

# 多元函数极值问题（矩阵求导）

- 1、多元函数导数的定义
- 2、最速下降方向
- 3、牛顿法
- 4、线性回归、岭回归、Logistic回归

## 1、基本概念

设 $f: R^n \rightarrow R$ 为一个 $n$ 元一阶可微函数 $y = f(x_1, \dots, x_n)$ ，定义其梯度向量为

$$\nabla f(x) (\nabla_x f(x)) = (\partial_{x_1} f, \partial_{x_2} f, \dots, \partial_{x_n} f) \quad \text{其中 } x = (x_1, x_2, \dots, x_n)$$

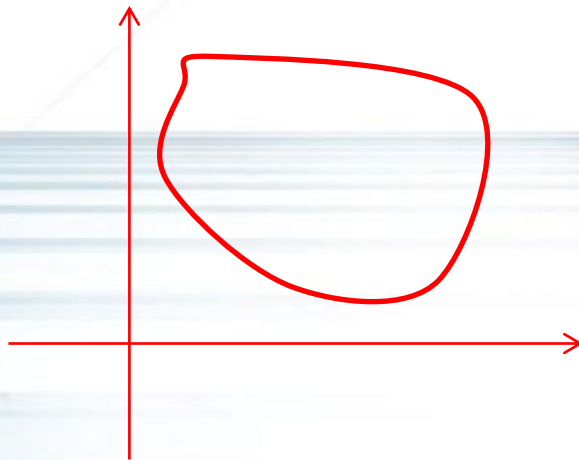
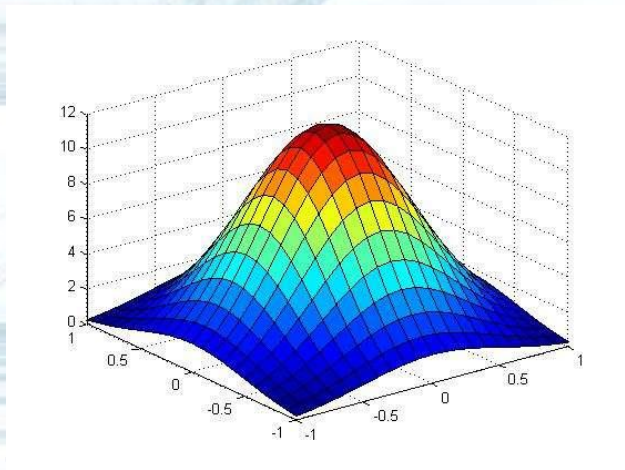
二阶导数 (Hessian矩阵)

$$Hf = [a_{ij}]_{n \times n} = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}_{n \times n}$$

其中 $a_{ij} = \partial_{x_i x_j} f$

注：由于 $\partial_{x_i x_j} f = \partial_{x_j x_i} f$ ，所以 $a_{ij} = a_{ji}$ ，所以 $Hf = Hf^T$ ，即 $Hf$ 是实对称矩阵

# 最速下降法（梯度下降）



# 最速下降法（梯度下降）

**方向导数：** 设 $\mathbf{u}$ (单位向量)是 $R^n$ 中的一个方向， $n$ 元函数 $f(x_0)$ 沿 $\mathbf{u}$ 方向的斜率，我们称之为 $f(x_0)$ 在 $\mathbf{u}$ 方向的方向导数，计算公式为：

$$\partial_{\mathbf{u}}f(\mathbf{x}_0)=\mathbf{u}\nabla_{\mathbf{x}}f(\mathbf{x}_0)^T \text{ ( 向量}\mathbf{u}\text{和梯度向量的内积 )}$$

我们希望找到使函数 $f$ 下降最快的方向 $\mathbf{u}$ ，由内积公式得：

$$\mathbf{u} \cdot \nabla_{\mathbf{x}}f(\mathbf{x}_0)=|\mathbf{u}||\nabla_{\mathbf{x}}f(\mathbf{x}_0)|\cos\theta$$

$|\mathbf{u}|=1$ ， $|\nabla_{\mathbf{x}}f(\mathbf{x}_0)|$ 与 $\mathbf{u}$ 无关。因此，当 $\theta=0$ 时，方向导数大于0，且取最大值；当 $\theta=\pi$ 时，方向导数小于0，且取最小值。

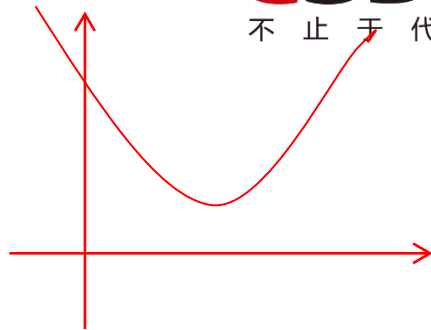
结论：对于函数 $f(x)$ 中的任意一点 $x \in R^n$ ，沿着和梯度向量一样的方向时，函数递增最快；沿着和梯度向量相反方向时，函数递减最快。因此最速下降方向为 $-\nabla f(x)$

最速下降建议点为

$$x_{t+1} = x_t - \epsilon \nabla f(x_t)$$

一元函数极值判别法：

导数分析法：若 $f'(x_0) = 0$ ，且 $f''(x_0) > (<)0$ ，则 $f(x)$ 在 $x = x_0$ 处取极小（大）值



$$\begin{aligned}\text{泰勒公式法：} f(x) &= f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(\xi)(x - x_0)^2 \\ &= f(x_0) + \frac{1}{2}f''(\xi)(x - x_0)^2 \quad \xi \in (x, x_0) \text{ or } (x_0, x)\end{aligned}$$

当 $f''(x_0) > 0$ 时，由导数的保号性可知，在 $x_0$ 很小的邻域内有 $f''(\xi) > 0$ ，因此

在 $x_0$ 很小的邻域内有 $f(x) > f(x_0)$ ，即函数 $f(x)$ 在 $x = x_0$ 处取极小值

# 多元函数的泰勒展开

$$f(\mathbf{x}) \approx f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T \nabla f(\mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^T Hf(\mathbf{x}_0) (\mathbf{x} - \mathbf{x}_0) \quad \mathbf{x}, \mathbf{x}_0 \in R^n$$

如果 $\mathbf{x}_0$ 点满足 $\nabla f(\mathbf{x}_0) = \mathbf{0}$ (零向量), 并且  $(\mathbf{x} - \mathbf{x}_0)^T Hf(\mathbf{x}_0) (\mathbf{x} - \mathbf{x}_0) > (<) 0$ 。  
那么, 函数  $f(\mathbf{x})$  在  $\mathbf{x} = \mathbf{x}_0$  处取极小(大)值。

# ▶ 正定（半正定）矩阵

$A$ 是一个 $n$ 阶对称矩阵，即 $A = A^T$  ( $a_{ij} = a_{ji}$ )

设 $n$ 维向量 $\mathbf{x} = (x_1, x_2, \dots, x_n)$ ，定义二次型（多项式）

$$f(x_1, x_2, \dots, x_n) = \mathbf{x}A\mathbf{x}^T = \sum_{i=1}^n a_{ii}x_i^2 + \sum_{i,j} a_{ij}x_ix_j$$

为 $A$ 对应的二次型（多项式）。对称方阵 $A$ 称为二次型对应的矩阵

**定义1：**若任意的 $\mathbf{x} \neq \mathbf{0}$ ，都有 $f(x_1, x_2, \dots, x_n) = \mathbf{x}A\mathbf{x}^T > (\geq) 0$ ，则称该二次型为正定（半正定）二次型，对于的矩阵 $A$ 为正定（半正定）矩阵。

**定义2：**若任意的 $\mathbf{x} \neq \mathbf{0}$ ，都有 $f(x_1, x_2, \dots, x_n) = \mathbf{x}A\mathbf{x}^T < (\leq) 0$ ，则称该二次型为负定（半负定）二次型，对于的矩阵 $A$ 为负定（半负定）矩阵。



结论：如果在 $x = x_0$ 处，有 $\nabla f(x_0) = \mathbf{0}$  (零向量)，我们称 $x_0$ 为 $f(x)$ 的驻点。

如果 $Hf(x_0)$ 正定矩阵， $f(x)$ 在 $x = x_0$ 处是一个局部极小值

如果 $Hf(x_0)$ 负定矩阵， $f(x)$ 在 $x = x_0$ 处是一个局部极大值

如果 $Hf(x_0)$ 不定矩阵， $f(x)$ 在 $x = x_0$ 处不取极值

二元函数求极值：

$$\text{step1、} \quad \begin{cases} f'_x(x, y) = 0 \\ f'_y(x, y) = 0 \end{cases} \Rightarrow \begin{cases} x = x_0 \\ y = y_0 \end{cases}$$

$$\text{Step2、} \quad A = f''_{xx}(x_0, y_0) \quad B = f''_{xy}(x_0, y_0) \quad C = f''_{yy}(x_0, y_0)$$

Step3、 当 $B^2 - AC < 0$ 且 $A < (>) 0$ ， $f(x_0, y_0)$ 是极大(小)值

当 $B^2 - AC > 0$ 时， $f(x_0, y_0)$ 不取极值

当 $B^2 - AC = 0$ 时，无法判断

## 正定矩阵的判断法：

**引理：**对称方阵一定可正交对角化，即

任意的对称矩阵A，必然存在一个正交矩阵Q，使得

$$Q^T A Q = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \lambda_n \end{bmatrix} \quad A = Q \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \lambda_n \end{bmatrix} Q^T$$

$$\begin{aligned} f(x_1, x_2, \dots, x_n) &= \mathbf{x} A \mathbf{x}^T = \mathbf{x} (Q \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) Q^T) \mathbf{x}^T \\ &= (\mathbf{x} Q) \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) (\mathbf{x} Q)^T \end{aligned}$$

记  $\mathbf{x}Q = \mathbf{y} = (y_1, y_2, \dots, y_n)$

$$f(x_1, x_2, \dots, x_n) = \mathbf{y} \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & \dots & 0 & \lambda_n \end{bmatrix} \mathbf{y}^T = \lambda_1 y_1^2 + \lambda_2 y_2^2 + \dots + \lambda_n y_n^2$$

**定理1**：二次型  $f(x_1, x_2, \dots, x_n) = \mathbf{x}A\mathbf{x}^T$  正定（半正定） $\Leftrightarrow A$  的每个特征值  $\lambda_i > (\geq) 0$

**定理 2:** 对称矩阵正定的充分必要条件是各阶顺序主子式为正数，即

$$a_{11} > 0, \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} > 0, \dots \begin{vmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{vmatrix} > 0.$$

$$A = \begin{pmatrix} 6 & 2 & -4 \\ 2 & 1 & -2 \\ -4 & -2 & 5 \end{pmatrix}, \quad \text{各阶顺序主子式为}$$

$$|6| = 6 > 0, \quad \begin{vmatrix} 6 & 2 \\ 2 & 1 \end{vmatrix} = 2 > 0, \quad \begin{vmatrix} 6 & 2 & -4 \\ 2 & 1 & -2 \\ -4 & -2 & 5 \end{vmatrix} = 2 > 0,$$

所以  $f$  是正定二次型.

一维模型： 观测到样本点 $(x_i, y_i)$  ( $i = 1, 2, \dots, m$ ) , 试图找到合适的 $a, b$  , 使得

$$f(x_i) = ax_i + b , \text{ 并且 } f(x_i) \approx y_i$$

$$\begin{aligned} (a^*, b^*) &= \operatorname{argmin}_{(a,b)} \sum_{i=1}^m (f(x_i) - y_i)^2 \quad \left( \operatorname{argmin}_{(a,b)} \|f(\mathbf{x}) - \mathbf{y}\|_{L^2} \right) \\ &= \operatorname{argmin}_{(a,b)} \sum_{i=1}^m (ax_i + b - y_i)^2 \end{aligned}$$

$$\text{令 } E(a, b) = \sum_{i=1}^m (ax_i + b - y_i)^2$$

$$\nabla E(a, b) = \left( \frac{\partial E(a, b)}{\partial a}, \frac{\partial E(a, b)}{\partial b} \right) = (0, 0)$$

$$a = \frac{\sum_{i=1}^m y_i (x_i - \bar{x})}{\sum_{i=1}^m x_i^2 - \frac{1}{n} \left( \sum_{i=1}^m x_i \right)^2}$$

$$b = \frac{1}{n} \sum_{i=1}^m (y_i - ax_i)$$

$$\text{其中 } \bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$$

高维模型：观测到样本点 $(x_i, y_i)$  ( $i = 1, 2, \dots, m$ )。其中，

$x_i \in R^n, y_i \in R$  试图找到合适 $a \in R^n, b \in R$ ，使得

$$f(x_i) = ax_i + b, \text{ 并且 } f(x_i) \approx y_i$$

注： $f(x_i) = ax_i + b = a_1x_{i1} + a_2x_{i2} + \dots + a_nx_{in} + b$

$$\text{令 } A = \begin{bmatrix} x_{11} & \dots & x_{1n} & 1 \\ x_{21} & \dots & x_{2n} & 1 \\ \dots & \dots & \dots & \vdots \\ x_{m1} & \dots & x_{mn} & 1 \end{bmatrix} \quad \omega = (a, b)^T \in R^{n+1} \quad y = (y_1, y_2, \dots, y_m)^T$$

$$\sum_{i=1}^m (f(x_i) - y_i)^2 = \|A\omega - y\|_2^2 = (y - A\omega)^T (y - A\omega)$$



$$\boldsymbol{\omega}^* = \underset{\boldsymbol{\omega}}{\operatorname{argmin}} (\boldsymbol{y} - A\boldsymbol{\omega})^T (\boldsymbol{y} - A\boldsymbol{\omega})$$

$$\text{令 } E(\boldsymbol{\omega}) = (\boldsymbol{y} - A\boldsymbol{\omega})^T (\boldsymbol{y} - A\boldsymbol{\omega})$$

$$\nabla E(\boldsymbol{\omega}) = 2A^T(A\boldsymbol{\omega} - \boldsymbol{y}) = \boldsymbol{\theta}$$

$$A^T A \boldsymbol{\omega} - A^T \boldsymbol{y} = \boldsymbol{\theta}$$

$$\boldsymbol{\omega} = (A^T A)^{-1} A^T \boldsymbol{y}$$

## 用SVD处理岭回归：

如果矩阵 $A^T A$ 不可逆，则用岭回归代替线性回归。

线性回归： $\hat{\mathbf{y}} = A\boldsymbol{\omega} = A(A^T A)^{-1}A^T \mathbf{y}$

岭回归： $\hat{\mathbf{y}} = A(A^T A + \lambda I)^{-1}A^T \mathbf{y}$

性质：损失无偏性，增加稳定性，从而得到较高的计算精度。

对矩阵 $A$ 进行SVD,  $A = U\Lambda V^T$  其中 $U, V$ 是正交矩阵

岭回归:  $\hat{y} = A(A^T A + \lambda I_n)^{-1} A^T \mathbf{y}$

$$= U\Lambda V^T (V\Lambda^T U^T U\Lambda V^T + \lambda V I_n V^T)^{-1} V\Lambda^T U^T \mathbf{y}$$

$$= U\Lambda (\Lambda^T \Lambda + \lambda I_n)^{-1} \Lambda^T U^T \mathbf{y}$$

$$= \sum_{i=1}^r U_i \frac{\sigma_i^2}{\sigma_i^2 + \lambda} U_i \mathbf{y}$$

## Logistic回归:

模型： $\mathbf{x} = (x_1, x_2, \dots, x_n)$ ，考虑变量的线性组合函数 $g(\mathbf{x}) = \omega_0 + \omega_1 x_1 + \dots + \omega_n x_n$   
( LR分类器的权重 )。概率模型满足sigmoid函数

$$P(y = 1|\mathbf{x}) = f(\mathbf{x}) = \frac{1}{1 + e^{-g(\mathbf{x})}}$$

$$P(y = 0|\mathbf{x}) = f(\mathbf{x}) = 1 - \frac{1}{1 + e^{-g(\mathbf{x})}}$$

观测到m个样本数据 $(\mathbf{x}^1, y_1), (\mathbf{x}^2, y_2), \dots, (\mathbf{x}^m, y_m)$ ，极大似然函数为：

$$L(\boldsymbol{\omega}) = \prod_{i=1}^m \left( f(\mathbf{x}^i) \right)^{y_i} \left( 1 - f(\mathbf{x}^i) \right)^{1-y_i}$$

左右取对数得：

$$\ln L(\omega) = \sum_{i=1}^m y_i \ln(f(\mathbf{x}^i)) + (1 - y_i) \ln(1 - f(\mathbf{x}^i))$$

计算梯度，并令 $\nabla \ln L(\omega) = \theta$ 。得到方程组

$$\frac{\partial \ln L(\omega)}{\partial \omega_k} = \sum_{i=1}^m x_k^i [y_i - f(\mathbf{x}^i)] = 0$$

# THANK YOU



AI100