

## 第4章 超越梯度下降

第3章介绍了梯度下降法。梯度下降法基于函数局部一阶特性。一阶近似是粗糙的，这种粗糙带来了一些问题。本章将介绍函数在局部的二阶特性。基于二阶特性分析函数在局部的性质。

本章首先回顾一些矩阵的相关知识，之后介绍如何在局部对函数进行二阶近似。有了函数的二阶近似就可以确定驻点的类型：极小点、极大点或者鞍点。之后本章介绍对原始梯度下降法的一些改进，这些改进有助于提高收敛速度，防止震荡或发散，规避局部极小。

最后，本章介绍两个基于函数二阶特性的优化算法：牛顿法和共轭方向法。然后介绍用牛顿法训练逻辑回归模型。二阶算法虽然不常用在神经网络和深度学习的训练中。阅读完本章，读者应该对函数的局部形态有更深刻的理解。

### 4.1 矩阵

首先回顾一下矩阵。这不是一个关于矩阵的全面介绍，例如行列式这个概念就没有出现。本节只介绍一下后文讨论中用得上的相关知识。

#### 4.1.1 矩阵基础

矩阵是实数构成的2维阵列。以一个 $3 \times 3$ 矩阵为例：

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} = (\mathbf{a}_{*1} \quad \mathbf{a}_{*2} \quad \mathbf{a}_{*3}) = \begin{pmatrix} \mathbf{a}_{1*}^T \\ \mathbf{a}_{2*}^T \\ \mathbf{a}_{3*}^T \end{pmatrix} \quad (4.1)$$

式(4.1)囊括了本书用到的对矩阵的各种表示。本书用大写粗斜体字母表示矩阵，例如 $\mathbf{A}$ 。 $a_{ij}$ 是实数，是矩阵 $\mathbf{A}$ 的第 $i$ 行、第 $j$ 列元素。 $\mathbf{a}_{*j}$ 是矩阵的第 $j$ 列，它是一个列向量：

$$\mathbf{a}_{*j} = \begin{pmatrix} a_{1j} \\ a_{2j} \\ a_{3j} \end{pmatrix} \quad (4.2)$$

$\mathbf{a}_{i*}$ 是矩阵的第 $i$ 行，它是一个行向量：

$$\mathbf{a}_{i*} = \begin{pmatrix} a_{i1} \\ a_{i2} \\ a_{i3} \end{pmatrix} \quad (4.3)$$

式(4.1)中对 $\mathbf{a}_{i*}$ 进行了转置，以表示一行。矩阵的行数和列数不一定相等，可以是 $m \times n$ ， $m \neq n$ 。表示成：

$$\mathbf{A}_{m \times n} = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix} \quad (4.4)$$

一般可省略下标 $m \times n$ 。两个相同形状的矩阵可以相加：

$$\mathbf{A} + \mathbf{B} = \begin{pmatrix} a_{11} + b_{11} & \cdots & a_{1n} + b_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} + b_{m1} & \cdots & a_{mn} + b_{mn} \end{pmatrix} \quad (4.5)$$

矩阵相加就是把相应元素相加。可以用实数（标量）乘一个矩阵：

$$k\mathbf{A} = \begin{pmatrix} ka_{11} & \cdots & ka_{1n} \\ \vdots & \ddots & \vdots \\ ka_{m1} & \cdots & ka_{mn} \end{pmatrix} \quad (4.6)$$

$-\mathbf{A}$ 就是 $(-1)\mathbf{A}$ 。显然有 $\mathbf{A} - \mathbf{A} = \mathbf{A} + (-\mathbf{A}) = \mathbf{O}$ 。 $\mathbf{O}$ 是所有元素都为0的矩阵——零矩阵。矩阵 $\mathbf{A}$ 的转置定义为：

$$\mathbf{A}^T = \begin{pmatrix} a_{11} & \cdots & a_{m1} \\ \vdots & \ddots & \vdots \\ a_{1n} & \cdots & a_{mn} \end{pmatrix} = (\mathbf{a}_{1*} \quad \mathbf{a}_{2*} \quad \mathbf{a}_{3*}) = \begin{pmatrix} \mathbf{a}_{*1}^T \\ \mathbf{a}_{*2}^T \\ \mathbf{a}_{*3}^T \end{pmatrix} \quad (4.7)$$

$\mathbf{A}^T$ 把 $\mathbf{A}$ 的行当做列，列当做行。如果 $\mathbf{A}$ 是 $m \times n$ 的，那么 $\mathbf{A}^T$ 就是 $n \times m$ 的。如果矩阵 $\mathbf{A}$ 是 $m \times n$ 的，它可以与一个 $n$ 维向量 $\mathbf{x}$ 相乘：

$$\mathbf{A}\mathbf{x} = \sum_{j=1}^n x_j \mathbf{a}_{*j} \quad (4.8)$$

矩阵 $\mathbf{A}$ 乘向量 $\mathbf{x}$ ，使用 $\mathbf{x}$ 的元素对矩阵的列进行线性组合。所以 $\mathbf{A}$ 的列数和 $\mathbf{x}$ 的维数必须相同。得到的结果是一个 $m$ 维向量。容易看出 $\mathbf{Ax}$ 的第 $i$ 个元素是 $\sum_{j=1}^n x_j a_{ij} = \mathbf{a}_{i*}^T \mathbf{x}$ ，即 $\mathbf{A}$ 的第 $i$ 行与 $\mathbf{x}$ 的内积。

有了矩阵和向量相乘的定义，就可以定义矩阵与矩阵相乘：

$$\mathbf{AB} = (\mathbf{Ab}_{*1} \quad \mathbf{Ab}_{*2} \quad \cdots \quad \mathbf{Ab}_{*k}) \quad (4.9)$$

$\mathbf{A}$ 与 $\mathbf{B}$ 的乘积是矩阵 $\mathbf{AB}$ 。 $\mathbf{AB}$ 的第 $j$ 列是 $\mathbf{A}$ 与 $\mathbf{B}$ 的第 $j$ 列 $\mathbf{b}_{*j}$ 的乘积。如果 $\mathbf{A}$ 是 $m \times n$ 的，那么 $\mathbf{b}_{*j}$ 必须是 $n$ 维向量，即 $\mathbf{B}$ 必须为 $n$ 行。 $\mathbf{B}$ 的列数任意，例如 $k$ 。所以要能够与 $m \times n$ 的 $\mathbf{A}$ 相乘， $\mathbf{B}$ 的形状必须是 $n \times k$ ， $k$ 任意。结果 $\mathbf{AB}$ 的形状是 $m \times k$ 。 $\mathbf{AB}$ 的第 $i$ 行、第 $j$ 列元素是：

$$\mathbf{a}_{i*}^T \mathbf{b}_{*j} = \sum_{s=1}^n a_{is} b_{sj} \quad (4.10)$$

仅从形状上看 $\mathbf{B}$ 与 $\mathbf{A}$ 不一定能够相乘，因为 $k$ 不一定等于 $m$ 。就算 $k = m$ ， $\mathbf{BA}$ 也不一定等于 $\mathbf{AB}$ 。即矩阵乘法不满足交换律。一个反例就可以证明这一点。这里不再赘述。

矩阵的乘法满足结合率：

$$(\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC}) \quad (4.11)$$

矩阵乘法对加法满足结合律：

$$\mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{AB} + \mathbf{AC}, \quad (\mathbf{A} + \mathbf{B})\mathbf{C} = \mathbf{AC} + \mathbf{BC} \quad (4.12)$$

矩阵乘法对数乘满足：

$$\mathbf{A}(k\mathbf{B}) = (k\mathbf{A})\mathbf{B} = k(\mathbf{AB}) \quad (4.13)$$

矩阵的数乘满足分配率：

$$k(\mathbf{A} + \mathbf{B}) = k\mathbf{A} + k\mathbf{B}, \quad (k + h)\mathbf{A} = k\mathbf{A} + h\mathbf{A} \quad (4.14)$$

矩阵乘积的转置是：

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T \quad (4.15)$$

上述几个规则的证明很简单，只需要检查一下矩阵元素的表达式。向量 $\mathbf{x}$ 的转置 $\mathbf{x}^T$ 可以乘一个矩阵 $\mathbf{A}$ ：

$$\mathbf{x}^T \mathbf{A} = (\mathbf{A}^T \mathbf{x})^T \quad (4.16)$$

把列向量和行向量分别看作 $n \times 1$ 和 $1 \times n$ 的矩阵，则矩阵与向量的乘法也满足上述几个规则。

行数和列数相同的矩阵是方阵。方阵 $\mathbf{A}$ 的对角线元素之和称为它的迹（trace）：

$$\text{tr}(\mathbf{A}) = \sum_{i=1}^n a_{ii} \quad (4.17)$$

方阵 $\mathbf{A}$ 和 $\mathbf{B}$ 的乘积 $\mathbf{AB}$ 的迹等于 $\mathbf{BA}$ 的迹，因为：

$$\text{tr}(\mathbf{AB}) = \sum_{i=1}^n \mathbf{a}_{i*} \mathbf{b}_{*i} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} b_{ji} = \sum_{j=1}^n \sum_{i=1}^n b_{ji} a_{ij} = \sum_{j=1}^n \mathbf{b}_{j*} \mathbf{a}_{*j} = \text{tr}(\mathbf{BA}) \quad (4.18)$$

如果一个 $n \times n$ 的方阵的对角线元素为 1，其余元素都是 0，那么它是单位阵：

$$\mathbf{I}_{n \times n} = \begin{pmatrix} 1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1 \end{pmatrix} \quad (4.19)$$

容易验证对于任何矩阵 $\mathbf{A}_{m \times n}$ ， $\mathbf{I}_{m \times m} \mathbf{A}_{m \times n} = \mathbf{A}_{m \times n} \mathbf{I}_{n \times n} = \mathbf{A}_{m \times n}$ 。在上下文很清晰时一般省略 $\mathbf{I}$ 的下标。

### 4.1.2 矩阵的逆

令 $\mathbf{A}$ 是 $n \times n$ 方阵，如果存在 $n \times n$ 方阵 $\mathbf{A}^{-1}$ 满足：

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I} \quad (4.20)$$

则称 $\mathbf{A}$ 是可逆的。 $\mathbf{A}^{-1}$ 是 $\mathbf{A}$ 的逆矩阵。 $\mathbf{A}$ 的逆矩阵是唯一的。因为假如任何一个矩阵 $\mathbf{B}$ 是 $\mathbf{A}$ 的逆矩阵，根据定义有：

$$\mathbf{B} = \mathbf{I}\mathbf{B} = \mathbf{A}^{-1}\mathbf{A}\mathbf{B} = \mathbf{A}^{-1} \quad (4.21)$$

如果 $\mathbf{A}$ 可逆则 $\mathbf{A}$ 的列线性独立。因为假如 $\mathbf{a}_{*j=1\dots n}$ 线性相关，则存在一组不全为 0 的系数 $w_1, w_2, \dots, w_n$ ，使得 $\sum_{i=1}^n w_i \mathbf{a}_{*i} = \mathbf{0}$ 。即存在向量 $\mathbf{w} = (w_1 \ \dots \ w_n)^T \neq \mathbf{0}$ 使：

$$\mathbf{A}\mathbf{w} = \mathbf{0} \quad (4.22)$$

因为 $\mathbf{A}$ 可逆，存在 $\mathbf{A}^{-1}$ ：

$$\mathbf{w} = \mathbf{A}^{-1}\mathbf{A}\mathbf{w} = \mathbf{0} \quad (4.23)$$

这与 $\mathbf{w} \neq \mathbf{0}$ 矛盾。所以可逆矩阵 $\mathbf{A}$ 的列 $\mathbf{a}_{*j=1\dots n}$ 一定线性独立。如果方阵 $\mathbf{A}$ 的逆矩阵是 $\mathbf{A}^T$ ，则称 $\mathbf{A}$ 为正交矩阵：

$$\mathbf{A}\mathbf{A}^T = \mathbf{A}^T\mathbf{A} = \mathbf{I} \quad (4.24)$$

从 (4.24) 可以看出 $\mathbf{A}$ 的列 $\mathbf{a}_{*j=1\dots n}$ 是单位向量且两两正交：

$$\begin{cases} \mathbf{a}_{*i}^T \mathbf{a}_{*j} = 0, & i \neq j \\ \mathbf{a}_{*i}^T \mathbf{a}_{*j} = 1, & i = j \end{cases} \quad (4.25)$$

也就是说，正交矩阵 $\mathbf{A}$ 的列都是单位向量， $\|\mathbf{a}_{*j=1\dots n}\| = 1$ 。任意两列是正交的（夹角为 $\pi/2$ ）。因为 $\mathbf{A}$ 可逆，所以 $\mathbf{a}_{*j=1\dots n}$ 线性独立，是 $n$ 维线性空间 $\mathbb{R}^n$ 的一组基。因为 $\mathbf{a}_{*j=1\dots n}$ 两两正交，还是单位向量，所以它们被称为 $\mathbb{R}^n$ 的一组标准正交基。

如果一组向量 $\mathbf{x}_{i=1\dots n}$ 是线性独立的，可以通过施密特正交化过程构造一组正交的向量 $\mathbf{x}'_{i=1\dots n}$ 。

构造过程是，首先令  $\mathbf{x}'_1 = \mathbf{x}_1$ 。然后令：

$$\mathbf{x}'_2 = \mathbf{x}_2 - \frac{\mathbf{x}_2^T \mathbf{x}'_1}{\mathbf{x}'_1^T \mathbf{x}'_1} \mathbf{x}'_1 = \mathbf{x}_2 - \frac{\|\mathbf{x}_2\| \cdot \cos \theta_{21}}{\|\mathbf{x}'_1\|} \mathbf{x}'_1 \quad (4.26)$$

$\theta_{21}$  是  $\mathbf{x}_2$  与  $\mathbf{x}'_1$  的夹角。 $\mathbf{x}'_2$  是  $\mathbf{x}_2$  减去  $\mathbf{x}_2$  向  $\mathbf{x}'_1$  的投影。 $\mathbf{x}_2$  与  $\mathbf{x}'_1 = \mathbf{x}_1$  线性独立，它不是  $\mathbf{x}'_1$  的数乘，所以  $\mathbf{x}'_2$  不是零向量。而且容易验证  $\mathbf{x}'_2$  与  $\mathbf{x}'_1$  正交。再令：

$$\mathbf{x}'_3 = \mathbf{x}_3 - \frac{\mathbf{x}_3^T \mathbf{x}'_1}{\mathbf{x}'_1^T \mathbf{x}'_1} \mathbf{x}'_1 - \frac{\mathbf{x}_3^T \mathbf{x}'_2}{\mathbf{x}'_2^T \mathbf{x}'_2} \mathbf{x}'_2 = \mathbf{x}_3 - \frac{\|\mathbf{x}_3\| \cdot \cos \theta_{31}}{\|\mathbf{x}'_1\|} \mathbf{x}'_1 - \frac{\|\mathbf{x}_3\| \cdot \cos \theta_{32}}{\|\mathbf{x}'_2\|} \mathbf{x}'_2 \quad (4.27)$$

$\theta_{31}$  是  $\mathbf{x}_3$  与  $\mathbf{x}'_1$  的夹角， $\theta_{32}$  是  $\mathbf{x}_3$  与  $\mathbf{x}'_2$  的夹角。 $\mathbf{x}'_3$  是  $\mathbf{x}_3$  减去  $\mathbf{x}_3$  向  $\mathbf{x}'_1, \mathbf{x}'_2$  张成空间的投影。如果  $\mathbf{x}'_3 = \mathbf{0}$ ，那么  $\mathbf{x}_3$  可以被  $\mathbf{x}'_1, \mathbf{x}'_2$  线性表出，也就可以被  $\mathbf{x}_1, \mathbf{x}_2$  线性表出，这与  $\mathbf{x}_{i=1 \dots n}$  线性独立矛盾。故  $\mathbf{x}'_3 \neq \mathbf{0}$ 。容易验证  $\mathbf{x}'_3$  正交于  $\mathbf{x}'_1, \mathbf{x}'_2$ 。此过程继续下去，最终可构造一组正交向量  $\mathbf{x}'_{i=1 \dots n}$ 。这就是施密特正交化过程。将  $\mathbf{x}'_{i=1 \dots n}$  的每一个向量除以各自的模，缩放到长度为 1，就得到了一组正交的单位向量。

### 4.1.3 特征值与特征向量

特征值和特征向量的概念不局限于方阵，但本书主要关注方阵。用方阵  $\mathbf{A}$  乘向量  $\mathbf{x}$  是在  $\mathbb{R}^n$  中进行一个变换，将  $\mathbf{x}$  变换成  $\mathbf{Ax}$ 。例如矩阵：

$$\mathbf{R} = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix} \quad (4.28)$$

用  $\mathbf{R}$  乘向量  $\mathbf{x}$  等于将  $\mathbf{x}$  逆时针旋转  $\theta$  度。这可以自行验证。任何方阵  $\mathbf{A}$  也改变不了零向量  $\mathbf{0}$ ，因为  $\mathbf{A}\mathbf{0} \equiv \mathbf{0}$ 。如果对于某非零向量  $\mathbf{v} \neq \mathbf{0}$ ， $\mathbf{A}$  只能改变  $\mathbf{v}$  的长度而不能改变其方向，即存在某个标量（可以为 0） $\lambda$ ，有：

$$\mathbf{Av} = \lambda \mathbf{v} \quad (4.29)$$

则称  $\lambda$  是  $\mathbf{A}$  的特征值， $\mathbf{v}$  是  $\mathbf{A}$  的对于  $\lambda$  的特征向量。同一个特征向量不可能对应两个特征值。假如  $\lambda_1 \neq \lambda_2$  都有  $\mathbf{v}$  是其特征向量：

$$(\lambda_1 - \lambda_2)\mathbf{v} = \lambda_1\mathbf{v} - \lambda_2\mathbf{v} = \mathbf{A}\mathbf{v} - \mathbf{A}\mathbf{v} = \mathbf{0} \quad (4.30)$$

$\lambda_1 - \lambda_2 \neq 0$  且  $\mathbf{v} \neq \mathbf{0}$ ，所以式 (4.30) 是不可能的。但是同一个特征值可以对应多个特征向量。如果  $\mathbf{v}$  是  $\lambda$  对应的特征向量，容易验证  $k\mathbf{v}$  也是  $\lambda$  对应的特征向量。线性独立的两个向量也有可能是同一个特征值对应的特征向量。

假如  $\lambda$  对应的特征向量  $\mathbf{v}$  和  $\mathbf{w}$  是线性独立的，即谁也不是另一个的数乘。那么  $k\mathbf{v} + l\mathbf{w}$  也是  $\lambda$  对应的特征向量。这也很容易验证。如果特征值  $\lambda$  共有  $k$  个线性独立的特征向量，由它们线性组合而得的向量也是  $\lambda$  的特征向量。这  $k$  个线性独立的特征向量张成的  $k$  维线性空间称为  $\lambda$  对应的特征空间，其中所有向量都是  $\lambda$  的特征向量。

将式 (4.29) 变形。如果  $\lambda$  是  $\mathbf{A}$  的特征值，它必须满足对某个  $\mathbf{v} \neq \mathbf{0}$ ，有：

$$(\mathbf{A} - \lambda\mathbf{I})\mathbf{v} = \mathbf{0} \quad (4.31)$$

因为  $\mathbf{v} \neq \mathbf{0}$ ，所以  $\mathbf{A} - \lambda\mathbf{I}$  的列线性相关。求  $\mathbf{A}$  的特征值和特征向量，就是求满足方程 (4.31) 的  $\lambda$  和  $\mathbf{v}$ 。若要  $\mathbf{A} - \lambda\mathbf{I}$  的列线性相关，则  $\mathbf{A} - \lambda\mathbf{I}$  的行列式  $|\mathbf{A} - \lambda\mathbf{I}|$  等于 0。本书没有涉及行列式，因为行列式与本书主线关系不大，加进来会影响流畅性。读者可以查阅任何一种线性代数教材。 $|\mathbf{A} - \lambda\mathbf{I}| = 0$  是  $\lambda$  的  $n$  次方程。它有  $n$  个根（包括重根和复数根）。即  $\mathbf{A}$  有  $n$  个特征值（包括重复的以及复特征值）。求得了  $\lambda$ ，就可以再求它对应的特征向量。

矩阵  $\mathbf{A}$  属于不同特征值的特征向量是线性独立的。现在证明这一点。 $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$  是  $k$  个不同特征值  $\lambda_1, \lambda_2, \dots, \lambda_k$  对应的特征向量。如果它们线性相关，则其中某一个  $\mathbf{v}_s$  可以被其他  $\mathbf{v}_{i \neq s}$  线性表出：

$$\mathbf{v}_s = \sum_{i \neq s} a_i \mathbf{v}_i \quad (4.32)$$

因为  $\mathbf{v}_s$  是  $\mathbf{A}$  的特征向量，所以它不是零向量。那么  $a_{i \neq s}$  一定不全为 0。另外根据特征值和特征向量的定义：

$$\lambda_s \mathbf{v}_s = \mathbf{A} \mathbf{v}_s = \sum_{i \neq s} a_i \mathbf{A} \mathbf{v}_i = \sum_{i \neq s} a_i \lambda_i \mathbf{v}_i \quad (4.33)$$

如果  $\lambda_s = 0$ ，那么  $\lambda_{i \neq s} \neq 0$ 。再加上  $a_{i \neq s}$  不全为 0，说明  $\mathbf{v}_{i \neq s}$  线性相关。在  $\lambda_s = 0$  情况下我们将问题规模减小了 1。如果  $\lambda_s \neq 0$ ，有：

$$\mathbf{v}_s = \sum_{i \neq s} a_i \frac{\lambda_i}{\lambda_s} \mathbf{v}_i \quad (4.34)$$

于是式 (4.34) 等于式 (4.32)，所以有：

$$\sum_{i \neq s} \left( a_i \frac{\lambda_i}{\lambda_s} \mathbf{v}_i - a_i \mathbf{v}_i \right) = \sum_{i \neq s} a_i \left( \frac{\lambda_i}{\lambda_s} - 1 \right) \mathbf{v}_i = \mathbf{0} \quad (4.35)$$

因为都是不同的特征值，所以  $\lambda_{i \neq s} / \lambda_s - 1 \neq 0$ 。再加上  $a_{i \neq s}$  不全为 0，说明  $\mathbf{v}_{i \neq s}$  线性相关。在  $\lambda_s \neq 0$  情况下我们也将问题规模减小了 1。这个过程持续下去，最终将只剩下两个向量  $\mathbf{v}_i$  和  $\mathbf{v}_j$ 。他们分属不同的特征值  $\lambda_i$  和  $\lambda_j$ 。且  $\mathbf{v}_i$  和  $\mathbf{v}_j$  线性相关，其中一个是另一个的数乘。不妨假设  $\mathbf{v}_i = k \mathbf{v}_j$ ，则  $\mathbf{v}_i$  也是  $\lambda_j$  的特征向量。之前已经证明，一个向量不可能同时属于两个不同特征值。这就推翻了最早的假设，证明了  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$  线性独立。

#### 4.1.4 对称矩阵的谱分解

如果方阵  $\mathbf{A}$  满足：

$$\mathbf{A} = \mathbf{A}^T \quad (4.36)$$

如果  $\mathbf{A}$  的元素都是实数，则它是一个实矩阵。实矩阵的特征值都是实数，特征向量是实向量。为了证明这个结论，我们需要暂时离开实数域。

复数  $\lambda = a + bi$  的共轭是  $\bar{\lambda} = a - bi$ 。

$$\lambda \bar{\lambda} = (a + bi)(a - bi) = a^2 + b^2 \geq 0 \quad (4.37)$$

只有当  $a = b = 0$ ，即  $\lambda = 0$  时，才有  $\lambda \bar{\lambda} = 0$ 。否则  $\lambda \bar{\lambda} > 0$ 。对于两个复数  $\lambda = a + bi$  和  $\xi = c + di$ ，有：

$$\lambda \xi = (ac - bd) + (bc + ad)i = \bar{\lambda} \bar{\xi} \quad (4.38)$$

把复矩阵  $\mathbf{A}$  的元素全都取共轭就得到  $\mathbf{A}$  的共轭  $\bar{\mathbf{A}}$ 。如果（复数） $\lambda$  和（复向量） $\mathbf{v}$  是  $\mathbf{A}$  的特征值及对应特征向量，由式 (4.38) 容易看出： $\bar{\lambda}$  和  $\bar{\mathbf{v}}$  是  $\bar{\mathbf{A}}$  的特征值及对应特征向量。



$$\bar{A}\bar{v} = Av = \lambda v = \bar{\lambda}\bar{v} \quad (4.39)$$

因为 $A$ 是实对称矩阵，有 $A = \bar{A} = A^T = \bar{A}^T$ ，所以有：

$$\lambda \bar{v}^T v = \bar{v}^T \lambda v = \bar{v}^T A v = \bar{v}^T \bar{A}^T v = (\bar{A}\bar{v})^T v = (\bar{\lambda}\bar{v})^T v = \bar{\lambda} \bar{v}^T v \quad (4.40)$$

因为 $x \neq 0$ ，根据式(4.37)  $\bar{v}^T v > 0$ 。所以 $\lambda = \bar{\lambda}$ ，即 $\lambda$ 是实数。 $A$ 是实矩阵， $\lambda$ 是实数，所以 $v$ 一定是实向量。这就证明了实对称矩阵的特征值都是实数，特征向量是实向量。后文谈到矩阵都是实矩阵。

所以实对称矩阵 $A$ 有 $n$ 个实特征值（可重复）。有一个结论我们不加证明：如果 $\lambda$ 是 $A$ 的 $k$ 重特征值（方程 $|A - \lambda I| = 0$ 的 $k$ 重根），则 $\lambda$ 对应的特征空间是 $k$ 维，即对于 $\lambda$ 能找到 $k$ 个线性独立的特征向量。

对于对称矩阵 $A$ ，求得它的 $n$ 个特征值 $\lambda_1, \lambda_2, \dots, \lambda_n$ ，从大到小排列。对每个特征值找到一个它的单位特征向量。如果某个特征值是 $k$ 重，那么找到它的 $k$ 个线性独立的特征向量，经过施密特正交化过程，再缩放到长度为1。这样，一共找到 $n$ 个单位特征向量 $v_1, v_2, \dots, v_n$ 。可以证明这 $n$ 个特征向量是线性独立的。因为假如它们线性相关，则存在一组不全为0的参数 $a_1, a_2, \dots, a_n$ 满足：

$$\sum_{i=1}^n a_i v_i = 0 \quad (4.41)$$

如果 $n$ 个特征值共有 $s$ 种不同取值，把属于同一个取值的特征向量归到一起：

$$\sum_{i=1}^s \sum_{j \in I(i)} a_j v_j = \sum_{i=1}^s v'_i = 0 \quad (4.42)$$

$I(i)$ 是属于第 $i$ 个排重特征值（一共 $s$ 排重特征值）的特征向量下标集合。式(4.42)中每一个 $v'_i = \sum_{j \in I(i)} a_j v_j$ 是对第 $i$ 个排重特征值的特征向量的线性组合，它仍然是第 $i$ 个排重特征值的特征向量。不可能所有 $v'_i$ 都是零向量。因为存在非0的 $a_j$ 且 $\sum_{j \in I(i)} a_j v_j = 0$ ，这与 $v_{j \in I(i)}$ 线性独立产生矛盾。那些非零的 $v'_i$ 加在一起是零向量，又与属于不同特征值的特征向量线性独立矛盾。所以 $v_1, v_2, \dots, v_n$ 是线性独立的。

对称矩阵 $A$ 属于不同特征值的特征向量是正交的。如果 $\lambda_i$ 和 $\lambda_j$ 是 $A$ 的两个不同特征值， $v_i$ 和 $v_j$ 是它们各自的特征向量，有：

$$\lambda_j \mathbf{v}_i^T \mathbf{v}_j = \mathbf{v}_i^T \mathbf{A} \mathbf{v}_j = \mathbf{v}_i^T \mathbf{A}^T \mathbf{v}_j = \mathbf{v}_j^T \mathbf{A} \mathbf{v}_i = \lambda_i \mathbf{v}_j^T \mathbf{v}_i = \lambda_i \mathbf{v}_i^T \mathbf{v}_j \quad (4.43)$$

式 (4.43) 多次用到了特征值/特征向量的定义,  $\mathbf{A}^T = \mathbf{A}$  以及标量的转置还是该实数本身。注意观察式 (4.43) 那些计算结果是标量。

根据式 (4.43), 有:

$$(\lambda_i - \lambda_j) \mathbf{v}_i^T \mathbf{v}_j = 0 \quad (4.44)$$

因为  $\lambda_i \neq \lambda_j$ , 所以必有  $\mathbf{v}_i^T \mathbf{v}_j = 0$ 。

$\mathbf{A}$  的属于同一个特征值的特征向量已经经过施密特正交化, 它们彼此正交。刚刚证明了属于  $\mathbf{A}$  不同特征值的特征向量是正交的, 所以  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$  是线性独立、两两正交的单位向量。用它们作为列, 构造矩阵:

$$\mathbf{V}^T = (\mathbf{v}_1 \quad \cdots \quad \mathbf{v}_n) \quad (4.45)$$

$\mathbf{V}^T$  是正交矩阵。用  $\mathbf{A}$  的特征值  $\lambda_1, \lambda_2, \dots, \lambda_n$  构造对角矩阵:

$$\mathbf{\Lambda} = \begin{pmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_n \end{pmatrix} \quad (4.46)$$

$\mathbf{\Lambda}$  的  $n$  个对角线元素是从大到小排列的特征值。其余元素是 0。有:

$$\mathbf{V}^T \mathbf{A} = (\mathbf{v}_1 \quad \cdots \quad \mathbf{v}_n) \begin{pmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_n \end{pmatrix} = (\lambda_1 \mathbf{v}_1 \quad \cdots \quad \lambda_n \mathbf{v}_n) = \mathbf{A} \mathbf{V}^T \quad (4.47)$$

因为  $\mathbf{V} \mathbf{V}^T = \mathbf{I}$ , 所以有:

$$\mathbf{A} = \mathbf{V}^T \mathbf{\Lambda} \mathbf{V} \quad (4.48)$$

这就是对称矩阵的谱分解。

#### 4.1.5 二次型

式(4.48)中 $\mathbf{V}^T$ 的列，也就是 $\mathbf{V}$ 的行是 $\mathbb{R}^n$ 的一组标准正交基。所以对于任何向量 $\mathbf{x}$ ， $\mathbf{V}\mathbf{x}$ 的每一个元素是 $\mathbf{x}$ 向 $n$ 个标准基向量的投影长度，即 $\mathbf{x}$ 在这组基下的坐标，写作 $\mathbf{V}\mathbf{x} = (c_1 \ \cdots \ c_n)^T$ 。这时候有：

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T \mathbf{V}^T \mathbf{\Lambda} \mathbf{V} \mathbf{x} = (\mathbf{V} \mathbf{x})^T \begin{pmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_n \end{pmatrix} \mathbf{V} \mathbf{x} = \sum_{i=1}^n \lambda_i c_i^2 \quad (4.49)$$

式(4.49)说明 $\mathbf{x}^T \mathbf{A} \mathbf{x}$ 等于向量 $\mathbf{x}$ 在坐标系 $\mathbf{V}^T$ 下各个坐标值的平方乘以对应特征值 $\lambda_i$ 再加和。

对于某个对称矩阵 $\mathbf{A}$ ：

$$q(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} \quad (4.50)$$

称作一个二次型。