

1 基础知识

2 凸函数的定义与性质

3 保凸的运算

4 凸函数的推广

梯度

定义 (梯度)

给定函数 $f: \mathbb{R}^n \rightarrow \mathbb{R}$, 且 f 在点 x 的一个邻域内有意义, 若存在向量 $g \in \mathbb{R}^n$ 满足

$$\lim_{p \rightarrow 0} \frac{f(x+p) - f(x) - g^T p}{\|p\|} = 0,$$

其中 $\|\cdot\|$ 是任意的向量范数, 就称 f 在点 x 处可微 (或Fréchet 可微). 此时 g 称为 f 在点 x 处的梯度, 记作 $\nabla f(x)$. 如果对区域 D 上的每一个点 x 都有 $\nabla f(x)$ 存在, 则称 f 在 D 上可微.

若 f 在点 x 处的梯度存在, 在定义式中令 $p = \varepsilon e_i$, e_i 是第 i 个分量为 1 的单位向量, 可知 $\nabla f(x)$ 的第 i 个分量为 $\frac{\partial f(x)}{\partial x_i}$. 因此,

$$\nabla f(x) = \left[\frac{\partial f(x)}{\partial x_1}, \frac{\partial f(x)}{\partial x_2}, \dots, \frac{\partial f(x)}{\partial x_n} \right]^T.$$

海瑟矩阵

定义 (海瑟矩阵)

如果函数 $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ 在点 x 处的二阶偏导数 $\frac{\partial^2 f(x)}{\partial x_i \partial x_j}$ $i, j = 1, 2, \dots, n$ 都存在, 则

$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_3} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_3} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} & \frac{\partial^2 f(x)}{\partial x_n \partial x_3} & \cdots & \frac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix}$$

称为 f 在点 x 处的海瑟矩阵.

当 $\nabla^2 f(x)$ 在区域 D 上的每个点 x 处都存在时, 称 f 在 D 上二阶可微.
若 $\nabla^2 f(x)$ 在 D 上还连续, 则称 f 在 D 上二阶连续可微, 可以证明此时海瑟矩阵是一个对称矩阵.

矩阵变量函数的导数

多元函数梯度的定义可以推广到变量是矩阵的情形. 对于以 $m \times n$ 矩阵 X 为自变量的函数 $f(X)$, 若存在矩阵 $G \in \mathbb{R}^{m \times n}$ 满足

$$\lim_{V \rightarrow 0} \frac{f(X+V) - f(X) - \langle G, V \rangle}{\|V\|} = 0,$$

其中 $\|\cdot\|$ 是任意矩阵范数, 就称矩阵变量函数 f 在 X 处 **Fréchet** 可微, 称 G 为 f 在 **Fréchet** 可微意义下的梯度. 类似于向量情形, 矩阵变量函数 $f(X)$ 的梯度可以用其偏导数表示为

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_{11}} & \frac{\partial f}{\partial x_{12}} & \cdots & \frac{\partial f}{\partial x_{1n}} \\ \frac{\partial f}{\partial x_{21}} & \frac{\partial f}{\partial x_{22}} & \cdots & \frac{\partial f}{\partial x_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_{m1}} & \frac{\partial f}{\partial x_{m2}} & \cdots & \frac{\partial f}{\partial x_{mn}} \end{bmatrix}.$$

其中 $\frac{\partial f}{\partial x_{ij}}$ 表示 f 关于 x_{ij} 的偏导数.

矩阵变量函数的导数

在实际应用中，矩阵Fréchet可微的定义和使用往往比较繁琐，为此我们需要介绍另一种定义——Gâteaux可微。

定义 (Gâteaux 可微)

设 $f(X)$ 为矩阵变量函数，如果对任意方向 $V \in \mathbb{R}^{m \times n}$ ，存在矩阵 $G \in \mathbb{R}^{m \times n}$ 满足

$$\lim_{t \rightarrow 0} \frac{f(X + tV) - f(X) - t \langle G, V \rangle}{t} = 0,$$

则称 f 关于 X 是Gâteaux可微的。满足上式的 G 称为 f 在 X 处在Gâteaux可微意义下的梯度。

可以证明，当 f 是Fréchet可微函数时， f 也是Gâteaux可微的，且这两种意义下的梯度相等。

矩阵变函数的导数

- 线性函数: $f(X) = \text{tr}(AX^T B)$, 其中 $A \in \mathbb{R}^{p \times n}$, $B \in \mathbb{R}^{m \times p}$, $X \in \mathbb{R}^{m \times n}$
对任意方向 $V \in \mathbb{R}^{m \times n}$ 以及 $t \in \mathbb{R}$, 有

$$\begin{aligned}\lim_{t \rightarrow 0} \frac{f(X + tV) - f(X)}{t} &= \lim_{t \rightarrow 0} \frac{\text{tr}(A(X + tV)^T B) - \text{tr}(AX^T B)}{t} \\ &= \text{tr}(AV^T B) = \langle BA, V \rangle.\end{aligned}$$

因此, $\nabla f(X) = BA$.

- 二次函数: $f(X, Y) = \frac{1}{2} \|XY - A\|_F^2$, 其中 $(X, Y) \in \mathbb{R}^{m \times p} \times \mathbb{R}^{p \times n}$
对变量 Y , 取任意方向 V 以及充分小的 $t \in \mathbb{R}$, 有

$$\begin{aligned}f(X, Y + tV) - f(X, Y) &= \frac{1}{2} \|X(Y + tV) - A\|_F^2 - \frac{1}{2} \|XY - A\|_F^2 \\ &= \langle tXV, XY - A \rangle + \frac{1}{2} t^2 \|XV\|_F^2 \\ &= t \langle V, X^T (XY - A) \rangle + \mathcal{O}(t^2).\end{aligned}$$

由定义可知 $\frac{\partial f}{\partial Y} = X^T (XY - A)$.

对变量 X , 同理可得 $\frac{\partial f}{\partial X} = (XY - A)Y^T$.

矩阵变量函数的导数

- *ln-det* 函数: $f(X) = \ln(\det(X))$, $X \in \mathcal{S}_{++}^n$, 给定 $X \succ 0$, 对任意方向 $V \in \mathcal{S}^n$ 以及 $t \in \mathbb{R}$, 我们有

$$\begin{aligned} & f(X + tV) - f(X) \\ &= \ln(\det(X + tV)) - \ln(\det(X)) \\ &= \ln(\det(X^{1/2}(I + tX^{-1/2}VX^{-1/2})X^{1/2})) - \ln(\det(X)) \\ &= \ln(\det(I + tX^{-1/2}VX^{-1/2})). \end{aligned}$$

由于 $X^{-1/2}VX^{-1/2}$ 是对称矩阵, 所以它可以正交对角化, 不妨设它的特征值为 $\lambda_1, \lambda_2, \dots, \lambda_n$, 则

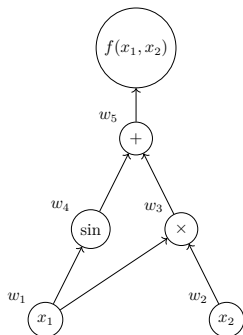
$$\begin{aligned} \ln(\det(I + tX^{-1/2}VX^{-1/2})) &= \ln \prod_{i=1}^n (1 + t\lambda_i) \\ &= \sum_{i=1}^n \ln(1 + t\lambda_i) = \sum_{i=1}^n t\lambda_i + \mathcal{O}(t^2) \\ &= t \operatorname{tr}(X^{-1/2}VX^{-1/2}) + \mathcal{O}(t^2) \\ &= t \langle (X^{-1})^T, V \rangle + \mathcal{O}(t^2). \end{aligned}$$

因此, 我们得到结论 $\nabla f(X) = (X^{-1})^T$.

自动微分

自动微分是使用计算机计算导数的算法。在神经网络中，损失函数 $f(x)$ 是由很多个简单函数复合而成的函数，根据复合函数的链式法则，可以通过每个简单函数的导数的乘积来计算对于各层变量的导数。

我们先考虑一个简单的例子 $f(x_1, x_2) = x_1 x_2 + \sin x_1$ 。计算该函数的过程可以用下图来表示。



$$w_1 = x_1$$

$$w_2 = x_2$$

$$w_3 = w_1 w_2$$

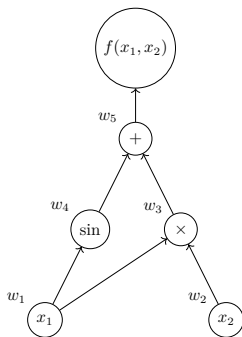
$$w_4 = \sin w_1$$

$$w_5 = w_3 + w_4$$

Figure: 函数 $f(x_1, x_2)$ 的计算过程

自动微分

在计算图中, w_1 和 w_2 为自变量, w_3 和 w_4 为中间变量, w_5 代表最终的目标函数值. 容易看出, 函数 f 计算过程中涉及的所有变量和它们之间的依赖关系构成了一个有向图: 每个变量 w_i 代表着图中的一个节点, 变量的依赖关系为该图的边. 如果有一条从节点 w_i 指向 w_j 的边, 我们称 w_i 为 w_j 的父节点, w_j 为 w_i 的子节点. 一个节点的值由其所有的父节点的值确定. 则称从父节点的值推子节点值的计算流为前向传播.



自动微分

利用计算导数的链式法则，我们可以依次计算

$$\frac{\partial f}{\partial w_5} = 1,$$

$$\frac{\partial f}{\partial w_4} = \frac{\partial f}{\partial w_5} \frac{\partial w_5}{\partial w_4} = 1,$$

$$\frac{\partial f}{\partial w_3} = \frac{\partial f}{\partial w_5} \frac{\partial w_5}{\partial w_3} = 1,$$

$$\frac{\partial f}{\partial w_2} = \frac{\partial f}{\partial w_3} \frac{\partial w_3}{\partial w_2} = w_1 = x_1,$$

$$\frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial w_3} \frac{\partial w_3}{\partial w_1} + \frac{\partial f}{\partial w_4} \frac{\partial w_4}{\partial w_1} = w_2 + \cos w_1 = \cos x_1 + x_2.$$

通过这种方式，就求得了导数

$$\frac{\partial f}{\partial x_1} = \cos x_1 + x_2, \quad \frac{\partial f}{\partial x_2} = x_1.$$

自动微分

自动微分有两种方式：前向模式和后向模式。在前向模式中，根据计算图，可以依次计算每个中间变量的取值及其对父变量的偏导数值。通过链式法则，可以复合得到每个中间变量对自变量的导数值。直至传播到最后一个子节点时，就得到了最终的目标函数值以及目标函数关于自变量的梯度值。

不同于前向模式，后向模式的节点求值和导数计算不是同时进行的。它是先利用前向模式计算各个节点的值，然后再根据计算图逆向计算对函数 f 关于各个中间变量的偏导数。

$$\frac{\partial f}{\partial w_i} = \sum_{w_j \text{ 是 } w_i \text{ 的子节点}} \frac{\partial f}{\partial w_j} \frac{\partial w_j}{\partial w_i}.$$

对于前向模式而言，后向模式的梯度的计算复杂度更低。具体地，后向模式的梯度计算代价至多为函数值计算代价的5倍，但是前向模式的计算代价可能多达函数值计算代价的 n （ n 为自变量维数）倍。因此对于神经网络中的优化问题，其自动微分采用的是后向模式。

广义实值函数与适当函数

定义 (广义实值函数)

令 $\overline{\mathbb{R}} \stackrel{\text{def}}{=} \mathbb{R} \cup \{\pm\infty\}$ 为广义实数空间, 则映射 $f: \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ 称为广义实值函数.

和数学分析一样, 我们规定

$$-\infty < a < +\infty, \quad \forall a \in \mathbb{R}$$

$$(+\infty) + (+\infty) = +\infty, \quad +\infty + a = +\infty, \quad \forall a \in \mathbb{R}.$$

定义 (适当函数)

给定广义实值函数 f 和非空集合 \mathcal{X} . 如果存在 $x \in \mathcal{X}$ 使得 $f(x) < +\infty$, 并且对任意的 $x \in \mathcal{X}$, 都有 $f(x) > -\infty$, 那么称函数 f 关于集合 \mathcal{X} 是适当的.

概括来说, 适当函数 f 的特点是“至少有一处取值不为正无穷”, 以及“处处取值不为负无穷”.

下水平集与上方图

定义 (α -下水平集)

对于广义实值函数 $f: \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$,

$$C_\alpha = \{x \mid f(x) \leq \alpha\}$$

称为 f 的 α -下水平集.

定义 (上方图)

对于广义实值函数 $f: \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$,

$$\text{epi } f = \{(x, t) \in \mathbb{R}^{n+1} \mid f(x) \leq t\}$$

称为 f 的上方图.

闭函数

定义 (闭函数)

设 $f: \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ 为广义实值函数, 若 $\text{epi } f$ 为闭集, 则称 f 为闭函数.

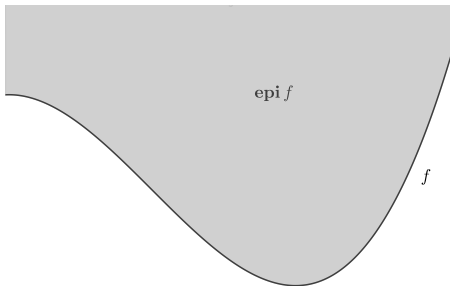


Figure: 函数 f 和其上方图 $\text{epi } f$

下半连续函数

定义 (下半连续函数)

设广义实值函数 $f: \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$, 若对任意的 $x \in \mathbb{R}^n$, 有

$$\liminf_{y \rightarrow x} f(y) \geq f(x),$$

则 $f(x)$ 为下半连续函数.

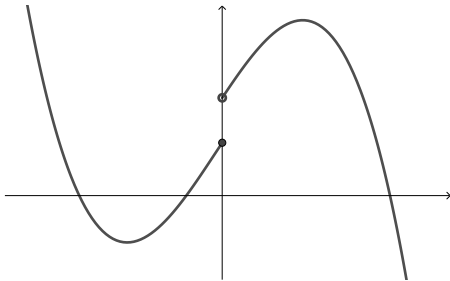


Figure: 下半连续函数 $f(x)$

闭函数与下半连续函数

虽然表面上看这两种函数的定义方式截然不同，但闭函数和下半连续函数是等价的。

定理

设广义实值函数 $f: \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ ，则以下命题等价：

- 1 $f(x)$ 的任意 α -下水平集都是闭集；
- 2 $f(x)$ 是下半连续的；
- 3 $f(x)$ 是闭函数。

闭函数与下半连续函数

闭（下半连续）函数间的简单运算会保持原有性质：

- 加法：若 f 与 g 均为适当的闭（下半连续）函数，并且 $\text{dom } f \cap \text{dom } g \neq \emptyset$ ，则 $f + g$ 也是闭（下半连续）函数。其中适当函数的条件是为了避免出现未定式 $(-\infty) + (+\infty)$ 的情况；
- 仿射映射的复合：若 f 为闭（下半连续）函数，则 $f(Ax + b)$ 也为闭（下半连续）函数；
- 取上确界：若每一个函数 f_α 均为闭（下半连续）函数，则 $\sup_\alpha f_\alpha(x)$ 也为闭（下半连续）函数。

提纲

1 基础知识

2 凸函数的定义与性质

3 保凸的运算

4 凸函数的推广

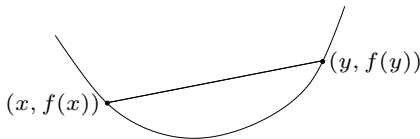
凸函数的定义

定义 (凸函数)

$f: \mathbb{R}^n \rightarrow \mathbb{R}$ 为适当函数, 如果 $\text{dom } f$ 是凸集, 且

$$f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y)$$

对所有 $x, y \in \text{dom } f$, $0 \leq \theta \leq 1$ 都成立, 则称 f 是凸函数



- 若 f 是凸函数, 则 $-f$ 是凹函数
- 若对所有 $x, y \in \text{dom } f$, $x \neq y$, $0 < \theta < 1$, 有

$$f(\theta x + (1 - \theta)y) < \theta f(x) + (1 - \theta)f(y)$$

则称 f 是严格凸函数

一元凸函数的例子

凸函数:

- 仿射函数: 对任意 $a, b \in \mathbb{R}$, $ax + b$ 是 \mathbb{R} 上的凸函数
- 指数函数: 对任意 $a \in \mathbb{R}$, e^{ax} 是 \mathbb{R} 上的凸函数
- 幂函数: 对 $\alpha \geq 1$ 或 $\alpha \leq 0$, x^α 是 \mathbb{R}_{++} 上的凸函数
- 绝对值的幂: 对 $p \geq 1$, $|x|^p$ 是 \mathbb{R} 上的凸函数
- 负熵: $x \log x$ 是 \mathbb{R}_{++} 上的凸函数

凹函数:

- 仿射函数: 对任意 $a, b \in \mathbb{R}$, $ax + b$ 是 \mathbb{R} 上的凹函数
- 幂函数: 对 $0 \leq \alpha \leq 1$, x^α 是 \mathbb{R}_{++} 上的凹函数
- 对数函数: $\log x$ 是 \mathbb{R}_{++} 上的凹函数

多元凸函数的例子

所有的仿射函数既是凸函数，又是凹函数。所有的范数都是凸函数。

欧氏空间 \mathbb{R}^n 中的例子

- 仿射函数: $f(x) = a^T x + b$
- 范数: $\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$ ($p \geq 1$) ; 特别地, $\|x\|_\infty = \max_k |x_k|$

矩阵空间 $\mathbb{R}^{m \times n}$ 中的例子

- 仿射函数:

$$f(X) = \text{tr}(A^T X) + b = \sum_{i=1}^m \sum_{j=1}^n A_{ij} X_{ij} + b$$

- 谱范数:

$$f(X) = \|X\|_2 = \sigma_{\max}(X) = (\lambda_{\max}(X^T X))^{1/2}$$

强凸函数

- 定义1: 若存在常数 $m > 0$, 使得

$$g(x) = f(x) - \frac{m}{2}\|x\|^2$$

为凸函数, 则称 $f(x)$ 为**强凸函数**, 其中 m 为**强凸参数**. 为了方便我们也称 $f(x)$ 为 m -强凸函数.

- 定义2: 若存在常数 $m > 0$, 使得对任意 $x, y \in \text{dom}f$ 以及 $\theta \in (0, 1)$, 有

$$f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y) - \frac{m}{2}\theta(1 - \theta)\|x - y\|^2,$$

则称 $f(x)$ 为强凸函数, 其中 m 为强凸参数.

- 设 f 为强凸函数且存在最小值, 则 f 的最小值点唯一.

凸函数判定定理

凸函数的一个最基本的判定方式是：先将其限制在任意直线上，然后判断对应的一维函数是否是凸的。

定理

$f: \mathbb{R}^n \rightarrow \mathbb{R}$ 是凸函数，当且仅当对每个 $x \in \text{dom } f, v \in \mathbb{R}^n$ ，函数 $g: \mathbb{R} \rightarrow \mathbb{R}$,

$$g(t) = f(x + tv), \quad \text{dom } g = \{t | x + tv \in \text{dom } f\}$$

是关于 t 的凸函数

例： $f(X) = -\log \det X$ 是凸函数，其中 $\text{dom } f = \mathbb{S}_{++}^n$ 。任取 $X \succ 0$ 以及方向 $V \in \mathbb{S}^n$ ，将 f 限制在直线 $X + tV$ (t 满足 $X + tV \succ 0$) 上，那么

$$\begin{aligned} g(t) &= -\log \det(X + tV) = -\log \det X - \log \det(I + tX^{-1/2}VX^{-1/2}) \\ &= -\log \det X - \sum_{i=1}^n \log(1 + t\lambda_i) \end{aligned}$$

其中 λ_i 是 $X^{-1/2}VX^{-1/2}$ 第 i 个特征值

对每个 $X \succ 0$ 以及方向 V ， g 关于 t 是凸的，因此 f 是凸的

凸函数判定定理

Proof.

必要性: 设 $f(x)$ 是凸函数, 要证 $g(t) = f(x + tv)$ 是凸函数. 先说明 $\text{dom}g$ 是凸集. 对任意的 $t_1, t_2 \in \text{dom}g$ 以及 $\theta \in (0, 1)$,

$$x + t_1v \in \text{dom}f, x + t_2v \in \text{dom}f$$

由 $\text{dom}f$ 是凸集可知 $x + (\theta t_1 + (1 - \theta)t_2)v \in \text{dom}f$,
这说明 $\theta t_1 + (1 - \theta)t_2 \in \text{dom}g$, 即 $\text{dom}g$ 是凸集. 此外, 我们有

$$\begin{aligned} g(\theta t_1 + (1 - \theta)t_2) &= f(x + (\theta t_1 + (1 - \theta)t_2)v) \\ &= f(\theta(x + t_1v) + (1 - \theta)(x + t_2v)) \\ &\leq \theta f(x + t_1v) + (1 - \theta)f(x + t_2v) \\ &= \theta g(t_1) + (1 - \theta)g(t_2). \end{aligned}$$

结合以上两点得到函数 $g(t)$ 是凸函数.

凸函数判定定理

Proof.

充分性:先说明 $\text{dom}f$ 是凸集, 取 $v = y - x$, 以及 $t_1 = 0, t_2 = 1$, 由 $\text{dom}g$ 是凸集可知 $\theta \cdot 0 + (1 - \theta) \cdot 1 \in \text{dom}g$, 即 $\theta x + (1 - \theta)y \in \text{dom}f$, 这说明 $\text{dom}f$ 是凸集. 再根据 $g(t) = f(x + tv)$ 的凸性, 我们有

$$\begin{aligned}g(1 - \theta) &= g(\theta t_1 + (1 - \theta)t_2) \\&\leq \theta g(t_1) + (1 - \theta)g(t_2) \\&= \theta g(0) + (1 - \theta)g(1) \\&= \theta f(x) + (1 - \theta)f(y).\end{aligned}$$

而等式左边有

$$g(1 - \theta) = f(x + (1 - \theta)(y - x)) = f(\theta x + (1 - \theta)y),$$

这说明 $f(x)$ 是凸函数.

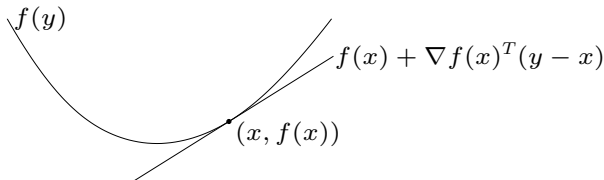


一阶条件

定理

一阶条件：对于定义在凸集上的可微函数 f ， f 是凸函数当且仅当

$$f(y) \geq f(x) + \nabla f(x)^T(y - x) \quad \forall x, y \in \text{dom } f$$



几何直观： f 的一阶逼近始终在 f 的图像下方

一阶条件

Proof.

必要性：设 f 是凸函数，则对于任意的 $x, y \in \text{dom} f$ 以及 $t \in (0, 1)$ ，有

$$tf(y) + (1 - t)f(x) \geq f(x + t(y - x)).$$

将上式移项，两边同时除以 t ，注意 $t > 0$ ，则

$$f(y) - f(x) \geq \frac{f(x + t(y - x)) - f(x)}{t}.$$

令 $t \rightarrow 0$ ，由极限保号性可得

$$f(y) - f(x) \geq \lim_{t \rightarrow 0} \frac{f(x + t(y - x)) - f(x)}{t} = \nabla f(x)^T (y - x).$$

这里最后一个等式成立是由于方向导数的性质。

一阶条件

Proof.

充分性：对任意的 $x, y \in \text{dom}f$ 以及任意的 $t \in (0, 1)$ ，定义 $z = tx + (1 - t)y$ ，应用两次一阶条件我们有

$$f(x) \geq f(z) + \nabla f(z)^T(x - z),$$

$$f(y) \geq f(z) + \nabla f(z)^T(y - z).$$

将上述第一个不等式两边同时乘 t ，第二个不等式两边同时乘 $1 - t$ ，相加得

$$tf(x) + (1 - t)f(y) \geq f(z) + 0.$$

这正是凸函数的定义，因此充分性成立。 □

梯度单调性

定理

设 f 为可微函数，则 f 为凸函数当且仅当 $\text{dom}f$ 为凸集且 ∇f 为单调映射，

$$(\nabla f(x) - \nabla f(y))^T(x - y) \geq 0, \quad \forall x, y \in \text{dom}f.$$

Proof.

必要性：若 f 可微且为凸函数，根据一阶条件，我们有

$$f(y) \geq f(x) + \nabla f(x)^T(y - x),$$

$$f(x) \geq f(y) + \nabla f(y)^T(x - y).$$

将两式不等号左右两边相加即可得到结论。

梯度单调性

Proof.

充分性：若 ∇f 为单调映射，构造一元辅助函数

$$g(t) = f(x + t(y - x)), \quad g'(t) = \nabla f(x + t(y - x))^T (y - x)$$

由 ∇f 的单调性可知 $g'(t) \geq g'(0), \forall t \geq 0$. 因此

$$\begin{aligned} f(y) = g(1) &= g(0) + \int_0^1 g'(t) dt \\ &\geq g(0) + g'(0) = f(x) + \nabla f(x)^T (y - x). \end{aligned}$$

□

定理

函数 $f(x)$ 为凸函数当且仅当其上方图 $\text{epi}f$ 是凸集.

Proof.

必要性: 若 f 为凸函数, 则对任意 $(x_1, y_1), (x_2, y_2) \in \text{epi}f, t \in [0, 1]$,

$$ty_1 + (1 - t)y_2 \geq tf(x_1) + (1 - t)f(x_2) \geq f(tx_1 + (1 - t)x_2),$$

故 $(tx_1 + (1 - t)x_2, ty_1 + (1 - t)y_2) \in \text{epi}f, t \in [0, 1]$.

充分性: 若 $\text{epi}f$ 是凸集, 则对任意 $x_1, x_2 \in \text{dom } f, t \in [0, 1]$,

$$(tx_1 + (1 - t)x_2, tf(x_1) + (1 - t)f(x_2)) \in \text{epi}f \Rightarrow \\ f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2).$$



二阶条件

定理

二阶条件：设 f 为定义在凸集上的二阶连续可微函数

- f 是凸函数当且仅当

$$\nabla^2 f(x) \succeq 0 \quad \forall x \in \text{dom } f$$

- 如果 $\nabla^2 f(x) \succ 0 \quad \forall x \in \text{dom } f$ ，则 f 是严格凸函数

例：二次函数 $f(x) = (1/2)x^T P x + q^T x + r$ (其中 $P \in \mathbb{S}^n$)

$$\nabla f(x) = P x + q, \quad \nabla^2 f(x) = P$$

f 是凸函数当且仅当 $P \succeq 0$

二阶条件

Proof.

必要性：反设 $f(x)$ 在点 x 处的海瑟矩阵 $\nabla^2 f(x) \not\geq 0$ ，即存在非零向量 $v \in \mathbb{R}^n$ 使得 $v^T \nabla^2 f(x) v < 0$ 。根据佩亚诺（Peano）余项的泰勒展开，

$$f(x + tv) = f(x) + t \nabla f(x)^T v + \frac{t^2}{2} v^T \nabla^2 f(x) v + o(t^2).$$

移项后等式两边同时除以 t^2 ，

$$\frac{f(x + tv) - f(x) - t \nabla f(x)^T v}{t^2} = \frac{1}{2} v^T \nabla^2 f(x) v + o(1).$$

当 t 充分小时，

$$\frac{f(x + tv) - f(x) - t \nabla f(x)^T v}{t^2} < 0,$$

这显然和一阶条件矛盾，因此必有 $\nabla^2 f(x) \succeq 0$ 成立。

二阶条件

Proof.

充分性：设 $f(x)$ 满足二阶条件 $\nabla^2 f(x) \succeq 0$ ，对任意 $x, y \in \text{dom} f$ ，根据泰勒展开，

$$f(y) = f(x) + \nabla f(x)^T(y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(x + t(y - x))(y - x),$$

其中 $t \in (0, 1)$ 是和 x, y 有关的常数。由半正定性可知对任意 $x, y \in \text{dom} f$ 有

$$f(y) \geq f(x) + \nabla f(x)^T(y - x).$$

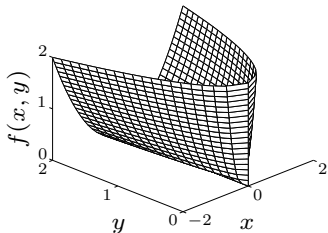
由凸函数判定的一阶条件知 f 为凸函数。进一步，若 $\nabla^2 f(x) > 0$ ，上式中不等号严格成立（ $x \neq y$ ）。利用一阶条件的充分性的证明过程可得 $f(x)$ 为严格凸函数。 □

二阶条件的应用

最小二乘函数: $f(x) = \|Ax - b\|_2^2$

$$\nabla f(x) = 2A^T(Ax - b), \quad \nabla^2 f(x) = 2A^T A$$

对任意 A , f 都是凸函数



quadratic-over-linear 函数: $f(x, y) = x^2/y$

$$\nabla^2 f(x, y) = \frac{2}{y^3} \begin{bmatrix} y \\ -x \end{bmatrix} \begin{bmatrix} y \\ -x \end{bmatrix}^T \succeq 0$$

是区域 $\{(x, y) \mid y > 0\}$ 上的凸函数

log-sum-exp函数: $f(x) = \log \sum_{k=1}^n \exp x_k$ 是凸函数

$$\nabla^2 f(x) = \frac{1}{\mathbf{1}^T z} \text{diag}(z) - \frac{1}{(\mathbf{1}^T z)^2} z z^T \quad (z_k = \exp x_k)$$

要证明 $\nabla^2 f(x) \succeq 0$, 我们只需证明对任意 v , $v^T \nabla^2 f(x) v \geq 0$, 即

$$v^T \nabla^2 f(x) v = \frac{(\sum_k z_k v_k^2)(\sum_k z_k) - (\sum_k v_k z_k)^2}{(\sum_k z_k)^2} \geq 0$$

由柯西不等式, 得 $(\sum_k v_k z_k)^2 \leq (\sum_k z_k v_k^2)(\sum_k z_k)$, 因此 f 是凸函数

几何平均: $f(x) = (\prod_{k=1}^n x_k)^{1/n}$ ($x \in \mathbb{R}_{++}^n$) 是凹函数

Jensen不等式

基础Jensen不等式: 设 f 是凸函数, 则对于 $0 \leq \theta \leq 1$,

$$f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y)$$

概率Jensen 不等式: 设 f 是凸函数, 则对任意随机变量 z

$$f(\mathbf{E}z) \leq \mathbf{E}f(z)$$

基础Jensen不等式可以视为概率Jensen 不等式在两点分布下的特殊情况

$$\text{prob}(z = x) = \theta, \quad \text{prob}(z = y) = 1 - \theta$$

提纲

1 基础知识

2 凸函数的定义与性质

3 保凸的运算

4 凸函数的推广

保凸的运算

验证一个函数 f 是凸函数的方法：

- ① 用定义验证（通常将函数限制在一条直线上）
- ② 利用一阶条件、二阶条件
- ③ 直接研究 f 的上方图 $\text{epi } f$
- ④ 说明 f 可由简单的凸函数通过一些保凸的运算得到
 - 非负加权和
 - 与仿射函数的复合
 - 逐点取最大值
 - 与标量、向量函数的复合
 - 取下确界
 - 透视函数

非负加权和与仿射函数的复合

非负数乘: 若 f 是凸函数, 则 αf 是凸函数, 其中 $\alpha \geq 0$.

求和: 若 f_1, f_2 是凸函数, 则 $f_1 + f_2$ 是凸函数.

与仿射函数的复合: 若 f 是凸函数, 则 $f(Ax + b)$ 是凸函数.

例子

- 线性不等式的对数障碍函数

$$f(x) = -\sum_{i=1}^m \log(b_i - a_i^T x), \quad \text{dom } f = \{x | a_i^T x < b_i, i = 1, \dots, m\}$$

- 仿射函数的(任意)范数: $f(x) = \|Ax + b\|$

逐点取最大值

若 f_1, \dots, f_m 是凸函数, 则 $f(x) = \max\{f_1(x), \dots, f_m(x)\}$ 是凸函数

例子

- 分段线性函数: $f(x) = \max_{i=1, \dots, m}(a_i^T x + b_i)$ 是凸函数
- $x \in \mathbb{R}^n$ 的前 r 个最大分量之和:

$$f(x) = x_{[1]} + x_{[2]} + \dots + x_{[r]}$$

是凸函数($x_{[i]}$ 为 x 的从大到小排列的第 i 个分量)

事实上, $f(x)$ 可以写成如下多个线性函数取最大值的形式:

$$f(x) = \max\{x_{i_1} + x_{i_2} + \dots + x_{i_r} \mid 1 \leq i_1 < i_2 < \dots < i_r \leq n\}$$

逐点取上界

若对每个 $y \in \mathcal{A}$, $f(x, y)$ 是关于 x 的凸函数, 则

$$g(x) = \sup_{y \in \mathcal{A}} f(x, y)$$

是凸函数

例子

- 集合 C 的支撑函数: $S_C(x) = \sup_{y \in C} y^T x$ 是凸函数
- 集合 C 点到给定点 x 的最远距离:

$$f(x) = \sup_{y \in C} \|x - y\|$$

- 对称矩阵 $X \in \mathbb{S}^n$ 的最大特征值

$$\lambda_{\max}(X) = \sup_{\|y\|_2=1} y^T X y$$

与标量函数的复合

给定函数 $g: \mathbb{R}^n \rightarrow \mathbb{R}$ and $h: \mathbb{R} \rightarrow \mathbb{R}$,

$$f(x) = h(g(x))$$

若 g 是凸函数, h 是凸函数, \tilde{h} 单调不减
 g 是凹函数, h 是凸函数, \tilde{h} 单调不增, 那么 f 是凸函数

- 对 $n = 1$, g, h 均可微的情形, 我们给出简证

$$f''(x) = h''(g(x))g'(x)^2 + h'(g(x))g''(x)$$

- 注意: 必须是 \tilde{h} 满足单调不减/不增的条件; 如果仅是 h 满足单调不减/不增的条件, 存在反例

推论

- 如果 g 是凸函数, 则 $\exp g(x)$ 是凸函数
- 如果 g 是正值凹函数, 则 $1/g(x)$ 是凸函数

与向量函数的复合

给定函数 $g : \mathbb{R}^n \rightarrow \mathbb{R}^k$ and $h : \mathbb{R}^k \rightarrow \mathbb{R}$:

$$f(x) = h(g(x)) = h(g_1(x), g_2(x), \dots, g_k(x))$$

若 g_i 是凸函数, h 是凸函数, \tilde{h} 关于每个分量单调不减, 那么 f 是凸函数
 g_i 是凹函数, h 是凸函数, \tilde{h} 关于每个分量单调不增

对 $n = 1$, g, h 均可微的情形, 我们给出简证

$$f''(x) = g'(x)^T \nabla^2 h(g(x)) g'(x) + \nabla h(g(x))^T g''(x)$$

推论

- 如果 g_i 是正值凹函数, 则 $\sum_{i=1}^m \log g_i(x)$ 是凹函数
- 如果 g_i 是凸函数, 则 $\log \sum_{i=1}^m \exp g_i(x)$ 是凸函数

取下确界

若 $f(x, y)$ 关于 (x, y) 整体是凸函数, C 是凸集, 则

$$g(x) = \inf_{y \in C} f(x, y)$$

是凸函数

例子

- 考虑函数 $f(x, y) = x^T A x + 2x^T B y + y^T C y$, 海瑟矩阵满足

$$\begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \succeq 0, \quad C \succ 0,$$

则 $f(x, y)$ 为凸函数. 对 y 求最小值得

$$g(x) = \inf_y f(x, y) = x^T (A - B C^{-1} B^T) x,$$

因此 g 是凸函数. 进一步地, A 的Schur 补 $A - B C^{-1} B^T \succeq 0$

- 点 x 到凸集 S 的距离 $\text{dist}(x, S) = \inf_{y \in S} \|x - y\|$ 是凸函数

透视函数

定义 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ 的透视函数 $g: \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}$,

$$g(x, t) = tf(x/t), \quad \text{dom } g = \{(x, t) | x/t \in \text{dom } f, t > 0\}$$

若 f 是凸函数, 则 g 是凸函数.

例子

- $f(x) = x^T x$ 是凸函数, 因此 $g(x, t) = x^T x/t$ 是区域 $\{(x, t) | t > 0\}$ 上的凸函数
- $f(x) = -\log x$ 是凸函数, 因此相对熵函数 $g(x, t) = t \log t - t \log x$ 是 \mathbb{R}_{++}^2 上的凸函数
- 若 f 是凸函数, 那么

$$g(x) = (c^T x + d)f((Ax + b)/(c^T x + d))$$

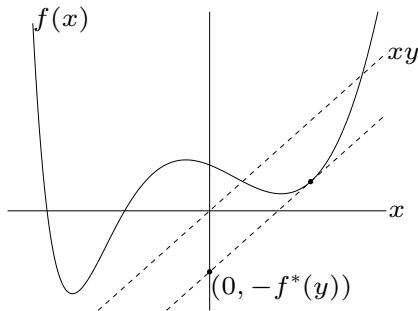
是区域 $\{x | c^T x + d > 0, (Ax + b)/(c^T x + d) \in \text{dom } f\}$ 上的凸函数

共轭函数

适当函数 f 的共轭函数定义为

$$f^*(y) = \sup_{x \in \text{dom } f} (y^T x - f(x))$$

- f^* 恒为凸函数，无论 f 是否是凸函数



例子

- 负对数 $f(x) = -\log x$

$$\begin{aligned} f^*(y) &= \sup_{x>0} (xy + \log x) \\ &= \begin{cases} -1 - \log(-y) & y < 0 \\ \infty & \text{其他} \end{cases} \end{aligned}$$

- 强凸二次函数 $f(x) = (1/2)x^T Qx$, $Q \in \mathbb{S}_{++}^n$

$$\begin{aligned} f^*(y) &= \sup_x (y^T x - (1/2)x^T Qx) \\ &= \frac{1}{2} y^T Q^{-1} y \end{aligned}$$

提纲

1 基础知识

2 凸函数的定义与性质

3 保凸的运算

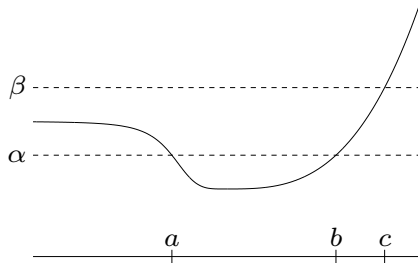
4 凸函数的推广

拟凸函数

$f: \mathbb{R}^n \rightarrow \mathbb{R}$ 称为拟凸的, 如果 $\text{dom } f$ 是凸集, 并且下水平集

$$S_\alpha = \{x \in \text{dom } f | f(x) \leq \alpha\}$$

对任意 α 都是凸的



- 若 f 是拟凸的, 则称 $-f$ 是拟凹的
- 若 f 既是拟凸的, 又是拟凹的, 则称 f 是拟线性的

拟凸、凹函数的例子

- $\sqrt{|x|}$ 是 \mathbb{R} 上的拟凸函数
- $\text{ceil}(x) = \inf\{z \in \mathbb{Z} | z \geq x\}$ 是拟线性的
- $\log x$ 是 \mathbb{R}_{++} 上的拟线性函数
- $f(x_1, x_2) = x_1 x_2$ 是 \mathbb{R}_{++}^2 上的拟凹函数
- 分式线性函数

$$f(x) = \frac{a^T x + b}{c^T x + d}, \quad \text{dom } f = \{x | c^T x + d > 0\}$$

是拟线性的

- 距离比值函数

$$f(x) = \frac{\|x - a\|_2}{\|x - b\|_2}, \quad \text{dom } f = \{x | \|x - a\|_2 \leq \|x - b\|_2\}$$

是拟凸的

例：内部回报率(IRR)

- 现金流 $x = (x_0, \dots, x_n)$; x_i 是 i 时段支付的现金
- 假定 $x_0 < 0, x_0 + x_1 + \dots + x_n > 0$
- 现值(present value)的表达式为

$$\text{PV}(x, r) = \sum_{i=0}^n (1+r)^{-i} x_i$$

其中 r 为利率

- 内部回报率(IRR)是最小的使得 $\text{PV}(x, r) = 0$ 的利率 r

$$\text{IRR}(x) = \inf\{r \geq 0 \mid \text{PV}(x, r) = 0\}$$

IRR 是拟凹的，因为它的上水平集是开半空间的交

$$\text{IRR}(x) \geq R \iff \sum_{i=0}^n (1+r)^{-i} x_i > 0 \text{ for } 0 \leq r < R$$

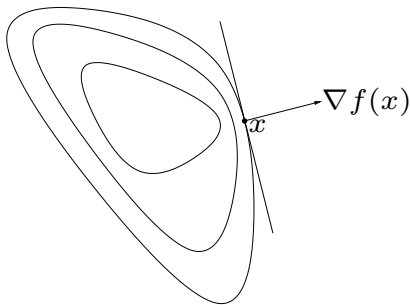
拟凸函数的性质

类Jensen不等式: 对拟凸函数 f

$$0 \leq \theta \leq 1 \implies f(\theta x + (1 - \theta)y) \leq \max\{f(x), f(y)\}$$

一阶条件: 定义在凸集上的可微函数 f 是拟凸的, 当且仅当

$$f(y) \leq f(x) \implies \nabla f(x)^T (y - x) \leq 0$$



注: 拟凸函数的和不一定是拟凸函数

对数凸函数

如果正值函数 f 满足 $\log f$ 是凸函数, 则 f 称为对数凸函数, 即

$$f(\theta x + (1 - \theta)y) \leq f(x)^\theta f(y)^{1-\theta} \quad \text{for } 0 \leq \theta \leq 1.$$

如果 $\log f$ 是凹函数, 则 f 称为对数凹函数,

- 幂函数: 当 $a \leq 0$ 时, x^a 是 \mathbb{R}_{++} 上的对数凸函数; 当 $a \geq 0$, x^a 是 \mathbb{R}_{++} 上的对数凹函数
- 许多常见的概率密度函数是对数凹函数, 例如正态分布

$$f(x) = \frac{1}{\sqrt{(2\pi)^n \det \Sigma}} e^{-\frac{1}{2}(x-\bar{x})^T \Sigma^{-1}(x-\bar{x})}$$

- 高斯分布的累计分布函数 Φ 是对数凹函数

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-u^2/2} du$$

对数凸、凹函数的性质

- 定义在凸集上的二阶可微函数 f 是对数凹的，当且仅当

$$f(x)\nabla^2 f(x) \preceq \nabla f(x)\nabla f(x)^T$$

对任意 $x \in \text{dom } f$ 成立

- 对数凹函数的乘积仍为对数凹函数
- 对数凹函数的和不一定为对数凹函数
- 若 $f: \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ 是对数凹函数，那么

$$g(x) = \int f(x, y) dy$$

是对数凹函数

对数凹函数的积分

- 对数凹函数 f, g 的卷积 $f * g$ 是对数凹函数

$$(f * g)(x) = \int f(x - y)g(y)dy$$

- 若 $C \subseteq \mathbb{R}^n$ 是凸集，并且随机变量 y 的概率密度函数是对数凹函数，则

$$f(x) = \text{prob}(x + y \in C)$$

是对数凹函数

证明： $f(x)$ 可表示为两个对数凹函数乘积的积分

$$f(x) = \int g(x + y)p(y)dy, \quad g(u) = \begin{cases} 1 & u \in C \\ 0 & u \notin C, \end{cases}$$

其中 p 是 y 的概率密度函数

例：生成函数

$$Y(x) = \text{prob}(x + w \in S)$$

- $x \in \mathbb{R}^n$: 产品的标称参数(nominal parameter)
- $w \in \mathbb{R}^n$: 制成品的参数是随机变量
- S : 接受集

若 S 是凸集，并且随机变量 w 的概率密度函数是对数凹函数，则

- Y 是拟凹函数
- 生成区域 $\{x | Y(x) \geq \alpha\}$ 是凸集

广义不等式意义下的凸函数

$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ 称为 K -凸函数: 如果 $\text{dom } f$ 是凸集, 并且

$$f(\theta x + (1 - \theta)y) \preceq_K \theta f(x) + (1 - \theta)f(y)$$

对任意 $x, y \in \text{dom } f$, $0 \leq \theta \leq 1$ 成立

例子 $f: \mathbb{S}^m \rightarrow \mathbb{S}^m, f(X) = X^2$ 是 \mathbb{S}_+^m -凸函数

证明: 对固定的 $z \in \mathbb{R}^m$, $z^T X^2 z = \|Xz\|_2^2$ 关于 X 是凸函数, 即

$$z^T (\theta X + (1 - \theta)Y)^2 z \leq \theta z^T X^2 z + (1 - \theta)z^T Y^2 z$$

对任意 $X, Y \in \mathbb{S}^m$, $0 \leq \theta \leq 1$ 成立

因此 $(\theta X + (1 - \theta)Y)^2 \preceq \theta X^2 + (1 - \theta)Y^2$