

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 3×3 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，它是一个列向量：

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

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 $\mathbf{w} = (w_1 \ \cdots \ w_n)^T \neq \mathbf{0}$ 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：

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

θ_{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：

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

θ_{31} 是 \mathbf{x}_3 与 \mathbf{x}'_1 的夹角， θ_{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 $\mathbf{x}'_3 = 0$ ，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 $\mathbf{x}'_3 \neq 0$ 。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 θ 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） λ ，有：

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

则称 λ 是 \mathbf{A} 的特征值， \mathbf{v} 是 \mathbf{A} 的对于 λ 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 λ 对应的特征向量，容易验证 $k\mathbf{v}$ 也是 λ 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 λ 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 λ 对应的特征向量。这也很容易验证。如果特征值 λ 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 λ 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 λ 对应的特征空间，其中所有向量都是 λ 的特征向量。

将式 (4.29) 变形。如果 λ 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 λ 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 λ 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 λ ，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 λ_i 和 λ_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 λ_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） λ 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。

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

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

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

因为 $\mathbf{x} \neq \mathbf{0}$ ，根据式 (4.37) $\bar{\mathbf{v}}^T\mathbf{v} > 0$ 。所以 $\lambda = \bar{\lambda}$ ，即 λ 是实数。 \mathbf{A} 是实矩阵， λ 是实数，所以 \mathbf{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 \mathbf{A} 有 n 个实特征值（可重复）。有一个结论我们不加证明：如果 λ 是 \mathbf{A} 的 k 重特征值（方程 $|\mathbf{A} - \lambda\mathbf{I}| = 0$ 的 k 重根），则 λ 对应的特征空间是 k 维，即对于 λ 能找到 k 个线性独立的特征向量。

对于对称矩阵 \mathbf{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 $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ 。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：

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

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：

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

$I(i)$ 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 $\mathbf{v}'_i = \sum_{j \in I(i)} a_j \mathbf{v}_j$ 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 \mathbf{v}'_i 都是零向量。因为存在非 0 的 a_j 且 $\sum_{j \in I(i)} a_j \mathbf{v}_j = \mathbf{0}$ ，这与 $\mathbf{v}_{j \in I(i)}$ 线性独立产生矛盾。那些非零的 \mathbf{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 $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ 是线性独立的。

对称矩阵 \mathbf{A} 属于不同特征值的特征向量是正交的。如果 λ_i 和 λ_j 是 \mathbf{A} 的两个不同特征值， \mathbf{v}_i 和 \mathbf{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:

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

\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:

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

$\mathbf{\Lambda}$ 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:

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

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

对于某个方阵 \mathbf{A} （不一定对称）：

$$q(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} \quad (4.49)$$

称作一个二次型。例如：

$$q(\mathbf{x}) = (x_1 \ x_2) \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = x_1^2 + 5x_1x_2 + 4x_2^2 \quad (4.50)$$

二次型只包含二次项，不包含一次项、常数（零次）项以及更高次项。这种情况称为是齐次的。从（4.50）可以看出， $q(\mathbf{x})$ 也可以写成：

$$q(\mathbf{x}) = (x_1 \ x_2) \begin{pmatrix} 1 & \frac{5}{2} \\ \frac{5}{2} & 4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = x_1^2 + 5x_1x_2 + 4x_2^2 \quad (4.51)$$

将“交叉”项（二元情况下就是 x_1x_2 ）的系数对半分，任何一个二次型都能写成关于对称矩阵的二次型。

一个对称矩阵 \mathbf{A} ，如果对于任何向量 \mathbf{x} 都有：

$$q(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0 \quad (4.52)$$

则称 \mathbf{A} 是半正定的（positive semidefinite）。如果不等号是严格的（>），则 \mathbf{A} 是正定的（positive definite）。相应地还有半负定（negative semidefinite）和负定（negative definite）。半正定矩阵的所有特征值都不为负，因为假如 \mathbf{A} 一个负数特征值 $\lambda_i < 0$ ，那么 \mathbf{A} 可以分解为：

$$\mathbf{A} = \mathbf{V}^T \mathbf{\Lambda} \mathbf{V} \quad (4.52)$$

因为 \mathbf{V}^T 是正交矩阵，根据正交矩阵的定义， \mathbf{V} 也是正交矩阵。 \mathbf{V} 的列是 \mathbb{R}^n 的基，即存在 \mathbf{x} 使

得 $\mathbf{e}_i = \mathbf{V}\mathbf{x}$ 。那么有：

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = (\mathbf{V}\mathbf{x})^T \mathbf{\Lambda} \mathbf{V}\mathbf{x} = \mathbf{e}_i^T \mathbf{\Lambda} \mathbf{e}_i = \lambda_i < 0 \quad (4.53)$$

这与 \mathbf{A} 的半正定性矛盾。所以半正定矩阵的所有特征值都非负。同样可以证明正定矩阵的所有特征值都为正。半负定矩阵的所有特征值都非正以及负定矩阵的所有特征值都为负。

半正定矩阵可以分解成两个矩阵的乘积：

$$\begin{aligned} \mathbf{A} = \mathbf{V}^T \mathbf{\Lambda} \mathbf{V} &= \mathbf{V}^T \begin{pmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_n \end{pmatrix} \mathbf{V} = \mathbf{V}^T \begin{pmatrix} \sqrt{\lambda_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sqrt{\lambda_n} \end{pmatrix} \begin{pmatrix} \sqrt{\lambda_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sqrt{\lambda_n} \end{pmatrix} \mathbf{V} = \\ &= \mathbf{V}^T \left(\mathbf{\Lambda}^{\frac{1}{2}} \right)^T \mathbf{\Lambda}^{\frac{1}{2}} \mathbf{V} = \left(\mathbf{\Lambda}^{\frac{1}{2}} \mathbf{V} \right)^T \mathbf{\Lambda}^{\frac{1}{2}} \mathbf{V} \end{aligned} \quad (4.54)$$

如果限定 \mathbf{x} 是单位向量，则半正定矩阵 \mathbf{A} 的二次型的最大值是 \mathbf{A} 的最大特征值 λ_1 ，达到这个最大值的 \mathbf{x} 是 λ_1 的某一个单位特征向量：

$$\max_{\|\mathbf{x}\|=1} \mathbf{x}^T \mathbf{A} \mathbf{x} = \lambda_1 \quad (4.55)$$

\mathbf{V} 是正交矩阵， \mathbf{x} 是任意向量，审视 $\mathbf{V}\mathbf{x}$ 的模：

$$\|\mathbf{V}\mathbf{x}\|^2 = (\mathbf{V}\mathbf{x})^T \mathbf{V}\mathbf{x} = \mathbf{x}^T \mathbf{V}^T \mathbf{V} \mathbf{x} = \mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|^2 \quad (4.66)$$

即用正交矩阵去变换任意向量，不改变向量的长度。从几何角度理解，正交矩阵的行是标准正交基。用正交矩阵乘一个向量相当于将该向量投影在新的标准坐标系中，所以长度不发生改变。

令 $\mathbf{\Lambda}$ 是对半正定 \mathbf{A} 谱分解得到的对角阵，对于任意单位向量 $\mathbf{x}' = \mathbf{V}\mathbf{x}$ ，有：

$$(\mathbf{x}')^T \mathbf{\Lambda} \mathbf{x}' = \sum_{i=1}^n \lambda_i (x'_i)^2 \quad (4.67)$$

显然在 $\sum_{i=1}^n (x'_i)^2 = 1$ 限制下，当 $x'_1 = 1$ 且 $x'_{i \neq 1} = 0$ ，也就是 $\mathbf{x}' = \mathbf{e}_1$ 时式（4.67）达到最大。如果 $\mathbf{V}\mathbf{x} = \mathbf{e}_1$ ，则 $\mathbf{x} = \mathbf{V}^T \mathbf{e}_1 = \mathbf{v}_1$ 。所以当 \mathbf{x} 是 \mathbf{A} 的最大特征值 λ_1 对应的某个单位特征向量 \mathbf{v}_1 时，

$\mathbf{x}^T \mathbf{A} \mathbf{x}$ 取得最大值 λ_1 。

问 $\|\mathbf{x}\| = 1$ 时次大的 $\mathbf{x}^T \mathbf{A} \mathbf{x}$ 值没有意义。二次型是连续的，当 \mathbf{x} 无限接近 \mathbf{v}_1 ， $\mathbf{x}^T \mathbf{A} \mathbf{x}$ 可以无限接近 λ_1 。我们加上一个限制条件： \mathbf{x} 须与 \mathbf{v}_1 正交。有结论：

$$\max_{\|\mathbf{x}\|=1, \mathbf{x} \perp \mathbf{v}_1} \mathbf{x}^T \mathbf{A} \mathbf{x} = \lambda_2 \quad (4.68)$$

在模为 1 并正交于 \mathbf{v}_1 约束下， $\mathbf{x}^T \mathbf{A} \mathbf{x}$ 的最大值是 \mathbf{A} 的第二大特征值 λ_2 。 \mathbf{x} 取 λ_2 的某个单位特征向量时达到此值。证明这个结论，回顾式 (4.67)，正交于 \mathbf{v}_1 的约束使 \mathbf{x}' 的第一个元素为 0。那么令式 (4.67) 最大就是把所有“量”都分配给 \mathbf{x}' 的第 2 个元素，因为它对应的系数是第二大特征值 λ_2 。所以 $\mathbf{V} \mathbf{x} = \mathbf{e}_2$ ，即 $\mathbf{x} = \mathbf{V}^T \mathbf{e}_2 = \mathbf{v}_2$ 。

总结一下：半正定矩阵 \mathbf{A} 的最大特征值 λ_1 对应的单位特征向量 \mathbf{v}_1 ，在 $\|\mathbf{x}\| = 1$ 的约束下使 $\mathbf{x}^T \mathbf{A} \mathbf{x}$ 达到最大值 λ_1 ； \mathbf{A} 的第二大特征值 λ_2 对应的单位特征向量 \mathbf{v}_2 ，在 $\|\mathbf{x}\| = 1$ 且正交于 \mathbf{v}_1 的约束下使 $\mathbf{x}^T \mathbf{A} \mathbf{x}$ 达到最大值 λ_2 ； \mathbf{A} 的第三大特征值 λ_3 对应的单位特征向量 \mathbf{v}_3 ，在 $\|\mathbf{x}\| = 1$ 且正交于 \mathbf{v}_1 和 \mathbf{v}_2 的约束下使 $\mathbf{x}^T \mathbf{A} \mathbf{x}$ 达到最大值 λ_3 。依次类推。

这个结论在后文介绍主成分分析和导致梯度下降发生震荡的损失函数病态形态时都有用处。

4.2 多元函数局部二阶特性与驻点的类型

第 3 章讲到：梯度为零向量的点——驻点，是局部极小/大点的必要条件。但是紧靠一阶梯度信息，无法判断驻点的具体类型。要做到这一点，就要分析函数在驻点附近的二阶特性。

4.2.1 赫森矩阵

n 元函数 $f(\mathbf{x})$ 在 \mathbf{x} 点的赫森（Hessian）矩阵是 $n \times n$ 矩阵 $\mathbf{H}(\mathbf{x})$ ：

$$\mathbf{H}(\mathbf{x}) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{pmatrix} \quad (4.69)$$

$\frac{\partial^2 f}{\partial x_j \partial x_i}$ 表示先求 f 对 x_i 求偏导，再对 x_j 求偏导，即 f 的二阶偏导。这些二阶导数是在 \mathbf{x} 点计算，

应记作 $\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x})$ 。文中为了清晰省略自变量 \mathbf{x} 。有一个结论本书不给出证明：如果 $\frac{\partial^2 f}{\partial x_j \partial x_i}$ 和

$\frac{\partial^2 f}{\partial x_i \partial x_j}$ 都是连续的，则它们相等。在本书考虑的问题中这种连续性要求都是满足的。所以 $\mathbf{H}(\mathbf{x})$

是一个对称矩阵。

赫森矩阵 $\mathbf{H}(\mathbf{x})$ 的第一行是 $\left(\frac{\partial f}{\partial x_1 \partial x_1} \quad \cdots \quad \frac{\partial f}{\partial x_n \partial x_1}\right)$ 。它的 n 个元素是 $\frac{\partial f}{\partial x_1}$ 对 x_1, x_2, \dots, x_n 的偏导。这

个行向量是函数 $\frac{\partial f}{\partial x_1}$ 的梯度的转置 $\left(\nabla \frac{\partial f}{\partial x_1}\right)^T$ 。

4.2.2 二阶泰勒展开

前文我们知道函数 $f(\mathbf{x})$ 在 \mathbf{x} 沿着方向 \mathbf{d} 的方向导数是 $\nabla f(\mathbf{x})^T \mathbf{d}$ 。当限制在一个方向上，可以把 $f(\mathbf{x})$ 看做一元函数，方向导数就是该一元函数的导数。那么当然也存在沿着方向 \mathbf{d} 的二阶导数，即导数的导数。首先我们需要求 $\nabla f^T \mathbf{d}$ 的梯度：

$$\frac{\partial \nabla f(\mathbf{x})^T \mathbf{d}}{\partial x_i} = \frac{\partial \sum_{j=1}^n d_j \frac{\partial f}{\partial x_j}}{\partial x_i} = \sum_{j=1}^n d_j \frac{\partial f}{\partial x_i \partial x_j} = \mathbf{h}_{*i}^T \mathbf{d} \quad (4.70)$$

其中 \mathbf{h}_{*i} 是 $\mathbf{H}(\mathbf{x})$ 的第 i 列。 $\nabla f(\mathbf{x})^T \mathbf{d}$ 对 x_i 的偏导数是 $\mathbf{H}(\mathbf{x})$ 的第 i 列与 \mathbf{d} 的内积，所以 $\nabla f(\mathbf{x})^T \mathbf{d}$ 在 \mathbf{x} 的梯度是：

$$\nabla(\nabla f(\mathbf{x})^T \mathbf{d}) = \begin{pmatrix} \frac{\partial \nabla f(\mathbf{x})^T \mathbf{d}}{\partial x_1} \\ \vdots \\ \frac{\partial \nabla f(\mathbf{x})^T \mathbf{d}}{\partial x_n} \end{pmatrix} = \begin{pmatrix} \mathbf{h}_{*1}^T \mathbf{d} \\ \vdots \\ \mathbf{h}_{*n}^T \mathbf{d} \end{pmatrix} = \mathbf{H}(\mathbf{x})^T \mathbf{d} \quad (4.71)$$

所以 $\nabla f(\mathbf{x})^T \mathbf{d}$ 沿 \mathbf{d} 方向的方向导，也就是 $f(\mathbf{x})$ 沿 \mathbf{d} 方向二阶导是：

$$\left(\nabla(\nabla f(\mathbf{x})^T \mathbf{d})\right)^T \mathbf{d} = (\mathbf{H}(\mathbf{x})^T \mathbf{d})^T \mathbf{d} = \mathbf{d}^T \mathbf{H}(\mathbf{x}) \mathbf{d} \quad (4.72)$$

结论是 $f(\mathbf{x})$ 在 \mathbf{x} 沿着方向 \mathbf{d} 的二阶导数是二次型 $\mathbf{d}^T \mathbf{H}(\mathbf{x}) \mathbf{d}$ 。回忆一下一元函数的二阶泰勒展开（一元微积分是本书的先导知识，故涉及一元微积分的结论直接使用）：

$$f(x+h) = f(x) + f'(x)h + \frac{f''(x)h^2}{2} + \mathcal{R}(h) \quad (4.73)$$

余项 $\mathcal{R}(\mathbf{h})$ 是 h^2 的高阶无穷小:

$$\lim_{h \rightarrow 0} \frac{\mathcal{R}(h)}{h^2} = 0 \quad (4.74)$$

二阶泰勒展开可以扩展到多维:

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + \frac{\mathbf{h}^T \mathbf{H}(\mathbf{x}) \mathbf{h}}{2} + \mathcal{R}(\mathbf{h}) \quad (4.75)$$

$\nabla f(\mathbf{x})$ 是 f 在 \mathbf{x} 的梯度, $\mathbf{H}(\mathbf{x})$ 是 f 在 \mathbf{x} 的赫森矩阵。 $\mathcal{R}(\mathbf{h})$ 是 $\|\mathbf{h}\|^2$ 的高阶无穷小:

$$\lim_{\|\mathbf{h}\| \rightarrow 0} \frac{\mathcal{R}(\mathbf{h})}{\|\mathbf{h}\|^2} = 0 \quad (4.76)$$

因为任何一个 $f(\mathbf{x} + \mathbf{h})$ 都可以看作是自变量从 \mathbf{x} 出发, 沿着 $\mathbf{d} = \mathbf{h}/\|\mathbf{h}\|$ 方向变化了 $\|\mathbf{h}\|$ 距离。将 $f(\mathbf{x} + \mathbf{h})$ 视为关于一个标量 t 的一元函数 $g(t) = f(\mathbf{x} + t \mathbf{h}/\|\mathbf{h}\|)$, 利用一元函数的二阶泰勒展开, 有:

$$f(\mathbf{x} + \mathbf{h}) = g(\|\mathbf{h}\|) = g(0) + g'(0)\|\mathbf{h}\| + \frac{g''(0)}{2}\|\mathbf{h}\|^2 + \mathcal{R} = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + \frac{\mathbf{h}^T \mathbf{H}(\mathbf{x}) \mathbf{h}}{2} + \mathcal{R} \quad (4.77)$$

当变化距离 $\|\mathbf{h}\|$ 趋近于 0 时, 余项 \mathcal{R} 也趋近于 0, 且比距离的平方消失得更快 ($\|\mathbf{h}\|^2$ 的高阶无穷小)。这就证明了式 (4.75)。

式 (4.75) 右边刨除余项的部分:

$$q(\mathbf{h}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + \frac{\mathbf{h}^T \mathbf{H}(\mathbf{x}) \mathbf{h}}{2} \quad (4.78)$$

$q(\mathbf{h})$ 是关于向量 \mathbf{h} 的二次函数。注意 $q(\mathbf{h})$ 不是二次型。二次型只包含二次项, 而 $q(\mathbf{h})$ 包含常数 (零次) 项、一次项以及二次项。二次函数的典型图像如图 4-1。

图 4-1 二次函数的典型图像

$q(\mathbf{h})$ 的图像过 $(0, f(\mathbf{x}))$ 。将 $q(\mathbf{h})$ 的图像平移，使 $(0, f(\mathbf{x}))$ 移动到 $(\mathbf{x}, f(\mathbf{x}))$ 。平移后的图像是 $f(\mathbf{x})$ 在 \mathbf{x} 附近的二阶近似。 $q(\mathbf{h})$ 的前两项 $f(\mathbf{x}) + \nabla f^T \mathbf{h}$ 是仿射函数。它的全部信息包含在 f 的梯度中。它是 $f(\mathbf{x})$ 的一阶近似。所以，二阶泰勒展开是对原函数精细到二阶的近似。误差信息包含在余项中。函数的局部一阶特性包含在梯度中，局部二阶特性包含在赫森矩阵中。如图 4-2 所示。

图 4-2 函数的局部二阶近似

4.2.3 驻点的类型

驻点到底是局部极小点，局部极大点或鞍点，这些信息一定程度上包含在函数的局部二阶特性中。假如 \mathbf{x}^* 是 $f(\mathbf{x})$ 的驻点。令 $\mathbf{H}(\mathbf{x})$ 是 $f(\mathbf{x})$ 在 \mathbf{x}^* 的赫森矩阵。因为有 $\nabla f(\mathbf{x}^*) = \mathbf{0}$ ，根据二阶泰勒展开，离开 \mathbf{x}^* 的变化量为 \mathbf{h} 的点的函数值是：

$$f(\mathbf{x}^* + \mathbf{h}) = f(\mathbf{x}^*) + \frac{\mathbf{h}^T \mathbf{H}(\mathbf{x}) \mathbf{h}}{2} + \mathcal{R}(\mathbf{h}) \quad (4.79)$$

令 $\mathbf{d} = \mathbf{h}/\|\mathbf{h}\|$ 是与 \mathbf{h} 同方向的单位向量。假如 $\mathbf{H}(\mathbf{x})$ 是正定的，有 $\mathbf{d}^T \mathbf{H}(\mathbf{x}) \mathbf{d} > 0$ 。 $\mathcal{R}(\mathbf{h})$ 是 $\|\mathbf{h}\|^2$ 的高阶无穷小。所以当 $\|\mathbf{h}\| \rightarrow 0$ ，式（4.79）右边后两项有：

$$\lim_{\|\mathbf{h}\| \rightarrow 0} \frac{\mathbf{h}^T \mathbf{H}(\mathbf{x}) \mathbf{h}/2 + \mathcal{R}(\mathbf{h})}{\|\mathbf{h}\|^2} = \mathbf{d}^T \mathbf{H}(\mathbf{x}) \mathbf{d} > 0 \quad (4.80)$$

所以当 $\mathbf{x} + \mathbf{h}$ 足够靠近 \mathbf{x} 时，有 $f(\mathbf{x}^* + \mathbf{h}) > f(\mathbf{x}^*)$ 。 \mathbf{x}^* 是局部极小点。注意，赫森矩阵必须为正定。如果仅仅是半正定，那么式（4.90）的极限只能是大于等于 0，那么自变量顺着某个方向可以从负侧接近 0，不能保证 \mathbf{x}^* 局部极小。类似地，如果 $\mathbf{H}(\mathbf{x})$ 是负定的， \mathbf{x}^* 是局部极大点。

正定矩阵特征值都为正，负定矩阵所有特征值都为负。如果特征值有正有负，则矩阵是不定的。如果 $\mathbf{H}(\mathbf{x})$ 不定，则沿着正特征值的特征方向二阶导数为正，函数向上翘；沿着负特征值方向，函数向下弯。此时在任何领域内，都有函数值大于 $f(\mathbf{x}^*)$ 的点，也有函数值小于 $f(\mathbf{x}^*)$ 的点， \mathbf{x}^* 是一个鞍点。如图 4-3 所示。

图 4-3 驻点的类型与赫森矩阵的关系

总结如下：

- 驻点 \mathbf{x}^* 的赫森矩阵 $\mathbf{H}(\mathbf{x})$ 正定， \mathbf{x}^* 是局部极小点；
- 驻点 \mathbf{x}^* 的赫森矩阵 $\mathbf{H}(\mathbf{x})$ 负定， \mathbf{x}^* 是局部极大点；
- 驻点 \mathbf{x}^* 的赫森矩阵 $\mathbf{H}(\mathbf{x})$ 不定， \mathbf{x}^* 是鞍点；
- 驻点 \mathbf{x}^* 的赫森矩阵 $\mathbf{H}(\mathbf{x})$ 半正定但非正定，或半正定但非负定， \mathbf{x}^* 的性质无法确定；

如果不考虑最后一种情况，那么只有当赫森矩阵是正定的，也就是其特征值全部大于 0 时，驻点是局部极小点。在没有任何先验知识时可以假设特征值大于 0 的概率是 0.5 且彼此独立。那么如果模型有 n 个参数，则损失函数的赫森矩阵的特征值全部大于 0 的概率是 $1/2^n$ 。这是一个极小的概率。所以在例如神经网络或深度学习这种参数很多的模型中，局部极小点应该是很罕见的。

函数的局部二阶特性提供了比一阶特性更多的信息。但如果赫森矩阵只是“半定”的，则仍然没有足够的信息完全判断驻点性质，这时就需要更精确的近似。

在梯度下降法中，梯度反方向 $-\nabla f$ 虽然是下降最快的方向，但如果知道 $\nabla f(\mathbf{x})^T \mathbf{H}(\mathbf{x}) \nabla f(\mathbf{x}) / \|\nabla f(\mathbf{x})\|^2 > 0$ 且绝对值很大，则沿着 $-\nabla f(\mathbf{x})$ 方向，（方向）导数上升且升得很快。这时步长如果较大，则有可能一步之后函数值不降反升。如果在选择前进方向时不仅参考梯度 $\nabla f(\mathbf{x})$ ，也参考赫森矩阵 $\mathbf{H}(\mathbf{x})$ ，优化算法能做出更优的选择。后文将要讲解的牛顿法和共轭方向法就是如此。

4.2.4 条件数与病态峡谷

令 $q(x_1, x_2)$ 是一个二元二次函数：

$$q(x_1, x_2) = a + b_1 x_1 + b_2 x_2 + \frac{c_{11}}{2} x_1^2 + \frac{c}{2} x_1 x_2 + \frac{c_{22}}{2} x_2^2 = a + \mathbf{b}^T \mathbf{x} + \frac{1}{2} \mathbf{x}^T \mathbf{C} \mathbf{x} \quad (4.81)$$

其中 $\mathbf{x} = (x_1 \ x_2)^T$ ， $\mathbf{b} = (b_1 \ b_2)^T$ ， $\mathbf{C} = \begin{pmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{pmatrix}$ ， $c_{12} = c_{21} = c/2$ 。考察 $q(x_1, x_2)$ 的梯度和赫森矩阵：

$$\nabla q = \begin{pmatrix} \frac{\partial q}{\partial x_1} \\ \frac{\partial q}{\partial x_2} \end{pmatrix} = \begin{pmatrix} b_1 + c_{11}x_1 + c_{12}x_2 \\ b_2 + c_{21}x_1 + c_{22}x_2 \end{pmatrix} = \mathbf{b} + \mathbf{C} \mathbf{x} \quad (4.82)$$

q 的赫森矩阵是：

$$\mathbf{H} = \begin{pmatrix} \frac{\partial q}{\partial x_1 \partial x_1} & \frac{\partial q}{\partial x_2 \partial x_1} \\ \frac{\partial q}{\partial x_1 \partial x_2} & \frac{\partial q}{\partial x_2 \partial x_2} \end{pmatrix} = \begin{pmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{pmatrix} = \mathbf{C} \quad (4.83)$$

q 的二阶泰勒展开就是它自身，余项为 0。即二次函数的二阶泰勒展开就精确地是它自身。二次函数的赫森矩阵是常矩阵 \mathbf{C} ，它的二阶特性是处处相同的。

现在假设 q 的 \mathbf{C} 是正定的，它的所有特征值为正。容易验证 \mathbf{C} 可逆。 \mathbf{C} 的逆矩阵是：

$$\mathbf{C}^{-1} = \mathbf{V}^T \mathbf{\Lambda}^{-1} \mathbf{V} = \mathbf{V}^T \begin{pmatrix} \frac{1}{\lambda_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\lambda_n} \end{pmatrix} \mathbf{V} \quad (4.84)$$

$\mathbf{V}^T \mathbf{\Lambda} \mathbf{V}$ 是 \mathbf{C} 的谱分解。

因为 \mathbf{C} 可逆，所以 q 的唯一驻点是 $\mathbf{x}^* = -\mathbf{C}^{-1} \mathbf{b}$ 是一个局部极小点。 q 在驻点 \mathbf{x}^* 上梯度为 0，则 q 沿任意方向的方向导数为 0。但是因为 \mathbf{C} 正定，所以 q 在 \mathbf{x}^* 沿任意方向的二阶导是 $\mathbf{d}^T \mathbf{C} \mathbf{d} > 0$ ，也就是沿任意方向的方向导数持续增大，过 \mathbf{x}^* 之前为负，过 \mathbf{x}^* 之后为正。也就是沿任意方向， q 都是先下降，在 \mathbf{x}^* 降到最低然后上升。 q 的图像呈碗状。因为 \mathbf{x}^* 是唯一的局部极小点，所以它是 q 的全局最小点。

q 的图形沿任意方向都是两端上升的抛物线，但上升的速度不一样。二阶导较大的点，方向导变化较快，沿该方向的抛物线上升更快；二阶导较小的点，方向导变化较慢，沿该方向的抛物线上升更慢。

使二阶导 $\mathbf{d}^T \mathbf{C} \mathbf{d}$ 最大的单位向量 \mathbf{d} 是 \mathbf{C} 的最大特征值 λ_1 对应的特征向量 \mathbf{v}_1 。在此方向上 q 具有最大的二阶导 λ_1 。垂直于 \mathbf{v}_1 且使二阶导 $\mathbf{d}^T \mathbf{C} \mathbf{d}$ 最大的单位向量是 \mathbf{C} 的次大特征值 λ_2 （在二元函数情况下就是最小特征值）。在此方向上 q 具有最大的二阶导 λ_2 。

称 λ_1/λ_2 是 \mathbf{C} 的条件数。条件数越大，最大特征值与最小特征值之比越大，在最大特征值的特征方向上函数图像的弯曲程度比垂直于该方向上函数图像的弯曲程度就更加剧烈。所以二次函数的赫森矩阵的条件数越大，则函数图像成越狭长的峡谷。这时候称二次函数的图像是病态的。反之若 $\lambda_1/\lambda_2 = 1$ ，则函数图像超各个方向的弯曲程度是相同的，图像成完美的碗状。如图 4-4 所示。

图 4-4 条件数与二次函数图像的狭长程度

狭长的峡谷称为病态的，它对梯度下降算法尤其不利。在第 3 章中已经看到在这种情况下，不合适的步长会引发震荡甚至不收敛。现在我们从某一个角度提供一个洞见。如果 \mathbf{C} 的条件数是 1，则 $\lambda_1 = \lambda_2$ ，令它们都等于 λ 。则有：

$$\mathbf{C}^{-1} = \mathbf{V}^T \mathbf{\Lambda}^{-1} \mathbf{V} = \mathbf{V}^T \begin{pmatrix} \frac{1}{\lambda} & 0 \\ 0 & \frac{1}{\lambda} \end{pmatrix} \mathbf{V} = \frac{1}{\lambda} \mathbf{V}^T \mathbf{V} = \frac{1}{\lambda} \mathbf{I} \quad (4.85)$$

那么 q 在的驻点 \mathbf{x}^* 是：

$$\mathbf{x}^* = -\mathbf{C}^{-1} \mathbf{b} = -\frac{1}{\lambda} \mathbf{b} \quad (4.86)$$

从任意点 \mathbf{x} 指向驻点 \mathbf{x}^* 的“箭头”是：

$$\mathbf{x}^* - \mathbf{x} = -\frac{1}{\lambda} \mathbf{b} - \mathbf{x} \quad (4.87)$$

而 q 在 \mathbf{x} 的负梯度是：

$$-\nabla q(\mathbf{x}) = -(\mathbf{b} + \mathbf{C}\mathbf{x}) = -(\mathbf{b} + \lambda\mathbf{x}) = \lambda(\mathbf{x}^* - \mathbf{x}) \quad (4.88)$$

式 (4.88) 的含义是，在任一点 \mathbf{x} ，其梯度正指向朝着 \mathbf{x}^* 的方向。于是每一步梯度下降都是朝着正确的方向前进。当然如果步长过大，有可能一步跨过 \mathbf{x}^* ，发生震荡。步长更大时，仍有可能发生不收敛。

现实中损失函数作为参数值的函数不是二次的，但是在驻点附近二次函数是很好的近似。本节分析的二次函数的特性，可以近似地表示任何函数在局部的特性。

4.3 基于函数二阶特性的优化

如果参考函数二阶特性信息，则优化过程会更有效率。但是高维情况下赫森矩阵的计算量是巨大的，这妨碍了二阶优化算法在神经网络或深度学习中的应用。人们提出了一些节省计算量的近似二阶算法，本书不对其进行讨论。本节介绍两个基于赫森矩阵的优化算法——牛顿法和共轭方向法，并探讨它们的几何意义。

4.3.1 牛顿法

牛顿法的思想是：在某点附近对函数进行二阶泰勒展开，之后求函数的近似二次函数的驻点，以该点为下一个点，迭代地进行此步骤。根据式(4.75)，离开 \mathbf{x} 一个变化量 \mathbf{h} 的函数值 $f(\mathbf{x} + \mathbf{h})$ 的二次近似函数是：

$$f(\mathbf{x} + \mathbf{h}) \approx q(\mathbf{h}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + \frac{\mathbf{h}^T \mathbf{H}(\mathbf{x}) \mathbf{h}}{2} \quad (4.89)$$

$q(\mathbf{h})$ 的梯度是：

$$\nabla q(\mathbf{h}) = \nabla f(\mathbf{x}) + \mathbf{H}(\mathbf{x}) \mathbf{h} \quad (4.90)$$

如果 $\mathbf{H}(\mathbf{x})$ 的特征值都非 0，则利用式(4.84)可以构造它的逆矩阵。特征值取某个特定值，例如 0 的概率很低，也就是绝大部分情况下都可以认为 $\mathbf{H}(\mathbf{x})$ 是可逆的。令 $\nabla q(\mathbf{h}) = \mathbf{0}$ ，解得 $q(\mathbf{h})$ 的驻点是：

$$\mathbf{h}^* = -\mathbf{H}(\mathbf{x})^{-1} \nabla f(\mathbf{x}) \quad (4.91)$$

所以，牛顿法对自变量 \mathbf{x} 的迭代更新量就是 $-\mathbf{H}(\mathbf{x})^{-1} \nabla f(\mathbf{x})$ 。牛顿法伪代码如下：

$\mathbf{x} \leftarrow$ randomly initialized

while $\|\nabla f(\mathbf{x})\| \geq \varepsilon$:

$\mathbf{x} \leftarrow \mathbf{x} - \mathbf{H}(\mathbf{x})^{-1} \nabla f(\mathbf{x})$

与梯度下降法一样， $\varepsilon > 0$ 是一个预设的阈值，当 $\|\nabla f(\mathbf{x})\| < \varepsilon$ 时认为 $\nabla f(\mathbf{x})$ 已经足够接近零向量，算法停止。也可以采用循环次数达到预设的最大值，或者函数值的下降幅度小于阈值等停止标准。

梯度下降法，向着反梯度方向 $-\nabla f(\mathbf{x})$ 更新自变量。而牛顿法将 $-\nabla f(\mathbf{x})$ 乘上了赫森矩阵的逆矩阵 $\mathbf{H}(\mathbf{x})^{-1}$ 。假设 $\mathbf{H}(\mathbf{x})$ 是正定的，根据对称矩阵的谱分解， $\mathbf{H}(\mathbf{x})^{-1}$ 可以分解成：

$$\mathbf{H}(\mathbf{x})^{-1} = \mathbf{V}^T \begin{pmatrix} \frac{1}{\lambda_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\lambda_n} \end{pmatrix} \mathbf{V} \quad (4.92)$$

\mathbf{V}^T 是由 $\mathbf{H}(\mathbf{x})$ 的特征向量构成的正交矩阵。 $\lambda_1, \lambda_2, \dots, \lambda_n$ 是从大到小排列的 $\mathbf{H}(\mathbf{x})$ 的 n 个特征值（可重）。牛顿算法的更新量可以写成：

$$-\mathbf{H}(\mathbf{x})^{-1} \nabla f(\mathbf{x}) = -\mathbf{V}^T \mathbf{\Lambda} \mathbf{V} \nabla f(\mathbf{x}) = -\mathbf{V}^T \begin{pmatrix} \frac{1}{\lambda_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\lambda_n} \end{pmatrix} \mathbf{V} \nabla f(\mathbf{x}) \quad (4.93)$$

\mathbf{V}^T 的列（ \mathbf{V} 的行）是 \mathbb{R}^n 的一组标准正交基，也就是一个坐标系。该坐标系的第一个轴沿着 $f(\mathbf{x})$ 二阶导最大的方向；第二个轴沿着垂直于第一个轴且 $f(\mathbf{x})$ 二阶导第二大的方向。向量 $\mathbf{V} \nabla f(\mathbf{x})$ 相当于对 $\nabla f(\mathbf{x})$ 做了一次旋转。如果忽略 $\mathbf{\Lambda}$ 矩阵， $\mathbf{V}^T \mathbf{V} \nabla f(\mathbf{x}) = \nabla f(\mathbf{x})$ 是将 $\nabla f(\mathbf{x})$ 旋转过去又旋转回来。

加上 $\mathbf{\Lambda}$ 的作用后，相当于用 $1/\lambda_i$ 惩罚 $\mathbf{H}(\mathbf{x})$ 的各个特征方向上 $\nabla f(\mathbf{x})$ 的分量。特征值 λ_i 越大，其对应的特征方向上 $f(\mathbf{x})$ 的二阶导越大。 $1/\lambda_i$ 越小，则对 $\nabla f(\mathbf{x})$ 在该特征方向上的分量惩罚越严厉。还记得狭长峡谷中梯度下降法的震荡发生的原因，就是二阶导较大的方向主宰了梯度方向，导致算法没有想着理想的局部极小点方向前进。牛顿法通过对二阶导最大的方向进行适当惩罚，缓解了这种现象。如图 4-5 所示。

图 4-5 牛顿法相当于对二阶导过大的方向进行惩罚

如果函数本身是二次函数，它的二阶泰勒展开就是它自身。牛顿法可以直接定位到二次近似的驻点。也就是说对于二次函数，牛顿法一步迭代就能寻找到它的唯一驻点——全局最小点。

4.3.2 共轭方向法

对于一个二次函数 $f(\mathbf{x})$ ，它在任意点的赫森矩阵是 \mathbf{H} 。假设 \mathbf{H} 正定。对于任意两个单位向量 \mathbf{d}_1 和 \mathbf{d}_2 ，如果满足：

$$\mathbf{d}_1^T \mathbf{H} \mathbf{d}_2 = 0 \quad (4.94)$$

则称 \mathbf{d}_1 与 \mathbf{d}_2 关于 \mathbf{H} 共轭。令自变量从某一个点 \mathbf{x} 开始，沿着 \mathbf{d}_1 运动，可以构造一个一元二次

函数：

$$g(h) = f(\mathbf{x} + h\mathbf{d}_1) = f(\mathbf{x}) + h\nabla f(\mathbf{x})^T \mathbf{d}_1 + \frac{h^2}{2} \mathbf{d}_1^T \mathbf{H} \mathbf{d}_1 \quad (4.95)$$

一元函数 $g(h)$ 的导数是：

$$g'(h) = \nabla f(\mathbf{x})^T \mathbf{d}_1 + h\mathbf{d}_1^T \mathbf{H} \mathbf{d}_1 \quad (4.96)$$

一元函数 $g(h)$ 的二阶导数是：

$$g''(h) = \mathbf{d}_1^T \mathbf{H} \mathbf{d}_1 > 0 \quad (4.97)$$

因为 \mathbf{H} 正定，所以二次函数 $g(h)$ 有全局最小点：

$$h^* = -\frac{\nabla f(\mathbf{x})^T \mathbf{d}_1}{\mathbf{d}_1^T \mathbf{H} \mathbf{d}_1} \quad (4.98)$$

因为 h^* 是 $g(h)$ 的全局最小点，所以 $g'(h^*) = 0$ ，即 $f(\mathbf{x})$ 在 $\mathbf{x} + h^*\mathbf{d}_1$ 沿 \mathbf{d}_1 的方向导数为 0。

$$\nabla_{\mathbf{d}_1} f(\mathbf{x} + h^*\mathbf{d}_1) = \nabla f(\mathbf{x} + h^*\mathbf{d}_1)^T \mathbf{d}_1 = 0 \quad (4.99)$$

\mathbf{d}_1 与 \mathbf{d}_2 关于 \mathbf{H} 共轭。

令自变量从 $\mathbf{x} + h^*\mathbf{d}_1$ 运动一个变化量 \mathbf{h} 。利用二阶泰勒展开， $f(\mathbf{x} + h^*\mathbf{d}_1 + \mathbf{h})$ 作为 \mathbf{h} 的方程是：

$$f(\mathbf{x} + h^*\mathbf{d}_1 + \mathbf{h}) = f(\mathbf{x} + h^*\mathbf{d}_1) + \nabla f(\mathbf{x} + h^*\mathbf{d}_1)^T \mathbf{h} + \frac{1}{2} \mathbf{h}^T \mathbf{H} \mathbf{h} \quad (4.100)$$

$f(\mathbf{x} + h^*\mathbf{d}_1 + \mathbf{h})$ 对 \mathbf{h} 的梯度是：

$$\nabla f(\mathbf{x} + h^* \mathbf{d}_1 + \mathbf{h}) = \nabla f(\mathbf{x} + h^* \mathbf{d}_1) + \mathbf{H} \mathbf{h} \quad (4.101)$$

$f(\mathbf{x} + h^* \mathbf{d}_1 + \mathbf{h})$ 沿 \mathbf{d}_1 的方向导数是:

$$\nabla_{\mathbf{d}_1} f(\mathbf{x} + h^* \mathbf{d}_1 + \mathbf{h}) = \nabla f(\mathbf{x} + h^* \mathbf{d}_1 + \mathbf{h})^T \mathbf{d}_1 = \nabla f(\mathbf{x} + h^* \mathbf{d}_1)^T \mathbf{d}_1 + \mathbf{h}^T \mathbf{H} \mathbf{d}_1 \quad (4.102)$$

如果变化量 \mathbf{h} 是沿着 \mathbf{d}_2 方向, 即 $\mathbf{h} = h \mathbf{d}_2$ 。根据式(4.99)以及 \mathbf{d}_1 与 \mathbf{d}_2 关于 \mathbf{H} 共轭, 有:

$$\nabla_{\mathbf{d}_2} f(\mathbf{x} + h^* \mathbf{d}_1 + h \mathbf{d}_2) = \nabla f(\mathbf{x} + h^* \mathbf{d}_1)^T \mathbf{d}_1 + h \mathbf{d}_1^T \mathbf{H} \mathbf{d}_2 = 0 \quad (4.103)$$

也就是说, 自变量从 \mathbf{x} 出发沿着