

21

Surfaces and Positive Definiteness

曲面和正定性

代数、微积分、几何、线性代数的结合体



神将一切几何化。

God ever geometrizes.

—— 柏拉图 (Plato) | 古希腊哲学家 | 424/423 ~ 348/347 BC



- ◀ matplotlib.pyplot.contour() 绘制等高线图
- ◀ matplotlib.pyplot.contourf() 绘制填充等高线图
- ◀ matplotlib.pyplot.scatter() 绘制散点图
- ◀ numpy.arange() 在给定间隔内返回均匀间隔的值
- ◀ numpy.array() 创建 array 数据类型
- ◀ numpy.cos() 余弦函数
- ◀ numpy.linalg.cholesky() Cholesky 分解函数
- ◀ numpy.linspace() 产生连续均匀向量数值
- ◀ numpy.meshgrid() 生成网格化数据
- ◀ numpy.multiply() 向量或矩阵逐项乘积
- ◀ numpy.roots() 多项式求根
- ◀ numpy.sin() 正弦函数
- ◀ numpy.sqrt() 平方根
- ◀ sympy.abc import x 定义符号变量 x
- ◀ sympy.diff() 求解符号导数和偏导解析式
- ◀ sympy.Eq() 定义符号等式
- ◀ sympy.evalf() 将符号解析式中未知量替换为具体数值
- ◀ sympy.plot_implicit() 绘制隐函数方程
- ◀ sympy.symbols() 定义符号变量

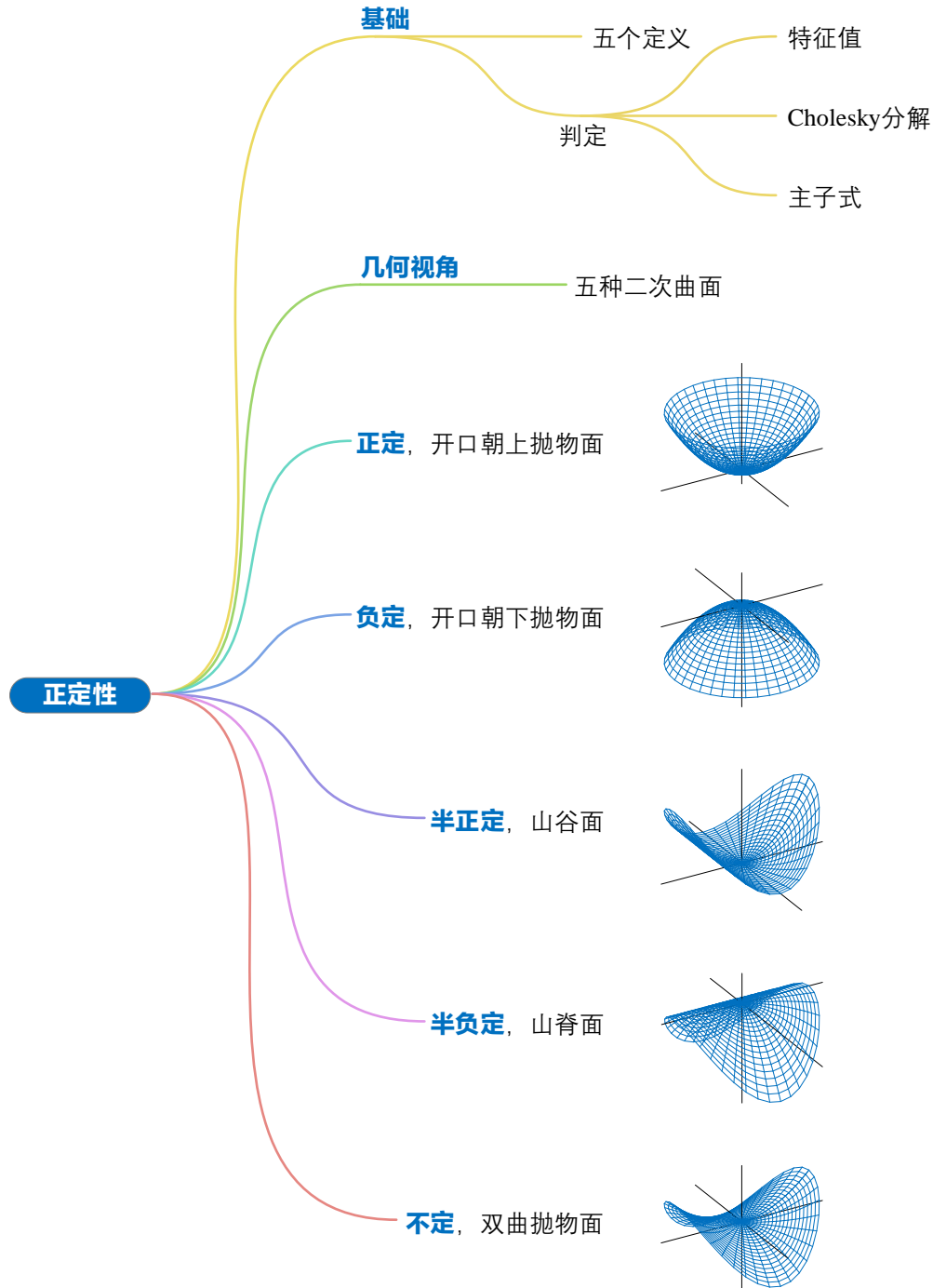
本 PDF 文件为作者草稿，发布目的为方便读者在移动终端学习，终稿内容以清华大学出版社纸质出版物为准。

版权归清华大学出版社所有，请勿商用，引用请注明出处。

代码及 PDF 文件下载：<https://github.com/Visualize-ML>

本书配套微课视频均发布在 B 站——生姜 DrGinger：<https://space.bilibili.com/513194466>

欢迎大家批评指教，本书专属邮箱：jiang.visualize.ml@gmail.com



本 PDF 文件为作者草稿，发布目的为方便读者在移动终端学习，终稿内容以清华大学出版社纸质出版物为准。

版权归清华大学出版社所有，请勿商用，引用请注明出处。

代码及 PDF 文件下载：<https://github.com/Visualize-ML>

本书配套微课视频均发布在 B 站——生姜 DrGinger：<https://space.bilibili.com/513194466>

欢迎大家批评指教，本书专属邮箱：jiang.visualize.ml@gmail.com

21.1 正定性

正定性 (positive definiteness) 是优化问题经常出现线性代数概念。本章结合三维曲面，特别是二次曲面 (quadratic surface)，和大家聊一聊正定性及其应用。

矩阵正定性分为如下五种情况。

当 $\mathbf{x} \neq \mathbf{0}$ (\mathbf{x} 为非零列向量) 时，如果满足：

$$\mathbf{x}^T \mathbf{A} \mathbf{x} > 0 \quad (1)$$

矩阵 \mathbf{A} 为**正定矩阵** (positive definite matrix)。

当 $\mathbf{x} \neq \mathbf{0}$ 时，

$$\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0 \quad (2)$$

矩阵 \mathbf{A} 为**半正定矩阵** (positive semi-definite matrix)。

当 $\mathbf{x} \neq \mathbf{0}$ 时，

$$\mathbf{x}^T \mathbf{A} \mathbf{x} < 0 \quad (3)$$

矩阵 \mathbf{A} 为**负定矩阵** (negative definite matrix)。

当 $\mathbf{x} \neq \mathbf{0}$ 时，

$$\mathbf{x}^T \mathbf{A} \mathbf{x} \leq 0 \quad (4)$$

矩阵 \mathbf{A} 为**半负定矩阵** (negative semi-definite matrix)。

矩阵 \mathbf{A} 不属于以上任何一种情况， \mathbf{A} 为**不定矩阵** (indefinite matrix)。

判定正定矩阵

判断矩阵是否为正定矩阵，本书主要采用如下两种方法：

- ◀ 若矩阵为对称矩阵，并且所有特征值为正，则矩阵为正定矩阵；
- ◀ 若矩阵可以进行 Cholesky 分解，则矩阵为正定矩阵。



Bk4_Ch21_01.py 介绍如何使用 Cholesky 分解判定矩阵是否为正定矩阵。

```
# Bk4_Ch21_01.py
import numpy as np
```

```
def is_pos_def(A):
    if np.array_equal(A, A.T):
        try:
            np.linalg.cholesky(A)
            return True
        except np.linalg.LinAlgError:
            return False
    else:
        return False

A = np.array([[1,0],
              [0,0]])

print(is_pos_def(A))
```

Cholesky 分解

如果矩阵 A 为正定矩阵，对 A 进行 Cholesky 分解，得到：

$$A = R^T R \quad (5)$$

利用 (5)，将 $x^T A x$ 写成如下形式：

$$x^T A x = x^T R^T R x = (R x)^T R x = \|R x\|^2 \quad (6)$$

R 中列向量线性无关，若 x 为非零向量，则 $R x \neq 0$ ，因此上式 $x^T A x > 0$ 。

特征值分解

对称矩阵 A 进行特征值分解得到：

$$A = V \Lambda V^T \quad (7)$$

将 (7) 代入 $x^T A x$ ，得到：

$$\begin{aligned} x^T A x &= x^T V \Lambda V^T x \\ &= \begin{pmatrix} V^T x \\ z \end{pmatrix}^T \Lambda \begin{pmatrix} V^T x \\ z \end{pmatrix} \end{aligned} \quad (8)$$

令：

$$z = V^T x \quad (9)$$

(8) 可以写成：

$$\begin{aligned} x^T A x &= z^T \Lambda z \\ &= \lambda_1 z_1^2 + \lambda_2 z_2^2 + \cdots + \lambda_D z_D^2 = \sum_{j=1}^D \lambda_j z_j^2 \end{aligned} \quad (10)$$

当上式中特征值均为正数，除非 z_1, z_2, \dots, z_D 均为 0 (即 z 为零向量)，否则上式大于 0。如果 x 和 z 存在 $x = Vz$ 这个等式关系，如果 z 为非零向量， x 也是非零向量。

若矩阵 A 为负定矩阵，则 A 的特征值均为负值。矩阵 A 为半正定矩阵，则矩阵 A 特征值为正值或 0。矩阵 A 为半负定矩阵，则矩阵特征值为负值或 0。

这一节介绍了正定性相关性质，但是要深入理解这个概念，还需要借助几何视角。

21.2 几何视角看正定性

给定如下 2×2 对称矩阵 A ：

$$A = \begin{bmatrix} a & b \\ b & c \end{bmatrix} \quad (11)$$

构造如下二元函数 $y = f(x_1, x_2)$ ：

$$\begin{aligned} y &= f(x_1, x_2) = \mathbf{x}^T A \mathbf{x} \\ &= \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} a & b \\ b & c \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &= ax_1^2 + 2bx_1x_2 + cx_2^2 \end{aligned} \quad (12)$$

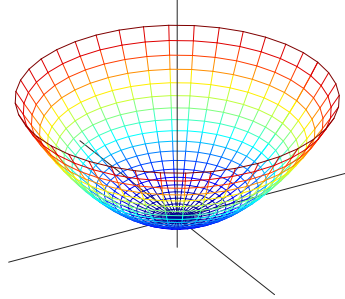
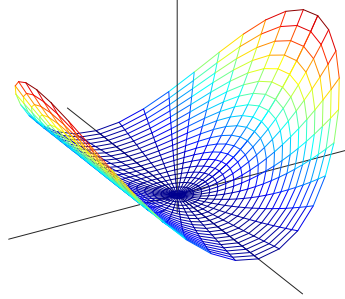
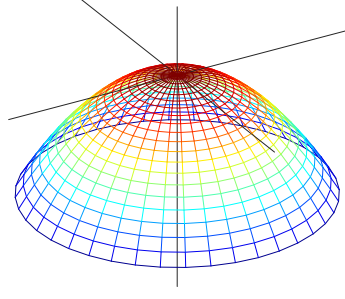
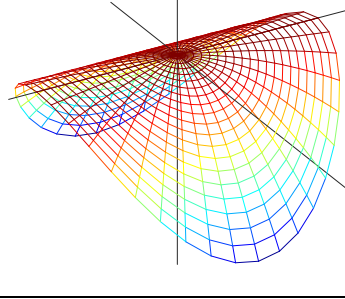
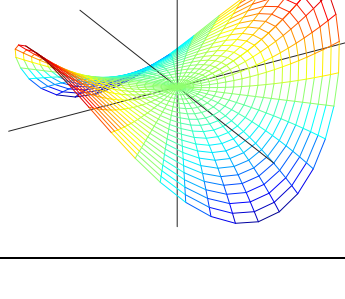
在三维正交空间中，当矩阵 $A_{2 \times 2}$ 正定性不同时， $y = f(x_1, x_2)$ 对应曲面展现出不同的形状：

- ◀ 当 $A_{2 \times 2}$ 为正定矩阵时， $y = f(x_1, x_2)$ 对应开口向上抛物面；
- ◀ 当 $A_{2 \times 2}$ 为半正定矩阵时， $y = f(x_1, x_2)$ 对应山谷面；
- ◀ 当 $A_{2 \times 2}$ 为负定矩阵时， $y = f(x_1, x_2)$ 对应开口向下抛物面；
- ◀ 当 $A_{2 \times 2}$ 为半负定矩阵时， $y = f(x_1, x_2)$ 对应山脊面；
- ◀ 当 $A_{2 \times 2}$ 不定时， $y = f(x_1, x_2)$ 为马鞍面，也叫做双曲抛物面。

表 1 总结了矩阵 A 不同正定性条件下对应的不同曲面形状。本章以下六节就按表中形状顺序展开。

表 1. 正定性的几何意义

$A_{D \times D}$	特征值	形状
------------------	-----	----

$A_{D \times D}$ 为正定矩阵 $\mathbf{x}^T A \mathbf{x} > 0, \mathbf{x} \neq \mathbf{0}$	D 个特征值均为正值	
$A_{D \times D}$ 为半正定矩阵，秩为 r $\mathbf{x}^T A \mathbf{x} \geq 0, \mathbf{x} \neq \mathbf{0}$	r 个正特征值， $D - r$ 个特征值为 0	
$A_{D \times D}$ 为负定矩阵 $\mathbf{x}^T A \mathbf{x} < 0$	D 个特征值均为负值	
$A_{D \times D}$ 为半负定矩阵，秩为 r $\mathbf{x}^T A \mathbf{x} \leq 0$	r 个负特征值， $D - r$ 个特征值为 0	
$A_{D \times D}$ 为不定矩阵	特征值符号正负不定	

21.3 开口朝上抛物面：正定

正圆

先来看一个单位矩阵的例子。矩阵 A 为 2×2 单位矩阵：

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (13)$$

构造如下二元函数 $y = f(x_1, x_2)$ ：

$$\begin{aligned} y &= f(x_1, x_2) = \mathbf{x}^T A \mathbf{x} \\ &= \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &= x_1^2 + x_2^2 \end{aligned} \quad (14)$$

观察上式，容易发现只有当 $x_1 = 0$ 且 $x_2 = 0$ 时， $y = f(x_1, x_2) = 0$ 。

容易求得 A 特征值分别为 $\lambda_1 = 1$ 和 $\lambda_2 = 1$ ，对应特征向量分别为：

$$\mathbf{v}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad (15)$$

计算矩阵 A 的秩， $\text{rank}(A) = 2$ 。

图 1 (a) 所示为 $y = f(x_1, x_2)$ 曲面。在该曲面边缘 A 、 B 和 C 放置小球，小球都会朝着曲面最低点滚动。

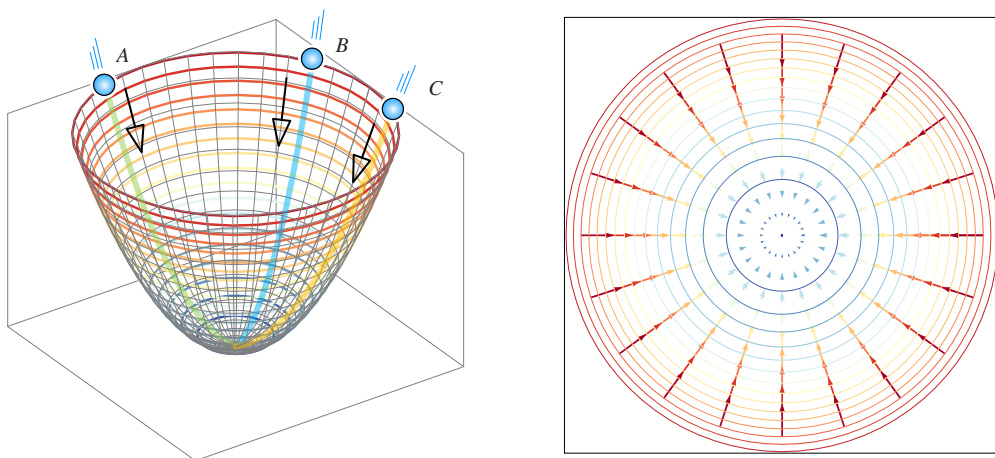


图 1. 正定矩阵曲面和梯度下降，正圆抛物面

(14) 的梯度下降向量为：

$$-\nabla f(\mathbf{x}) = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix} \quad (16)$$

图 1 (b) 展示曲面等高线为正圆，和不同位置的梯度下降向量。注意，图中给出的是梯度下降向量（下山方向），方向和梯度向量（上山方向）正好相反。

而梯度向量为 0 的点，就是 $y = f(x_1, x_2)$ 两个偏导均为 0 的点。本系列丛书《数学要素》介绍过，(0, 0) 这个点被称作驻点。通过图 1，我们可以很容易判断 (0, 0) 就是二元函数最小值点。

正椭圆

再看一个 2×2 正定矩阵例子。矩阵 A 具体值如下：

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \quad (17)$$

同样，构造二元函数 $y = f(x_1, x_2)$ ，具体如下：

$$\begin{aligned} y = f(x_1, x_2) &= \mathbf{x}^T A \mathbf{x} \\ &= \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &= x_1^2 + 2x_2^2 \end{aligned} \quad (18)$$

同样，只有 $x_1 = 0$ 且 $x_2 = 0$ 时， $y = f(x_1, x_2) = 0$ 。图 2 所示为 (18) 对应开口向上正椭圆抛物面。

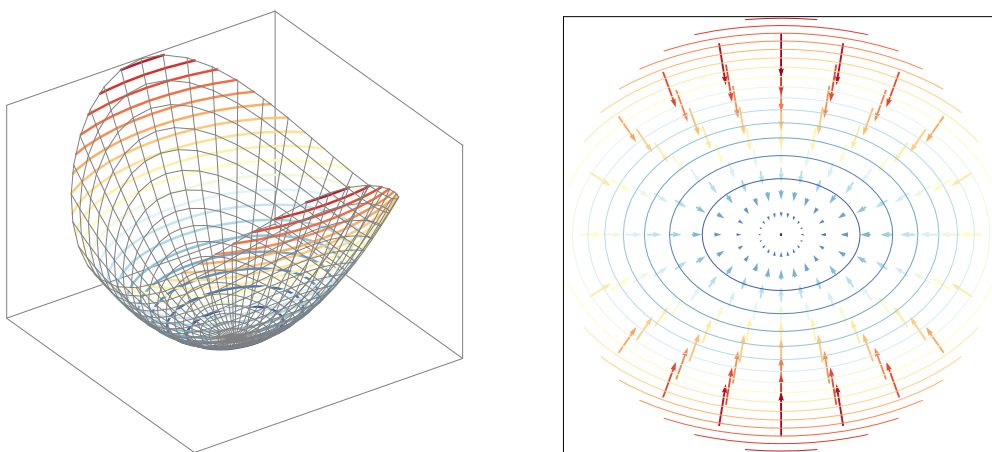


图 2. 正定矩阵曲面和梯度下降，正椭圆抛物面

容易求得 A 特征值分别为 $\lambda_1 = 1$ 和 $\lambda_2 = 2$ ，对应特征向量分别为：

$$\mathbf{v}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad (19)$$

(14) 的梯度下降向量为：

$$-\nabla f(\mathbf{x}) = - \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} -2x_1 \\ -4x_2 \end{bmatrix} \quad (20)$$

梯度向量为 0 的点 (0, 0) 也是函数的最小值点。

旋转椭圆

本节前两个例子对应的曲面的等高线分别是正圆和正椭圆，下面再看一个旋转椭圆情况。 \mathbf{A} 矩阵具体如下：

$$\mathbf{A} = \begin{bmatrix} 1.5 & 0.5 \\ 0.5 & 1.5 \end{bmatrix} \quad (21)$$

构造函数 $y = f(x_1, x_2)$ ：

$$\begin{aligned} y &= f(x_1, x_2) = \mathbf{x}^T \mathbf{A} \mathbf{x} \\ &= \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 1.5 & 0.5 \\ 0.5 & 1.5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &= 1.5x_1^2 + x_1x_2 + 1.5x_2^2 \end{aligned} \quad (22)$$

同样，只有当 $x_1 = 0$ 且 $x_2 = 0$ 时， $y = f(x_1, x_2) = 0$ 。

经过计算得到 \mathbf{A} 特征值也是 $\lambda_1 = 1$ 和 $\lambda_2 = 2$ ；这两个特征值对应特征向量分别为：

$$\mathbf{v}_1 = \begin{bmatrix} -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix} \quad (23)$$

(22) 梯度下降向量为：

$$-\nabla f(\mathbf{x}) = - \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} -3x_1 - x_2 \\ -x_1 - 3x_2 \end{bmatrix} \quad (24)$$

$y = f(x_1, x_2)$ 曲面对应图像如图 3。

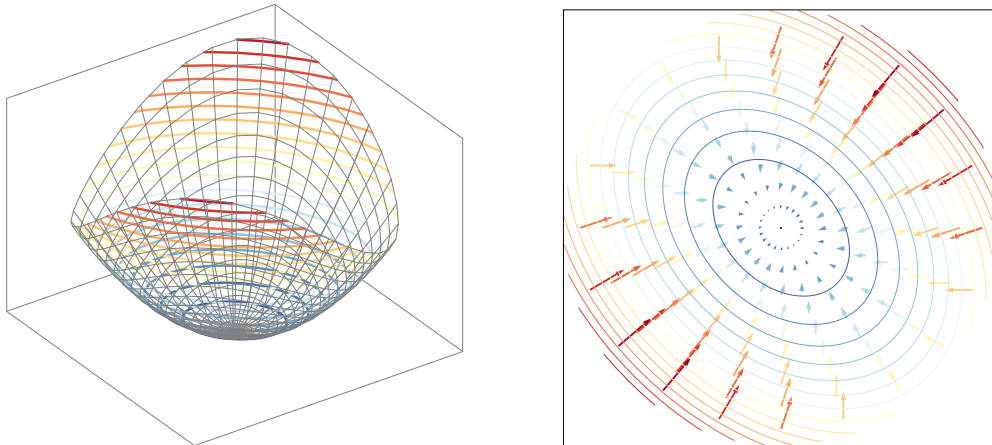
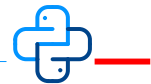


图 3. 正定矩阵曲面和梯度下降，开口向上旋转椭圆抛物面



Bk4_Ch21_02.py 绘制图 1、图 2、图 3，此外请大家修改绘制本章其他图像。

```
# Bk4_Ch21_02.py
import sympy
import numpy as np
import matplotlib.pyplot as plt

def mesh_circ(c1, c2, r, num):
    theta = np.arange(0, 2*np.pi+np.pi/num, np.pi/num)
    r = np.arange(0, r, r/num)
    theta, r = np.meshgrid(theta, r)
    xx1 = np.cos(theta)*r + c1
    xx2 = np.sin(theta)*r + c2
    return xx1, xx2

#define symbolic vars, function
x1, x2 = sympy.symbols('x1 x2')

A = np.array([[1.5, 0.5],
              [0.5, 1.5]])

x = np.array([[x1, x2]]).T

f_x = x.T@A@x
f_x = f_x[0][0]
print(f_x)

#take the gradient symbolically
grad_f = [sympy.diff(f_x, var) for var in (x1, x2)]
print(grad_f)

f_x_fcn = sympy.lambdify([x1, x2], f_x)

#turn into a bivariate lambda for numpy
grad_fcn = sympy.lambdify([x1, x2], grad_f)
```

本 PDF 文件为作者草稿，发布目的为方便读者在移动终端学习，终稿内容以清华大学出版社纸质出版物为准。

版权归清华大学出版社所有，请勿商用，引用请注明出处。

代码及 PDF 文件下载：<https://github.com/Visualize-ML>

本书配套微课视频均发布在 B 站——生姜 DrGinger：<https://space.bilibili.com/513194466>

欢迎大家批评指教，本书专属邮箱：jiang.visualize.ml@gmail.com

```

xx1, xx2 = mesh_circ(0, 0, 4, 20)

# coarse mesh
xx1_, xx2_ = mesh_circ(0, 0, 4, 10)
V = grad_fcn(xx1_, xx2_)
V_z = np.ones_like(V[1]);

if isinstance(V[1], int):
    V[1] = np.zeros_like(V[0])

elif isinstance(V[0], int):
    V[0] = np.zeros_like(V[1])

ff_x = f_x_fcn(xx1, xx2)

color_array = np.sqrt(V[0]**2 + V[1]**2)
l_3D_vectors = np.sqrt(V[0]**2 + V[1]**2 + V_z**2)

# 3D visualization
ax = plt.figure().add_subplot(projection='3d')
ax.plot_wireframe(xx1, xx2, ff_x, rstride=1,
                  cstride=1, color = [0.5, 0.5, 0.5],
                  linewidth = 0.2)
ax.contour3D(xx1, xx2, ff_x, 20, cmap = 'RdYlBu_r')

ax.xaxis.set_ticks([])
ax.yaxis.set_ticks([])
ax.zaxis.set_ticks([])
plt.xlim(xx1.min(), xx1.max())
plt.ylim(xx2.min(), xx2.max())
ax.set_proj_type('ortho')
ax.view_init(30, -125)
ax.set_xlabel('$x_1$')
ax.set_ylabel('$x_2$')
ax.set_zlabel('$f(x_1, x_2)$')
plt.tight_layout()

color_array = np.sqrt(V[0]**2 + V[1]**2)

# 2D visualization
fig, ax = plt.subplots()
plt.quiver(xx1_, xx2_, -V[0], -V[1], color_array,
           angles='xy', scale_units='xy',
           edgecolor='none', alpha=0.8, cmap = 'RdYlBu_r')

plt.contour(xx1, xx2, ff_x, 20, cmap = 'RdYlBu_r')
plt.show()
ax.set_aspect('equal')
ax.xaxis.set_ticks([])
ax.yaxis.set_ticks([])
ax.set_xlabel('$x_1$')
ax.set_ylabel('$x_2$')
plt.tight_layout()

```

21.4 山谷面：半正定

下面来聊一聊半正定矩阵情况。举个例子，矩阵 A 取值如下：

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \quad (25)$$

本 PDF 文件为作者草稿，发布目的为方便读者在移动终端学习，终稿内容以清华大学出版社纸质出版物为准。

版权归清华大学出版社所有，请勿商用，引用请注明出处。

代码及 PDF 文件下载：<https://github.com/Visualize-ML>

本书配套微课视频均发布在 B 站——生姜 DrGinger：<https://space.bilibili.com/513194466>

欢迎大家批评指教，本书专属邮箱：jiang.visualize.ml@gmail.com

容易判定 $\text{rank}(\mathbf{A}) = 1$ ；构造如下二元函数 $y = f(x_1, x_2)$ ：

$$\begin{aligned} y = f(x_1, x_2) &= \mathbf{x}^T \mathbf{A} \mathbf{x} \\ &= \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &= x_1^2 \end{aligned} \quad (26)$$

$x_1 = 0$ 时，不管 x_2 取任何值，上式为 0。

图 4 展示 $y = f(x_1, x_2)$ 对应曲面。观察该图容易发现，除了位于纵轴上点以外任意点处放置一个小球，小球都会滚动到山谷面谷底。

(26) 梯度下降向量为：

$$-\nabla f(\mathbf{x}) = -\begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} -2x_1 \\ 0 \end{bmatrix} \quad (27)$$

谷底位置对应一条直线，这条直线上每一点处梯度向量均为 0 向量，它们都是函数 $y = f(x_1, x_2)$ 极小值。

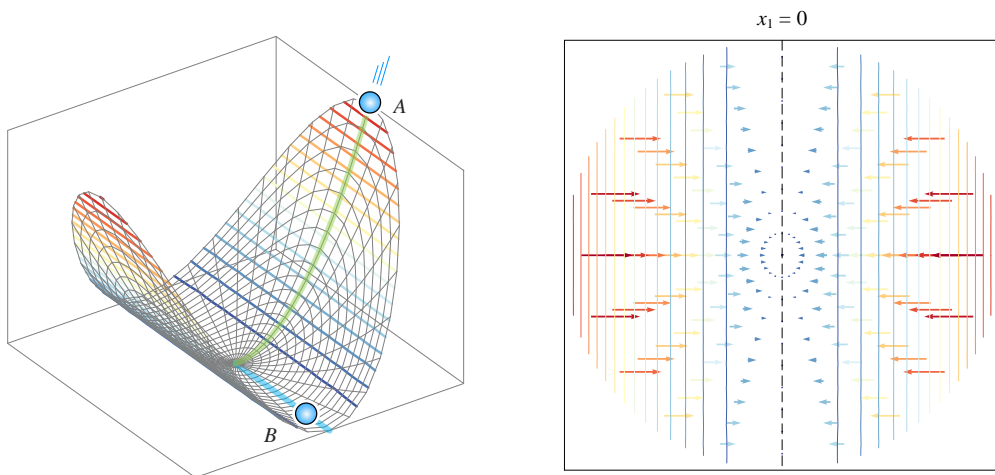


图 4. 半正定矩阵对应曲面

旋转山谷面

下式中矩阵 \mathbf{A} 也是半正定矩阵：

$$\mathbf{A} = \begin{bmatrix} 0.5 & -0.5 \\ -0.5 & 0.5 \end{bmatrix} \quad (28)$$

构造函数 $y = f(x_1, x_2)$ ：

$$\begin{aligned}
 y &= f(x_1, x_2) = \mathbf{x}^T \mathbf{A} \mathbf{x} \\
 &= \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 0.5 & -0.5 \\ -0.5 & 0.5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\
 &= 0.5x_1^2 - x_1x_2 + 0.5x_2^2
 \end{aligned} \tag{29}$$

(29) 配方得到：

$$f(x_1, x_2) = 0.5x_1^2 - x_1x_2 + 0.5x_2^2 = \frac{1}{2}(x_1 - x_2)^2 \tag{30}$$

容易发现，任何满足 $x_1 = x_2$ 的点，都会使得 $y = f(x_1, x_2)$ 为 0。

(29) 中矩阵 \mathbf{A} 特征值为 $\lambda_1 = 0$ 和 $\lambda_2 = 1$ ，对应特征向量如下：

$$\mathbf{v}_1 = \begin{bmatrix} -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix} \tag{31}$$

图 5 展示旋转山谷面。同样，小球沿图 5 中 \mathbf{v}_1 (特征值为 0 对应特征向量) 方向运动，函数值没有任何变化。这条直线上的点都是 (30) 二元函数极小值点。

(30) 梯度下降向量为：

$$-\nabla f(\mathbf{x}) = - \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} -x_1 + x_2 \\ x_1 - x_2 \end{bmatrix} \tag{32}$$

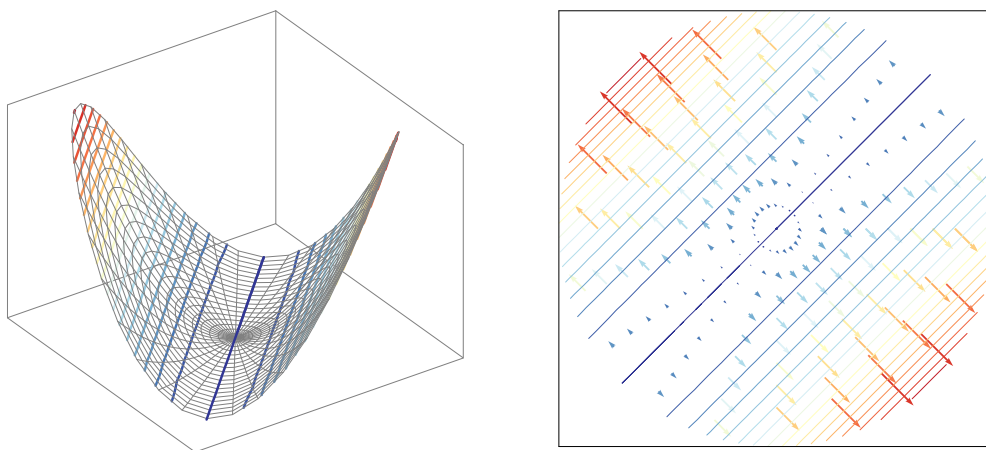


图 5. 旋转山谷面

21.5 开口朝下抛物面：负定

最简单的负定矩阵是单位矩阵取负，即 $-I$ 。 $-I$ 的特征值都为 -1 。

下面也用 2×2 矩阵讨论负定。如下 A 为负定矩阵：

$$A = \begin{bmatrix} -1 & 0 \\ 0 & -2 \end{bmatrix} \quad (33)$$

构造如下二元函数 $y = f(x_1, x_2)$ ：

$$\begin{aligned} y = f(x_1, x_2) &= \mathbf{x}^T A \mathbf{x} \\ &= \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} -1 & 0 \\ 0 & -2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &= -x_1^2 - 2x_2^2 \end{aligned} \quad (34)$$

观察上式，容易发现只有当 $x_1 = 0$ 且 $x_2 = 0$ 时， $y = f(x_1, x_2) = 0$ 。

很容易求得 A 特征值分别为 $\lambda_1 = -2$ 和 $\lambda_2 = -1$ ，对应特征向量分别为：

$$\mathbf{v}_1 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad (35)$$

图6展示负定矩阵对应曲面，容易发现 $y = f(x_1, x_2)$ 对应曲面为凹面。在曲面最大值处放置一个小球，小球处于不稳定平衡状态。受到轻微扰动后，小球沿着任意方向运动，都会下落。

(34)中 $y = f(x_1, x_2)$ 梯度下降向量为：

$$-\nabla f(\mathbf{x}) = - \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2x_1 \\ 4x_2 \end{bmatrix} \quad (36)$$

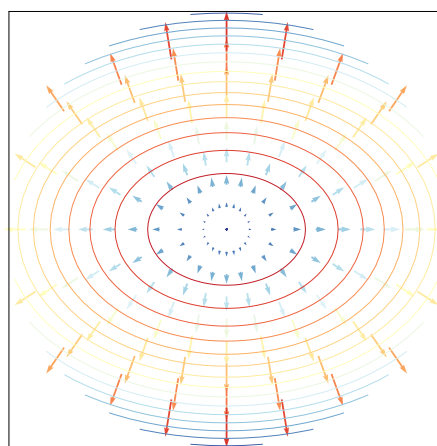
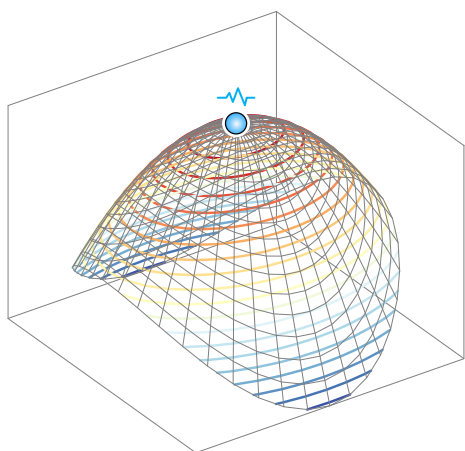


图 6. 负定矩阵对应曲面

21.6 山脊面：半负定

下面看一个半负定矩阵例子，矩阵 A 取值如下：

$$A = \begin{bmatrix} 0 & 0 \\ 0 & -1 \end{bmatrix} \quad (37)$$

构造 $y = f(x_1, x_2)$ ：

$$\begin{aligned} y = f(x_1, x_2) &= \mathbf{x}^T A \mathbf{x} \\ &= \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &= -x_2^2 \end{aligned} \quad (38)$$

$x_2 = 0$, x_1 为任意值，上式为 0。矩阵 A 的秩为 1, $\text{rank}(A) = 1$ 。

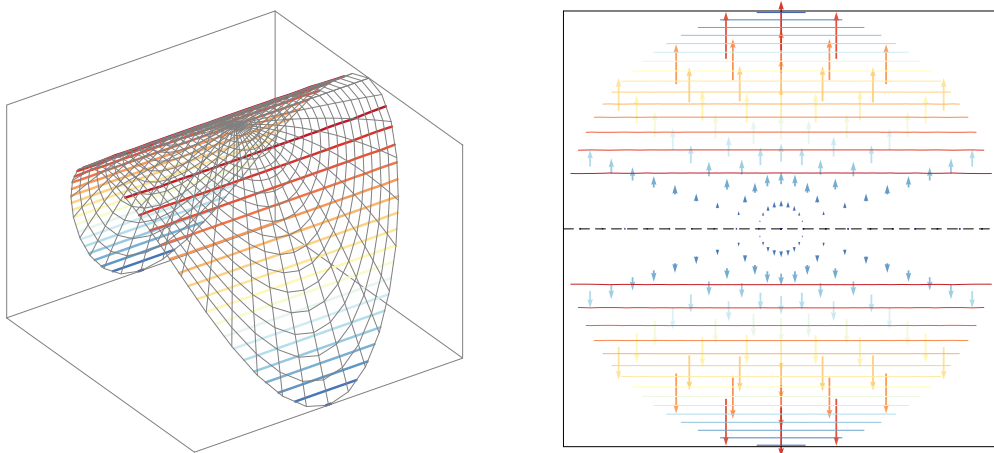


图 7. 半负定矩阵对应山脊面

图 7 展示半负定矩阵对应山脊面，发现曲面有无数个极大值。在任意极大值（山脊方向）处放置一个小球，受到扰动后，小球会沿着曲面滚下。沿着山脊方向运动，函数值没有任何变化。

(38) 梯度下降向量为：

$$-\nabla f(\mathbf{x}) = -\begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 0 \\ -2x_2 \end{bmatrix} \quad (39)$$

21.7 双曲抛物面：不定

本节最后聊一下不定矩阵情况。举个例子， A 为：

$$A = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} \quad (40)$$

构造函数 $y = f(x_1, x_2)$ ：

$$\begin{aligned} y = f(x_1, x_2) &= \mathbf{x}^T A \mathbf{x} \\ &= \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &= x_1^2 - x_2^2 \end{aligned} \quad (41)$$

求得矩阵 A 对应特征值为 $\lambda_1 = -1$ 和 $\lambda_2 = 1$ ，对应特征向量如下：

$$\mathbf{v}_1 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad (42)$$

图 8 展示 $y = f(x_1, x_2)$ 对应曲面。

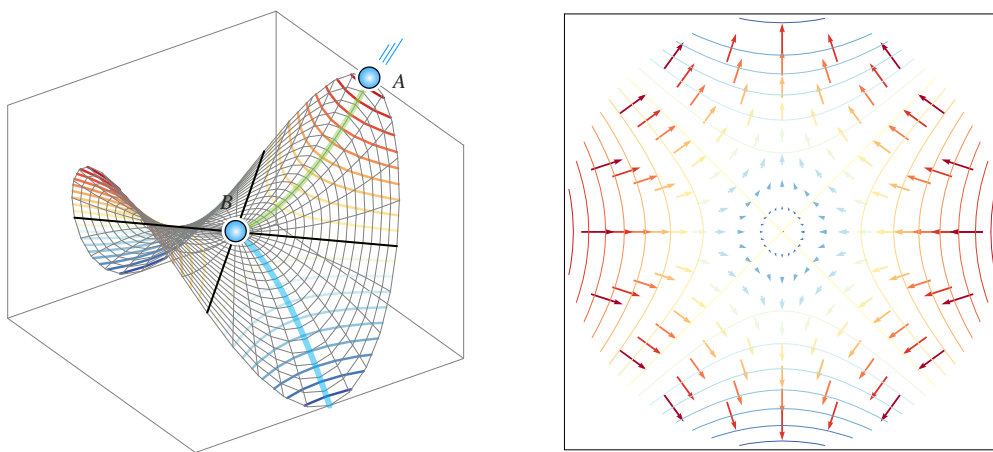


图 8. 不定矩阵对应曲面，马鞍面

当 y 不为零时，曲面对应等高线为双曲线。当 y 为零时，曲面对应等高线是两条在 x_1x_2 平面内直线 (图 8 (a) 中深色直线)，这两条直线即双曲线渐近线。

图 8 告诉我们，曲面边缘不同位置放置小球会有完全不同结果。 A 点处松手小球会向中心方向滚动， B 点处小球受到轻微扰动后会朝远离中心方向滚动。

$y = f(x_1, x_2)$ 梯度下降向量为：

$$-\nabla f(\mathbf{x}) = -\begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} -2x_1 \\ 2x_2 \end{bmatrix} \quad (43)$$

图 8 所示马鞍面中心既不是极小值点 (如图 1 曲面)，也不是极大值点 (如图 6 曲面)；图 8 中马鞍面中心点被称之为鞍点 (saddle point)。另外，沿着图 8 中深色轨道运动，小球高度没有任何变化。

旋转双曲抛物面

图 8 中马鞍面顺时针旋转 45° 得到图 9 曲面。图 9 曲面对应矩阵 A 如下：

$$A = \begin{bmatrix} 0 & -1 \\ -1 & 0 \end{bmatrix} \quad (44)$$

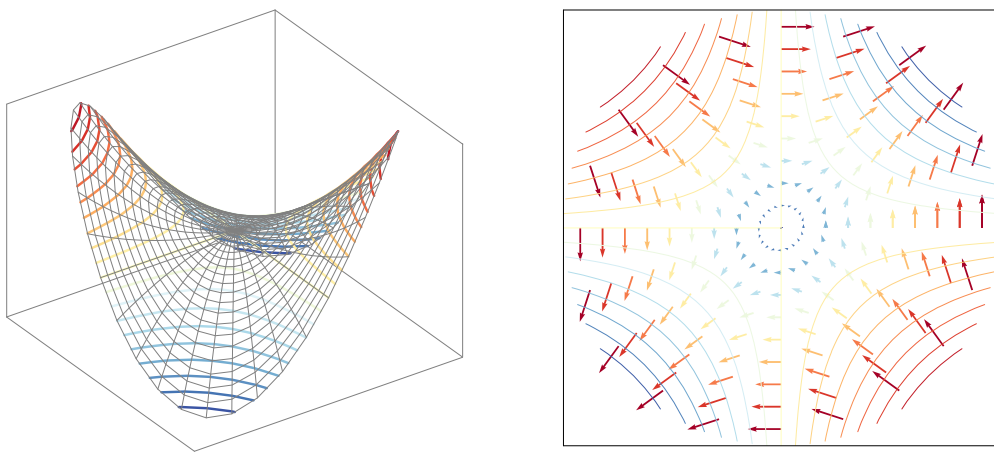


图 9. 不定矩阵对应曲面，旋转马鞍面

构造如下二元函数 $y = f(x_1, x_2)$ ：

$$\begin{aligned}
 y = f(x_1, x_2) &= \mathbf{x}^T \mathbf{A} \mathbf{x} \\
 &= \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 0 & -1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\
 &= -2x_1x_2
 \end{aligned} \tag{45}$$

在 $y = f(x_1, x_2)$ 为非零定值时，发现上式相当于反比例函数。

(45) 的梯度下降向量为：

$$-\nabla f(\mathbf{x}) = - \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2x_2 \\ 2x_1 \end{bmatrix} \tag{46}$$

21.8 多极值曲面：局部正定性

判定二元函数极值点

本系列丛书在《数学要素》一册介绍过如何判定二元函数 $y = f(x_1, x_2)$ 的极值。对于 $y = f(x_1, x_2)$ ，一阶偏导数 $f_{x_1}(x_1, x_2) = 0$ 和 $f_{x_2}(x_1, x_2) = 0$ 同时成立的点 (x_1, x_2) 为二元函数 $f(x_1, x_2)$ 的驻点；如图 10 所示，驻点可以是极大值、极小值或鞍点。

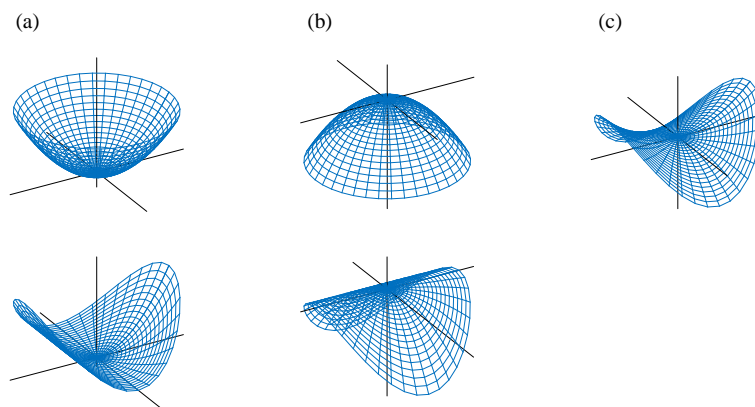


图 10. 二元函数驻点的三种情况

当时，我们聊过为了进一步判定驻点到底是极大值、极小值或是鞍点，我们需要知道二元函数 $f(x_1, x_2)$ 二阶偏导。如果 $f(x_1, x_2)$ 在 (a, b) 邻域内连续， $f(x_1, x_2)$ 二阶偏导连续。令，

$$A = f_{x_1x_1}, \quad B = f_{x_1x_2}, \quad C = f_{x_2x_2} \tag{47}$$

$f(a, b)$ 是否为极值点可以通过如下条件判断：

- a) $AC - B^2 > 0$ 存在极值，且当 $A < 0$ 有极大值， $A > 0$ 时有极小值；
- b) $AC - B^2 < 0$ 没有极值；
- c) $AC - B^2 = 0$ ，可能有极值，也可能没有极值，需要进一步讨论。

当时我们留了一个问题， $AC - B^2$ 这个表达值的含义到底是什么？本节就来回答这个问题。

(12) 中函数的黑塞矩阵 (Hessian matrix) 为：

$$\mathbf{H} = \frac{\partial(\mathbf{x}^T \mathbf{A} \mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}^T} = 2\mathbf{A} = 2 \begin{bmatrix} a & b \\ b & c \end{bmatrix} \quad (48)$$

\mathbf{A} 的行列式值为：

$$|\mathbf{A}| = ac - b^2 \quad (49)$$

相信大家已经在上式中看到和 $AC - B^2$ 一样的形式。

注意，对于二元函数， \mathbf{A} 的形状为 2×2 。 \mathbf{A} 为正定或负定时， \mathbf{A} 的两个特征值同号，因此 \mathbf{A} 的行列式值都大于 0。而 a 的正负则决定了开口方向，也就是决定了 \mathbf{A} 是正定还是负定，因此决定了极大值或极小值。

再进一步， a 实际上是 \mathbf{A} 的一阶主子式。这引出了，判定正定的另一个方法。 \mathbf{A} 正定的充分必要条件为 \mathbf{A} 的顺序主子式全大于零。

举个例子

继续采用《数学要素》一书中反复出现的多极值曲面的例子。

图 11 为曲面平面等高线。图中 \times 对应的位置为梯度向量为 $\mathbf{0}$ 。观察图中等高线不难发现，I、II、III 点为极大值点，其中 I 为最大值点。IV、V、VI 为极小值点，其中 IV 为最小值点。VII、VIII、IX 是鞍点。

图 12 给出的是二元函数的梯度向量图。极大值点处，梯度向量汇聚；极小值点处，梯度向量发散。

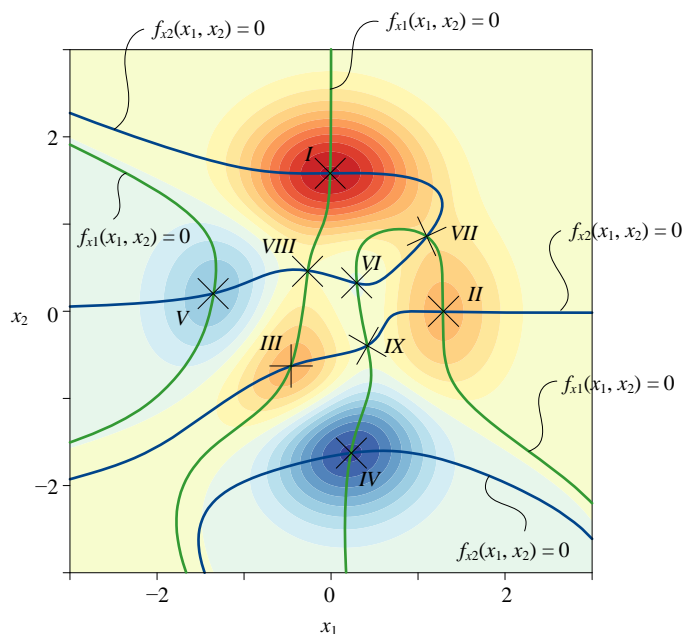


图 11. $f_{x1}(x_1, x_2) = 0$ 和 $f_{x2}(x_1, x_2) = 0$ 同时投影在 $f(x_1, x_2)$ 曲面填充等高线，来自本系列丛书《数学要素》

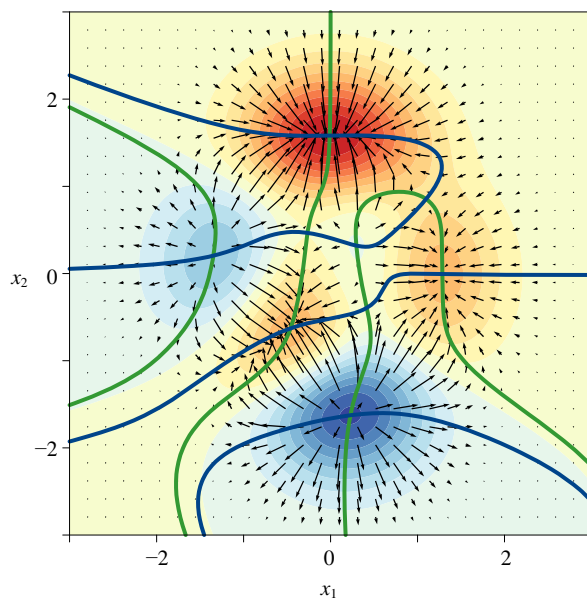


图 12. $f(x_1, x_2)$ 梯度向量图

图 13 所示为二次函数黑塞矩阵行列式值对应的等高线图，阴影圈出来的六个点对应行列式值为正，因此它们是要考察的极值点。图 13 中虚线为行列式值为 0 对应位置。梯度向量为零点没有出现在虚线位置处。

根据图 14 所示一阶主子式对应等高线，可以进一步判定极值点为极大值或极小值点。得出的结论和图 11 一致。

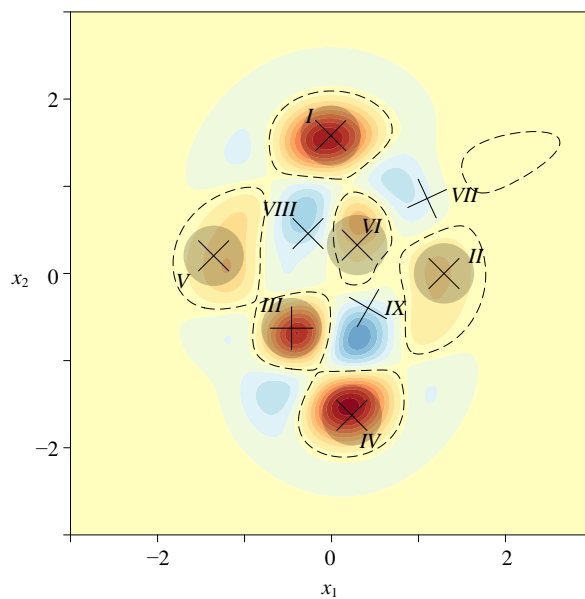


图 13. 黑塞矩阵行列式值

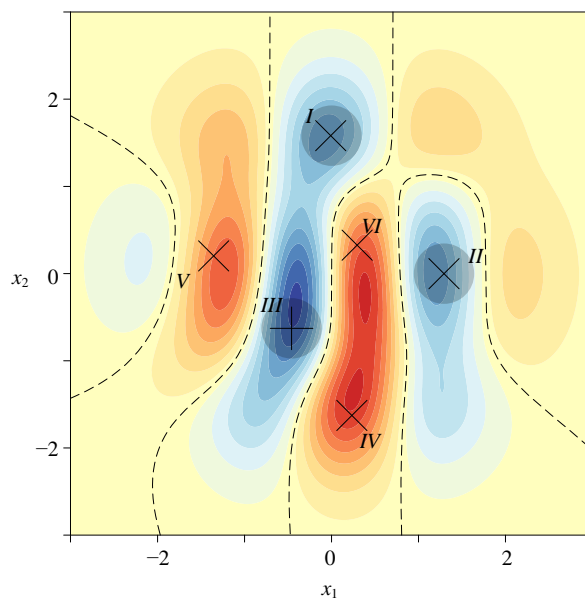


图 14. 一阶主子式正负

更一般情况

对于多元函数 $f(\mathbf{x})$ ，利用本书前文介绍的二次逼近可以写成：

$$\begin{aligned}
f(\mathbf{x}) &\approx f(\mathbf{x}_p) + \nabla f(\mathbf{x}_p)^\top (\mathbf{x} - \mathbf{x}_p) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_p)^\top \nabla^2 f(\mathbf{x}_p) (\mathbf{x} - \mathbf{x}_p) \\
&= f(\mathbf{x}_p) + \nabla f(\mathbf{x}_p)^\top \Delta \mathbf{x} + \frac{1}{2} \Delta \mathbf{x}^\top \nabla^2 f(\mathbf{x}_p) \Delta \mathbf{x} \\
&= f(\mathbf{x}_p) + \nabla f(\mathbf{x}_p)^\top \Delta \mathbf{x} + \frac{1}{2} \Delta \mathbf{x}^\top \mathbf{H} \Delta \mathbf{x}
\end{aligned} \tag{50}$$

其中 \mathbf{x}_p 为展开点。

假设 \mathbf{x}_p 处存在梯度向量，且梯度向量为 $\mathbf{0}$ 。

当 $\mathbf{x} \rightarrow \mathbf{x}_p$ 时， $\nabla f(\mathbf{x}_p)^\top \Delta \mathbf{x} \rightarrow 0$ 。但是如果在 \mathbf{x}_p 点处黑塞矩阵 \mathbf{H} 为正定， $\frac{1}{2} \Delta \mathbf{x}^\top \mathbf{H} \Delta \mathbf{x}$ 为正。这意味着：

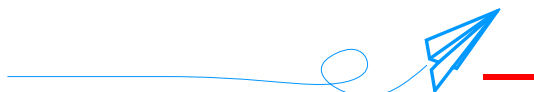
$$f(\mathbf{x}) = f(\mathbf{x}_p) + \underbrace{\nabla f(\mathbf{x}_p)^\top \Delta \mathbf{x}}_{\rightarrow 0} + \underbrace{\frac{1}{2} \Delta \mathbf{x}^\top \mathbf{H} \Delta \mathbf{x}}_{+} > f(\mathbf{x}_p) \tag{51}$$

我们称 \mathbf{x}_p 局部正定，对应 \mathbf{x}_p 为极小值点。这个判断也适用于半正定情况，不过要将上式的 $>$ 改为 \geq 。

同理，如果在 \mathbf{x}_p 点处黑塞矩阵 \mathbf{H} 为负定， $\frac{1}{2} \Delta \mathbf{x}^\top \mathbf{H} \Delta \mathbf{x}$ 为负，因此：

$$f(\mathbf{x}) = f(\mathbf{x}_p) + \underbrace{\nabla f(\mathbf{x}_p)^\top \Delta \mathbf{x}}_{\rightarrow 0} + \underbrace{\frac{1}{2} \Delta \mathbf{x}^\top \mathbf{H} \Delta \mathbf{x}}_{-} < f(\mathbf{x}_p) \tag{52}$$

我们称 \mathbf{x}_p 局部负定，对应 \mathbf{x}_p 为极大值点。如上判断也适用于半负定情况，同样将上式的 $<$ 改为 \leq 。



本章把曲面、梯度向量、正定性、极值这几个重要的概念有机的联系起来。这个例子告诉我们几何视角对于理解线性代数概念至关重要。请大家再次回顾图 15 给出的五种情况，相信大家会觉得正定性变得极容易理解。

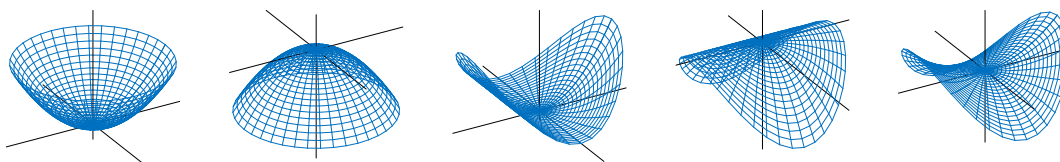


图 15. 总结本章重要内容的五副图