frac{\lambda_i}{\lambda_1})^k \boldsymbol{\xi}_i\}} \end{aligned}$$

可见

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \frac{\alpha_1 \boldsymbol{\xi}_1}{\max\{\alpha_1 \boldsymbol{\xi}_1\}} = \frac{\boldsymbol{\xi}_1}{\max\{\boldsymbol{\xi}_1\}} = \boldsymbol{\xi}_1^0.$$

又

$$\begin{aligned} \mathbf{v}^{(k)} &= \mathbf{A}\mathbf{x}^{(k-1)} = \frac{\mathbf{A}^k \mathbf{x}^{(0)}}{m_{k-1} \cdots m_1} = \frac{\mathbf{A}^k \mathbf{x}^{(0)}}{\max\{\mathbf{A}^{(k-1)} \mathbf{x}^{(0)}\}} \\ &= \frac{\lambda_1^k [\alpha_1 \boldsymbol{\xi}_1 + \sum_{i=2}^n \alpha_i (\frac{\lambda_i}{\lambda_1})^k \boldsymbol{\xi}_i]}{\lambda_1^{k-1} \max\{[\alpha_1 \boldsymbol{\xi}_1 + \sum_{i=2}^n \alpha_i (\frac{\lambda_i}{\lambda_1})^{k-1} \boldsymbol{\xi}_i]\}} \end{aligned}$$

注意到 m_k 是 $\mathbf{v}^{(k)}$ 模的最大的分量, 既有

$$m_k = \max\{\mathbf{v}^{(k)}\} = \lambda_1 \frac{\max\{\alpha_1 \boldsymbol{\xi}_1 + \sum_{i=2}^n \alpha_i (\frac{\lambda_i}{\lambda_1})^k \boldsymbol{\xi}_i\}}{\max\{\alpha_1 \boldsymbol{\xi}_1 + \sum_{i=2}^n \alpha_i (\frac{\lambda_i}{\lambda_1})^{k-1} \boldsymbol{\xi}_i\}}$$

从而 $\lim_{k \rightarrow \infty} m_k = \lambda_1$ 成立. 证毕. □

幂法的收敛速度与比值 $|\lambda_2/\lambda_1|$ 的大小有关, $|\lambda_2/\lambda_1|$ 越小, 收敛速度越快, 当此比值接近于1时, 收敛速度是非常缓慢的。因此可以对原矩阵作原点位移, 令

$$\mathbf{B} = \mathbf{A} - \alpha \mathbf{I}$$

式中 α 为参数。选择此参数可使矩阵 B 的上述比值更小, 以加快幂法的收敛速度。设矩阵 A 的特征值为 $\lambda_1, \lambda_2, \dots, \lambda_n$, 对应的特征向量为 $\xi_1, \xi_2, \dots, \xi_n$, 则矩阵 B 的特征值为 $\lambda_1 - \alpha, \lambda_2 - \alpha, \dots, \lambda_n - \alpha$, B 的特征向量与 A 的特征向量相同。假设原点位移后, B 的特征值 $\lambda_1 - \alpha$ 仍为模最大的特征值, 选择 α 的目的是使

$$\max_{2 \leq i \leq n} \frac{\lambda_i - \alpha}{\lambda_1 - \alpha} < \frac{\lambda_2}{\lambda_1}$$

适当地选择 α 可使幂法的收敛速度得到加速。此时 $m_k \rightarrow \lambda_1 - \alpha, m_k + \alpha \rightarrow \lambda_1$, 而 $\mathbf{x}^{(K)}$ 仍然收敛于 A 的特征向量 ξ_1^0 。这种加速方法叫做原点位移法。

在实际计算中, 由于矩阵的特征值分布情况事先一般是不知道的, 参数 α 的选取存在困难, 因为 α 的选取要保证 $\lambda_1 - \alpha$ 仍然是矩阵 $B (= A - \alpha I)$ 模最大的特征值, 故原点位移法是很难实现的。但是在反幂法中, 原点位移参数 α 是很容易选取的, 因此, 带原点位移的反幂法已经成为改进特征值和特征向量精度的标准算法。

5.3.3 反幂法

设 A 可逆, 则对 A 的逆阵 A^{-1} 施以幂法称为反幂法。由于 $A\xi_i = \lambda_i\xi_i$ 时, 成立 $A^{-1}\xi_i = \lambda_i^{-1}\xi_i$ 。因此, 若 $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_{n-1}| > |\lambda_n|$, 则 λ_n^{-1} 是 A^{-1} 按模最大的特征值, 此时按反幂法, 必有

$$m_k \rightarrow \lambda_n^{-1}, \mathbf{x}^{(k)} \rightarrow \xi_n^0$$

其收敛率为 $|\lambda_n/\lambda_{n-1}|$ 。任取初始向量 $\mathbf{x}^{(0)}$, 构造向量序列

$$\mathbf{x}^{(k+1)} = A^{-1}\mathbf{x}^{(k)}, k = 0, 1, 2, \dots$$

按幂法计算即可, 但用上述式子计算, 首先要求 A^{-1} , 这比较麻烦而且效率不高, 实际计算中通常用解方程组的办法, 即用

$$A\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}, k = 0, 1, 2, \dots$$

求 $\mathbf{x}^{(k+1)}$ 。为防止计算机溢出, 实际计算时所用公式为

$$\begin{aligned} \mathbf{v}^{(k)} &= \mathbf{x}^{(k)} / \max(\mathbf{x}^{(k)}), \\ A\mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)}, \end{aligned} \quad k = 0, 1, 2, \dots$$

式中: $\max(\mathbf{x}^{(k)})$ 为 $\mathbf{x}^{(k)}$ 模最大的分量。

若 A 的特征值是 λ , 则 $\lambda - \alpha$ 是 $A - \alpha I$ 的特征值。因此反幂法可以用于已知矩阵的近似特征值为 α 时, 求矩阵的特征向量并且提高特征值精度。

此时, 可以用原点位移法来加速迭代过程, 于是上式相应为

$$A\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}, k = 0, 1, 2, \dots$$

求 $\mathbf{x}^{(k+1)}$ 。为防止计算机溢出, 实际计算时所用公式为

$$\begin{aligned} \mathbf{v}^{(k)} &= \mathbf{x}^{(k)} / \max(\mathbf{x}^{(k)}), \\ (A - \alpha I)\mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)}, \end{aligned} \quad k = 0, 1, 2, \dots$$

算法（反幂法）

- (1) 选取初值 $\mathbf{x}^{(0)}$, 近似值 α , 误差限 ϵ , 最大迭代次数 N 。记 m_0 为 $\mathbf{x}^{(0)}$ 中按模最大的分量, $\mathbf{v}^{(0)} = \mathbf{x}^{(0)} / m_0$ 。置 $k = 0$
- (2) 解方程组 $(\mathbf{A} - \alpha \mathbf{I})\mathbf{x}^{(k+1)} = \mathbf{v}^{(k+1)}$ 。
- (3) 记 m_{k+1} 为 $\mathbf{x}^{(k+1)}$ 中按模最大的分量, $\mathbf{v}^{(k+1)} = \mathbf{x}^{(k+1)} / m_{k+1}$ 。
- (4) 若 $|m_{k+1}^{-1} - m_k^{-1}| < \epsilon$, 则置 $\lambda = m_{k+1}^{-1} + \alpha$, 输出 λ 和 $\mathbf{x}^{(k+1)}$, 停止计算; 否则, 转 (5)
- (5) 若 $k < N$, 置 $k = k + 1$, 转 (2), 否则输出计算失败信息, 停止计算。

5.3.4 特征值计算的应用: Pagerank 网页排名

接下来我们将介绍一个用于网页排名的算法——Pagerank。

(1) **问题背景** 互联网 (Internet) 的使用已经深入到人们的日常生活中, 其巨大的信息量和强大的功能给生产、生活带来了很大的便利。随着网络的信息量越来越庞大, 如何有效地搜索出用户真正需要的信息变得十分重要。自 1998 年搜索引擎网站 Google 创立以来, 网络搜索引擎成为解决上述问题的重要手段。

1998 年, 美国斯坦福大学的博士生 Larry Page 和 Sergey Brin 创立了 Google 公司, 他们的核心技术就是通过 Pagerank 技术对海量的网页进行重要性分析。该技术利用网页相互链接的关系对网页进行组织, 确定出每个网页的重要级别 (Pagerank)。当用户进行搜索时, Google 找出符合搜索要求的网页, 并按它们的 Pagerank 大小依次列出。这样, 用户一般显示结果的第一页或者前几页就能找到真正有用的结果。

形象地解释, Pagerank 技术的基本原理是: 如果网页 A 链接到网页 B, 则认为“网页 A 投了网页 B”一票, 而且如果网页 A 是级别高的网页, 则网页 B 的级别也相应地高。

(2) **数学问题建模** 假设 n 是 Internet 中所有可访问网页的数目, 此数值非常大, 再 2010 年已接近 100 亿。定义 $n \times n$ 的网页连接矩阵 $\mathbf{G} = (g_{ij}) \in \mathbb{R}^{n \times n}$, 若从网页 j 有一个链接到网页 i , 则 $g_{ij} = 1$, 否则 $g_{ij} = 0$ 。矩阵 \mathbf{G} 有如下特点:

- (1) \mathbf{G} 矩阵是大规模稀疏矩阵;
- (2) 第 j 列非零元素的位置表示了从网页 j 链接出去的所有网页;
- (3) 第 i 列非零元素的位置表示了所有链接到网页 i 的网页;
- (4) \mathbf{G} 中非零元的数目为整个 Internet 中存在的超链接的数量;
- (5) 记 \mathbf{G} 矩阵行元素之和 $r_i = \sum_j g_{ij}$, 它表示第 i 个网页的“入度”;
- (6) 记 \mathbf{G} 矩阵列元素之和 $c_j = \sum_i g_{ij}$, 它表示第 j 个网页的“出度”。

要计算 PageRank, 可假设一个随机上网“冲浪”的过程, 即每次看完当前网页后, 有两种选择:

- (1) 在当前网页中随机选一个超链接进入下一个网页;
- (2) 随机地新开一个网页;

设 p 为选择当前网页上链接的概率 (比如 $p = 0.85$), 则 $1 - p$ 为不选当前网页的链接而随机打开一个网页的概率。若当前网页是网页 j , 则如何计算下一步浏览到达网页 i 的概率 (网页 j 到 i 的转移概率)? 它有两种可能性:

(1) 若网页 i 在网页 j 的链接上, 其概率为 $p \cdot 1/c_j + (1 - p) \cdot 1/n$;

(2) 若网页 i 不在网页 j 的链接上, 其概率为 $(1 - p) \cdot 1/n$ 。

由于网页 i 是否在网页 j 的链接上由 g_{ij} 决定, 网页 j 到 i 的转移概率为

$$a_{ij} = g_{ij} \left[p \cdot \frac{1}{c_j} + (1 - p) \cdot \frac{1}{n} \right] + (1 - g_{ij}) \left[(1 - p) \cdot \frac{1}{n} \right] = \frac{pg_{ij}}{c_j} + \frac{1 - p}{n}$$

应注意的是, 若 $c_j = 0$ 意味着 $g_{ij} = 0$, 上式改为 $a_{ij} = 1/n$ 。任意两个网页之间的转移概率形成了一个转移矩阵 $A = (a_{ij})_{n \times n}$ 。设矩阵 D 为各个网页出度的倒数 (若没有出度, 设为 1) 构成的 n 阶对角阵, e 为全是 1 的 n 维向量, 则

$$A = pGD + \frac{1 - p}{n}ee^T.$$

这在数学上称为马尔可夫过程。若这样的随机“冲浪”一直进行下去, 某个网页被访问的极限概率就是它的 PageRank。

设 $x_i^{(k)}, i = 1, 2, \dots, n$ 表示某时刻 k 浏览网页 i 的概率 ($\sum x_i^{(k)} = 1$), 向量 $x^{(k)}$ 表示当前时刻浏览个网页的概率分布。那么下一时刻浏览到网页 i 的概率为 $\sum_{j=1}^n a_{ij}x_j^{(k)}$, 此时浏览个网页的概率分布为 $x^{(k+1)} = Ax^{(k)}$ 。

当这个过程无限进行下去, 达到极限情况, 即网页访问概率 $x^{(k)}$ 收敛到一个极限值, 这个极限向量 x 为个网页的 PageRank, 他满足 $Ax = x$, 且 $\sum_{i=1}^n x_i = 1$ 。

总结一下, 我们要求解的问题是在给定 $n \times n$ 的网页连接矩阵 G , 以及选择当前网页链接的概率 p 时, 计算特征值 1 对应的特征向量 x

$$\begin{cases} Ax = x \\ \sum_{i=1}^n x_i = 1 \end{cases}$$

易知 $\|A\|_1 = 1$, 所以 $\rho(A) \leq 1$ 。又考虑矩阵 $L = I - A$, 容易验证它各列元素和均为 0, 则 L 为奇异矩阵, 所以 $\det(I - A) = 0$, 1 是 A 的特征值且为主特征值。更进一步, 用圆盘定理考察矩阵 A^T 的特征值分布, 图 (a) 显示了第 j 个圆盘 $D_j (j = 1, 2, \dots, n)$, 显然其圆心 $a_{jj} > 0$, 半径 r_j 满足 $a_{jj} + r_j = 1$, 因此除了 1 这一点, 圆盘上任何一点到圆心的距离 (即复数的模) 都小于 1。这就说明, 1 是矩阵 A^T 和 A 的唯一主特征值。对于实际的大规模稀疏矩阵 A , 幂法是求其主特征向量的可靠的、唯一的选择。

网页的 PageRank 完全由所有网页的超链接结构所决定, 隔一段时间重新算一次 PageRank 以反映互联网的发展变化, 此时将上一次计算的结果作为幂法的迭代初值可提高收敛速度。由于迭代向量以及矩阵 A 的物理意义。在使用幂法时并不需要对向量进行规格化, 而且不需要形成矩阵 A 。通过遍历整个网页的数据库, 根据网页间超链接关系即可得到 $Ax^{(k)}$ 的结果。

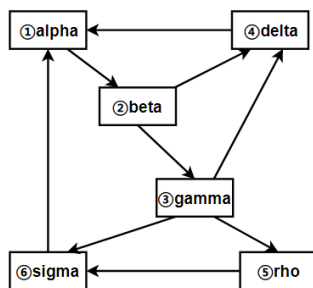


图 5.3: 网页链接关系

用一个只有 6 个网页的微型网络作为例子,其网页链接关系如图5.3所示。通过下述 MATLAB 命令可生成矩阵 G

```
 $i = [2\ 3\ 4\ 4\ 5\ 6\ 1\ 6\ 1];$ 
```

```
 $j = [1\ 2\ 2\ 3\ 3\ 3\ 4\ 5\ 6];$ 
```

```
 $n = 6;$ 
```

```
 $G = \text{aparse}(i, j, 1, n, n);$ 
```

再使用下述命令得到矩阵 A

```
 $c = \text{full}(\text{sum}(G));$ 
```

```
 $D = \text{spdiags}(\text{ones}(1, c), 0, n, n);$ 
```

```
 $e = \text{ones}(n, 1);$ 
```

```
 $p = .85; \text{delta} = (1 - p)/n;$ 
```

```
 $A = p * G * D + \text{delta} * e * e';$ 
```

得到的矩阵 A 为

$$\begin{pmatrix} 0.025 & 0.025 & 0.025 & 0.875 & 0.025 & 0.875 \\ 0.875 & 0.025 & 0.025 & 0.025 & 0.025 & 0.025 \\ 0.025 & 0.45 & 0.025 & 0.025 & 0.025 & 0.025 \\ 0.025 & 0.45 & 0.3083 & 0.025 & 0.025 & 0.025 \\ 0.025 & 0.025 & 0.3083 & 0.025 & 0.025 & 0.025 \\ 0.025 & 0.025 & 0.3083 & 0.025 & 0.875 & 0.025 \end{pmatrix}$$

使用幂方法可求出其主特征向量, 即 PageRank 为

$$\mathbf{x} = [0.2675\ 0.2524\ 0.1323\ 0.1697\ 0.0625\ 0.1156]^T$$

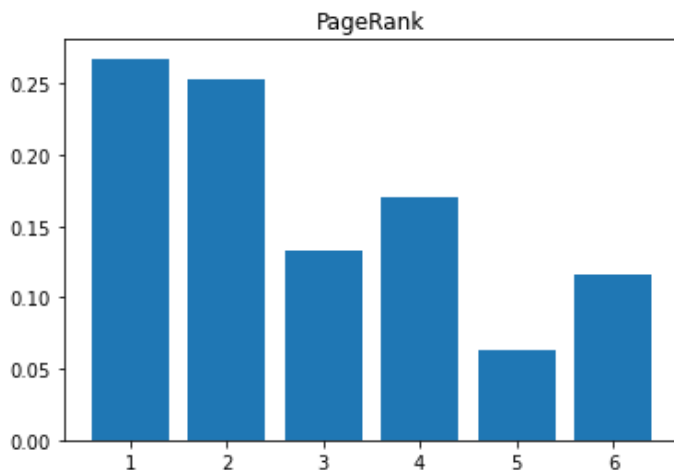


图 5.4: 网页的级别高低

使用 MATLAB 的 bar 命令, 将 x 的各分量显示如图 5.4 所示, 从中看出各个网页的级别高低, 虽然链接数目一样, 但是网页 alpha 1 的链接比 delta 4 和 sigma 6 都高, 而由于高级别的 alpha 1 链接在 beta 2 上面, 受到 alpha 1 的影响, beta 2 的级别第二高。

5.4 阅读材料

本章介绍了数值线性代数三大核心主题内容, 包括线性方程的求解、最小二乘问题和特征值的求解。数据科学中的很多问题最终都归结为线性方程的求解, 因此这一章主要介绍线性方程组的类型和解的结构, 引入线性方程组和最小二乘问题的求解方法, 并讨论解的敏感性, 这些内容将与后续优化问题求解、数据科学中的线性回归问题相联系。此外, 还介绍了大规模矩阵求解特征值的一些计算方法, 包括幂迭代法, 这已被广泛应用于数据科学中的搜索技术 pagerank 的矩阵特征值计算。此外, 用于对高维数据进行非线性降维和聚类的更现代的谱方法, 如 Isomap (Tenenbaum 等, 2000), Laplacian 特征映射 (Belkin 和 Niyogi, 2003), Hessian 特征映射 (Donoho 和 Grimes, 2003), 谱聚类 (Shi 和 Malik, 2000) 等, 每一个都需要计算正定核的特征向量和特征值, 这些核心计算通常由低秩矩阵近似技术 (Belabbas 和 Wolfe, 2009) 支持, 正如我们在 SVD 中遇到的那样。另外, 关于稠密数值线性代数可参考 (Golub 和 Van Loan, 1989), (Trefethen 和 Bau, 1997) 等。(Gill, Murray, 1981) 和 (Wright, 1997), (Wright 以及 Nocedal, 1999) 等书籍注重于数值优化问题的数值线性代数介绍。关于数值线性代数软件包, 可以参考 LAPACK, 其包括常规的稠密的线性代数算法的高质量实现。LAPACK 在基本线性代数子程序 (BLAS) 的基础上建成, 后者是基本的向量和矩阵运算的程序库, 可以很容易地根据具体的计算机结构的优点进行定制, 也可得到求解稀疏的线性方程组的一些源代码, 包括 SPOOLES, SuperLU, UMFPACK

以及 WSMP 等等, 这里提到的只是其中少数几个。

习题

习题 5.1. 设 $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}$, $b = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$ 用正则化方法求对应的 LS 问题的解。

习题 5.2. 设 $A = \begin{bmatrix} 1 & 3 & 1 & 1 \\ 2 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$, $b = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$ 求对应的 LS 问题的全部解。

习题 5.3. 设 $A \in \mathbb{R}^{m \times n}$ 且存在 $X \in \mathbb{R}^{n \times m}$ 使得对每一个 $b \in \mathbb{R}^m$, $x = Xb$ 均极小化 $\|Ax - b\|_2$. 证明 $AXA = A$ 和 $(AX)^T = AX$.

习题 5.4. 利用等式

$$\|A(x + \alpha w) - b\|_2^2 = \|Ax - b\|_2^2 + 2\alpha w^T A^T (Ax - b) + \alpha^2 \|Aw\|_2^2$$

证明: 如果 $x \in X_{LS}$, 那么 $A^T Ax = A^T b$

习题 5.5. 给定点集 $p_1, \dots, p_m \in \mathbb{R}^n$ 构成的 $m \times n$ 矩阵 $P = [p_1, \dots, p_m]$. 考虑问题

$$\min_X F(X) = \sum_{i=1}^m \|x_i - p_i\|_2^2 + \frac{\lambda}{2} \sum_{1 \leq i, j \leq m} \|x_i - x_j\|_2^2$$

其中 $\lambda \geq 0$ 为参数, 变量是一个 $m \times n$ 矩阵 $X = [x_1, \dots, x_m]$, 其中 $x_i \in \mathbb{R}^n$ 是 X 的第 i 列, $i = 1, \dots, m$. 上述问题尝试聚类点集 p_i , 第一项鼓励聚类中心 x_i 靠近对应的点 p_i , 第二项鼓励 x_i 们之间彼此靠近, 当 λ 增大的时候, 对应更高的组群影响。

1. 请说明这个问题属于最小二乘类问题。不需要明确阐述这个问题的形式。
2. 证明 $\frac{1}{2} \sum_{1 \leq i, j \leq m} \|x_i - x_j\|_2^2 = \text{trace} X H X^T$, 其中 $H = mI_m - \mathbf{1}\mathbf{1}^T$ 是一个 $m \times m$ 矩阵, I_m 是 $m \times m$ 单位矩阵, $\mathbf{1}$ 是 \mathbb{R}^n 中的单位向量。
3. 证明 H 是半正定的。
4. 证明函数 F 在矩阵 X 处的梯度是一个 $n \times m$ 矩阵, 为:

$$\nabla F(X) = 2(X - P + \lambda XH)$$

提示: 对于第二项, 找到函数的一阶展式, $\Delta \rightarrow \text{trace}((X + \Delta)H(X + \Delta)^T)$, 其中 $\Delta \in \mathbb{R}^{n, m}$ 。

5. 依据最小二乘问题的最优条件为目标函数的梯度为零。证明最优点集的形式为:

$$x_i = \frac{1}{m\lambda + 1} p_i + \frac{m\lambda}{m\lambda + 1} \hat{p}, i = 1, \dots, m,$$

其中 $\hat{p} = (1/m)(p_1 + \dots + p_m)$ 是给定点集的中心。

6. 阐述你的结果, 你认为这是聚类点集的一个好的模型么?

习题 5.6. 判断 $[1, 3, 4]$ 的转置是否在 A 的零空间中

$$A = \begin{bmatrix} 3 & 5 & -3 \\ 6 & -2 & 0 \\ -8 & 4 & 1 \end{bmatrix}$$

习题 5.7. 求矩阵

$$\begin{bmatrix} 5 & 21 & 19 \\ 13 & 23 & 2 \\ 8 & 14 & 1 \end{bmatrix}$$

的行空间和列空间

习题 5.8. 简答：阐述非负矩阵分解和主成分分析的相同点和不同点

参考文献

- [1] E.Anderson, Z.Bai, C.Bischof, S.Blackford, J.Demrnel, J.Dongarra, J.DuCroz, A.Greenbaum, S.Hammarling, A.McKenney, and D.Sorensen. LAPACK Users' Guide. Society for Industrial and Applied Mathematics, third edition, 1999. Available from www.netlib.org/lapack.
- [2] C.Ashcraft, D.Pierce, D.K.Wah, and J.Wu.The Reference Manual for SPOOLES Version 2.2: An Object Oriented Software Library for Solving Sparse Linear Systems of Equations, 1999. Available from www.netlib.org/linalg/spooles/spooles.2.2.html.
- [3] T.A.Davis.UMFPACK User Guide, 2003.Available from wv.cise.ufl.edu/research/sparse/umfpack.
- [4] J.W.Dexnmel.Applied Numerical Linear Algebra.Society for Industrial and Applied Mathematics, 1997.
- [5] I.S.Duff, A.M.Erismaim, and J.K.Reid.Direct Methods for Sparse Matrices.Clarendon Press, 1986.
- [6] J.W.Demxnel, J.R.Gilbert, and X.S.Li.SuperLU Users' Guide, 2003.Available from crd.lbl.gov/xiaoye/SuperLU.
- [7] I.S.Duff.The solution of augmented systems.In D.F.Griffiths and G.A.Watson, editors, Numerical Analysis 1993.Proceedings of the 15th Dundee Conference, pages 40- 55.Longman Scientific & Technical, 1993.
- [8] G.Golub and C.F.Van Loan.Matrix Computations.Johns Hopkins University Press, second edition, 1989.
- [9] A.George and J.W.-H.Liu.Computer sohstion of large sparse positive definite systems.Prentice-Hall, 1981.

- [10] Belkin, Mikhail, and Niyogi, Partha. 2003. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation*, 15(6), 1373–1396.
- [11] Tenenbaum, Joshua B, De Silva, Vin, and Langford, John C. 2000. A global geometric framework for nonlinear dimensionality reduction. *science*, 290(5500), 2319–2323.
- [12] Donoho, David L, and Grimes, Carrie. 2003. Hessian eigenmaps: Locally linear embedding techniques for high-dimensional data. *Proceedings of the National Academy of Sciences*, 100(10), 5591–5596.
- [13] Belabbas, Mohamed-Ali, and Wolfe, Patrick J. 2009. Spectral methods in machine learning and new strategies for very large datasets. *Proceedings of the National Academy of Sciences*, pnas – 0810600105.
- [14] Shi, Jianbo, and Malik, Jitendra. 2000. Normalized cuts and image segmentation. *IEEE Transactions on pattern analysis and machine intelligence*, 22(8), 888–905.
- [15] A.Gupta.WSMP: Watson Sparse Matrix Package.Part I —Direct Solution of Symmetric Sparse Systems.Part II —Direct Solution of General Sparse Systems, 2000.Available from www.cs.umn.edu/~agupta/wsmp.
- [16] N.J.Higham.Accuracy and Stability of Numerical Algorithms.Society for Industrial and Applied Mathematics, 1996,
- [17] P.E.Gill, W.Murray, and M.H.Wright.Practical Optimization, Academic Press, 1981.
- [18] J.Nocedal and S.J.Wright.Numerical OptimizationL Springer, 1999.
- [19] L.N.Trefethen and D.Bau, III.Numerical Linear Algebra.Society for Industrial and Applied Mathematics, 1997.
- [20] S.J.Wright.Primal-Dual Interior-Point Methods.Society for Industrial and Applied Mathematics, 1997.