



组合优化与凸优化(combinatorial optimization and convex optimization)

课堂讨论

刘绍辉

计算机科学与技术学院 哈尔滨工业大学

shliu@hit.edu.cn

2023年春季



简介

- <https://optimization-online.org/> : 类似arxiv.org的网站, 专门发表跟优化相关的电子论文, 2000年
- <https://www.coin-or.org/> : COmputational Infrastructure for Operations Research, Open Source for the operations research community, 2000年.
- <http://Neos-server.org> : NEOS Server: SOTA solvers for Numerical Optimization, a free internet-based for solving numerical optimization problems. 由威斯康星大学麦迪逊分校托管的求解器运行在由HTCondor软件支持的分布式高性能机器上; 远程求解器在亚利桑那州立大学、奥地利克拉根福大学和葡萄牙米尼奥大学的机器上运行。

优化求解应用广泛-涉及国民经济的各个行业

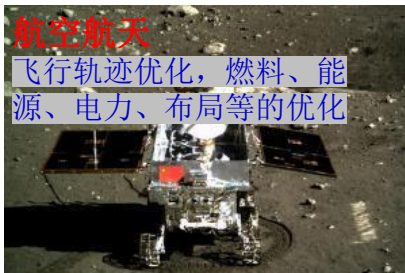
军事

后勤优化, 如资源调度、存储、管理、装备维修等



航空航天

飞行轨迹优化, 燃料、能源、电力、布局等的优化



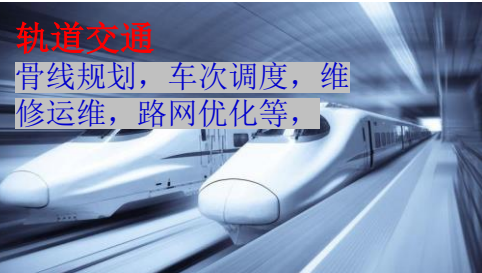
机械制造

布局优化, 路径优化、空间、材料等



轨道交通

骨干线规划, 车次调度, 维修运维, 路网优化等,



能源电网

电网并联, 机组运行, 石油管道布设, 传感定位, 故障恢复等,



物流

路径, 路网, 路线优化, 仓储选址与规划等



金融与经济

风控, 资产管理与配置, 供应链管理, 市场博弈, 资源配置等



智慧城市与现代生活

社交网络管理, 供水排水网络管控与优化, 突发事件决策与管理等



问题建模 + 优化求解是解决实际应用问题的关键



调度

规划

...

决策



优化求解器

- 科学技术
 - 科学-发现规律，总结规律
 - 技术-利用规律，解决问题
- 在国防、能源、制造、交通、金融、通信、计算机等各个领域出现了很多实际问题，对其建模后，很多都形式化为运筹优化问题，并指导各个领域和行业的精益化发展
- 实际问题通过规律和实际约束对其进行**数学建模**，然后对这些数学问题进行**求解**，这就是专门解决优化问题的求解器(Solver)
- 优化求解器主要以工业软件的形式为各个领域提供问题求解工具，能够将实际问题中的大规模复杂问题进行优化求解！



商用求解器vs开源求解器

- 数据驱动

数据获取 → 分析内在规律 → 建模与求解

- 数学规划/优化理论

- 问题建模与求解的核心所在
- 求解器：编程实现求解方法，给出最优解
- 决策部门根据最优解及其敏感性分析作出最优决策！

- 商用和开源

- Gurobi, IBM CPLEX, FICO XPRESS, Matlab, Excel (Frontline Solvers)
- SCIP (德国), MOSEK (丹麦), GLPK (俄罗斯)
- CVX-求解器建模平台 (斯坦福大学: Stephen Boyd, Yinyu Ye)

华为: OPTV

阿里: MDOPT

杉数: COPT

中科院: CMIP

上财: LEAVES



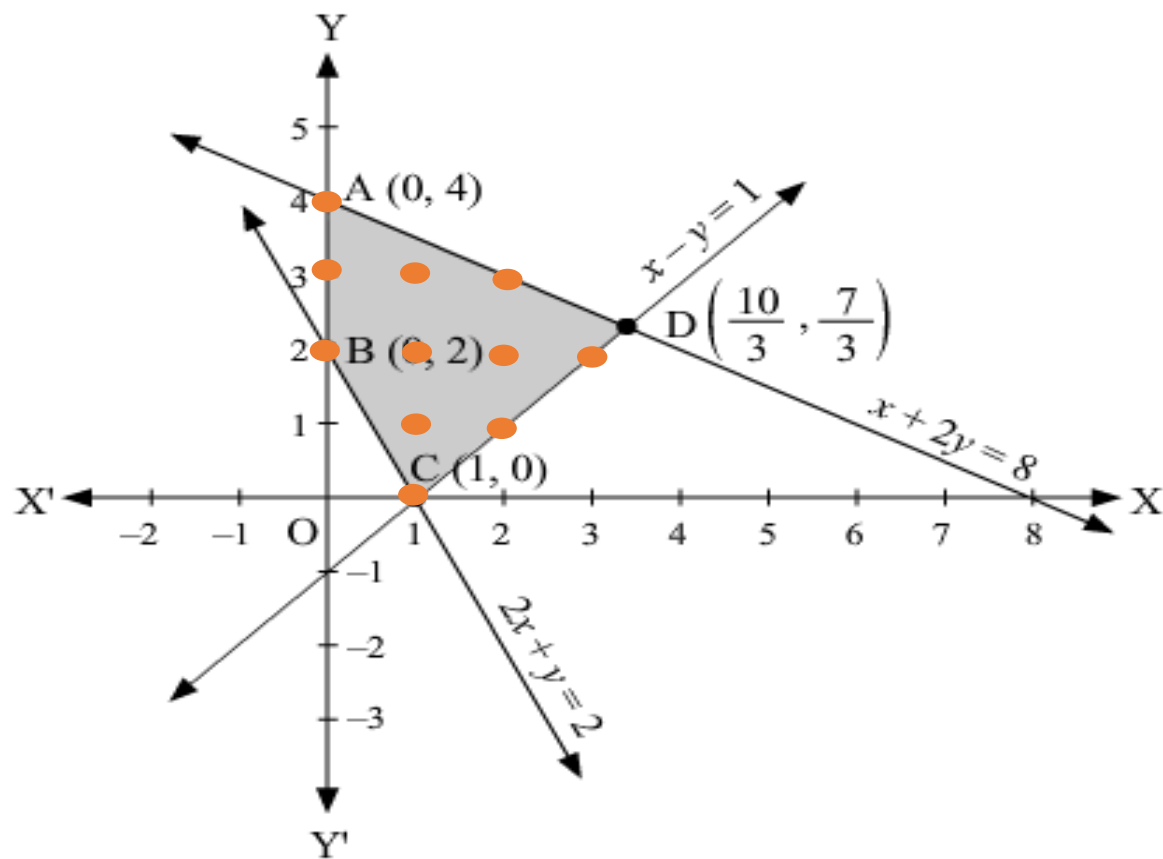
数学规划的分类

- 连续优化问题vs离散优化问题
 - COP: 约束集包含无穷多个元素, 并具有连续性特征, 利用微积分和凸函数性质来分析
 - DOP: 约束集是有限集, 调度、路径规划、匹配问题、整数规划, 利用组合数学和离散数学进行分析, 同时也利用COP相关的方法
- 无约束优化 (Unconstrained) vs约束优化
- 凸优化vs非凸优化
- 线性规划 (Linear Programming) vs非线性规划 (Nonlinear Programming)
- 网络优化问题: 同时具有连续优化和离散优化问题特点



线性规划LP-IP

- 可行域内找最优解
 - 规模大：数百万变量和约束
 - 稀疏性：稀疏矩阵计算
 - 数值：受计算机精度影响
- LP*
- 可行域是离散点
 - 经典的NP完全问题
 - 问题结构复杂多样
 - 调用LP求解
- IP*





非线性规划 (NLP)

- 最小二乘LS，线性规划LP，以及凸二次规划 (QP) 的统一和推广
- 其中凸规划 (Convex Programming) 是非线性规划的子集，近二十年来，凸规划已成为与LS、LP和QP一样无处不在的技术。
- 与成熟的LS，LP和QP只需要对计算有基本的理解即可不一样，一般的CP的广泛使用要求具有较高的专业知识水平，用户必须理解凸分析的基础知识既是合理的，也是不可避免的;但事实上，还需要更深入的理解。
- 用户必须找到一种方法将他的问题转化为众多有限的**标准形式**之一;或者，如果做不到这一点，开发一个**自定义求解器**。对于那些专注于特定应用而不是底层数学的面向应用的用户来说，这些要求构成了使用凸规划的巨大障碍，特别是如果还不能确定结果是否会比其他方法更好的话。



优化求解器

- LP, MIP, NLP在实际应用中的占比
 - LP: 15%
 - MIP: 79%
 - NLP: 7%
- 线性规划LP: 约5万行代码
- 整数规划IP: 100~200万行
- 非线性各个模块: 10万行以下



第一章、绪论

- 各行各业都或多或少的应用运筹学或最优化的思想

问题一：你觉得在哪些方面已经应用了这种思想？请从宏观和微观方面，举例说明。

问题二：结合你自己的研究方向？举例说明。

machine learning = representation + optimization + evaluation

P. Domingos, an AAAI Fellow and a Professor of University of Washington



第一章、绪论

- 线性系统 $Ax = b, A \in R^{m \times n}, m < n$
 - 无穷多解
 - 无解：为此假设 A 是满秩矩阵
- 实际工程考虑
 - 图像超分辨率，输入为模糊下采样后的图像 b , 矩阵 A 标识降质运算，目标是从输入 b 重构原始图像 x
 - $\min_x J(x) \text{ s.t. } b = Ax$
 - $J(x)$ 用来评估解的可取性，例如超分中可以是平滑性、分段平滑性等
 - 常见的 $J(x)$ 为 $\|x\|_2^2$, 此时定义拉格朗日函数为：
 - $L(x) = \|x\|_2^2 + \lambda^T (Ax - b)$, 求导数为0即可得： $\hat{x}_{opt} = -\frac{1}{2} A^T \lambda$
 - 代入约束 $Ax = b \Rightarrow A\hat{x}_{opt} = -\frac{1}{2} AA^T \lambda = b \Rightarrow \lambda = -2(AA^T)^{-1} b$
 - 最后得 $\hat{x}_{opt} = A^T (AA^T)^{-1} b = A^+ b$, 其中 A^+ 表示伪逆
 - 问题：如果 $J(x) = \|Bx\|_2^2$ (注意： $B^T B$ 可逆)，最优解=？

$$\hat{x}_{opt} = (B^T B)^{-1} A^T (A(B^T B)^{-1} A^T)^{-1} b$$



第一章、绪论

- 为何 l_2 使用广泛？
 - 具有唯一的解析形式的解，简单
 - 在反问题，信号表达等很多地方常见
 - 那是否 l_2 就是最好的选择呢？为什么？
- 如果不是，如何给出解？
 - l_p 范数($p \geq 1$): $\|x\|_p^p = \sum_i |x_i|^p$, 其中特殊情况为 $l_\infty = \max_i(|x_i|)$, $l_1 = \sum_i |x_i|$
 - 而其中 l_1 在一定条件下往往倾向于稀疏解，并且 $l_p(p \geq 1)$ 是凸的，因此很常用
 - 并且，此时该问题与线性规划问题等价！
 - l_0 范数：不是真正的范数，且 l_0 非凸，求解困难

为什么？



第一章、绪论

- 压缩感知 (CS: Compressive sensing)
 - 求解方程 $Ax = b, x \in R^n, b \in R^m, A \in R^{m \times n}, m \ll n$
 - 如接收到信号 b , 已知矩阵 A , 重构向量 x (无穷多), 先验: 稀疏性, 即 x 中有很多0元素
 - $m=128; n=256; A=\text{randn}(m,n); u=\text{sprandn}(n,1,0.1); b=A*u;$
 - 理论上, u 是如下 l_0 范数问题的最优解:
 - $\min_{x \in R^n} \|x\|_0, s. t. Ax = b$
 - 这是一个NP难的, 如何求解?
 - 定义 l_1 范数: $\|x\|_1 = \sum_{i=1}^n |x_i|$, 则在一定条件下, l_0 可等价转换为 l_1 问题:
 - $\min_{x \in R^n} \|x\|_1, s. t. Ax = b$
 - 能否使用 l_2 范数来替代 l_0 范数呢?



第一章、绪论

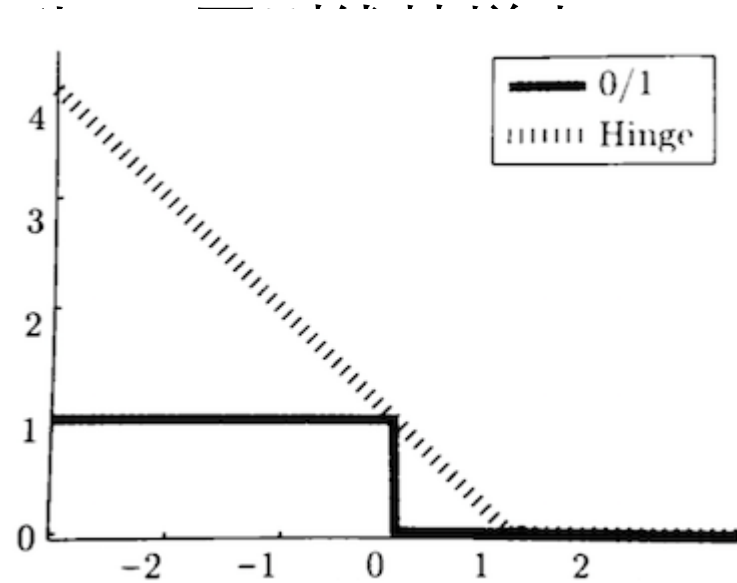
- SVM
 - 最大间隔算法与拉格朗日对偶
 - Hinge loss
- 二分类问题：标签 $\mathbf{y} = \pm 1$, 预测值 $\hat{\mathbf{y}} \in \mathbf{R}$
 - $\hat{\mathbf{y}} > +1 / \hat{\mathbf{y}} < -1$, 此时分类正确, 损失为0
 - $\hat{\mathbf{y}} \in (-1, 1)$, 分类不确定, 损失不为0
 - $\hat{\mathbf{y}} = 0$, 损失最大
- 对输出 \mathbf{y} , 则损失可写为
 - $l(\mathbf{y}) = \max\{0, 1 - \mathbf{y} \cdot \hat{\mathbf{y}}\}$
 - 这就是Hinge loss在二分类问题的变体, 可看做是双向Hinge loss.



第一章、绪论

- 单向hinge loss

- $y = \pm 1, y \geq 1, \text{loss}$

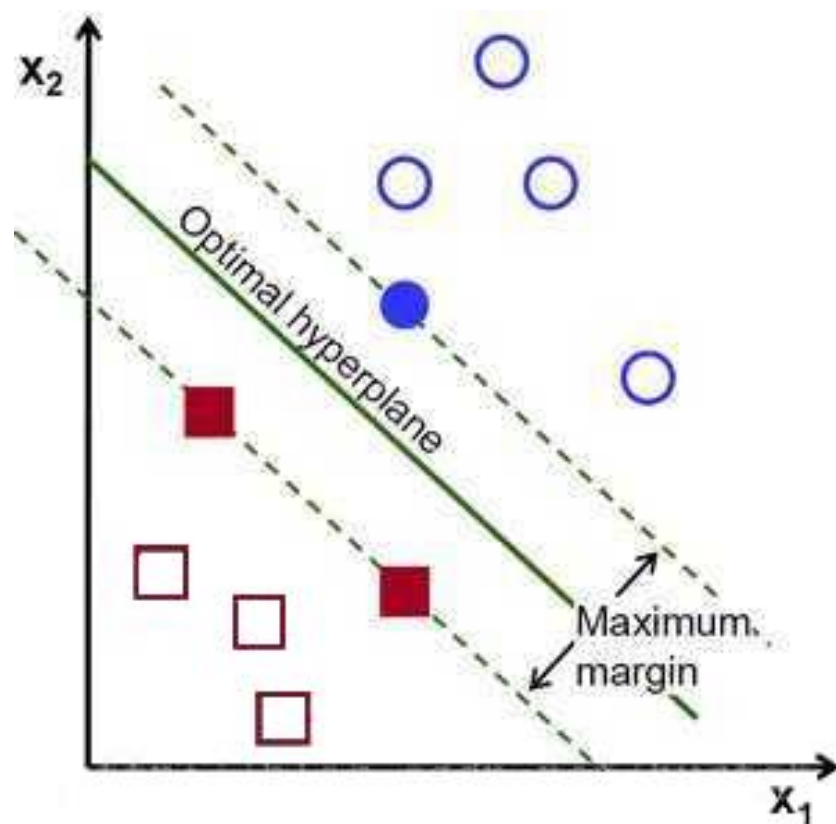


- SVM呢?



第一章、绪论

- 线性可分SVM的预测值 $\hat{y} = \mathbf{w}\mathbf{x} + \mathbf{b} \in \mathbf{R}$ ，如果分离超平面在如图，且支持向量与分割平面之间的距离为1
 - 每个 $\mathbf{y} = 1$ 的样本 $\hat{y} \geq 1$ ， $\mathbf{y} = -1$ 的样本 $\hat{y} \leq -1$ ，样本集上损失为0
 - 如果分类错误，则Hinge loss大于0。Hinge loss约束超平面作出调整。如果超平面距离支持向量的距离小于1，则Hinge loss大于0，且就算分离超平面满足最大间隔，Hinge loss仍大于0





第一章、绪论

- 预测值 $|\hat{y}|$ 越大，表明样本点越远离分离超平面，分类越容易。选择优化设计时，没必要关注离超平面很远的点，因此可通过对距离分离超平面的距离选择一个阈值，来去除这些样本！
 - $l(y) = \max\{0, 1 - y \cdot \hat{y}\}$ ，其中的1就是阈值，可看做是超参数，对预测值超过阈值的样本无需考虑
 - 考虑多分类
 - $l(y) = \max\{0, 1 + \max_{\hat{y} \neq y} w_{\hat{y}}x - w_yx\}$
 - $l(y) = \sum_{t \neq y} \max\{0, 1 + w_{\hat{y}}x - w_yx\}$
- 损失函数知道，如何求解？
 - $\frac{\partial l}{\partial w_i} = \begin{cases} -y \cdot x_i & y \cdot \hat{y} < 1 \\ 0 & \text{其它} \end{cases}$ ，有什么问题？



第一章、绪论

- 损失函数在 $\mathbf{y} \cdot \hat{\mathbf{y}} = 1$ 处不可导
 - 各种光滑变体

- 分段光滑[2]: $l(\hat{\mathbf{y}}) = \begin{cases} \frac{1}{2} - \mathbf{y} \cdot \hat{\mathbf{y}} & \mathbf{y} \cdot \hat{\mathbf{y}} \leq 0 \\ \frac{1}{2} (1 - \mathbf{y} \cdot \hat{\mathbf{y}})^2 & 0 < \mathbf{y} \cdot \hat{\mathbf{y}} \leq 1 \\ 0 & 1 \leq \mathbf{y} \cdot \hat{\mathbf{y}} \end{cases}$

- 平方光滑: $l_{\gamma}(\hat{\mathbf{y}}) = \begin{cases} \frac{1}{2\gamma} \max\{0, 1 - \mathbf{y} \cdot \hat{\mathbf{y}}\}^2 & \mathbf{y} \cdot \hat{\mathbf{y}} \geq 1 - \gamma \\ 1 - \frac{\gamma}{2} - \mathbf{y} \cdot \hat{\mathbf{y}} & \text{其它} \end{cases}$

- 其中更改的Huberloss是 $\gamma = 2$ 是平方光滑的特例 $L(\mathbf{y}, \hat{\mathbf{y}}) = 4l_2(\hat{\mathbf{y}})$

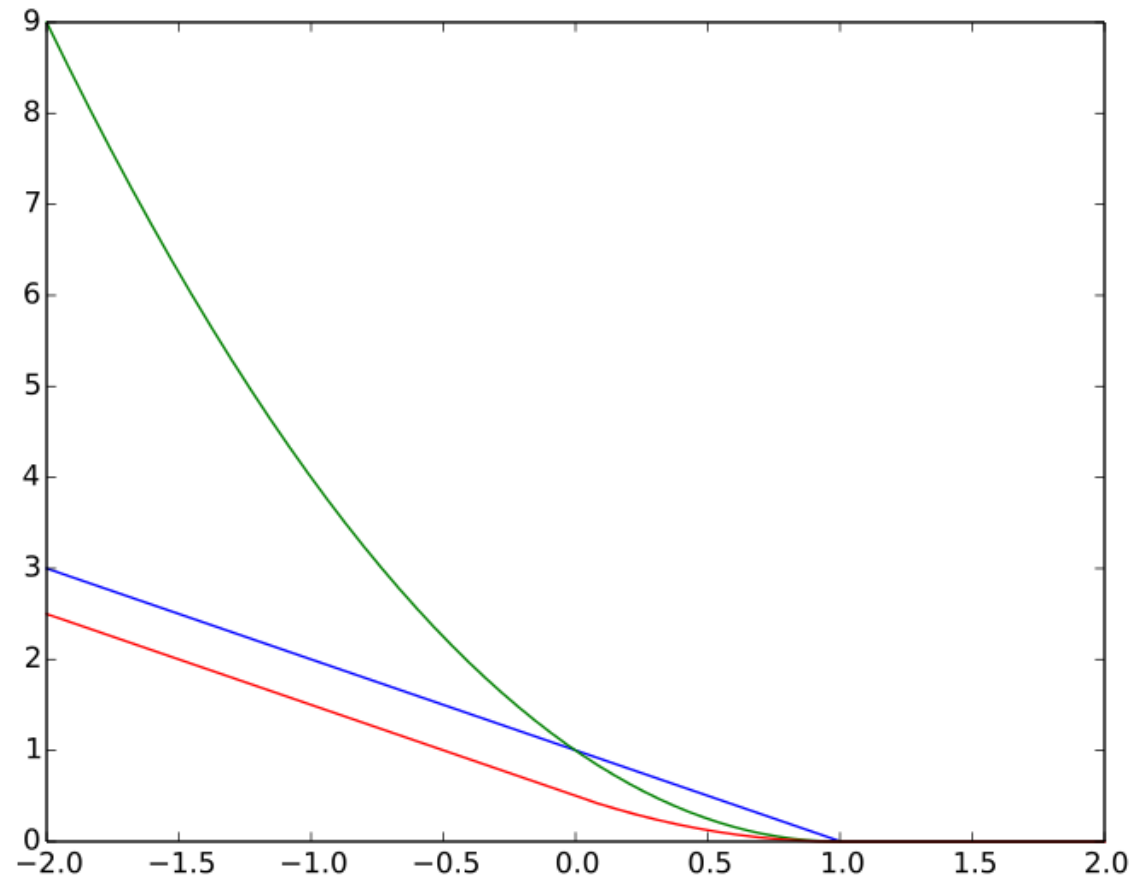
1. Zhang, Tong. Solving large scale linear prediction problems using stochastic gradient descent algorithms (PDF). ICML. 2004.

2. Rennie, Jason D. M.; Srebro, Nathan. Loss Functions for Preference Levels: Regression with Discrete Ordered Labels (PDF). Proc. IJCAI. 2005



第一章、绪论

- “普通变体”（蓝色），平方变体（绿色），以及 Rennie 和 Srebro 提出的分段：





第一章、绪论

- 典型问题

- $\min_{\mathbf{w} \in \mathbb{R}^n} \frac{1}{m} \sum_{i=1}^m l(\mathbf{p}(x_i; \mathbf{w}), y_i) + \lambda R(\mathbf{w})$

- $l(\mathbf{p}, y) = \frac{1}{2} (\mathbf{p} - y)^2$; $l(\mathbf{p}, y) = \log(1 + e^{-\mathbf{p}y})$; $l(\mathbf{p}, y) = \max\{0, 1 - \mathbf{p}y\}$

- $\mathbf{p}(x; \mathbf{w}) = \mathbf{w}^T \mathbf{x} - b$; $\mathbf{p}(x; W) = \phi(W_n \phi(W_{n-1} \cdots \phi(W_1 x) \cdots))$

- $R(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$; $R(\mathbf{w}) = \|\mathbf{w}\|_1$

- SVM: hinge loss, 线性分类函数, \mathbf{l}_2 正则化

- 正则化逻辑斯谛回归: logistic损失, 线性回归函数, \mathbf{l}_2 正则化

- 多层感知机: 平方损失, 前向反馈函数, $R(W)=0$

- LASSO问题: 平方损失, 线性回归函数, \mathbf{l}_1 正则化



第一章、绪论

- 典型问题

- $\min_{X \in \mathbb{R}^{m \times n}} \|X\|_*, s. t. X_{ij} = D_{ij}, \forall (i, j) \in \Omega$, 其中 Ω 表示观测值的位置

- 低秩表示问题

- $\min_{Z \in \mathbb{R}^{m \times n}, E \in \mathbb{R}^{m \times n}} \|Z\|_* + \lambda \|E\|_1, s. t. D = DZ + E$

- 为减少计算代价和存储空间，低秩矩阵可分解为两个小得多的矩阵乘积： $\mathbf{X} = \mathbf{UV}^T$, 此时上述矩阵补全问题可形式化为如下的非凸问题

- $\min_{U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}} \frac{1}{2} \sum_{(i,j) \in \Omega} \|U_i V_j^T - D_{ij}\|_F^2 + \frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2)$



第一章、绪论

- 优化分类
- $\min f_0(x), s. t. f_j(x) \leq (\geq, =) 0, j = 1, \dots, m, x \in Q$
- 可行集 $\mathcal{F} = \{x \in Q | f_j(x) \leq 0, j = 1, \dots, m\}$
- 约束问题: $\mathcal{F} \subset \mathbb{R}^n$; 无约束问题: $\mathcal{F} = \mathbb{R}^n$; 光滑问题: $f_j(\cdot)$ 都可微; 非光滑问题: $f_k(\cdot)$ 不可微; 线性约束问题: $f_j(x) = \sum_{i=1}^n a_j^{(i)} x^{(i)} + b_j \equiv \langle a_j, x \rangle + b_j, j = 1, \dots, m$ 都是仿射的, 如果 $f_0(x)$ 也是仿射的, 则实线性优化问题; 如果目标函数 $f_0(x)$ 是二次的, 则为二次优化问题; 如果所有约束都是二次函数, 则为二次约束的二次优化问题



第一章、绪论

- 优化问题也可分为可行的，和不可行的
- 如果 $\mathcal{F} \neq \emptyset$,则称问题是可行的
- 如果存在某个 $\mathbf{x} \in \mathbf{Q}$ ，使得所有不等式约束和等式严格成立，则称为严格可行的
- 一般说来，优化问题应该是不可解的（？）
 - 的确，从我们的现实生活经验来看，很难相信存在一种能够解决世界上所有问题的万能工具。
- 根据函数的性质称为零阶，一阶和二阶
 - 零阶：只利用函数值
 - 一阶：利用函数值和导数值
 - 二阶：利用函数值、导数值和Hessian矩阵值



第一章、绪论

- 优化(Optimization)无所不在，包罗万象，但在国内的应用跟国外还有一些差距

问题一：在宏观方面，你觉的哪些方面能够采用这些技术来提高效率？

问题二：优化技术应用的难点在哪？



第一章、绪论

- 优化或运筹学是解决实际应用的问题，这显然需要建模

问题一：建模的基本流程是怎么样的？

问题二：数学建模的难点在哪？

问题三：你听说过幸存者定律，墨菲定律（Murphy's Law）？

Anything that can go wrong will go wrong



第一章、绪论

- 最优化(Optimization) 是数学当中比较流行的词汇

问题一：你听说过次梯度，梯度投影，拉格朗日对偶，强对偶，弱对偶等名词吗？

问题二：听说过秩1校正、秩2校正吗？



第一章、绪论

- 最优化(Optimization) 是数学当中比较流行的词汇

李普希兹连续，梯度李普希兹连续，变量拆分，块坐标下降 (block coordinate descent), Q-收敛速度，超线性、线性、次线性、二次收敛，R-收敛速度，优化算法复杂度（

$$N(\epsilon) = \mathcal{O}\left(\frac{1}{\epsilon^2}\right)$$



第一章、绪论

- 作为最基本的线性规划

问题一：听说过丹齐格，敏感性分析吗？

问题二：为什么要学线性规划？



第一章、绪论

- 作为最优化迭代的基本形式

$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \mathbf{p}_k$, 其中 α 为步长, \mathbf{p}_k 为方向,

问题一：你知道几种步长和方向的求法？

问题二：你知道线搜索和信任域方法的区别吗？



第一章、绪论

- 优化算法中一般涉及泰勒逼近：

$$f(x + p) = f(x) + \nabla f(x + tp)^T p, t \in (0, 1)$$

$$f(x + p) = f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x + tp) p, t \in (0, 1)$$

问题一：连续可微性为什么很重要？

问题二：迭代求解中，如果保证目标函数值下降，是否可以求得极小值？



第一章、绪论

- 优化算法中一般需要求极小值：

$$f(x + p) = f(x) + \nabla f(x + tp)^T p, t \in (0, 1)$$

$$f(x + p) = f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x + tp) p, t \in (0, 1)$$

问题一：连续可微性为什么很重要？

问题二：迭代求解中，如果保证目标函数值下降，是否可以求得极小值？



第二章 线性规划

- 集合的概念
- 开集, 闭集, 紧集
- 每个有界并单调非递减或非递增序列都收敛。每个单调非递减序列 $\{x_k\}$ 要么有界, 要么 $x_k \rightarrow \infty$
- 集合 $X \subset \mathbf{R}^n$, x 为集合 X 的封闭点(closure), 如果存在序列 $\{x_k\} \subset X$ 收敛于 x 。集合 X 的闭包: $cl(X)$ 为所有闭点的集合。
- 闭集: 闭包与其本身相等
- 开集: 闭集的补集
- 紧集: 闭且有界



第二章 线性规划

- 设 $L \subset \mathbf{R}^n$ 为子空间，其余子空间或正交补记为：

$$L^\perp = \{\mathbf{x} \in \mathbf{R}^n | \mathbf{x}^T \mathbf{y} = 0, \forall \mathbf{y} \in L\}$$

- 显然，任意 $\mathbf{z} \in \mathbf{R}^n$ ，存在唯一分解 $\mathbf{z} = \mathbf{x} + \mathbf{y}$ ，使得 $\mathbf{x} \in L, \mathbf{y} \in L^\perp$. 称 \mathbf{x}, \mathbf{y} 分别为 \mathbf{z} 在子空间 L , L^\perp 上的投影，这时有： $\|\mathbf{z}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2$ ，这就是投影定理
- 问题1. 这个定理内容是否理解？
- 问题2. 你是否听说过投影定理 (Projection Theorem)



第二章 线性规划

- 设 L 为 R^n 为子空间, 那么, $\forall z \in R^n, \exists$ 唯一 $x \in L, y \in L^\perp$, 使得 $z = x + y$, 且 x 为问题:

$$\begin{cases} \min ||z - u|| \\ s. t. u \in L \end{cases}$$

的唯一解, 而问题的最优值为 $||y||$

- 问题1. 这个定理跟之前的投影定理是否有不同?
- 问题2. 子空间 L 是否是凸的?
- 问题3. 参看教材17页的定理1-2.



第二章 线性规划

- 你知道对线性规划影响最大的三种算法是什么算法吗？
 - 单纯形法 (simplex algorithms)
 - 内点法 (interior point algorithm)
 - 椭球法 (ellipsoid method)
- 迄今未发现单纯形法的任何变种具有多项式运行时间
- 椭球法已证明对LP问题具有多项式时间算法，但其效率低下不适合实际使用
- 但内点法（最坏情况具有指数运行时间），和单纯形法远比椭球法更有效，因而实际广泛使用



第二章 线性规划

- 设 L 为 R^n 的非空闭凸子集, 那么, $\forall z \in R^n, \exists$ 唯一的向量 $x^* \in L$, 且 x^* 为问题:

$$\begin{cases} \min ||z - x|| \\ s.t. x \in L \end{cases}$$

的唯一解, 称其为向量 z 在集合 L 上的投影。并且, 向量 x^* 是 z 在 L 上的投影当前仅当 (if and only if) :

$$(z - x^*)'(x - x^*) \leq 0, \forall x \in L \text{ ----- } (*)$$

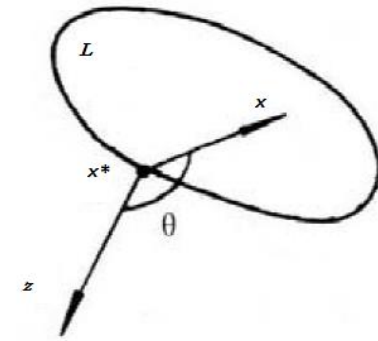
- 问题1. 这种描述之前听说过吗?

- 问题2. 是否会证明?

- 最小化 $||z - x||$ 等价于最小化凸可微函数: $f(x) = \frac{1}{2} ||z - x||^2$, 则 x^*
 $\nabla f(x^*)'(x - x^*) \geq 0, \forall x \in L$, 而这正好是上述(*)的内容!

- 问题3. 唯一性如何证明?

- 在集合 L 上最小化 f 等价于在紧集 $L \cap \{||z - x|| \leq ||z - w||\}, \forall w \in L$. 由Weierstrass定理 (连续函数 $f: R^n \rightarrow R$ 在 R^n 的任何紧子集上都有极小点), 存在一个最小化向量。若存在两个最小化向量, 则必定相等!



于



第二章 线性规划

- 凸函数的操作:
- $f: \mathbf{R}^m \rightarrow (-\infty, \infty]$, A 为 $m \times n$ 的矩阵, 令 $F: \mathbf{R}^n \rightarrow (-\infty, \infty]$, $F(\mathbf{x}) = f(A\mathbf{x})$, $\mathbf{x} \in \mathbf{R}^n$
- **问题1.** 如果 f 是凸函数, 则 F 是否是凸函数? 如果 f 是闭的, 则 F 也是闭的。如何证明?
- 问题2. 函数的和是否是与这些函数的凸性有关呢?
 - $F(\mathbf{x}) = \gamma_1 f_1(\mathbf{x}) + \cdots + \gamma_m f_m(\mathbf{x})$, $\mathbf{x} \in \mathbf{R}^n$, $\gamma_i > 0$
- 问题3. 多个函数的上确界呢?
 - $f(\mathbf{x}) = \sup_{i \in I} f_i(\mathbf{x})$, I 为任意索引集, $f_i: \mathbf{R}^n \rightarrow (-\infty, \infty]$
 - $f(\mathbf{x}) = \max_{i \in I} f_i(\mathbf{x})$

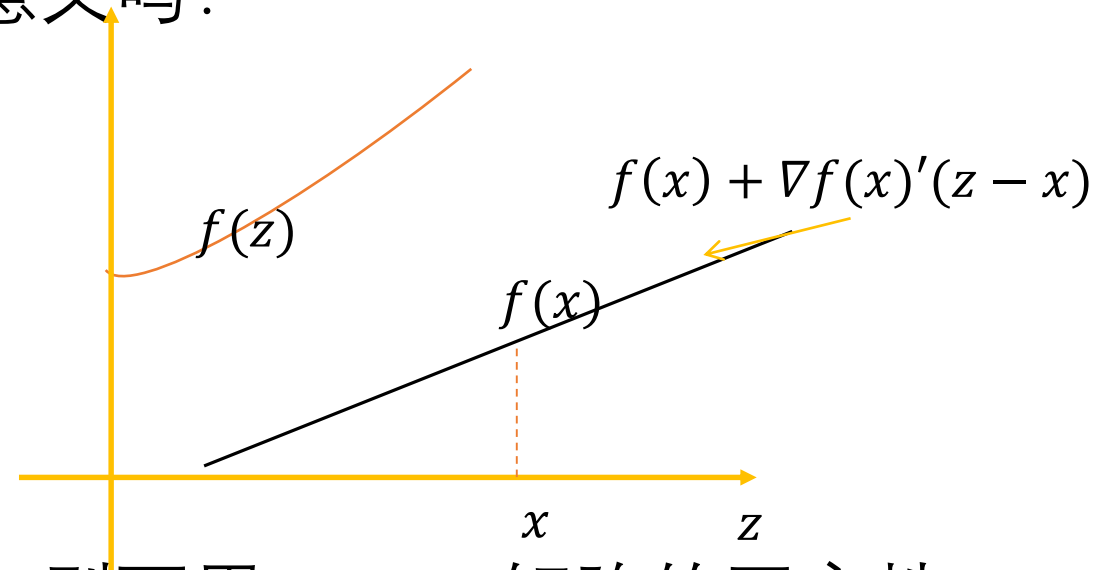


第二章 线性规划

- 凸函数的几何意义等价于
 - $f(z) \geq f(x) + \nabla f(x)'(z - x), \forall x, z, \in C, C$ 为非空凸子集
 - 严格凸的时候, 则等号没有
- 问题1.理解凸函数的几何意义吗?
- 问题2.如何证明上述结论?

反例:

$f(x) = x_1^2 - x_2^2, C = \{(x_1, 0) | x_1 \in R\}, f$ 是凸的, 但 $\nabla^2 f(x)$ 并不是正定的; 又如: $f(x) = x^4$, 严格凸, 但 $\nabla^2 f(x)$ 并不一定 > 0 , 例如在0点就为0.



- 问题3.如果函数二次可微, 则可用Hessian矩阵的正定性来判断: 如果为凸 $\Leftrightarrow \nabla^2 f(x) \geq 0$; 若 $\nabla^2 f(x) > 0 \Rightarrow f$ 为严格凸。问: 若 f 严格凸, 则 $\nabla^2 f(x) > 0$ 是否成立?

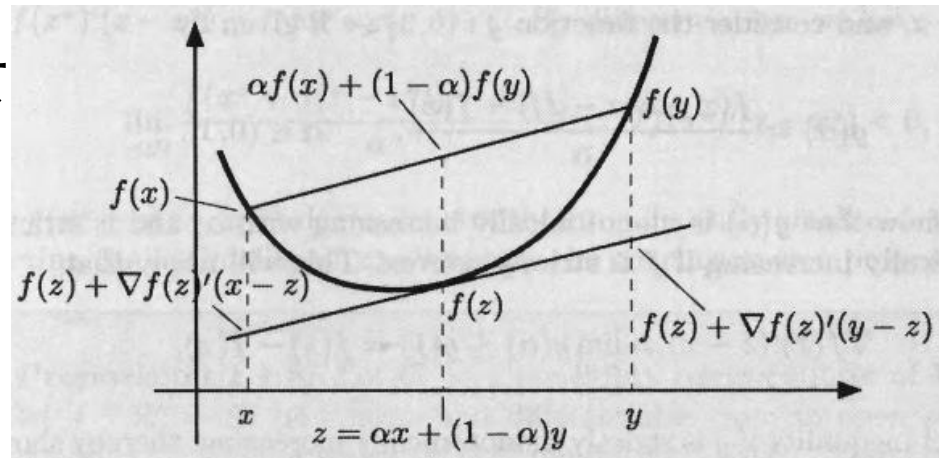


第二章 线性规划

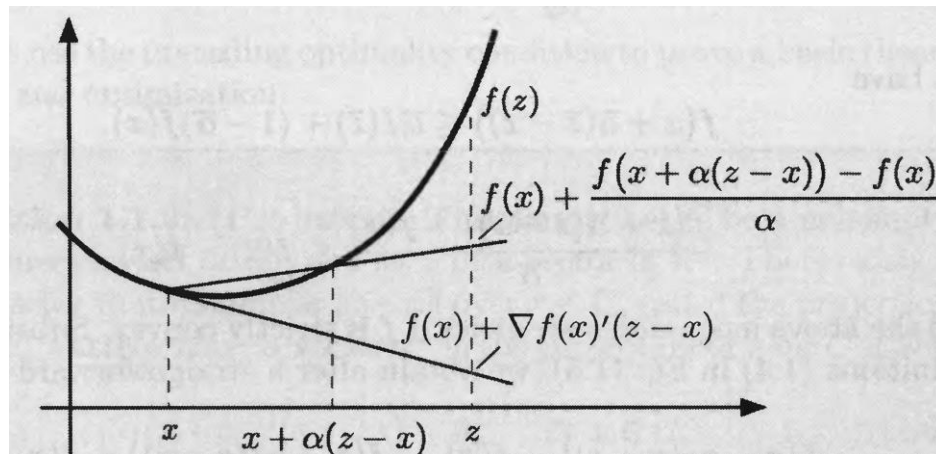
- 凸函数组合的几何意义
 - $x, y \in C, \alpha \in [0, 1]$, 设 $z = \alpha x + (1 - \alpha)y$, 则根据上页的不等式有:
 - $f(x) \geq f(z) + \nabla f(z)'(x - z), \text{--- -- (1)}$
 - $f(y) \geq f(z) + \nabla f(z)'(y - z), \text{--- -- (2)}$
 - $(1) * \alpha + (2) * (1 - \alpha) \Rightarrow \alpha f(x) + (1 - \alpha)f(y) \geq f(z) + \nabla f(z)'[\alpha(x - z) + (1 - \alpha)(y - z)] = f(z) + \nabla f(z)'[\alpha x + (1 - \alpha)y - z] = f(z)$
- 问题1. 上述定义的几何意义怎么表达? 能画出图形吗?
- 问题2. 当 α 变化时候, $f(x) + \frac{f(x + \alpha(z - x)) - f(x)}{\alpha}$ 有何意义? (参见下页图)

第二章 线性规划

- 问题1：解释下



- 问题2：



这表示什么意思？

考虑 $0 < \alpha_1 < \alpha_2 < 1$, 令 $\bar{\alpha} = \frac{\alpha_1}{\alpha_2}$, $\bar{z} = x +$

$\alpha_2(z - x)$, 则有:

$f(x + \bar{\alpha}(\bar{z} - x)) \leq \bar{\alpha}f(\bar{z}) + (1 - \bar{\alpha})f(x)$ 即:

$\frac{f(x + \bar{\alpha}(\bar{z} - x)) - f(x)}{\bar{\alpha}} \leq f(\bar{z}) - f(x)$, 将 \bar{z} 代入:

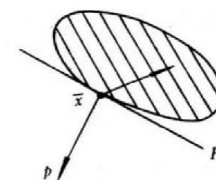
$\frac{f(x + \alpha_1(z - x)) - f(x)}{\alpha_1} \leq \frac{f(x + \alpha_2(z - x)) - f(x)}{\alpha_2}$,

令 $g(\alpha) = \frac{f(x + \alpha(z - x)) - f(x)}{\alpha}$, $\alpha \in (0, 1]$, 则其单调递增!



第二章 线性规划

- 分离定理
- 两个凸集可以用超平面分开！直观上很显然！但由此可以推导出Farkas定理. 介绍!
- 定义：对于非空集合 $C \subset \mathbf{R}^n$, $\bar{x} \in \partial C$ 表示边界上的点，若有：
 - $C \subset H^+ = \{x \in C | p^T(x - \bar{x}) \geq 0\}$ 或
 - $C \subset H^- = \{x \in C | p^T(x - \bar{x}) \leq 0\}$
 - 则称超平面 $H = \{x | p^T(x - \bar{x}) = 0\}$ 是 C 在 \bar{x} 处的支撑超平面。若 C 不包含在 H 中，则 H 称为 C 在 \bar{x} 的正常支撑超平面。
- 从上述凸集和性质，可以看出，凸集在每个边界点都有一个切平面支撑。
- 两个凸集 C_1, C_2 , 若 $p^T x \geq \alpha, \forall x \in C_1, p^T x \leq \alpha, \forall x \in C_2$ ，则称超平面 $H = \{x | p^T x = \alpha\}$ 分离 C_1, C_2 。若都是不等号，则称为严格分离。
- 分离定理：两个不相交凸集存在分离超平面，即存在非零向量 $p \in \mathbf{R}^n$, 使得：
$$p^T x_1 \leq p^T x_2, \forall x_1 \in cl(C_1), \forall x_2 \in cl(C_2)$$
- 并可经一步推广到强分离定理：
$$\inf\{p^T x | x \in C_1\} \geq \epsilon + \sup\{p^T x | x \in C_2\}, p \neq 0, \epsilon > 0$$





第二章 线性规划

- 典型的凸函数，凸函数由于具有很好的性质，在实际应用中颇为广泛
- 问题1. 函数 $f(x) = x^p$, 若 $x > 0, 0 < p < 1$, 函数 $f(x)$ 的凸凹性如何? 若 $x > 0, p > 1$ 呢?
- 问题2. 尝试证明: Minkowski不等式: $\|x + y\|_p \leq \|x\|_p + \|y\|_p, p \geq 1$
- 问题3. 令 $M_\alpha(x) = \left(\frac{x_1^\alpha + x_2^\alpha + \cdots + x_n^\alpha}{n} \right)^{\frac{1}{\alpha}}, \alpha \geq 0$, 则 $M_\alpha(x)$ 是 α 的减函数。
- 问题4. $f(x) = \ln x, x \in (0, +\infty)$ 是 $(0, +\infty)$ 上的凸函数还是凹函数? $f(x) = x \ln x (0 < x < +\infty)$ 呢?



第二章 线性规划

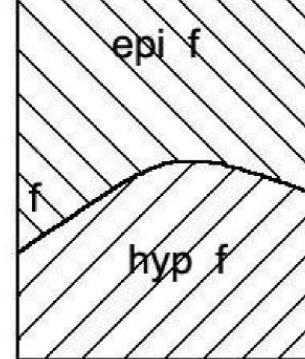
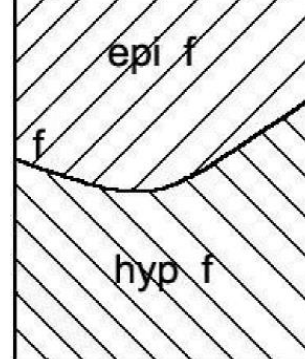
- 很多著名不等式都可以采用类似的办法

- Young不等式: $x, y > 0, p, q > 1, \frac{1}{p} + \frac{1}{q} = 1, x^{\frac{1}{p}} y^{\frac{1}{q}} \leq \frac{x}{p} + \frac{y}{q}$

- Holder不等式: $\alpha > 1, \beta > 1, \frac{1}{\alpha} + \frac{1}{\beta} = 1, a_i > 0, b_i > 0, i = 1, 2, \dots, n$, 则

有: $\sum_{i=1}^n a_i b_i \leq (\sum_{i=1}^n a_i^\alpha)^{\frac{1}{\alpha}} \left(\sum_{i=1}^n b_i^\beta \right)^{\frac{1}{\beta}}$, 且 a_i^α, b_i^β 成正比例时等号成立。

第二章 线性规划



- 水平集(level set)
 - 函数 $f(x)$ 在集合 S 上关于实数 α 的水平集定义为:
$$L(f, \alpha) = \{x \in S \mid f(x) \leq \alpha\}$$
 - 这里 $S \subset R^n, f: R^n \rightarrow R$. 并且函数 $f(x)$ 是凸集 $S \subset R^n$ 上的凸函数时, 水平集也是凸集
- 设 $f(x)$ 是凸集合 $S \subset R^n$ 上的凸函数, 则对 $\forall \alpha \in R^1$, 水平集 $L(f, \alpha)$ 是凸集。
- 函数 f 的上图或上镜图(Epigraph):
$$epi(f) = \left\{ \begin{pmatrix} x \\ \alpha \end{pmatrix} \mid f(x) \leq \alpha, x \in S, \alpha \in R^1 \right\}$$
- 问题1. 函数 $f(x)$ 在凸集 S 上是凸函数, 等价于 f 的上镜图 $epi(f)$ 是 R^{n+1} 上的凸集, 是否正确? (尝试证明)



第二章 线性规划

- 设 $f(x)$ 定义在凸集 $S \subseteq R^n$ 上, $x, y \in S, x \neq y$. 令 $\phi(t) = f(tx + (1-t)y), t \in [0, 1]$. 则 $f(x)$ 是 S 上的凸函数的充要条件为: 对 $\forall x, y \in S, x \neq y$, 一元函数 $\phi(t)$ 是 $[0, 1]$ 上的凸函数。
- 问题1. 上述论断是否成立, 为什么?



若: $C = \{(x_1, x_2) | x_1^2 + (x_2 - 1)^2 \leq 1\}$
问: $\text{cone}(C) = ?$

第二章 线性规划

- 非空集合 $C \subset \mathbb{R}^n$ 的元素的非负组合为:
 $x_i \in C, \alpha_i \geq 0, \sum_{i=1}^m \alpha_i x_i$, 若所有 α_i 为正的, 则称为正组合。
- 定义: 由集合 C 所形成的锥定义为: 集合 C 的元素所形成的所有非负组合集合, 表示为 $\text{cone}(C)$.
注意锥是凸的! 【所有包含集合 C 的凸集的交称为集合 C 的凸包: convex hull】
- 定理: 非空子集 $C \subset \mathbb{R}^n$, $\text{cone}(C)$ 中的任意非零向量都可以表示为 C 中线性无关向量的正组合. 并且锥中的任意向量都可以表示为 C 中不超过 $n + 1$ 个向量的凸组合!
- 问题1. 如何证明这条定理?
 - $x \neq 0, x \in \text{cone}(C)$, 令 m 是满足: $x = \sum_{i=1}^m \alpha_i x_i, \alpha_i > 0, x_i \in C, i = 1, \dots, m$ 的最小正整数。反证法! 若向量 x_i 线性相关, 则存在 $\lambda_i, \sum_{i=1}^m \lambda_i x_i$ 线性相关, 且至少一个 $\lambda_i > 0$. 则线性组合: $\sum_{i=1}^m (\alpha_i - \bar{\gamma} \lambda_i) x_i$, 其中 $\bar{\gamma} = \min_i \left\{ \frac{\alpha_i}{|\lambda_i|} \right\}$, 提供了一个 $m - 1$ 个元素的正组合, 从而矛盾。注: $\alpha_i - \gamma \lambda_i > 0 \Rightarrow$
$$\begin{cases} \gamma < \frac{\alpha_i}{\lambda_i}, \lambda_i > 0 \\ \gamma > \frac{\alpha_i}{\lambda_i}, \lambda_i < 0 \end{cases} \Rightarrow \text{取其交集!}$$
 - 将上述定理应用到子空间: $Y = \{(y, 1) | y \in C\}$ 上即可!

定理的作用是什么? 能将任意非凸集凸化为凸集, 而通过凸化凸函数的上镜图, 则可以将非凸函数转化为凸函数!



第二章 线性规划

- 都说凸集、凸函数有很好的性质，但实际上大部分集合和函数都不是凸的？
- 问题1.如何将非凸集合或函数，转换为凸集或凸函数呢？
- 问题2.听说过共轭函数吗？什么叫共轭函数？ Conjugate function: Fenchel conjugate, 如果 f 可微，又称为Legendre 变换
 - 令 $f: A \rightarrow R$ 是定义在子集 $A \subset R^n$ 上的一个函数。其共轭函数 $f^*: R^n \rightarrow R$ 定义为：
 - $f^*(y) = \sup_{x \in A} (y^T x - f(x)), y \in R^n$
 - 例子：假设 $f = ax + b, x, a, b \in R$, 则 $f^*(y) = \begin{cases} -b, & y = a \\ +\infty, & \text{otherwise} \end{cases}$



第二章 线性规划

- 共轭函数的计算

- $f(x) = -\log x, x \in R, x > 0$

- $f(x) = e^x, x \in R$

- $f(x) = x \log x, x \geq 0$

- $f(x) = \frac{1}{2} x^T A x, A > 0, \text{对称}, x \in R^n$

- $f(X) = \log \det(X^{-1}), X \text{对称正定矩阵}$

- $f(x) = \|x\|, \text{dom}(f) = R^n, \|\cdot\| \text{为任意范数}$

- 问题1. 如何计算共轭函数?

- 问题2. 从几何上如何理解共轭函数?

$$f^*(y) = -\log(-y) - 1, y \in R, y < 0$$

$$f^*(y) = y \log y - y, y \geq 0, 0 \log 0 = 0$$

$$f^*(y) = e^{y-1}, y \in R$$

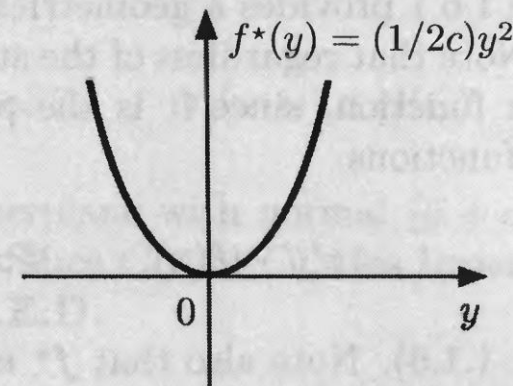
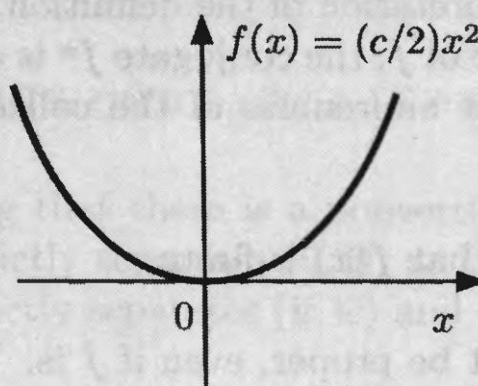
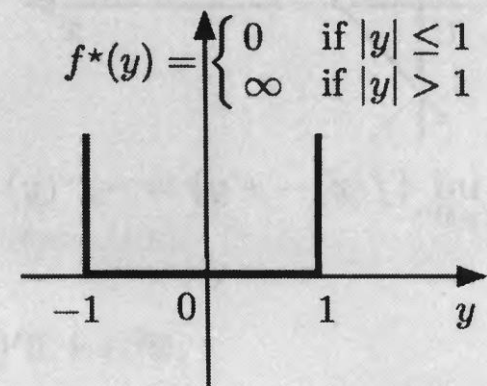
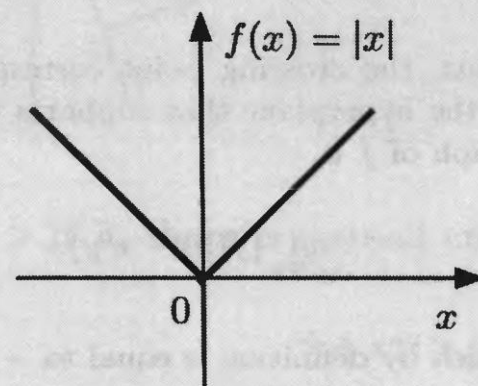
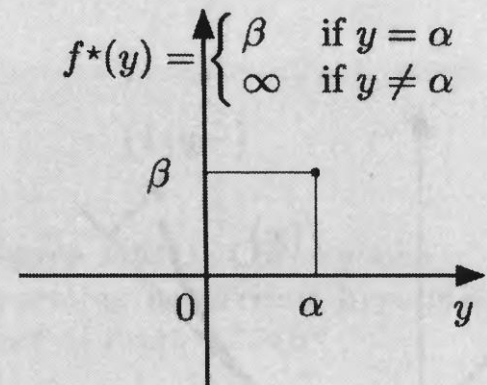
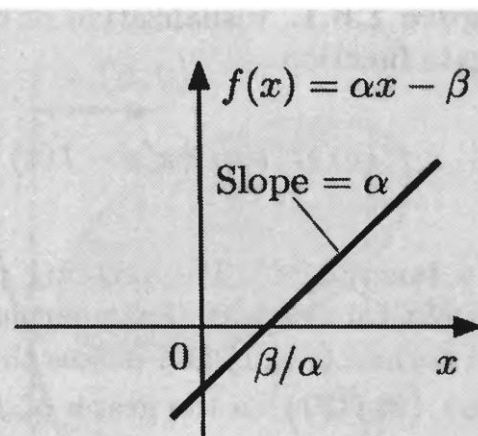
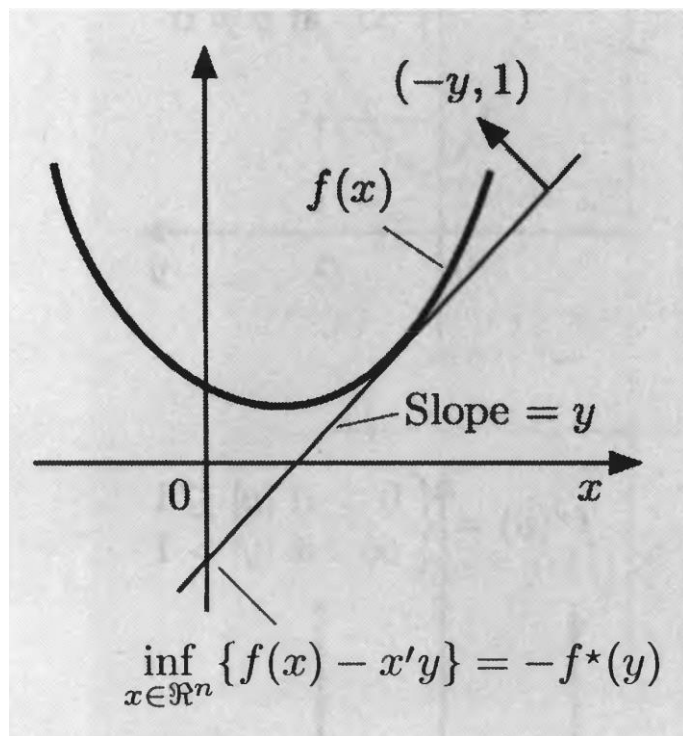
$$f^*(y) = -\frac{1}{2} y^T A^{-1} y$$

$$f(Y) = \log \det((-Y)^{-1}) - n$$



第二章 线性规划

- 几何直观解释





第二章 线性规划

- 线性规划和整数规划有什么联系？
 - $A \in \mathbb{Z}^{m \times n}$, 且 $b \in \mathbb{Z}^m, c \in \mathbb{Z}^n$, 求 $x \in \mathbb{Z}^n$, 使得 $Ax \leq b, \max cx$
- 几乎所有的组合优化问题都可以形式化表达为整数规划问题
 - 可行解表示为 $\{x: Ax \leq b, x \in \mathbb{Z}^n\}$, 集合 $P := \{x \in \mathbb{R}^n: Ax \leq b\}$ 为多面体, 令 $P_I := \{x \in \mathbb{R}^n: Ax \leq b\}$ 为 P 中整数向量的凸包(convex hull), 称为 P 的整数凸包
- 命题: 令 $P = \{x: Ax \leq b\}$ 是整数凸包非空的有理多面体, 令 c 为某一向量 (不一定是 有理数)。则 $\max\{cx: x \in P\}$ 有界 $\Leftrightarrow \max\{cx: x \in P_I\}$ 是有界的。
- 割平面法求解: $P \supset P' \supset P_I$, 期望 $\max\{cx: x \in P'\}$ 在整数向量上达到



第二章 线性规划

- Klee-Minty提出的一类LP问题，每个采用单纯形求解的迭代次数都是指数次迭代
- 如： $maximize \sum_{j=1}^n 10^{n-j} x_j$,
 - s.t. $2 \sum_{j=1}^{i-1} 10^{i-j} x_j + x_i \leq 100^{i-1}, i = 1, \dots, n$
 - $x_j \geq 0, j = 1, \dots, n$
- 对于 $n = 3$ ，其标准形式为：
- $maximize 100x_1 + 10x_2 + x_3$
 $x_1 \leq 1$
- S.t. $20x_1 + x_2 \leq 100$
 $200x_1 + 20x_2 + x_3 \leq 10000$
 $x_1 \geq 0, x_2 \geq 0, x_3 \geq 0$
- 一般来说，这类问题都需要顶点数-1次迭代。 $n = 50, 2^{50} - 1 \approx 10^{15}$, 每秒百万次计算能力也需要大约33年，但已经证明解具有固定折扣率的MDP问题是多项式的。



第三章非线性规划的数学模型

- 线性规划转换

- 基追踪问题

- $\min_{x \in \mathbb{R}^n} \|x\|_1, s. t. Ax = b$

- 法1.引入新变量 z_i , $\min_{(z \in \mathbb{R}^n)} \sum_{i=1}^n z_i, s. t. Ax = b, -z_i \leq x_i \leq z_i, i = 1, 2, \dots, n$

- 法2.引入 $u, v \in \mathbb{R}^n, u \geq 0, v \geq 0$,使得 $x = u - v$,则 $\min_{u, v \in \mathbb{R}^n} \sum_{i=1}^n (u + v), s. t. Au - Av = b, u, v \geq 0$

- 数据拟合

- 最小二乘, 最小 l_1 范数, 最小 l_∞ 范数模型

- $\min_{x \in \mathbb{R}^n} \|Ax - b\|_1, \min_{x \in \mathbb{R}^n} \|Ax - b\|_\infty$

- 引入变量 $y = Ax - b$,则转换为 $\min_{x, y \in \mathbb{R}^n} \|y\|_1, s. t. y = Ax - b$

- 引入 $t = \|Ax - b\|_\infty$, 等价转换为 $\min_{x \in \mathbb{R}^n, t \in \mathbb{R}} t, s. t. \|Ax - b\|_\infty \leq t \Leftrightarrow \min_{x \in \mathbb{R}^n, t \in \mathbb{R}} t, s. t. -tE \leq Ax - b \leq tE.$



第三章 非线性规划的数学模型

- 问题1.我们都说非线性规划，那你认为最简单的非线性规划是什么？
- 问题2.为什么二次规划很重要？ $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{x}^T \mathbf{b} + \mathbf{c}$,有什么实际物理意义吗？这里矩阵 \mathbf{A} 对称！
 - 实际上，这可以看作是一种能量的表达形式，大家回忆，在计算机中是不是很多情况都涉及能量？都是怎么来表示能量的？
 - 物理和工程中的许多问题都可以表达为能量函数的最小化！力学的基本原理就是保证能量最小化（熵增，稳态对应能量最小）
 - 最简单的能量函数就是二次函数！ $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} \pm \mathbf{x}^T \mathbf{b}$
- 问题3.在什么样的条件下， $f(\mathbf{x})$ 有全局最小值，并唯一呢？



第三章 非线性规划的数学模型

- 函数 $f(x) = \frac{1}{2}x^T Ax - x^T b + c$ 当矩阵 A 是对称正定的时候，函数 $f(x)$ 有唯一全局极小点.
- **结论**：二次函数 $f(x) = \frac{1}{2}x^T Ax - x^T b + c$ 中，如果 A 是对称正定的，则 $f(x)$ 的全局唯一极小点是线性方程 $Ax = b$ 的解， f 的极小值为：
$$f(A^{-1}b) = -\frac{1}{2}b^T A^{-1}b + c.$$
- 问题1. 如何证明上述结论？
 - A 对称正定，可逆，因此其特征值均为正，令 $x = A^{-1}b$ ，对 $\forall y \in R^n$ 计算
$$f(y) - f(x) = \frac{1}{2}y^T Ay - y^T b - \frac{1}{2}x^T Ax + x^T b = \frac{1}{2}y^T Ay - y^T Ax + \frac{1}{2}x^T Ax = \frac{1}{2}(y - x)^T A(y - x) \geq 0$$



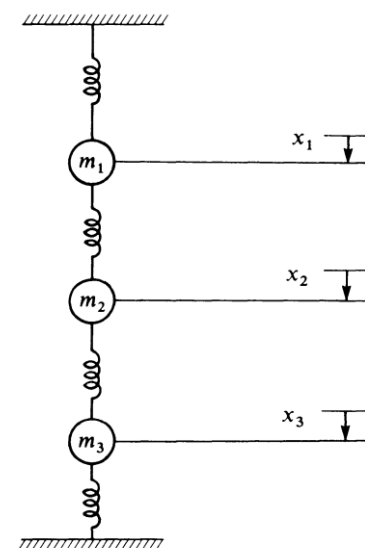
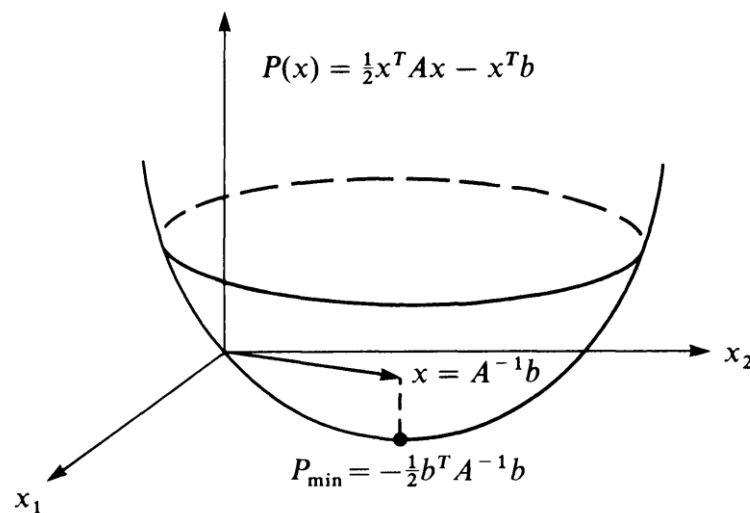
第三章 非线性规划的数学模型

- 上述二次函数最小值的求解等价于求解线性方程组 $A\mathbf{x} = \mathbf{b}$. 这表明也可以将 $A\mathbf{x} = \mathbf{b}$ 表述为二次函数求极值（对应泛函求极值，采用变分法求解）
- 问题 1. 假设二次函数为 $Q(x_1, x_2) = \frac{1}{2}(x_1^2 + x_2^2)$, s. t. $2x_1 - y_1 = 5$, 函数 Q 的最小值点是什么?
 - $(2, -1)$

第三章 非线性规划的数学模型

• 二次型的几何意义

- 最小势能：四个弹簧+3个质块，三个质块的位移 \mathbf{x} : x_1, x_2, x_3 应该是怎么样的？弹簧的延伸长度 \mathbf{e} ，弹簧上的作用力 \mathbf{y} ，以及节点上的外力 \mathbf{f} .
 - $A\mathbf{x} = \mathbf{e}, C\mathbf{e} = \mathbf{y}, A^T\mathbf{y} = \mathbf{f} \Rightarrow A^T C A \mathbf{x} = \mathbf{f}$, 且 $K = A^T C A$ 是对称正定的，物理中称为刚性矩阵(stiffness matrix)
 - 弹簧的势能可以表示为: $\frac{1}{2} \mathbf{x}^T A^T C A \mathbf{x}$, 质块的势能可以表示为: $-\mathbf{x}^T \mathbf{f}$, 关于位移 \mathbf{x} 的系统的整体势能即为: $P(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A^T C A \mathbf{x} - \mathbf{x}^T \mathbf{f}$





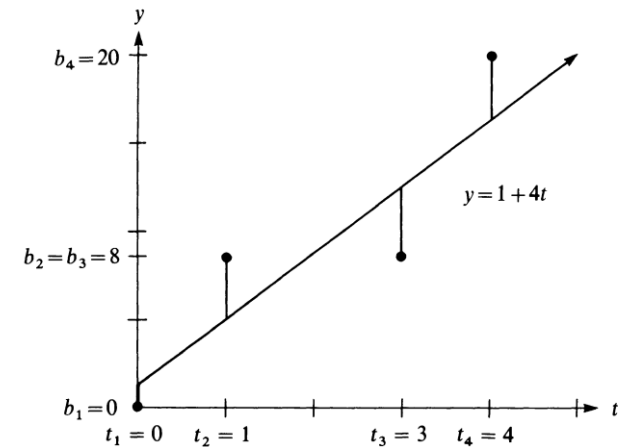
第三章 非线性规划的数学模型

- 最小二乘解 $Ax = b$

- m 个方程, n 个未知数, $m > n$. 超定方程!
- 通过右端项 b 可得: 矩阵 A 的列空间, 该列空间在 m 维子空间当中 (因为所有列都只有 m 个分量)
- 问题1. $Ax = b$ 什么时候有解?

- 若
$$\begin{cases} C + t_1 D = b_1 \\ C + t_2 D = b_2 \\ C + t_3 D = b_3 \\ C + t_4 D = b_4 \\ C + t_5 D = b_5 \end{cases} \Rightarrow \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ \mathbf{1} & \mathbf{2} \\ 1 & 3 \\ 1 & 4 \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \mathbf{b_3} \\ b_4 \\ b_5 \end{bmatrix}$$

- 问题2. 这个方程的几何意义是什么? 什么时候有解?





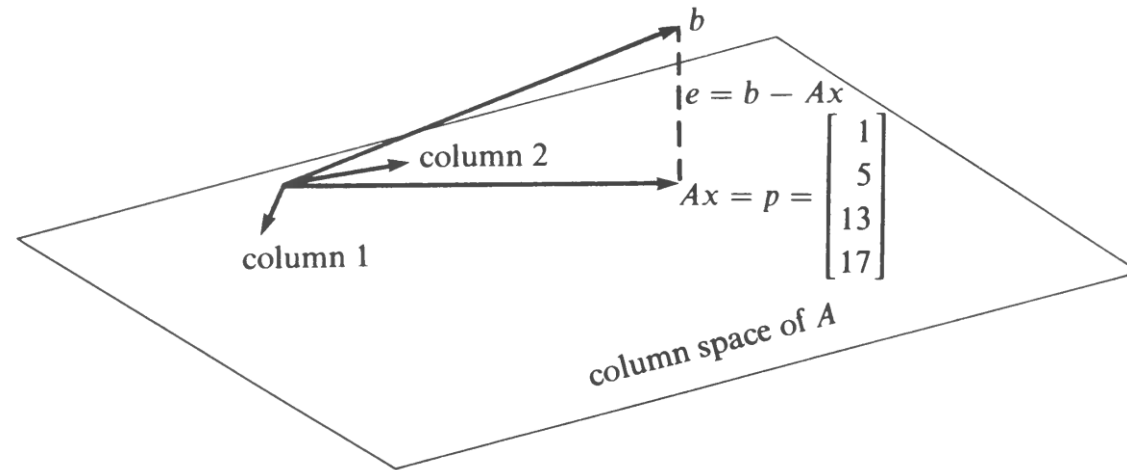
第三章 非线性规划的数学模型

- error误差 $\mathbf{e} = \mathbf{b} - \mathbf{Ax}$
 - 最好的直线就是使误差 \mathbf{e} 尽可能小的直线
 - $\min ||\mathbf{Ax} - \mathbf{b}||^2 = (\mathbf{Ax} - \mathbf{b})^T (\mathbf{Ax} - \mathbf{b})$
 - $||\mathbf{e}||^2 = \mathbf{e}^T \mathbf{e} = e_1^2 + e_2^2 + \dots + e_n^2$
 - 问题1. 哪个向量 \mathbf{x} 能最小化 $(\mathbf{Ax} - \mathbf{b})^T (\mathbf{Ax} - \mathbf{b}) = \mathbf{x}^T \mathbf{A}^T \mathbf{Ax} - \mathbf{x}^T \mathbf{A}^T \mathbf{b} - \mathbf{b}^T \mathbf{Ax} + \mathbf{b}^T \mathbf{b}$?
 - 最后可化简为 $P = \frac{1}{2} \mathbf{x}^T \mathbf{A}^T \mathbf{Ax} - \mathbf{x}^T \mathbf{A}^T \mathbf{b}$, 这跟之前的问题没有本质区别, 只是 $\mathbf{A} \rightarrow \mathbf{A}^T \mathbf{A}, \mathbf{b} \rightarrow \mathbf{A}^T \mathbf{b} \Rightarrow \mathbf{Ax} = \mathbf{b} \rightarrow \mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}$
 - 显然: $\min_x ||\mathbf{Ax} - \mathbf{b}||^2 \Rightarrow \mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}$, 其解为: $\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}$
 - 向量 $\mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$ 是方程 $\mathbf{Ax} = \mathbf{b}$ 的最小二乘解
 - 误差 $\mathbf{e} = \mathbf{b} - \mathbf{Ax}$ 不为0, 但与 \mathbf{A} 的每列向量的内积为0. $\mathbf{A}^T \mathbf{e} = \mathbf{0}$, 或 $\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}$, 表明误差 \mathbf{e} 垂直于 \mathbf{A} 的列空间, 而 \mathbf{b} 分解为



第三章 非线性规划的数学模型

- 而 \mathbf{b} 分解为
- $\mathbf{b} = \mathbf{Ax} +$



- 问题1.最近的点的几何意义?
 - \mathbf{p} 称为 \mathbf{b} 在列空间上的投影!
 - 从代数意义上: 如果 $\mathbf{Ax} = \mathbf{b}$ 无解, 则 $\mathbf{A}^T(\mathbf{Ax}) = \mathbf{A}^T\mathbf{b}$ 再解方程!
 - 这就是线性回归 (Linear Regression)
 - 预条件 (preconditioning)



第三章 非线性规划的数学模型

- 回归问题的表示
 - 给定不同时刻 t_1, \dots, t_m 的测量值 b_1, \dots, b_m ,最小化误差 $\|e\|^2$ 的直线 $y = C + Dt$ 由方程:
 - $A^T A x = A^T b$ 或者 $\begin{bmatrix} m & \sum t_i \\ \sum t_i & \sum t_i^2 \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} \sum b_i \\ \sum t_i b_i \end{bmatrix}$
 - 并且最佳拟合直线方程就是:
 - $y = \bar{b} + D(t - \bar{t}), D = \frac{\sum (t_i - \bar{t}) b_i}{\sum (t_i - \bar{t})^2}$,这条直线通过点中心点: (\bar{t}, \bar{b})
- 问题1.如果矩阵 A 的列不独立, 这时 $A^T A$ 就不是正定的, 不是正定的怎么办? $A^T A x = A^T b$
 - 伪逆(pseudoinverse).实际应用中一般列都是独立的



第三章 非线性规划的数学模型

- Q收敛速率
 - (α, β) 有关，在优化中可用其刻画迭代序列的收敛速度。
- 问题1. 假设有两个序列 $\{x_1^k\}, \{x_2^k\}$, 其 Q 阶和 Q 因子分别为 $(\alpha_1, \beta_1), (\alpha_2, \beta_2)$, 若 $\alpha_1 > \alpha_2$, 那个序列收敛快? ; 若 $\alpha_1 = \alpha_2, \beta_1 < \beta_2$, 这时候那个序列的收敛速度快?
- 一般主要就是 Q -线性, Q -超线性, 以及 Q -二次收敛, 通常算法的收敛速度为超线性或二次收敛, 则称为具有快速收敛速度
- 问题2. 拟牛顿方法的收敛速度属于哪类? 牛顿法的收敛速度属于哪类?



第三章 非线性规划的数学模型

- R收敛 (Root-convergence rate) 速率

- 比Q收敛要弱,令 $\{x^k\} \subset \mathbf{R}^n$ 为任意收敛到 x^* 的序列。令 $R_p = \begin{cases} \limsup_{k \rightarrow \infty} \|x^k - x^*\|^{\frac{1}{k}}, & \text{if } p = 1 \\ \limsup_{k \rightarrow \infty} \|x^k - x^*\|^{\frac{1}{p^k}}, & \text{if } p > 1 \end{cases}$

- 如果 $R_1 = 0$, 则 $\{x^k\}$ 称为 R -超线性收敛到 x^* .

- 如果 $0 < R_1 < 1$, 则 $\{x^k\}$ 称为 R -线性收敛到 x^* .

- 如果 $R_1 = 1$, 则 $\{x^k\}$ 称为 R -次线性收敛到 x^* (sublinearly)

- $R_2 = 0$, 则称为超平方收敛, $0 < R_2 < 1$, R -平方收敛, $R_2 \geq 1$, R -次平方收敛

- 问题1. 如果存在非负标量序列 $\{q^k\}$, 使得 $\|x^k - x^*\| \leq q^k, \forall k$, 并且 $\{q^k\}$ Q-线性收敛到0, 则序列 $\{x^k\}$ 也是 R -线性收敛的, 是否正确? 实际上, 在这种情况下, 根据 $\{q^k\}$ 的Q收敛特性对应R-收敛特性!

- 问题1. R -收敛速率依赖于R-阶 p 和 R -因子 R_p 。 p 越大, 则收敛速度越大? 若 R -阶 p 相同, 则 R -因子 R_p 越小收敛速度越快?



第三章 非线性规划的数学模型

- 令 $\mathcal{X} \subset \mathbf{R}^n, \min f(x), x \in \mathcal{X}$, 迭代法的基本方式
 - $x^{k+1} = x^k + \lambda_k p_k$
- 如果 x^* 为凸规划问题的最优解, 则表明从该点出发的任何方向都不是可行下降方向
 - $S_d(x) = \{d \in \mathbf{R}^n | d^T \nabla f(x) < 0\}$ = 下降方向集合
 - $S_f(x) = \{d \in \mathbf{R}^n | d = x' - x, \forall x' \in \mathcal{X}\}$ = 可行方向集合
- 上述最优性条件表示为:
 - $x^* \in \mathcal{X} \Leftrightarrow S_f(x^*) \cap S_d(x^*) = \emptyset$
- 实际上也等价于: $(x' - x)^T \nabla f(x) \geq 0, \forall x' \in \mathcal{X}$
- 由于 f 的凸性 $\Leftrightarrow (x - y)^T (\nabla f(x) - \nabla f(y)) \geq 0$, 也称凸函数的梯度算子为单调算子!

变分不等式



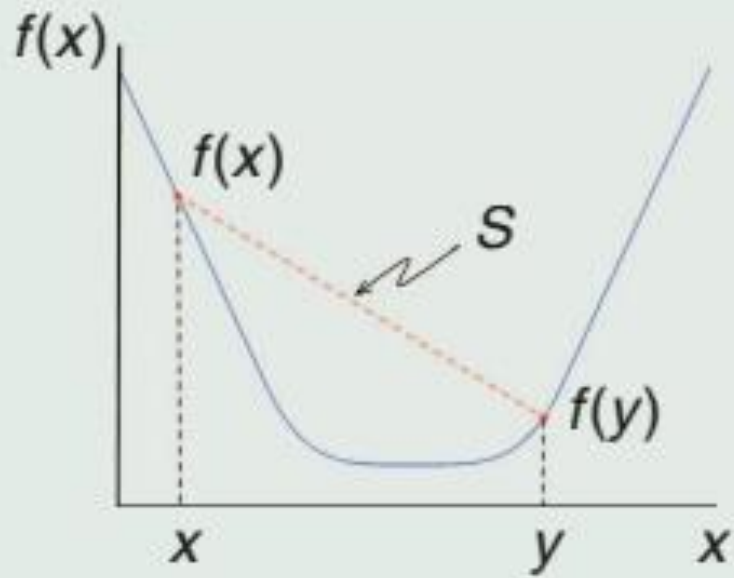
第三章 非线性规划的数学模型

- 变分不等式问题 (Variational Inequality Problem)
 - 给定闭凸集 $\mathcal{X} \subseteq \mathbb{R}^n$ 和映射 $F: \mathcal{X} \rightarrow \mathbb{R}^n$, 变分问题表示为 $VI(\mathcal{X}, F)$: 找到向量 $x^* \in \mathcal{X}$ 使得:
 - $(y - x^*)^T F(x^*) \geq 0, \forall y \in \mathcal{X}$
 - 即找到变分问题的一个解 x^* 。
- 注:
 - 如果某个函数 f , 其梯度 $\nabla f \rightarrow F$, 则显然有: $(y - x^*)^T \nabla f(x^*) \geq 0, \forall y \in \mathcal{X}$, 即满足 x^* 为函数 f 的极小值点!
 - 因此, 在一般意义下, 变分问题的解等价于寻找连续可微凸函数的极值
 - 如果函数 F 不能表示成某些函数的梯度形式, 则变分问题和最优化问题不同, 但 VI 问题包含最优化问题
 - 但如果变分问题函数 F 的 Jacobian 矩阵是对称的, 则 F 可表示为一个函数 f 的梯度形式, 例如 $F(x) = Ax + b$, A 是 $n \times n$ 的对称方阵, 则 $f(x) = \frac{1}{2}(x^T Ax + b^T x)$ 即满足 $F = \nabla f$

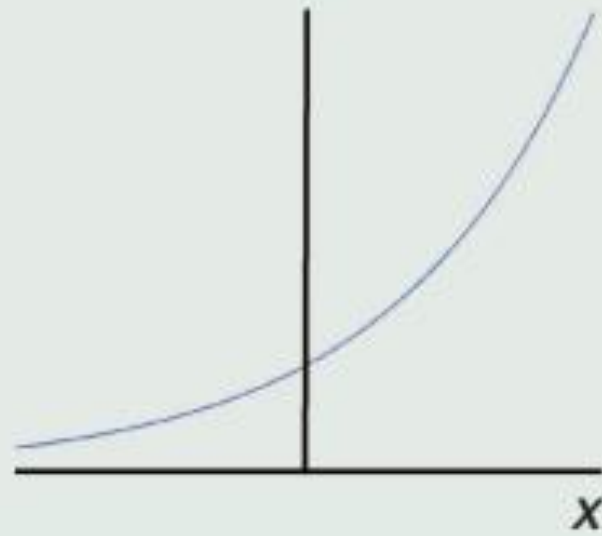


第三章 非线性规划的数学模型

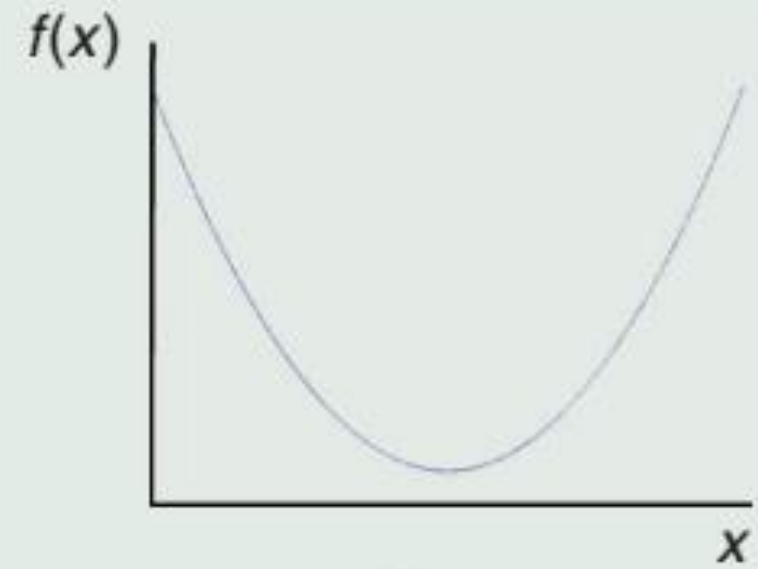
- 若 $\mathcal{X} = \mathbb{R}^n$: 方程组求解, 此时 $VI(\mathbb{R}^n, F)$ 等价于找到一个 $x^* \in \mathbb{R}^n$ 使得 $F(x^*) = 0$, 因为唯一一个向量 $F(x^*)$ 与所有 \mathbb{R}^n 中的向量成非钝角的向量就是 0 向量!
- 若 $\mathcal{X} = \mathbb{R}_+^n$: 非线性互补问题 (Nonlinear complementarity Problem-NCP), 找向量 $x^* \in \mathcal{X}$, 使得: $0 \leq x^* \perp F(x^*) \geq 0$, 等价于 $x_i^* F_i(x^*) = 0, \forall i = 1, 2, \dots, n$
- 变分不等式解的存在性和唯一性如何?
 - 可以从最优化问题的解所必须满足的条件出发去思考?



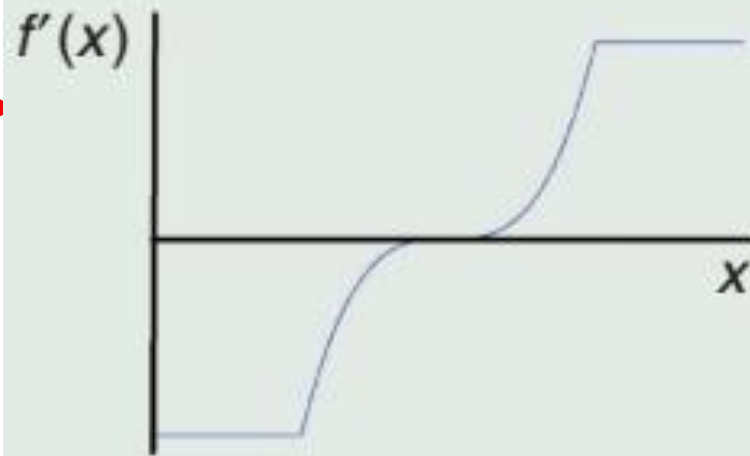
(a)



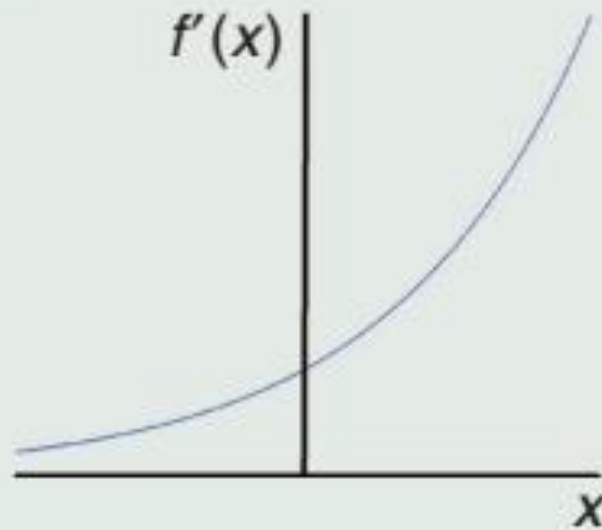
(b)



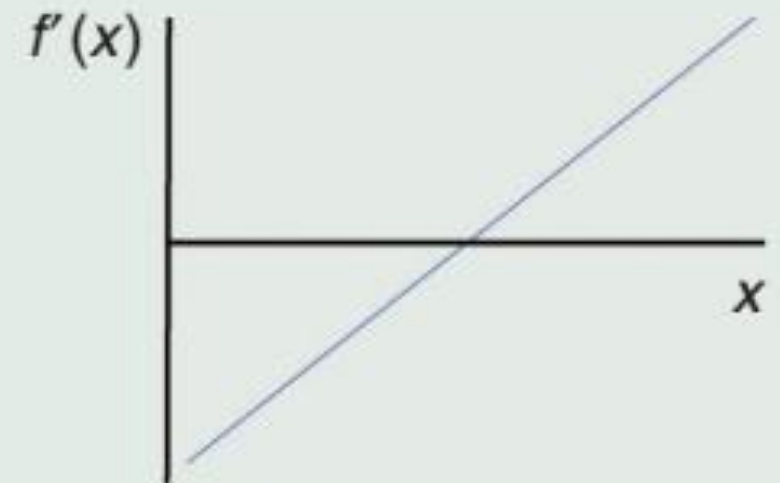
(c)



(d)



(e)



(f)



第三章 非线性规划的数学模型

- 强单调 \rightarrow 严格单调 \rightarrow 单调，反之不然
- 这种单调性和函数 F 的Jacobian矩阵的半正定性之间存在联系！
- 对于仿射函数, $F(x) = Ax + b, A_{n \times n}$ 为矩阵, b 为 n 维向量, 以下结果成立:
- $F(x) = Ax + b$ 是单调的, 当且仅当 A 是半正定的。而对正定矩阵而言, 严格单调和强单调等价。如果向量函数 F 是某个标量函数 f 的梯度 ∇f , 上述的单调性与函数 f 的凸性相关
 - f 是凸的 $\Leftrightarrow \nabla f$ 单调
 - f 是严格凸的 $\Leftrightarrow \nabla f$ 严格单调
 - f 是强凸的 $\Leftrightarrow \nabla f$ 强单调

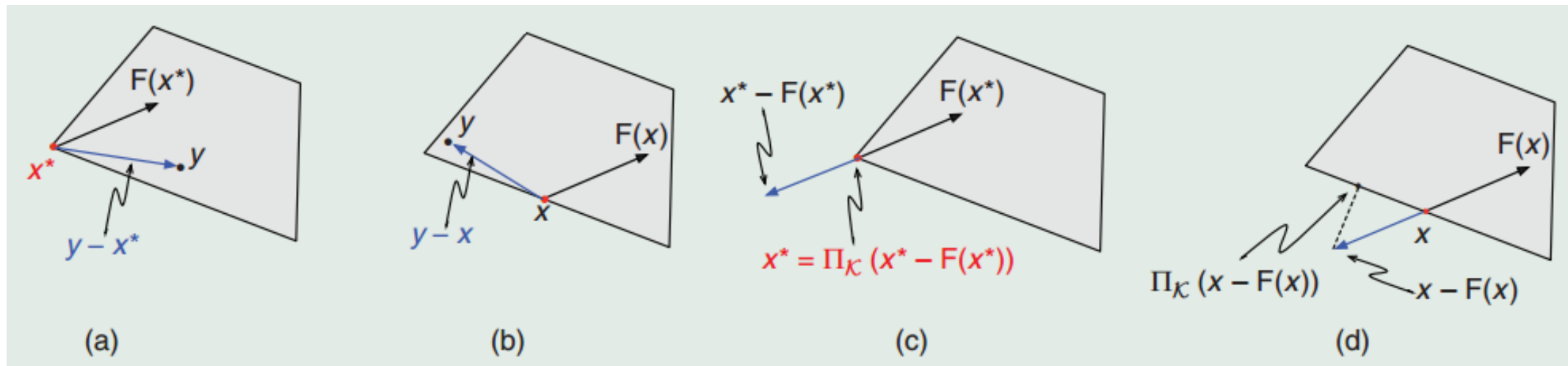


第三章 非线性规划的数学模型

- 基于上述结论，在不要要求集合 \mathcal{X} 的紧性情况下，变分问题 $VI(\mathcal{X}, F)$ 的解具有如下的性质（注意 F 在集合 \mathcal{X} 上的连续性）
 - 若 F 在 \mathcal{X} 上单调，则 $VI(\mathcal{X}, F)$ 的解集是闭且凸的
 - 若 F 在 \mathcal{X} 上严格单调，则 $VI(\mathcal{X}, F)$ 存在最多一个解
 - 若 F 在 \mathcal{X} 上强单调，则 $VI(\mathcal{X}, F)$ 存在唯一解
- 其中第2条表明，严格单调并不能保证其解存在，例如 $F(x) = e^x$ 严格单调，但 $VI(\mathbb{R}, e^x)$ 无解
- 同时，上述结果可以引导出凸优化解的存在性和唯一性
 - 例如，如果如果 f 强凸，则 $\min f(x), s.t. x \in \mathcal{X}$ 有唯一解，这与下述论述等价：如果 ∇f 强单调， $VI(\mathcal{X}, \nabla f)$ 存在唯一解

第三章 非线性规划的数学模型

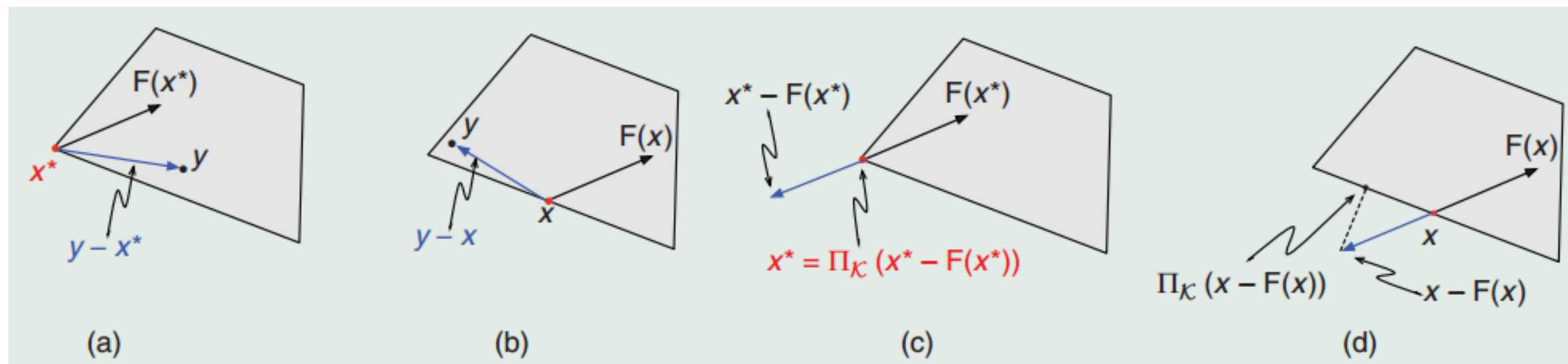
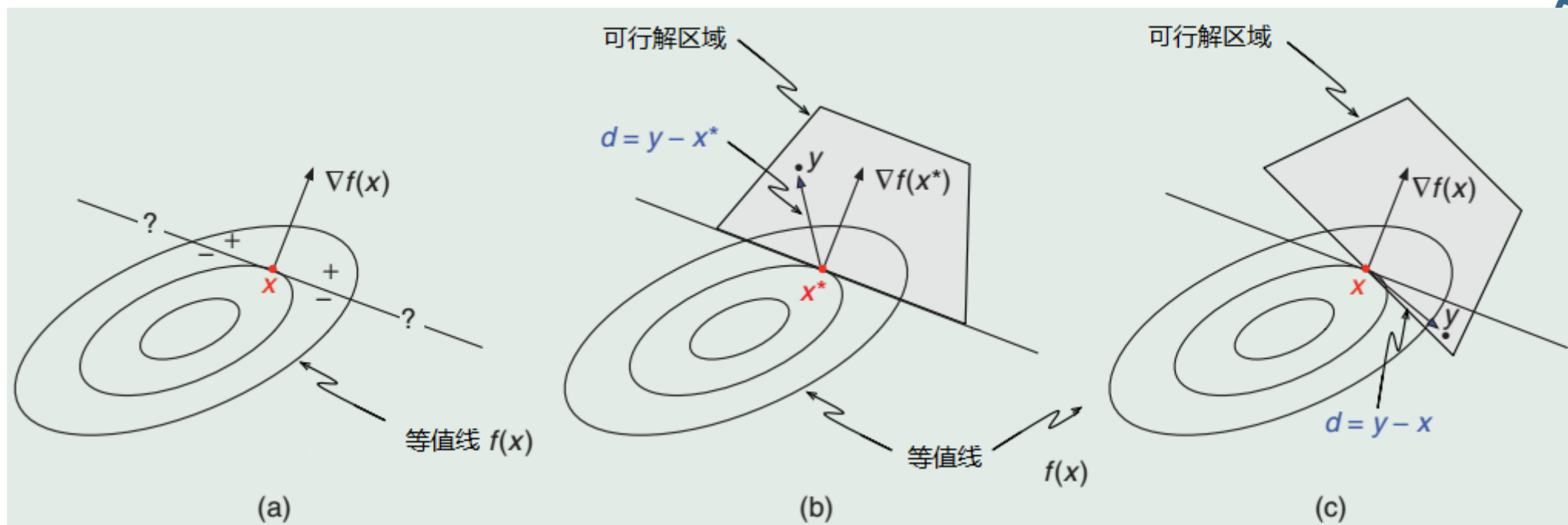
- 投影定理：向量 x_0 在闭凸集 \mathcal{X} 的欧几里德投影表示为 $\Pi_{\mathcal{X}}(x_0)$,表示 \mathcal{X} 中的与 x_0 的欧几里德范数意义上最近的唯一向量，采用优化的观点表示为： $\min_y ||y - x_0||^2, s.t. y \in \mathcal{X}$ 的解，目标函数是强凸的，因此其解存在且唯一。
- $VI(\mathcal{X}, F)$ 可以形式化为一个不动点问题：
- x^* 是变分问题 $VI(\mathcal{X}, F)$ 的解 $\Leftrightarrow x^* = \Pi_{\mathcal{X}}(x^* - F(x^*))$





第三章 非线性规划的数学模型

- 传统使用优化方法来求解的问题，博弈模型也越来越广泛的用来解决这类问题，尤其是当参与者之间的交互是不可忽略的情形下，单纯的中心化方法不再适合！
- 非合作博弈（Noncooperative games）
 - 纳什均衡（Nash Equilibrium problem: NEPs）和广义纳什均衡（GNEPs）
 - 前者只在目标函数级别上进行交互，后者会考虑对手的行为
- NEPs
 - 假设 Q 个玩家分别控制变量 $x_i \in \mathbb{R}^{n_i}$,其中用 $x \triangleq (x_1, x_2, \dots, x_Q)$, 用 $x_{-i} \triangleq (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_Q)$ 表示除第 i 个玩家外所有其他玩家对应的变量。给定其他玩家的策略 x_{-i} , 玩家 i 的目标是选择一个 $x_i \in Q_i$ 来最小化其支付函数 $f_i(x_i, x_{-i})$, 即
 - $\min_{x_i} f_i(x_i, x_{-i}), s.t. x_i \in Q_i$
 - 大概来看，一个NEP就是一个耦合的优化问题集合。假设 f_i 连续可微，单作为 x_i 的函数，是凸的集合 $Q_i \subseteq \mathbb{R}^{n_i}$ 是闭凸集。如果 $x_i \in Q_i$ 对所有玩家 i 都成立，则点 x 称为可行点。一个纯粗略NE或简单一个NEP的解是可行点 x^* ：
 - $f_i(Ix_i^*, x_{-i}^*) \leq f_i(x_i, x_{-i}^*), \forall x_i \in Q_i$, 对所有每一个玩家 $i = 1, 2, \dots, Q$ 成立。
 - 一个NE是具有以下性质的一个可行策略档 x^* ：如果所有玩家都遵循这个策略，则没有单个玩家能够从 x_i^* 的单边偏离中获益。





第三章 非线性规划的数学模型

- **变分表达**: 设 $\mathcal{X} \subset \mathbf{R}^n$ 是闭凸集, $\theta(x), \phi(x)$ 都是 $\mathbf{R}^n \rightarrow \mathbf{R}$ 上的凸函数, 如果 $\phi(x)$ 可微并且最优化问题: $\min\{\theta(x) + \phi(x) | x \in \mathcal{X}\}$ 有解, 则 $\tilde{x} \in \operatorname{argmin}\{\theta(x) + \phi(x) | x \in \mathcal{X}\}$ 的充要条件是:

$$\tilde{x} \in \mathcal{X}, \theta(x) - \theta(\tilde{x}) + (x - \tilde{x})^T \nabla \phi(\tilde{x}) \geq 0, \forall x \in \mathcal{X}$$

- 问题1. 如何证明?

- 必要性: x^* 存在, 令 $x_\alpha = (1 - \alpha)x^* + \alpha x, \forall \alpha \in (0, 1]$, 则有, $\frac{\theta(x_\alpha) - \theta(x^*)}{\alpha} + \frac{\phi(x_\alpha) - \phi(x^*)}{\alpha} \geq 0$
- $\theta(x)$ 是凸函数, $\theta(x_\alpha) \leq (1 - \alpha)\theta(x^*) + \alpha\theta(x)$, 从而: $\theta(x) - \theta(x^*) \geq \frac{\theta(x_\alpha) - \theta(x^*)}{\alpha}, \forall \alpha \in (0, 1]$, 因此
- $\theta(x) - \theta(x^*) + \frac{\phi(x_\alpha) - \phi(x^*)}{\alpha} \geq 0$, 由于 $\phi(x_\alpha) = \phi(x^* + \alpha(x - x^*))$, 若令 $\alpha \rightarrow 0_+$, 立即可得 $\theta(x) - \theta(\tilde{x}) + (x - \tilde{x})^T \nabla \phi(\tilde{x}) \geq 0, \forall x \in \mathcal{X}$
- 充分性: $\phi(x_\alpha) \leq (1 - \alpha)\phi(x^*) + \alpha\phi(x) \Rightarrow \phi(x_\alpha) - \phi(x^*) \leq \alpha(\phi(x) - \phi(x^*)) \Rightarrow \phi(x) - \phi(x^*) \geq \frac{\phi(x_\alpha) - \phi(x^*)}{\alpha} = \frac{\phi(x^* + \alpha(x - x^*)) - \phi(x^*)}{\alpha}, \forall \alpha \in (0, 1]$, 两边令 $\alpha \rightarrow 0_+ \Rightarrow \phi(x) - \phi(x^*) \geq \nabla \phi(x^*)^T (x - x^*) \Rightarrow \theta(x) - \theta(x^*) + \phi(x) - \phi(x^*) \geq 0$



第三章 非线性规划的数学模型

- 具有线性等式约束的凸优化问题求解
 - $\min\{\theta(x)|Ax = b, x \in \mathcal{X}\}$, 其中函数 $\theta(x)$ 为凸函数, $A_{m \times n}$
- 问题1. 该问题的拉格朗日乘子法的形式应该是什么?
 - $L(x, \lambda) = \theta(x) - \lambda^T(Ax - b)$
 - 点 (x^*, λ^*) 称为拉格朗日函数的鞍点, 如果满足下式:
$$L_{\lambda \in R^m}(x^*, \lambda) \leq L(x^*, \lambda^*) \leq L_{x \in \mathcal{X}}(x, \lambda^*)$$
- 该条件等价于如下条件:
$$\begin{cases} x^* \in \mathcal{X}, & \theta(x) - \theta(x^*) + (x - x^*)^T(-A^T \lambda^*) \geq 0, \forall x \in \mathcal{X} \\ \lambda^* \in R^m, & (\lambda - \lambda^*)^T(Ax^* - b) \geq 0, \forall \lambda \in R^m \end{cases}$$
- 若令 $w = \begin{pmatrix} x \\ \lambda \end{pmatrix}$, $F(w) = \begin{pmatrix} -A^T \lambda \\ Ax - b \end{pmatrix}$, $w \in \mathcal{X} \times R^m$



第三章 非线性规划的数学模型

- 此时最优性条件可表达为变分不等式:
- $\mathbf{w}^* \in \mathcal{X} \times \mathbf{R}^m, \boldsymbol{\theta}(\mathbf{x}) - \boldsymbol{\theta}(\mathbf{x}^*) + (\mathbf{w} - \mathbf{w}^*)^T \mathbf{F}(\mathbf{w}^*) \geq \mathbf{0}, \forall \mathbf{w} \in \mathcal{X} \times \mathbf{R}^m$
- 注意算子 \mathbf{F} 的单调性, 因为
 - $(\mathbf{w} - \tilde{\mathbf{w}})^T (\mathbf{F}(\mathbf{w}) - \mathbf{F}(\tilde{\mathbf{w}})) \geq \mathbf{0}$, 实际上
 - $(\mathbf{w} - \tilde{\mathbf{w}})^T (\mathbf{F}(\mathbf{w}) - \mathbf{F}(\tilde{\mathbf{w}})) = 0$

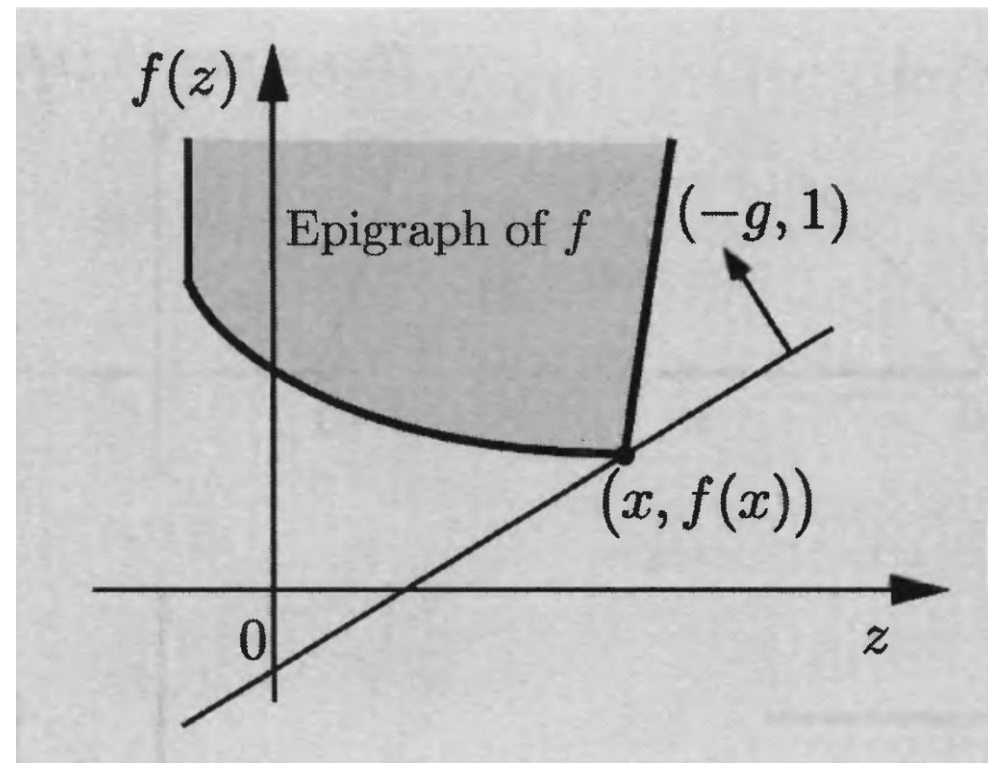
第三章 非线性规划的数学模型

- 次梯度(subgradient)
 - 梯度的替代项, 不可微的时候用次梯度
 - $f(z) \geq f(x) + g^T(z - x)$

g 是函数 f 在点 x 的次梯度等价于

$$\begin{aligned} f(z) - g^T z \\ \geq f(x) - g^T x \end{aligned}$$

在 $n + 1$ 维, 过点
 $(x, f(x))$, 法向量为
 $(-g, 1)$ 的超平面支撑函
数 f 的上镜图(epigraph)





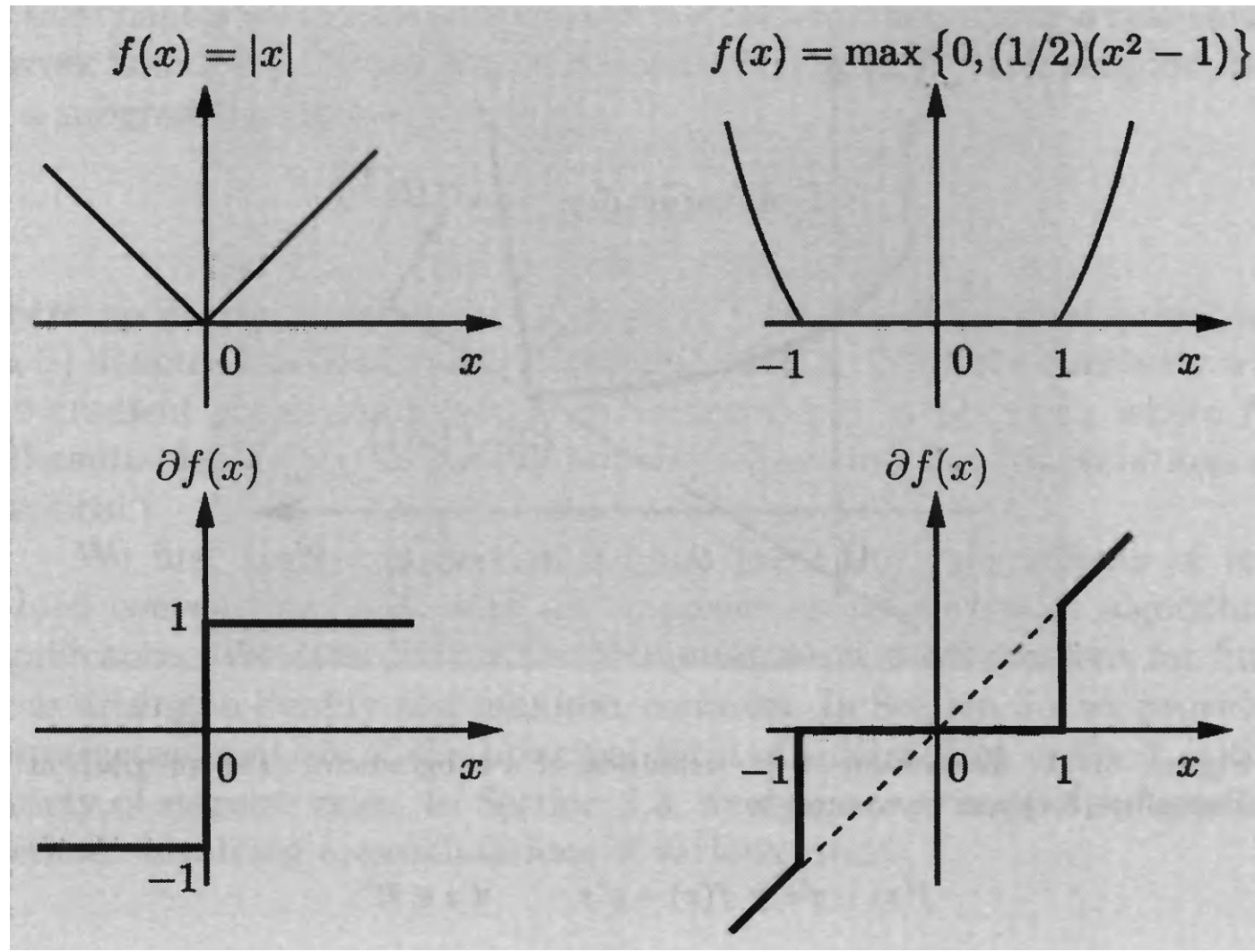
第三章 非线性规划的数学模型

- 次梯度的计算

向量 \mathbf{x} 是凸函数 f 在的极小值点等价于存在一个次梯度 g , 使得

$$g^T(z - x) \geq 0, \forall z \in \mathcal{X}$$

在 \mathbf{x} 点处的所有次梯度的集合称为该点的次微分 $\partial f(\mathbf{x})$





第三章 非线性规划的数学模型

- 临近点算法PPA (Proximal Point Algorithms)
 - 本章介绍了优化算法的一般迭代格式 $\mathbf{x}^{k+1} = \mathbf{x}^k + \lambda_k \mathbf{p}_k$
 - 多次强调在步长和方向上的变化!
- 问题1.你觉得除了步长和方向上有变化外,还能有哪些变化或扩展?
 - $\mathbf{x}^{k+1} = \mathbf{P}_\chi(\mathbf{x}^k - \lambda_k \mathbf{p}_k)$,其中 \mathbf{p}_k 可以是任何可行的方向, 梯度方向, 次梯度方向, λ_k 是正的步长, \mathbf{P}_χ 表示在可行集合上的欧式投影
 - 方向可以是次梯度, 此时对应次梯度方法
 - 若将次梯度放松, $f(\mathbf{z}) \geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{z} - \mathbf{x}) - \epsilon, \epsilon > 0$,则有 ϵ -次梯度方法



第三章 非线性规划的数学模型

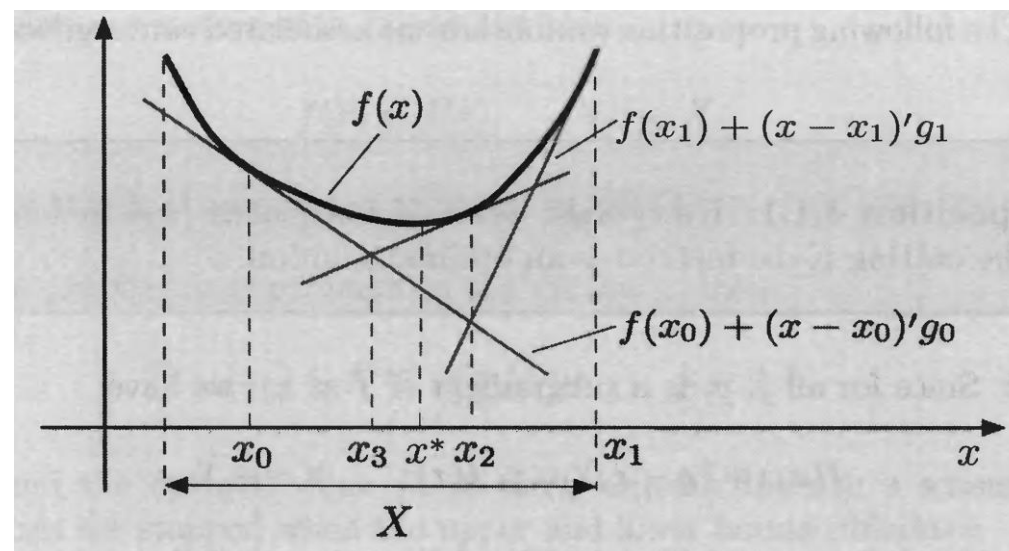
- 多边形近似算法 (Polyhedral Approximation Algorithms)
 - 在迭代算法中每次生成的 x^{k+1} 满足:
 - $x^{k+1} \in \operatorname{argmin}_{x \in \mathcal{X}_k} F_k(x)$
 - 其中 F_k 是近似 f 的一个多面体函数, \mathcal{X}_k 是近似 \mathcal{X} 的一个多面体集合
- 从而将原来问题变为逼近问题, 多面体结构比原问题易于求解, 从而解决原问题
 - 外部线性化方法-割平面法
 - 凸函数的求解中可以采用迭代点列的割平面 (支撑) 所围成的区域来近似



第三章 非线性规划的数学模型

- 外部线性化方法

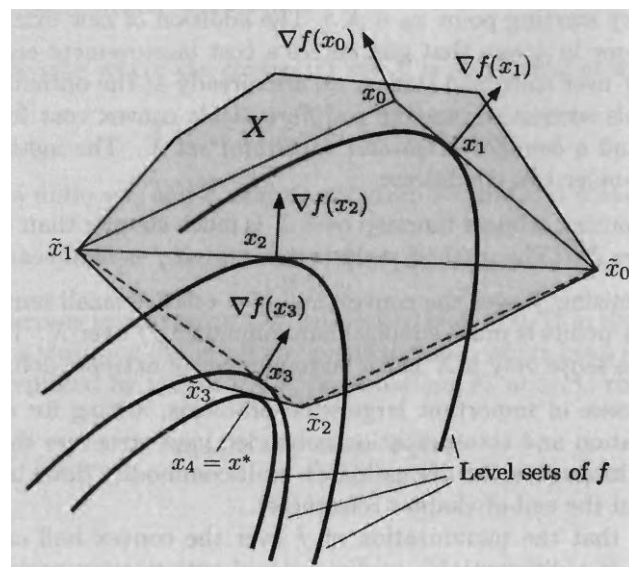
- 凸函数 $f: \mathbf{R}^n \rightarrow \mathbf{R}$, 初始点 \mathbf{x}^0 , 次梯度 $\mathbf{g}_0 \in \partial f(\mathbf{x}_0)$, 其迭代格式表示为:
- $\min F_k(\mathbf{x}), s. t. \mathbf{x} \in \mathcal{X}$
- f 采用多面体近似, 由已有 $\mathbf{x}^0, \mathbf{x}^1, \dots, \mathbf{x}^k$ 来进行近似, 且相关的次梯度为 $\mathbf{g}_i \in \partial f(\mathbf{x}_i)$, 例如
- $F_k(\mathbf{x}) = \max\{f(\mathbf{x}^0) + (\mathbf{x} - \mathbf{x}^0)^T \mathbf{g}_0, \dots, f(\mathbf{x}^k) + (\mathbf{x} - \mathbf{x}^k)^T \mathbf{g}_k\}$, 然后计算
- $\mathbf{x}^{k+1} \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} F_k(\mathbf{x})$





第三章 非线性规划的数学模型

- 最小化闭凸集上的凸函数 f
 - 通过有限集合 $X_k \subset \mathcal{X}$ 的凸包来逼近可行域 \mathcal{X} ,其中 X_k 包括 \mathcal{X} 的极点和任一个初始点 $x^0 \in \mathcal{X}$
- 基本步骤
 - 令初始点 x^0 构成初始集合 $X_0 = \{x^0\}$,求解如下问题生成 \tilde{x}^k 作为 \mathcal{X} 的极点:
 - $\min \nabla f(x^k)^T (x - x^k), s. t. x \in \mathcal{X}$
 - 然后将 \tilde{x}^k 加入到集合 X_k 中, $X_{k+1} = \{\tilde{x}^k\} \cup X_k$,求解下面的优化问题得到 x^{k+1} :
 - $\min f(x), s. t. x \in \text{conv}(X_{k+1})$

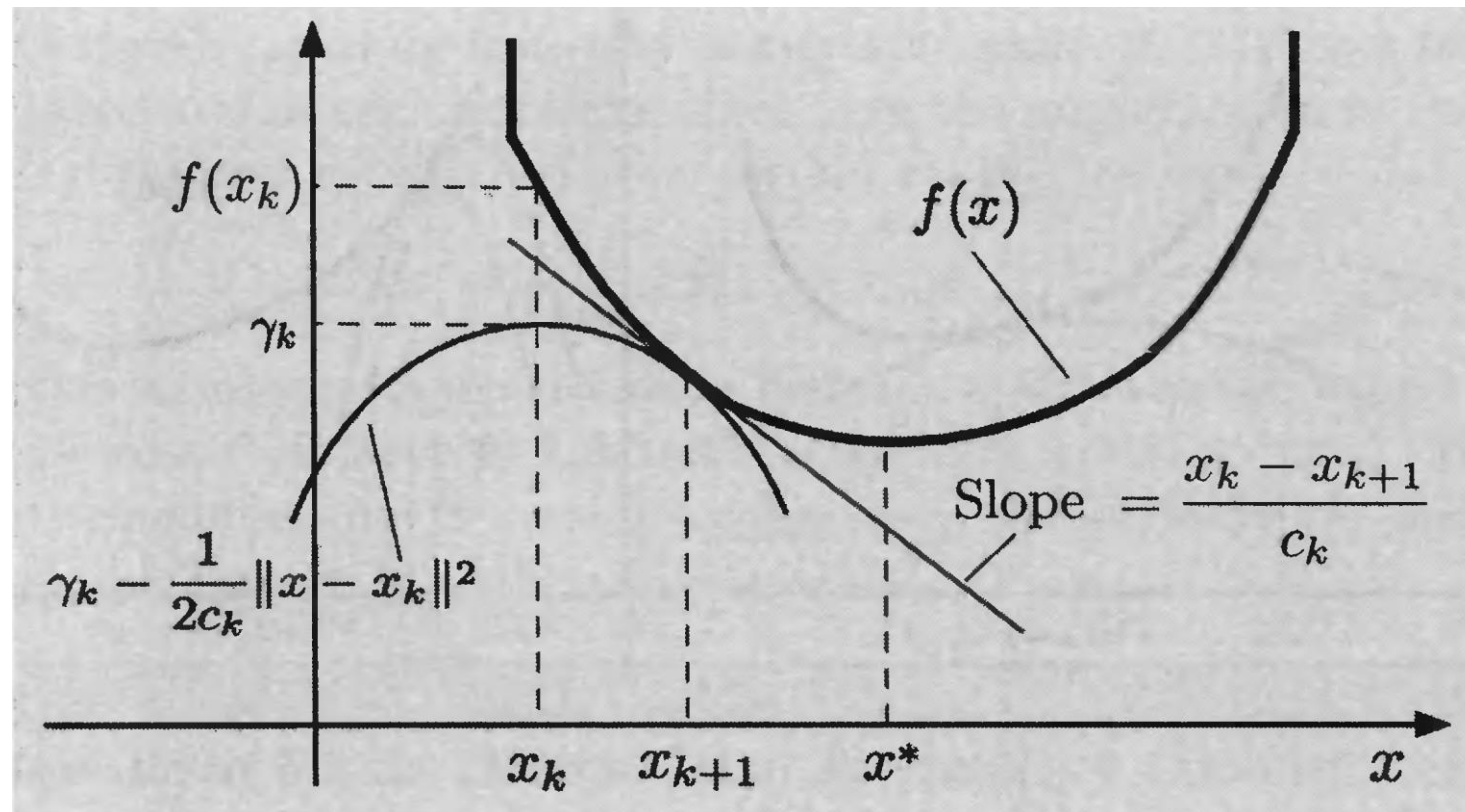




第三章 非线性规划的数学模型

- 添加正则项求解凸规划

- $x^{k+1} \in \operatorname{argmin}_{x \in R^n} \left\{ f(x) + \frac{1}{2c_k} \|x - x^k\|^2 \right\}$, 其中初始点任意, $c_k > 0$ 为标量。



第四章 一维搜索与无约束最优化



- 黑箱模型
 - 假设计算资源无限，约束集 \mathcal{X} 已知，目标函数 $f: \mathcal{X} \rightarrow \mathbb{R}$ 未知，但可以通过oracle查询
 - 零阶的oracle接受 $\mathbf{x} \in \mathcal{X}$ 作为输入，输出函数 f 在点 \mathbf{x} 处的函数值
 - 一阶的oracle接受 $\mathbf{x} \in \mathcal{X}$ 作为输入，输出函数 f 在点 \mathbf{x} 处的次梯度
 - 凸优化的oracle复杂性：必须经过多少对oracle的查询才足以找到凸函数的 ϵ -近似极小值
 - 能够导出一个完整的凸优化理论，获得匹配各种有趣凸函数子类的oracle复杂度上下界
 - 模型本身并不限制计算资源，允许对约束集的任何操作，但会注意算法的计算复杂性（即算法需要执行基本操作的数量）
 - 如果约束集合 \mathcal{X} 是未知的，并且只能通过分离oracle得到：给定 $\mathbf{x} \in \mathbb{R}$, $\Rightarrow \mathbf{x} \in \mathcal{X}$ 或者 $\mathbf{x} \notin \mathcal{X}$, 那么它输出 \mathbf{x} 和 \mathcal{X} 之间的分离超平面
 - 开发维度无关的oracle复杂性算法是可能的，对高维优化问题非常有意义
 - 在黑箱模型中开发的算法对oracle输出中的噪声具有鲁棒性，这对于随机优化特别有意义，并且与机器学习应用紧密相关
- 结构性优化，试图考虑约束集和目标函数的全局结构，如内点法

第四章 一维搜索与无约束最优化



• 通用迭代算法的复杂性

输入：初始点 x_0 和精度 $\epsilon > 0$

初始化：令 $k = 0, \psi_{-1} = \emptyset$, 这里 k 是迭代计数, ψ_k 是累积的信息集

主循环：

1. 在点 x_k 处调用Oracle \mathcal{O}
2. 更新信息集： $\psi_k = \psi_{k-1} \cup (x_k, \mathcal{O}(x_k))$
3. 将方法 \mathcal{M} 的规则应用于 ψ_k , 生成一个新点 x_{k+1}
4. 检验停止准则 \mathcal{T}_ϵ : 如果满足停止准则, 则输出 \bar{x} ; 否则置 $k := k + 1$, 转到第1步

引入两种度量准则来衡量算法 \mathcal{M} 求解优化问题 \mathcal{P} 的复杂度

解析复杂度 (Analytical Complexity): 为使问题 \mathcal{P} 达到精度 ϵ , 需要调用Oracle的次数

算术复杂度 (Arithmetical Complexity): 为使问题 \mathcal{P} 达到精度 ϵ , 需要的算术运算总量 (包括Oracle的调用计算量和算法 \mathcal{M} 的计算量)。



第四章 一维搜索与无约束最优化

- 了解优化方法的复杂性吗？
- 问题 \mathcal{P} : $\min_{\mathbf{x} \in \mathbf{B}_n} \mathbf{f}(\mathbf{x}), \mathbf{B}_n = \{\mathbf{x} \in \mathbb{R}^n | 0 \leq x_i \leq 1, i = 1 \cdots n\}$
 - 假设距离为 $l_\infty := \|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|$
 - 目标函数 $\mathbf{f}(\cdot): \mathbb{R}^n \rightarrow \mathbb{R}$ 是在 \mathbf{B}_n 上 Lipschitz 连续的: $|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})| \leq L \|\mathbf{x} - \mathbf{y}\|_{(\infty)}, \forall \mathbf{x}, \mathbf{y}, \in \mathbf{B}_n, L$ 为 Lipschitz 常数
- 假设采用均匀网格方法求解 \mathcal{P} , 其输入参数为整数 $p \geq 1$
 - $\mathbf{x}_\alpha = \left(\frac{2i_1-1}{2p}, \frac{2i_2-1}{2p}, \dots, \frac{2i_n-1}{2p} \right)^T, \alpha \equiv (i_1, \dots, i_n) \in \{1, \dots, p\}^n$
 - 在所有点 \mathbf{x}_α 上求具有最小目标函数值的点 $\bar{\mathbf{x}}$
 - 方法输出为 $(\bar{\mathbf{x}}, \mathbf{f}(\bar{\mathbf{x}}))$



第四章 一维搜索与无约束最优化

该算法在 B_n 内形成测试点的均匀网格，在网格上计算目标的最优值，并将此值作为问题 \mathcal{P} 的近似解。属于零阶迭代算法，现在看看其效率估计。

- **定理：**若 f^* 为全局最优解，则 $f(\bar{x}) - f^* \leq \frac{L}{2p}$
 - 对多索引 $\alpha \equiv (i_1, \dots, i_n)$ ，定义 $X_\alpha = \{x \in R^n: \|x - x_\alpha\|_\infty \leq \frac{1}{2p}\}$ ，显然 $\bigcup_{\alpha \in \{1, \dots, p\}^n} X_\alpha = B_n$ ，由 x^* 是全局解，存在多索引 α^* 使得 $x^* \in X_{\alpha^*}$ 。注意到 $\|x^* - x_{\alpha^*}\|_\infty \leq \frac{1}{2p}$ ，从而得证。
- **推论：**假设原问题变为：求 $\bar{x} \in B_n: f(\bar{x}) - f^* \leq \epsilon$ ，则有：上述问题的解析复杂性最多为 $\left(\left\lfloor \frac{L}{2\epsilon} \right\rfloor + 1\right)^n$ 。令 $p = \left\lfloor \frac{L}{2\epsilon} \right\rfloor + 1$ ，则 $p \geq \frac{L}{2\epsilon} \Rightarrow f(\bar{x}) - f^* \leq \frac{L}{2p} \leq \epsilon$ 。
- 确定了问题类的复杂度上界。存在的问题1.证明粗糙，实际性能可能会更好；2.不能确定是否算法就是解决问题的合理方法，可能存在更好的。**下界？**
- **定理：**对于 $\epsilon < \frac{1}{2}L$ ，问题的解析复杂度至少为 $\left\lfloor \frac{L}{2\epsilon} \right\rfloor^n$ 次调用oracle。（令 $p = \left\lfloor \frac{L}{2\epsilon} \right\rfloor (\geq 1)$ ）
- 对于上述均匀网格法的性能，将其效率估计值与下界进行比较 $\left(\left\lfloor \frac{L}{2\epsilon} \right\rfloor + 1\right)^n \Leftrightarrow \left\lfloor \frac{L}{2\epsilon} \right\rfloor^n$ ，如果 $\epsilon \leq O\left(\frac{L}{n}\right)$ ，则除了一个绝对常数乘子的意义下，下界和上界是一致的。这意味着，对于这个精度，算法对该问题类来说是最优的



第四章 一维搜索与无约束最优化

- 考虑上述问题参数为： $L = 2, n = 10, \epsilon = 0.01$. 问题规模非常小，且只要求适中的精度1%。
 - 该问题的复杂度下界是Oracle的 $\left\lfloor \frac{L}{2\epsilon} \right\rfloor^n$ 次调用，对于这个例子，看看具体值
 - 下界：Oracle 的 10^{20} 次调用
 - Oracle的复杂度：至少 n 次算术运算
 - 总体复杂度： 10^{21} 次算术运算
 - 处理器性能：每秒 10^6 次算术运算
 - 总时间： 10^{15} 秒
 - 一年：不超过 3.2×10^7 秒
 - 需要：31250000年，即使处理器达到 10^8 , $n=11$ 时仍然成立！
 - 与组合优化中的NP难问题的复杂度 比较，结果也令人沮丧，为找到精确解，最难的组合问题只需要 2^n 次算术运算！



第四章 一维搜索与无约束最优化

- 令 Q 是 R^n 中的一个子集, 用 $C_L^{k,p}(Q)$ 表示满足下面性质的函数类
 - 任意 $f \in C_L^{k,p}(Q)$ 在 Q 上是 k 次连续可微的
 - 其 p 阶倒数在 Q 上关于常数 L 是李普希兹连续的, 即对 $\forall x, y \in Q$, 都有 $\|\nabla^p f(x) - \nabla^p f(y)\| \leq L\|x - y\|$
 - 显然, $p \leq k$ 。如果 $q \geq k$, 则 $C_L^{q,p}(Q) \subseteq C_L^{k,p}(Q)$
 - 如果 $f_1 \in C_{L_1}^{k,p}(Q), f_2 \in C_{L_2}^{k,p}(Q), \alpha_1, \alpha_2 \in R$, 则对 $L_3 = |\alpha_1|L_1 + |\alpha_2|L_2$, 我们有 $\alpha_1 f_1 + \alpha_2 f_2 \in C_{L_3}^{k,p}(Q)$
- 定理: 函数 $f(\cdot) \in C_L^{2,1}(R^n) \subset C_L^{1,1}(R^n)$, 当且仅当 $\forall x \in R^n$, 我们有 $\|\nabla^2 f(x)\| \leq L$
 - 证明: 注意 $\forall x, y \in R^n, \nabla f(y) = \nabla f(x) + \int_0^1 \nabla^2 f(x + \tau(y-x))(y-x) d\tau$ 得证 \Rightarrow 。
 - $\Leftarrow \forall s \in R^n, \alpha > 0$, 有 $\|(\int_0^1 \nabla^2 f(x + \tau s) d\tau) \cdot s\| = \|\nabla f(x + \alpha s) - \nabla f(x)\| \leq \alpha L\|s\|$, 用 α 除这个等式, 同时令 $\alpha \downarrow 0$, 即得证。



第四章 一维搜索与无约束最优化

- 定理：令 $f \in C_L^{1,1}(R^n)$, 则 $\forall x, y \in R^n$, 我们有
 - $|f(y) - f(x) - \langle \nabla f(x), y - x \rangle| \leq \frac{L}{2} \|y - x\|^2$
- 由此, 可知, 对于任意 $f \in C_L^{1,1}(R^n)$, 取定点 $x_0 \in R^n$, 则可构造两个二次函数
 - $\phi_1(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle - \frac{L}{2} \|x - x_0\|^2$
 - $\phi_2(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle + \frac{L}{2} \|x - x_0\|^2$
- 则函数 f 的图像位于 ϕ_1 和 ϕ_2 之间, 即
 - $\phi_1(x) \leq f(x) \leq \phi_2(x), \forall x \in R^n$
- 定理：令 $f \in C_L^{2,2}(R^n), \forall x, y \in R^n$, 我们有
 - $\|\nabla f(y) - \nabla f(x) - \nabla^2 f(x)(y - x)\| \leq \frac{L}{2} \|y - x\|^2$
 - $|f(y) - f(x) - \langle \nabla f(x), y - x \rangle - \frac{1}{2} \langle \nabla^2 f(x)(y - x), y - x \rangle| \leq \frac{L}{6} \|y - x\|^3$
 - 证明略
- 推论：令 $f \in C_L^{2,2}(R^n)$ 和 $x, y \in R^n$, 满足 $\|y - x\| = r$, 则有
- $\nabla^2 f(x) - LrI_n \leq \nabla^2 f(y) \leq \nabla^2 f(x) + LrI_n$
- 证明：令 $G = \nabla^2 f(y) - \nabla^2 f(x)$, 因为 $f \in C_L^{2,2}(R^n)$, 从而 $\|G\| < Lr$, 因此矩阵 G 的特征值 $|\lambda_i| \leq Lr, i = 1, \dots, n$, 因此 $-LrI_n \leq G \leq LrI_n$



第四章 一维搜索与无约束最优化-梯度法收敛性分析

- 梯度法及其收敛性分析

- $\mathbf{x}_0 \in R^n$

- $\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{h}_k \nabla f(\mathbf{x}_k), k = 0, 1, \dots, \text{步长 } \mathbf{h}_k > 0$

- \mathbf{h}_k 的选取有很多变形

- $\{\mathbf{h}_k\}_{k=0}^{\infty}$: 如 $\mathbf{h}_k = h > 0, \mathbf{h}_k = \frac{h}{\sqrt{k+1}}$

- 全松弛 (精确步长) : $\mathbf{h}_k = \operatorname{argmin}_{h \geq 0} f(\mathbf{x}_k - h \nabla f(\mathbf{x}_k))$

- Armijo规则: 对 $h > 0$, 确定 $\mathbf{x}_{k+1} = \mathbf{x}_k - h \nabla f(\mathbf{x}_k)$, 满足

- $\alpha < \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}_{k+1} > \leq f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})$

- $\beta < \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}_{k+1} > \geq f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})$

- 其中, $0 < \alpha < \beta < 1$ 是一些固定参数.

- 看第三种策略, 确定 $\mathbf{x} \in R^n, \nabla f(\mathbf{x}) \neq \mathbf{0}$, 则只需研究单变量函数 $\phi(h) = f(\mathbf{x} - h \nabla f(\mathbf{x})), h \geq 0$



第四章 一维搜索与无约束最优化-梯度法收敛性分析

- 看第三种策略, 确定 $\mathbf{x} \in \mathbf{R}^n, \nabla f(\mathbf{x}) \neq \mathbf{0}$, 则只需研究单变量函数 $\phi(h) = f(\mathbf{x} - h\nabla f(\mathbf{x})), h \geq 0$
 - 则Armijo规则表明可接受的步长值对应于函数 ϕ 的图像的特定部分, 该部分介于两个线性函数的图像之间
 - $\phi_1(h) = f(\mathbf{x}) - \alpha h \|\nabla f(\mathbf{x})\|^2, \phi_2 = f(\mathbf{x}) - \beta h \|\nabla f(\mathbf{x})\|^2$
 - 注意 $\phi(0) = \phi_1(0) = \phi_2(0), \phi'(0) < \phi_2'(0) < \phi_1'(0) < 0$, 因此除非 ϕ 没有下界, 否则可接受的步长总是存在。
- 考虑问题 $\min_{\mathbf{R}^n} f(\mathbf{x})$, 满足 $f \in C_L^{1,1}(\mathbf{R}^n)$, 并假设函数 f 在 \mathbf{R}^n 有下界
 - 考虑 $\mathbf{y} = \mathbf{x} - h\nabla f(\mathbf{x})$, 此时 $f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2 = f(\mathbf{x}) - h \|\nabla f(\mathbf{x})\|^2 + \frac{h^2}{2} L \|\nabla f(\mathbf{x})\|^2 = f(\mathbf{x}) - h \left(1 - \frac{h}{2} L\right) \|\nabla f(\mathbf{x})\|^2$
 - 为了获得减少量的最优上界, 必须解 $\Delta(h) = -h \left(1 - \frac{h}{2} L\right) \rightarrow \min_h$, 计算其导数, 得最优步长必满足方程: $\Delta'(h) = hL - 1 = 0$. 因为 $\Delta''(h) = L > 0$, 因此 $h^* = \frac{1}{L}$ 就是 $\Delta(h)$ 的极小点。表明一步至少按 $f(\mathbf{y}) \leq f(\mathbf{x}) - \frac{1}{2L} \|\nabla f(\mathbf{x})\|^2$ 来降低目标函数值。



第四章 一维搜索与无约束最优化-梯度法收敛性分析

- 令 $x_{k+1} = x_k - h_k \nabla f(x_k)$, 则对于定步长策略, $h_k = h$, 我们有:

- $f(x_k) - f(x_{k+1}) \geq h \left(1 - \frac{1}{2} Lh\right) \|\nabla f(x_k)\|^2$

- 因此, 如果选择 $h_k = \frac{2\alpha}{L}$, 满足 $\alpha \in (0, 1)$, 则

- $f(x_k) - f(x_{k+1}) \geq \frac{2}{L} \alpha (1 - \alpha) \|\nabla f(x_k)\|^2$

- 当然最优选择为 $h_k = \frac{1}{L}$

- 对于全松弛策略, 我们有

- $f(x_k) - f(x_{k+1}) \geq \frac{1}{2L} \|\nabla f(x_k)\|^2$

- 最大的减少量不会比步长为 $h_k = \frac{1}{L}$ 的情形差

- 最后, 对于Armijo规则, 有

- $f(x_k) - f(x_{k+1}) \leq \beta < \nabla f(x_k), x_k - x_{k+1} > = \beta h_k \|\nabla f(x_k)\|^2$

- 根据下降法一步梯度推导:

- $f(x_k) - f(x_{k+1}) \geq h_k \left(1 - \frac{1}{2} Lh_k\right) \|\nabla f(x_k)\|^2$

- 因此 $h_k \geq \frac{2}{L} (1 - \beta)$, 根据Armijo第一条规则, 有



第四章 一维搜索与无约束最优化-梯度法收敛性分析

- $f(x_k) - f(x_{k+1}) \geq \alpha \langle \nabla f(x_k), x_k - x_{k+1} \rangle = \alpha h_k \|\nabla f(x_k)\|^2$
- 结合 $h_k \geq \frac{2}{L}(1 - \beta)$, 则可以得到
 - $f(x_k) - f(x_{k+1}) \geq \frac{2}{L} \alpha (1 - \beta) \|\nabla f(x_k)\|^2$
- 从而证明了, 所有条件下都满足
 - $f(x_k) - f(x_{k+1}) \geq \frac{\omega}{L} \|\nabla f(x_k)\|^2$, 其中 ω 是一个正常数
- 梯度法的性能如何?
 - 从 $f(x_k) - f(x_{k+1}) \geq \frac{\omega}{L} \|\nabla f(x_k)\|^2$ 能得到什么?
 - $\frac{\omega}{L} \sum_{k=0}^N \|\nabla f(x_k)\|^2 \leq f(x_0) - f(x_{N+1}) \leq f(x_0) - f^*$
 - f^* 为目标函数值的下界, 上式有界, 从而 $\lim_{k \rightarrow \infty} \|\nabla f(x_k)\| \rightarrow 0$



第四章 一维搜索与无约束最优化-梯度法收敛性分析

- 收敛率如何?
 - 定义 $g_N^* = \min_{0 \leq k \leq N} \|\nabla f(x_k)\|$, 根据 $\frac{\omega}{L} \sum_{k=0}^N \|\nabla f(x_k)\|^2 \leq f(x_0) - f(x_{N+1}) \leq f(x_0) - f^*$, 可得
 - $g_N^* \leq \frac{1}{\sqrt{N+1}} \left[\frac{1}{\omega} L (f(x_0) - f^*) \right]^{\frac{1}{2}}$
 - 这描述了序列 $\{g_N^*\}$ 收敛到0的速率
- 一般的非线性优化中, 只想找到接近优化问题的局部极小点, 但这个目标, 有时候对梯度法也实现不了。
- 例: 考虑下列函数: $f(x) = f(x_1, x_2) = \frac{1}{2}x_1^2 + \frac{1}{4}x_2^4 - \frac{1}{2}x_2^2$, 其梯度为 $\nabla f(x) = (x_1, x_2^3 - x_2)^T$, 其稳定点分别为 $x^{(1)} = (0, 0)^T, x^{(2)} = (0, -1)^T, x^{(3)} = (0, 1)^T$. 计算其Hessian矩阵
- $\nabla^2 f(x) = \begin{bmatrix} 1 & 0 \\ 0 & 3x_2^2 - 1 \end{bmatrix}$, 容易判断三个稳定点的极值情况. 其中 $x^{(1)}$ 为稳定点, 但非极值点, $f(x^{(1)}) = 0$, 对任意 $\epsilon > 0, f(x^{(1)} + \epsilon e_2) = \frac{\epsilon^4}{4} - \frac{\epsilon^2}{2} < 0$
- 此外, 以 $x_0 = (1, 0)$ 为梯度法的初始点, 则其迭代路径产生的序列收敛到 $x^{(1)} = (0, 0)^T$. 因此, 对于一阶无约束极小化方法, 若没有额外限制, 不能保证全局收敛到一个局部极小点, 只能靠近稳定点。



第四章 一维搜索与无约束最优化-梯度法收敛性分析

- 研究如下问题类：

- 模型：1. 无约束极小化，2. $f \in C_L^{1,1}(\mathbf{R}^n)$, 3. f^* 是 $f(\cdot)$ 的下界
- Oracle：一阶黑箱
- ϵ -最优解： $f(\bar{x}) \leq f(x_0), \|\nabla f(\bar{x})\| \leq \epsilon$

- 注意到 $g_N^* \leq \frac{1}{\sqrt{N+1}} \left[\frac{1}{\omega} L(f(x_0) - f^*) \right]^{\frac{1}{2}}$ 用于得到迭代次数的上界，这对找到梯度范数小的点很必要。为此，令 $g_N^* \leq \frac{1}{\sqrt{N+1}} \left[\frac{1}{\omega} L(f(x_0) - f^*) \right]^{\frac{1}{2}} \leq \epsilon \Rightarrow$ 如果 $N + 1 \geq \frac{L}{\omega \epsilon^2} (f(x_0) - f^*)$ ，则必然有 $g_N^* \leq \epsilon$ ，因此我们用 $g_N^* \leq \frac{1}{\sqrt{N+1}} \left[\frac{1}{\omega} L(f(x_0) - f^*) \right]^{\frac{1}{2}}$ 来作为该问题类的复杂度上界，注意这个上界比之前用均匀网格法的上界 $\left(\left\lfloor \frac{L}{2\epsilon} \right\rfloor + 1 \right)^n$ 更好，与 n 无关！但其复杂度下界还未知！



第四章 一维搜索与无约束最优化-梯度法收敛性分析

- 下面研究梯度法的局部收敛怎么描述!考虑无约束极小化问题
- $\min_{x \in R^n} f(x)$, 满足如下假设
 - $f \in C_M^{2,2}(R^n)$
 - f 有一个局部极小值点 $x^* \in R^n$, 该点处的Hessian矩阵正定
 - 该点处的Hessian矩阵, 知道其上下界 $0 < \mu \leq L < \infty$, 即
$$\mu I_n \leq \nabla^2 f(x^*) \leq L I_n$$
 - 初始点 x_0 足够接近 x^*
- 研究迭代过程: $x_{k+1} = x_k - h_k \nabla f(x_k)$, 因 $\nabla f(x^*) = 0$, 因此
- $\nabla f(x_k) = \nabla f(x_k) - \nabla f(x^*) = \int_0^1 \nabla^2 f(x^* + \tau(x_k - x^*)) (x_k - x^*) d\tau = G_k(x_k - x^*)$, 其中 $G_k = \int_0^1 \nabla^2 f(x^* + \tau(x_k - x^*)) d\tau$. 因此
- $x_{k+1} - x^* = x_k - x^* - h_k G_k(x_k - x^*) = (I_n - h_k G_k)(x_k - x^*)$
- 此时跟前面提到的收缩映射有关。设序列 $\{a_k\}$ 定义为
- $a_0 \in R^n, a_{k+1} = A_k a_k$, 其中 A_k 是 $n \times n$ 矩阵, $\forall k \geq 0, q \in (0, 1)$, 有 $\|A_k\| \leq 1 - q$, 这样估计序列 $\{a_k\}$ 的收敛速度



第四章 一维搜索与无约束最优化-梯度法收敛性分析

- 估计序列 $\{a_k\}$ 的收敛到0的速度
- $\|a_{k+1}\| \leq (1 - q)\|a_k\| \leq (1 - q)^{k+1}\|a_0\| \rightarrow 0$
- 现在估计 $\|I_n - h_k G_k\|$ 。令 $r_k = \|x_k - x^*\|$, 根据之前函数类的推论有: $\nabla^2 f(x^*) - \tau M r_k I_n \leq \nabla^2 f(x^* + \tau(x_k - x^*)) \leq \nabla^2 f(x^*) + \tau M r_k I_n$, 结合条件3
- $\left(\mu - \frac{r_k}{2} M\right) I_n \leq G_k \leq \left(L + \frac{r_k}{2} M\right) I_n$, 从而有
- $\left(1 - h_k \left(L + \frac{r_k}{2} M\right)\right) I_n \leq I_n - h_k G_k \leq \left(1 - h_k \left(\mu - \frac{r_k}{2} M\right)\right) I_n$, 得
- $\|I_n - h_k G_k\| \leq \max\{a_k(h_k), b_k(h_k)\}$,
- 其中 $a_k(h) = 1 - h \left(\mu - \frac{r_k}{2} M\right)$, $b_k(h) = h \left(L + \frac{r_k}{2} M\right) - 1$
- 注意到 $a_k(0) = 1$, $b_k(0) = -1$, 因此如果 $0 < r_k < \bar{r} = \frac{2\mu}{M}$, 则 $a_k(\cdot)$ 是一个严格递减函数, 对足够小的 h_k 可确保 $\|I_n - h_k G_k\| < 1$, 此时有 $r_{k+1} < r_k$



第四章 一维搜索与无约束最优化-梯度法收敛性分析

- 步长选择策略, 如 $\mathbf{h}_k = \frac{1}{L}$, 极小化 $\|\mathbf{I}_n - \mathbf{h}_k \mathbf{G}_k\| \leq \max\{a_k(\mathbf{h}_k), b_k(\mathbf{h}_k)\} \rightarrow \min_{\mathbf{h}}$ 的右端项, 假设 $r_0 < \bar{r}$, 我们利用最优策略序列得到序列 $\{\mathbf{x}_k\}$, 可以保证 $r_{k+1} < r_k < \bar{r}$. 进一步最优步长可从方程
- $a_k(\mathbf{h}) = b_k(\mathbf{h}) \Leftrightarrow 1 - h\left(\mu - \frac{r_k}{2}M\right) = h\left(L + \frac{r_k}{2}M\right) - 1$ 得到
- $\mathbf{h}_k^* = \frac{2}{L+\mu}$, 注意与M无关, 此时
- $r_{k+1} \leq \frac{(L-\mu)r_k}{L+\mu} + \frac{Mr_k^2}{L+\mu}$
- 现在估计迭代过程的收敛速度, 令 $q = \frac{2\mu}{L+\mu}$, $a_k = \frac{M}{L+\mu} r_k (< q)$, 则
- $a_{k+1} \leq (1-q)a_k + a_k^2 = a_k(1 + (a_k - q)) = \frac{a_k(1-(a_k-q)^2)}{1-(a_k-q)} \leq \frac{a_k}{1+q-a_k}$, 因此 $\frac{1}{a_{k+1}} \geq \frac{1+q}{a_k} - 1$



第四章 一维搜索与无约束最优化-梯度法收敛性分析

- 或者 $\frac{q}{a_{k+1}} - 1 \geq \frac{q(1+q)}{a_k} - q - q = (1+q) \left(\frac{q}{a_k} - 1 \right)$
- 所以, $\frac{q}{a_k} - 1 \geq (1+q)^k \left(\frac{q}{a_0} - q \right) = (1+q)^k \left(\frac{2\mu}{L+\mu} \cdot \frac{L+\mu}{r_0 M} - 1 \right) = (1+q)^k \left(\frac{\bar{r}}{r_0} - 1 \right)$
- 从而有 $a_k \leq \frac{qr_0}{r_0 + (1+q)^k(\bar{r} - r_0)} \leq \frac{qr_0}{\bar{r} - r_0} \left(\frac{1}{1+q} \right)^k$
- 定理: 设函数 $f(\cdot)$ 满足我们的假设, 且初始点 x_0 足够接近一个严格局部极小点 x^* , 即

$$r_0 = \|x_0 - x^*\| < \bar{r} = \frac{2\mu}{M}$$

则步长为 $h_k^* = \frac{2}{L+\mu}$ 的梯度法收敛如下:

- $\|x_k - x^*\| \leq \frac{\bar{r}r_0}{\bar{r} - r_0} \left(1 - \frac{2\mu}{L+3\mu} \right)^k$, 这种收敛速度称为线性收敛



第五章 无约束最优化方法及其收敛性分析-牛顿法

• 牛顿法

- 最开始用于单变量函数求根. 令 $\phi(\cdot): \mathbf{R} \rightarrow \mathbf{R}$, 考虑 $\phi(t^*) = 0$, 其原理由线性近似得到。假设知道距 t^* 足够近的 $t \in \mathbf{R}$, 注意到 $\phi(t + \Delta t) = \phi(t) + \phi'(t)\Delta t + o(|\Delta t|)$, 因此方程 $\phi(t + \Delta t) = 0$ 的解可用右端来近似为: $\phi(t) + \phi'(t)\Delta t = 0$, 在某些条件下, 我们希望增量 Δt 是最优增量 $\Delta t^* = t^* - t$ 的一个好的近似. 将其转化为算法得到: $t_{k+1} = t_k - \frac{\phi(t_k)}{\phi'(t_k)}$
- 该算法可自然推广到解非线性方程组的问题: $F(x) = 0, x \in \mathbf{R}^n, F(\cdot): \mathbf{R}^n \rightarrow \mathbf{R}^n$ 。此时定义一个增量 Δx 是下面线性方程组的解: $F(x) + F'(x)\Delta x = 0$ (牛顿系统)。如果 Jacobian 矩阵 $F'(x)$ 非退化, 我们可以计算增量 $\Delta x = -[F'(x)]^{-1}F(x)$, 相应的迭代方法如下:
- $x_{k+1} = x_k - [F'(x_k)]^{-1}F(x_k)$
- 最优, 由无约束优化的必要条件求解 $\nabla f(x) = 0$ 来代替求解 $\min f(x)$ (非退化情况), 为解必要条件, 使用标准的解非线性方程组的牛顿法, 这时, 牛顿系统为: $\nabla f(x) + \nabla^2 f(x)\Delta x = 0$
- 因此, 对优化问题, 牛顿法写为:
$$x_{k+1} = x_k - [\nabla^2 f(x_k)]^{-1} \nabla f(x_k)$$
- 注意, 也可以用二阶逼近来推导!

如何用二阶逼近
来推导牛顿法?



第五章 无约束最优化方法及其收敛性分析-牛顿法

- 牛顿法在严格局部极小点的一个邻域内，其收敛速度非常快
 - 约束1：若 $\nabla^2 f(\mathbf{x}_k)$ 退化，则牛顿法失败
 - 约束2：邻域内，否则可能不收敛
- 例：使用牛顿法求解 $\phi(t) = \frac{t}{\sqrt{1+t^2}}$ 的一个根。显然， $t^* = 0$ ，但 $\phi'(t) = \frac{1}{(1+t^2)^{\frac{3}{2}}}$ ，因此牛顿法的步骤如下：
 - $t_{k+1} = t_k - \frac{\phi(t_k)}{\phi'(t_k)} = t_k - \frac{t_k}{\sqrt{1+t_k^2}} \cdot (1+t_k^2)^{\frac{3}{2}} = -t_k^3$
 - 若 $|t_0| < 1$ ，则算法收敛且收敛速度很快
 - 点 $t_0 = \pm 1$ 为算法的震荡点
 - $|t_0| > 1$ ，则算法发散
- 为避免可能的发散，在实际中可以使用阻尼牛顿法：

$$\mathbf{x}_{k+1} = \mathbf{x}_k - h_k [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k), \text{ 其中 } h_k \text{ 为步长参数}$$



第五章 无约束最优化方法及其收敛性分析-牛顿法

- 注意，一般 \mathbf{h}_k 使用梯度法中的步长策略，但最后阶段，选择 $\mathbf{h}_k = \mathbf{1}$ 比较合理。
- 下面推导牛顿法的局部收敛速度
- 考虑问题 $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$
- 满足如下假设：
 - 1. $f \in C_M^{2,2}(\mathbb{R}^n)$
 - 2. 函数 f 存在一个局部极小点 \mathbf{x}^* ，该点处Hessian矩阵正定： $\nabla^2 f(\mathbf{x}^*) \succeq \mu \mathbf{I}_n, \mu > 0$
 - 3. 初始点 \mathbf{x}_0 足够接近 \mathbf{x}^*
- $\mathbf{x}_{k+1} = \mathbf{x}_k - [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$ 的收敛速度如何？



第五章 无约束最优化方法及其收敛性分析-牛顿法

推论: 令 $f \in C_L^{2,2}(R^n)$ 和 $x, y \in R^n$, 满足 $\|y - x\| = r$, 则有
 $\nabla^2 f(x) - LrI_n \leq \nabla^2 f(y) \leq \nabla^2 f(x) + LrI_n$

- 与梯度法类似
- $x_{k+1} - x^* = x_k - x^* - [\nabla^2 f(x_k)]^{-1} \nabla f(x_k) = x_k - x^* - [\nabla^2 f(x_k)]^{-1} \int_0^1 \nabla^2 f(x^* + \tau(x_k - x^*)) (x_k - x^*) d\tau = [\nabla^2 f(x_k)]^{-1} G_k (x_k - x^*)$
- 其中 $G_k = \int_0^1 (\nabla^2 f(x_k) - \nabla^2 f(x^* + \tau(x_k - x^*))) d\tau$
- 令 $r_k = \|x_k - x^*\|$, 则
- $\|G_k\| = \|\int_0^1 (\nabla^2 f(x_k) - \nabla^2 f(x^* + \tau(x_k - x^*))) d\tau\| \leq \int_0^1 \|(\nabla^2 f(x_k) - \nabla^2 f(x^* + \tau(x_k - x^*)))\| d\tau \leq \int_0^1 M(1 - \tau)r_k d\tau = \frac{r_k}{2} M$
- 因此, 由之前P71页推论我们有: $\nabla^2 f(x_k) \geq \nabla^2 f(x^*) - Mr_k I_n \geq (\mu - Mr_k) I_n$
- 因此, 若 $r_k < \frac{\mu}{M}$, 则 $\nabla^2 f(x_k)$ 正定, 且有 $[\nabla^2 f(x_k)]^{-1} \leq (\mu - Mr_k)^{-1}$
- 所以, 对于足够小的 $r_k (r_k \leq \frac{2\mu}{3M})$, 有 $r_{k+1} \leq \frac{Mr_k^2}{2(\mu - Mr_k)} (\leq r_k)$

二次收敛



第五章 无约束最优化方法及其收敛性分析-牛顿法

- 定理：令函数 $f(\cdot)$ 满足上述假设，假设初始点 \mathbf{x}_0 足够接近 \mathbf{x}^* ,即
$$\|\mathbf{x}_0 - \mathbf{x}^*\| \leq \bar{r} = \frac{2\mu}{3M}$$
- 则对任意 k 有 $\|\mathbf{x}_k - \mathbf{x}^*\| \leq \bar{r}$,牛顿法二次收敛，即：
- $$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \frac{M\|\mathbf{x}_k - \mathbf{x}^*\|^2}{2(\mu - M\|\mathbf{x}_k - \mathbf{x}^*\|)}$$
- 与梯度法的局部收敛率相比较，牛顿法显然收敛更快。
 - ! 牛顿法的二次收敛区域与梯度法的线性收敛区域几乎相同
 - 表明：标准推荐极小化过程的初始阶段使用梯度法来接近局部极小点，然后再利用牛顿法快速收敛到最优点！



第五章 无约束最优化方法及其收敛性分析-牛顿法

- 通过梯度法和牛顿法的收敛率，其与复杂度的界之间的对应关系，可知这些问题类的解析复杂度的上界是收敛率的反函数
 - 次线性速率。该速率由迭代计算器的幂函数来表示。例如，假设对于某算法可以证明其收敛率为 $r_k \leq \frac{c}{\sqrt{k}}$ ，此时对于相应的问题类，该算法的复杂度上界是 $\left(\frac{c}{\epsilon}\right)^2$
 - 次线性速率是比较慢的，就复杂度而言，最优值中每得到一位新的正确数字都需要经历与以前的总的工作量相当的迭代次数。注意常数 c 对相应的复杂度的上界影响很大
 - 线性速度。该速率是根据迭代计数器的指数函数给出的。例如， $r_k \leq c(1-q)^k \leq ce^{-qk}, 0 < q \leq 1$ ，注意其相应的复杂度上界是 $\frac{1}{q} \left(\ln c + \ln \frac{1}{\epsilon} \right)$
 - 这个速率很快：最优值中每得到一位新的正确数字需要经历大约固定的迭代次数。此外，复杂度估计对于常数 c 的依赖性非常弱
 - 二次收敛速率。该速率对迭代计数器都是双指数依赖的。例如， $r_{k+1} \leq cr_k^2$ ，相应复杂度估计依赖于所需精度的双对数： $\ln \ln \frac{1}{\epsilon}$
 - 这个收敛率非常快：每次迭代都会双倍增加最优值的正确数字。常数 c 仅对二次收敛的开始时刻很重要（ $cr_k < 1$ ）。例如，在 $cr_k \leq \frac{1}{2}$ 之后，我们可以保证一个较大的收敛速率 $r_{k+1} \leq \frac{1}{2}r_k$ ，不再依赖于常数 c 。



第五章 有约束最优化方法及其收敛性分析-牛顿法

- 非线性优化中的一阶方法

- 在梯度下降法中, 选取 $\mathbf{h}_k = \frac{1}{L}$, 一般其收敛速度为 $\mathcal{O}\left(\frac{1}{k}\right)$, 若目标函数为 μ 强凸函数(满足:

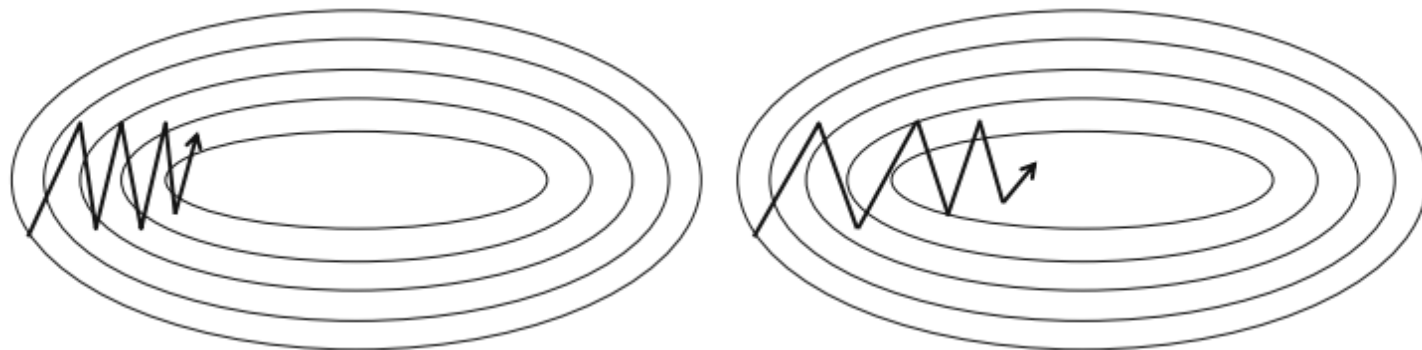
$$f(\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}) - \frac{\mu\alpha(1-\alpha)}{2} \|\mathbf{x} - \mathbf{y}\|^2), \text{ 则其收敛速率为 } \mathcal{O}\left(\left(1 - \frac{\mu}{L}\right)^k\right)$$

- 加速梯度法可将收敛速率提至 $\mathcal{O}\left(\frac{1}{k^2}\right)$ 和 $\mathcal{O}\left(\left(1 - \sqrt{\frac{\mu}{L}}\right)^k\right)$

- 加速梯度法 $\mathbf{x}_0 = \mathbf{x}_{-1} \in \mathbb{R}^n$

- $\mathbf{y}_k = \mathbf{x}_k + \beta_k(\mathbf{x}_k - \mathbf{x}_{k-1}), \mathbf{x}_{k+1} = \mathbf{y}_k - \frac{1}{L}\nabla f(\mathbf{y}_k), k = 0, 1, \dots,$
- 例如令 $\beta_k = \frac{k}{k+1}$

-





第五章 无约束最优化方法及其收敛性分析-牛顿法

- 二阶加速算法

$$\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 \leq L_2 \|x - y\|_2, \forall x, y \in \mathbb{R}^n$$

	Lower Bound	Upper Bound
$p = 2$	$\Omega \left(\left(\frac{L_2 \ x_0 - x^*\ ^3}{\epsilon} \right)^{\frac{2}{11}} \right)$ [Agarwal and Hazan (2018)]	$O \left(\left(\frac{L_2 \ x_0 - x^*\ ^3}{\epsilon} \right)^{\frac{1}{3}} \right)$ [Nesterov (2008)]
	$\Omega \left(\left(\frac{L_2 \ x_0 - x^*\ ^3}{\epsilon} \right)^{\frac{2}{7}} \right)$ [Arjevani et al. (2018)] [Nesterov (2018)]	$\tilde{O} \left(\left(\frac{L_2 \ x_0 - x^*\ ^3}{\epsilon} \right)^{\frac{2}{7}} \right)$ [Monteiro, Svaiter (2013)]

Key idea, second order approximation:

$$f(x) \approx f(x^i) + (x - x^i)^\top \nabla f(x^i) + \frac{1}{2} (x - x^i)^\top \nabla^2 f(x^i) (x - x^i)$$



第五

• 高阶

	下界	上界
$p = 1$	$\Omega\left(\left(\frac{L_1 x_0 - x^* ^2}{\epsilon}\right)^{\frac{1}{2}}\right)$ [Nemirovski,Yudin(1983)]	$\mathcal{O}\left(\left(\frac{L_1 x_0 - x^* ^2}{\epsilon}\right)^{\frac{1}{2}}\right)$ [Nesterov(1983)]
$p = 2$	$\Omega\left(\left(\frac{L_2 x_0 - x^* ^3}{\epsilon}\right)^{\frac{2}{7}}\right)$ [Arjevani et.al (2018)]	$\mathcal{O}\left(\left(\frac{L_2 x_0 - x^* ^3}{\epsilon}\right)^{\frac{2}{7}}\right)$ [Monteiro,Svaiter(2013)]
$p \geq 3$	$\Omega\left(\left(\frac{L_p x_0 - x^* ^{p+1}}{\epsilon}\right)^{\frac{2}{3p+1}}\right)$ [Arjevani et al(2018)] [Nesterov[2018]]	$\mathcal{O}\left(\left(\frac{L_p x_0 - x^* ^{p+1}}{\epsilon}\right)^{\frac{1}{p+1}}\right)$ [Baes(2009)] [Nesterov [2018]]
$p \geq 3$		$\mathcal{O}\left(\left(\frac{L_p x_0 - x^* ^{p+1}}{\epsilon}\right)^{\frac{2}{3p+1}}\right)$

。得到阶的问降

谢谢！