

#### 杨周旺

中国科学技术大学 数学科学学院

2023年2月

YZW (USTC)

1/467

æ

イロン イ理 とく ヨン イ ヨン

# Outline I

- Unconstrained Optimization
  - 2 Constrained Optimization
    - 二次规划
    - 非线性约束最优化
- 3 Convex Optimization
  - Convex Set and Convex Function
  - Convex Optimization and Algorithms
- 4 Sparse Optimization
  - Sparse Optimization Models
  - Sparse Optimization Algorithms
- 5 Optimization Methods for Machine Learning

► < Ξ ►</p>

# Outline II

- Typical Form of Problems
- Stochastic Algorithms
- Other Popular Methods



臣

(I)

- The course is devoted to the mathematical fundamentals of optimization and the practical algorithms of optimization.
- The course covers the topics of nonlinear continuous optimization, sparse optimization, and optimization methods for machine learning.

3 × 4 3 ×

Objectives of the course are

- to develop an understanding of the fundamentals of optimization;
- to learn how to analyze the widely used algorithms for optimization;
- to become familiar with the implementation of optimization algorithms.

∃ ▶ ∢ ∃ ▶

- Knowledge of Linear Algebra, Real Analysis, and Mathematics of Operations Research are very important for this course.
- Simultaneously, the ability to write computer programs of algorithms is also required.

글 제 제 글 제

- Unconstrained Optimization
- Constrained Optimization
- Convex Optimization
- Sparse Optimization
- Optimization Methods for Large-scale Machine Learning

3 × 4 3 ×

- 1 R. Fletcher. Practical Methods of Optimization (2nd Edition), John Wiley & Sons, 1987.
- 2 J. Nocedal and S. J. Wright. Numerical Optimization (2nd Edition), Springer, 2006.
- 3 S. Boyd and L. Vandenberghe. Convex Optimization, Cambridge University Press, 2004.
- 4 M. Elad. Sparse and Redundant Representations: From Theory to Applications in Signal and Image Processing. Springer, 2010.
- 5 L. Bottou, F.E. Curtis, J. Nocedal. Optimization methods for large-scale machine learning. SIAM Review, 60(2): 223-311, 2018.

- (1) Homework (10%)
- (2) Project (30%)
- (3) Final Exam (60%)

臣

イロト イヨト イヨト イヨト

# Outline I

- Unconstrained Optimization
  - 2 Constrained Optimization
    - 二次规划
    - 非线性约束最优化
  - 3 Convex Optimization
    - Convex Set and Convex Function
    - Convex Optimization and Algorithms
  - A Sparse Optimization
    - Sparse Optimization Models
    - Sparse Optimization Algorithms

## Optimization Methods for Machine Learning

YZW (USTC)

Image: A image: A

# Outline II

- Typical Form of Problems
- Stochastic Algorithms
- Other Popular Methods



æ

イロト 不得下 イヨト イヨト



#### 

- - 二次规划
  - 非线性约束最优化
- - Convex Set and Convex Function
  - Convex Optimization and Algorithms
- - Sparse Optimization Models
  - Sparse Optimization Algorithms
- - Typical Form of Problems
  - Stochastic Algorithms
  - Other Popular Methods YZW (USTC)

#### 无约束最优化问题

$$\min_{\mathsf{x}\in\mathbb{R}^n} \quad f(\mathsf{x}) \tag{1}$$

イロト 不得 トイヨト イヨト

其目标函数f是定义在 $\mathbb{R}$ "上的实值函数,决策变量×的可取值之集合是全空间  $\mathbb{R}$ ".

臣

梯度向量▽f(x)是函数f在点×处增加最快的方向,故它成为最优化时的 重要工具。实际上针对无约束最优化问题,大家所知的求解算法中大多 属于下面的梯度方法类。

#### GRADIENT (梯度法类)

- (0) 初始化:选取适当的初始点 $x^{(0)} \in \mathbb{R}^n$ , 令k := 0.
- (1) 计算搜索方向:利用适当的正定对称阵 $H_k$ 计算搜索方向向  $d^{(k)} := -H_k \nabla f(x^{(k)}).$  (如果 $\nabla f(x^{(k)}) = 0$ ,则结束计算)
- (2) 确定步长因子: 解一维最优化问题min  $f(x^{(k)} + \alpha d^{(k)})$ , 求出步 长 $\alpha = \alpha_k$ , 令 $x^{(k+1)} = x^{(k)} + \alpha_k d^{(k)}$ , k := k + 1, 回到第(1)步。

无约束最优化问题:  $\min_{x \in \mathbb{R}^n} f(x)$ 

$$f(\mathbf{x}) = f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^T (\mathbf{x} - \mathbf{x}^{(k)}) + O(\|\mathbf{x} - \mathbf{x}^{(k)}\|^2)$$
(2)

取负梯度方向

$$\mathsf{d}^{(k)} = -\nabla f(\mathsf{x}^{(k)}),$$

则当 $\alpha_k$ 足够小时,总能使

$$f(\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}) < f(\mathbf{x}^{(k)}).$$

æ

イロト イヨト イヨト イヨト

$$f(\mathbf{x}) = f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^{T} (\mathbf{x} - \mathbf{x}^{(k)}) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(k)})^{T} \nabla^{2} f(\mathbf{x}^{(k)}) (\mathbf{x} - \mathbf{x}^{(k)}) + O(\|\mathbf{x} - \mathbf{x}^{(k)}\|^{3})$$
(3)

取搜索方向

$$\mathsf{d}^{(k)} = -G_k^{-1}\nabla f(\mathsf{x}^{(k)}),$$

其中 $G_k = \nabla^2 f(\mathbf{x}^{(k)})$ 为函数 $f \mathbf{a} \mathbf{x}^{(k)}$ 点处的Hesse矩阵。

▲□▶ ▲圖▶ ▲ 圖▶ ▲ 圖▶ ― 圖 … のへで

### 在迭代格式中, 通过解一维最优化问题

$$\min_{\alpha \ge 0} \varphi(\alpha) = f(\mathsf{x}^{(k)} + \alpha \mathsf{d}^{(k)}) \tag{4}$$

确定步长因子的方法称为一维搜索(Line Search).



イロト イ団ト イヨト イヨト 二日

若以问题(4)的最优解为步长,此时称为**精确一维搜索**(Exact Line Search).

经常用到的精确一维搜索有黄金分割法和插值迭代法。即使说是精确一 维搜索,通过有限次计算求出问题(4)的严密解一般也是不可能的,实际 上在得到有足够精度的近似解时,就采用它作为步长。

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

在实际计算中,往往不是求解一维最优化问题(4),而是找出满足某些适 当条件的粗略近似解作为步长,此时称为**非精确一维搜索**(Inexact Line Search).

与精确一维搜索相比,在很多情况下采用非精确一维搜索可以提高整体 计算效率。

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

# 确定步长因子:一维搜索

设 $\bar{\alpha}_k$ 是使得

$$f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) = f(\mathbf{x}^{(k)})$$

的最小正数 $\alpha$ .

于是,我们将在区间[0, $\bar{\alpha}_k$ ]内求得满足适当条件的可接受的步长因子, 即 $\alpha \in [0, \bar{\alpha}_k]$ .



YZW (USTC)

20 / 467

∃ ▶ ∢ ∃ ▶

Goldstein(1965) conditions:

$$\varphi(\alpha) \le \varphi(0) + \rho \alpha \varphi'(0)$$
 (5)

$$\varphi(\alpha) \ge \varphi(0) + (1 - \rho)\alpha \varphi'(0) \tag{6}$$

其中 $\rho \in (0, 1/2)$ 是一个固定参数。

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <



#### Goldstein(1965) conditions:



YZW (USTC)

22 / 467

æ

イロト 不得下 イヨト イヨト

#### Wolfe(1968)-Powell(1976) conditions:

$$\varphi(\alpha) \le \varphi(0) + \rho \alpha \varphi'(0)$$
 (7)

$$\varphi'(\alpha) \ge \sigma \varphi'(0)$$
 (8)

イロト イポト イヨト イヨト 一日

其中 $\sigma \in (\rho, 1)$ 是另一个固定参数。



#### Wolfe(1968)-Powell(1976) conditions:



YZW (USTC)

24 / 467

• • = • • = •

#### 在很多实际算法中,式(8)常被强化的双边条件所取代

$$|\varphi'(\alpha)| \le -\sigma \varphi'(0)$$
 (9)

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

YZW (USTC)

#### 基于Wolfe-Powell准则的非精确一维搜索算法:

(0) 给定初始一维搜索区间[0,  $\bar{\alpha}$ ], 以及 $\rho \in (0, 1/2), \sigma \in (\rho, 1)$ . 计算 $\varphi_0 = \varphi(0) = f(\mathbf{x}^{(k)}), \varphi'_0 = \varphi'(0) = \nabla f(\mathbf{x}^{(k)})^T \mathsf{d}^{(k)}.$ 并令 $a_1 = 0, a_2 = \bar{\alpha}, \varphi_1 = \varphi_0, \varphi'_1 = \varphi'_0.$ 选取适当的 $\alpha \in (a_1, a_2).$ 

#### 基于Wolfe-Powell准则的非精确一维搜索算法:

(1) 计算 $\varphi = \varphi(\alpha) = f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$ . 若 $\varphi(\alpha) \le \varphi(0) + \rho \alpha \varphi'(0)$ , 则转到 第(2)步。否则,由 $\varphi_1, \varphi'_1, \varphi$ 构造两点二次插值多项式 $p^{(1)}(t)$ , 并得 其极小点

$$\hat{\alpha} = a_1 + rac{1}{2} rac{(a_1 - lpha)^2 arphi_1'}{(arphi_1 - arphi) - (a_1 - lpha) arphi_1'}.$$

于是置 $a_2 = \alpha, \alpha = \hat{\alpha},$ 重复第(1)步。

《曰》 《問》 《글》 《글》 \_ 글

基于Wolfe-Powell准则的非精确一维搜索算法:

(2) 计算 $\varphi' = \varphi'(\alpha) = \nabla f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})^T \mathbf{d}^{(k)}$ . 若 $\varphi'(\alpha) \ge \sigma \varphi'(0)$ , 则输 出 $\alpha_k = \alpha$ , 并停止搜索。否则, 由 $\varphi, \varphi', \varphi'_1$ 构造两点二次插值多项 式 $p^{(2)}(t)$ , 并得其极小点

$$\hat{\alpha} = \alpha - \frac{(\mathbf{a}_1 - \alpha)\varphi'}{\varphi'_1 - \varphi'}.$$

于是置 $a_1 = \alpha, \alpha = \hat{\alpha}, \varphi_1 = \varphi, \varphi'_1 = \varphi'$ , 返回第(1)步。

イロト イ理ト イヨト イヨト

# [思考题:请写出上述基于Wolfe-Powell准则的非精确一维搜索算法中插 值多项式 $p^{(1)}(t), p^{(2)}(t)$ 的具体表达式。]

イロト イタト イヨト イヨト 二日

从任意初始点出发,如果某迭代算法产生的点列的极限(聚点),在适 当假定下可保证恒为问题的最优解(或者稳定点),则称该迭代法具有 全局收敛性(Global Convergence).

与此相对,如果仅在解的附近选取初始点时,才可以保证所生成的点列收敛于该解,则称这样的迭代法有局部收敛性(Local Convergence).

イロト (周) (三) (三) (三)

为了证明迭代法的下降性,我们应尽量避免搜索方向与负梯度方向几乎 正交的情形,即要求d<sup>(k)</sup>偏离g<sup>(k)</sup> =  $\nabla f(x^{(k)})$ 的正交方向远一些。否则, g<sup>(k)<sup>T</sup></sup>d<sup>(k)</sup>接近于零,d<sup>(k)</sup>几乎不是下降方向。

为此,我们假设d<sup>(k)</sup>与–g<sup>(k)</sup>的夹角 $\theta_k$ 满足

$$\theta_k \le \frac{\pi}{2} - \mu, \ \forall k$$
(10)

< ロ > < 同 > < 三 > < 三 > 、

其中 $\mu > 0$ (与k无关)。

显然 $\theta_k \in [0, \pi/2)$ , 其定义为

$$\cos \theta_{k} = \frac{-g^{(k)}{}^{T} d^{(k)}}{\|g^{(k)}\| \|d^{(k)}\|} = \frac{-g^{(k)}{}^{T} s^{(k)}}{\|g^{(k)}\| \|s^{(k)}\|}$$
(11)  
$$\& \Xi s^{(k)} = \alpha_{k} d^{(k)} = x^{(k+1)} - x^{(k)}.$$

æ

イロト イヨト イヨト イヨト

下面给出各种步长准则下的下降算法的全局收敛性结论。

#### 全局收敛性定理:

设∇*f*(x)在水平集  $L(x^{(0)}) = \{x | f(x) \le f(x^{(0)})\}$  上存在且连续。下降算 法的搜索方向d<sup>(k)</sup>与−∇*f*(x<sup>(k)</sup>) 之间的夹角 $\theta_k$ 满足式(10), 其中步长 $\alpha_k$ 由 三种方法之一确定:

- (1) 精确一维搜索
- (2) Goldstein准则 (5),(6)
- (3) Wolfe-Powell准则 (7),(8)

那么,或者对某个k有 $\nabla f(x^{(k)}) = 0$ ,或者 $f(x^{(k)}) \rightarrow -\infty$ ,或者 $\nabla f(x^{(k)}) \rightarrow 0$ .

< ロ > < 同 > < 三 > < 三 > 、

### 全局收敛性证明: (只证明Wolfe-Powell准则的情形)

假设对所有的 $k, g^{(k)} = \nabla f(\mathbf{x}^{(k)}) \neq 0$ 和 $f(\mathbf{x}^{(k)})$ 有下界, 故 $f(\mathbf{x}^{(k)}) - f(\mathbf{x}^{(k+1)}) \rightarrow 0$ .由式(7)得, $-g^{(k)} \mathbf{x}^{(k)} \rightarrow 0$ .

(反证)若g<sup>(k)</sup>  $\rightarrow$  0不成立,那么存在 $\varepsilon$  > 0和子列{x<sup>(k)</sup>}<sub>k \in K</sub>使 得 $\|g^{(k)}\| \ge \varepsilon$ . 从而由

$$-\mathbf{g}^{(k)}\mathbf{s}^{(k)} = \|\mathbf{g}^{(k)}\|\|\mathbf{s}^{(k)}\|\cos\theta_k \ge \varepsilon\|\mathbf{s}^{(k)}\|\sin\mu$$

以及式(10)有 $\|s^{(k)}\| \to 0.$ 

< 口 > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

#### 全局收敛性证明(续):

又因为g(x) =  $\nabla f(x)$ 在 $L(x^{(0)})$ 上连续,所以

$$g^{(k+1)^{T}}s^{(k)} = g^{(k)^{T}}s^{(k)} + o(||s^{(k)}||) \frac{\psi}{g^{(k+1)^{T}}s^{(k)}} \rightarrow 1.$$
(12)

YZW (USTC)

▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ 臣 のへで



### 全局收敛性证明(续):

而这与Wolfe-Powell准则的式(8)

$$\frac{\mathsf{g}^{(k+1)}{}^{\mathsf{T}}\mathsf{s}^{(k)}}{\mathsf{g}^{(k)}{}^{\mathsf{T}}\mathsf{s}^{(k)}} \le \sigma < 1 \tag{13}$$

相矛盾。因此有 $g^{(k)} \rightarrow 0$ 

[思考题:请补充证明基于Goldstein准则的非精确一维搜索算法的全局收 敛性。]

イロト イポト イヨト イヨト 一日
#### 最速下降法取负梯度作为迭代算法的搜索方向,其迭代格式为

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{k}^{(k)}).$$

YZW (USTC)

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

#### 算法:

- (0) 选取初始点 $x^{(0)}$ ,设置终止误差 $\varepsilon > 0$ , 令k := 0.
- (1) 计算g<sup>(k)</sup> =  $\nabla f(\mathbf{x}^{(k)})$ . 若 $\|\mathbf{g}^{(k)}\| < \varepsilon$ , 则停止迭代并输出 $\mathbf{x}^{(k)}$ . 否则进行第(2)步。
- (2) 令 $d^{(k)} = -g^{(k)}$ ,并由一维搜索确定步长因子 $\alpha_k$ 使得

$$f(\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}) = \min_{\alpha > 0} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}).$$

(3) 迭代更新 $x^{(k+1)} = x^{(k)} + \alpha_k d^{(k)}$ , 置k := k + 1, 回到第(1)步。

< 口 > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

#### 最速下降法全局收敛性定理:

设 $f(x) \in C^1$ , 在最速下降法中采用(精确或非精确)一维搜索,则产生的迭代点列 $\{x^{(k)}\}$ 的每一个聚点都是驻点。

YZW (USTC)

39 / 467

A D N A (B) N A B N A B N B





YZW (USTC)

40 / 467

臣

#### 一般地,最速下降法只有线性收敛速度。

## 如下例子是一个非常著名的测试函数 (Rosenbrock function) $f(x_1, x_2) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2.$

YZW (USTC)

41 / 467

イロト イポト イヨト イヨト 二日

生顿法

设*f* (x)是二次可微实函数,在x<sup>(k)</sup>附近作二阶Taylor展开近似

$$f(\mathbf{x}^{(k)} + \mathbf{s}) \approx q^{(k)}(\mathbf{s}) = f(\mathbf{x}^{(k)}) + {\mathbf{g}^{(k)}}^{T}\mathbf{s} + \frac{1}{2}\mathbf{s}^{T}G_{k}\mathbf{s}$$
 (14)

其中g<sup>(k)</sup> =  $\nabla f(\mathbf{x}^{(k)}), G_k = \nabla^2 f(\mathbf{x}^{(k)}).$ 

将q<sup>(k)</sup>(s)极小化便得

$$s = -G_k^{-1}g^{(k)}.$$
 (15)

上式给出的搜索方向 $-G_k^{-1}g^{(k)}$ 称为牛顿方向(Newton Direction).

▲ロ▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ ▲ 国 ● のへで

生顿法

#### 在目标函数是正定二次函数

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^{\mathsf{T}} \mathbf{G} \mathbf{x} - \mathbf{c}^{\mathsf{T}} \mathbf{x}$$

的情况下(*G*为正定阵),对任意的x有 $\nabla^2 f(x) = G$ .

在第一次迭代里令 $H_0 = G^{-1}$ ,则有

 $d^{(0)} = -H_0 \nabla f(x^{(0)}) = -G^{-1}(Gx^{(0)} - c) = -(x^{(0)} - x^*).$ 

这里,  $x^* = G^{-1}$ c是问题的最优解。若 $x^{(0)} \neq x^*$ , 取步长 $\alpha_0 = 1$ , 于是得  $x^{(1)} = x^{(0)} + \alpha_0 d^{(0)} = x^*$ . 由此知道,不管初始点 $x^{(0)}$ 如何取,在一次迭 代后即可到达最优解 $x^*$ .

生顿法

根据以上事实,可以认为即使对于一般的非线性函数*f*(x),在迭代中令 搜索方向

$$d^{(k)} = -\nabla^2 f(x^{(k)})^{-1} \nabla f(x^{(k)})$$

也是较合适的。

特别地,步长 $\alpha_k \equiv 1$ 的迭代公式为

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{d}^{(k)} = \mathbf{x}^{(k)} - G_k^{-1} \mathbf{g}^{(k)}.$$
 (16)

这就是经典的牛顿迭代法

对于正定二次函数而言,牛顿法一步即可达到最优解。对于非二次函数,牛顿法并不能保证经有限次迭代求得最优解。但由于目标函数在极 小点附近可用二次函数较好地近似,故当初始点靠近极小点时,牛顿法 的收敛速度一般会很快。

可以证明牛顿法的局部收敛性和二阶收敛速率。

< 口 > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

#### 牛顿法收敛定理:

设 $f \in C^2$ ,  $x^{(k)}$ 充分靠近 $x^*$ , 其中 $\nabla f(x^*) = 0$ . 如果 $\nabla^2 f(x^*)$ 正定, 目标函数的Hesse矩阵G(x)满足Lipschitz条件,即存在 $\beta > 0$ 使得对所有(i, j)有

$$|G_{ij}(\mathsf{x}) - G_{ij}(\mathsf{y})| \le \beta \|\mathsf{x} - \mathsf{y}\|.$$
(17)

则对一切的k, 牛顿迭代(16)有定义, 所得序列 $\{x^{(k)}\}$ 收敛到 $x^*$ , 且具有二阶收敛速率。

< 口 > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >



#### 证明一:

记  $g(x) = \nabla f(x)$ , 因为 $f \in C^2$ , 我们有

$$g(x-h) = g(x) - G(x)h + O(||h||^2).$$

$$\mathbf{x} = \mathbf{x}^{(k)}, \mathbf{h} = \mathbf{h}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}^*$$
代入上式得

$$0 = g(x^*) = g(x^{(k)} - h^{(k)}) = g(x^{(k)}) - G(x^{(k)})h^{(k)} + O(||h^{(k)}||^2).$$
(18)

▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ 臣 のへで

牛顿法

#### 证明一 (续):

由于G(x)满足Lipschitz条件,易证 $[G(x^{(k)})]^{-1}$ 有界。方程(18)两边同时 乘以 $[G(x^{(k)})]^{-1}$ 得

$$0 = [G(x^{(k)})]^{-1}g(x^{(k)}) - h^{(k)} + O(||h^{(k)}||^2)$$
  
=  $x^* - (x^{(k)} - [G(x^{(k)})]^{-1}g(x^{(k)})) + O(||h^{(k)}||^2)$   
=  $x^* - x^{(k+1)} + O(||h^{(k)}||^2)$   
=  $-h^{(k+1)} + O(||h^{(k)}||^2)$ 

所以 $\|\mathbf{h}^{(k+1)}\| = O(\|\mathbf{h}^{(k)}\|^2)$ ,即牛顿迭代法具有二阶收敛速率。

48 / 467

A D N A (B) N A B N A B N B

牛顿法

#### 证明二:

对于牛顿迭代法,我们记

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - G_k^{-1} \mathbf{g}^{(k)} \triangleq \mathcal{A}(\mathbf{x}^{(k)}).$$
(19)

注意到g(x\*) = 0, G(x\*)正定(非奇异), 有A(x\*) = x\*.  
于是由x<sup>(k+1)</sup> - x\* = A(x<sup>(k)</sup>) - A(x\*) 得  
$$\|x^{(k+1)} - x^*\| = \|A(x^{(k)}) - A(x^*)\|$$
$$\leq \|A'(x^*)(x^{(k)} - x^*)\| + \frac{1}{2}\|A''(\bar{x})\|\|x^{(k)} - x^*\|^2,$$

其中<sup>x</sup>位于x<sup>(k)</sup>和x\*之间的线段上。

▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ 臣 のへで

证明二 (续): 显然  $\mathcal{A}'(\mathbf{x}) = [\mathbf{x} - G(\mathbf{x})^{-1}g(\mathbf{x})]' = -[G(\mathbf{x})^{-1}]'g(\mathbf{x})$ 所以 $\mathcal{A}'(\mathbf{x}^*) = 0$ . 从而有  $\|\mathbf{h}^{(k+1)}\| = \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \le \gamma \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 = \gamma \|\mathbf{h}^{(k)}\|^2$ 

其中常数 $\gamma$ 仅依赖于f(x)在 $x^*$ 附近的三阶导数。

▲ロ▶ ▲圖▶ ▲臣▶ ▲臣▶ ―臣 - 釣��

#### 在式(16)的牛顿迭代法里,如果选取的初始点×<sup>(0)</sup>不在解×\*的附近,那么 生成的点列{×<sup>(k)</sup>}未必收敛于最优解。

为保证算法的全局收敛性,有必要对牛顿法作某些改进。



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

生顿法

### 比如,在牛顿法中也可采用一维搜索来确定步长。 **阻尼牛顿法**:

- (0) 选取初始点 $x^{(0)}$ ,设置终止误差 $\varepsilon > 0$ , 令k := 0.
- (1) 计算g<sup>(k)</sup> =  $\nabla f(\mathbf{x}^{(k)})$ . 若 $\|\mathbf{g}^{(k)}\| < \varepsilon$ , 停止迭代并输出 $\mathbf{x}^{(k)}$ . 否则进行第(2)步。
- (2) 解线性方程组 $G_k d = -g^{(k)}$ , 求出牛顿方向 $d^{(k)}$ .
- (3) 采用一维搜索确定步长因子 $\alpha_k$ , 令 $x^{(k+1)} = x^{(k)} + \alpha_k d^{(k)}$ , 置k := k + 1, 回到第(1)步。

## 牛顿法面临的主要困难是Hesse矩阵 $G_k = \nabla^2 f(\mathbf{x}^{(k)})$ 不正定。这时二阶近 似模型不一定有极小点,即二次函数 $q^{(k)}(\mathbf{s})$ 是无界的。

为了克服这些困难,人们提出了很多修正措施。

A D N A (B) N A B N A B N B

Goldstein & Price (1967)

$$\mathsf{d}^{(k)} = \begin{cases} -G_k^{-1} \mathsf{g}^{(k)}, & \text{if } \cos \theta_k > \eta \\ \\ -\mathsf{g}^{(k)}, & \text{otherwise} \end{cases}$$



▲□▶ ▲圖▶ ▲厘▶ ▲厘▶ 二厘

(20)

#### Levenberg(1944), Marquardt(1963), Goldfeld et. al(1966)

$$(G_k + \mu_k I)d^{(k)} = -g^{(k)}$$
 (21)

イロト イヨト イヨト イヨト



æ



### 设×是函数f的一个不定点,若方向d满足

 $\mathsf{d}^T \nabla^2 f(\mathsf{x}) \mathsf{d} < 0,$ 

则称d为f在x处的负曲率方向。

当Hesse矩阵 $\nabla^2 f(\mathbf{x}^{(k)})$ 不正定时,负曲率方向法是修正牛顿法的另一种途径。

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ ― 臣 – ∽へ⊙

牛顿法的突出优点是局部收敛很快(具有二阶收敛速率), 但运用牛顿法需要计算二阶导,而且目标函数的Hesse矩阵 $\nabla^2 f(x^{(k)})$ 可 能非正定,甚至奇异。为了克服这些缺点, 人们提出了拟牛顿法。其基本思想是:用不含二阶导数的矩阵 $H_k$ 近似牛顿法中的Hesse矩阵的逆 $G(x^{(k)})^{-1}$ .

由构造近似矩阵的方法不同,将出现不同的拟牛顿法。

イロト (周) (三) (三) (三)

#### 回顾牛顿法的迭代

$$\begin{cases} G_k \mathsf{d} = -\mathsf{g}^{(k)} \\ \mathsf{x}^{(k+1)} = \mathsf{x}^{(k)} + \alpha_k \mathsf{d}^{(k)} \end{cases}$$

为了构造Hesse矩阵逆 $G_k^{-1}$ 的近似 $H_k$ , 我们先分析二阶导 $\nabla^2 f(\mathbf{x}^{(k)})$ 与一阶导 $\nabla f(\mathbf{x}^{(k)})$ 的关系。

设第k次迭代后得到 $x^{(k+1)}$ ,将目标函数f(x)在 $x^{(k+1)}$ 处二阶Taylor展开:

$$\begin{split} f(\mathbf{x}) &\approx \quad f(\mathbf{x}^{(k+1)}) + \nabla f(\mathbf{x}^{(k+1)})^T (\mathbf{x} - \mathbf{x}^{(k+1)}) \\ &+ \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(k+1)})^T \nabla^2 f(\mathbf{x}^{(k+1)}) (\mathbf{x} - \mathbf{x}^{(k+1)}), \end{split}$$

#### 进一步有

$$\nabla f(\mathbf{x}) \approx \nabla f(\mathbf{x}^{(k+1)}) + \nabla^2 f(\mathbf{x}^{(k+1)})(\mathbf{x} - \mathbf{x}^{(k+1)}),$$

于是令 $x = x^{(k)}$ 得

$$\nabla f(\mathbf{x}^{(k)}) \approx \nabla f(\mathbf{x}^{(k+1)}) + \nabla^2 f(\mathbf{x}^{(k+1)})(\mathbf{x}^{(k)} - \mathbf{x}^{(k+1)}).$$

YZW (USTC)

▲□▶ ▲圖▶ ▲圖▶ ▲圖▶ ▲圖 - のへの

记s<sup>(k)</sup> = x<sup>(k+1)</sup> - x<sup>(k)</sup>, y<sup>(k)</sup> = 
$$\nabla f(x^{(k+1)}) - \nabla f(x^{(k)})$$
, 则有  
 $\nabla^2 f(x^{(k+1)})s^{(k)} \approx y^{(k)} \text{ or } \nabla^2 f(x^{(k+1)})^{-1}y^{(k)} \approx s^{(k)}.$ 

这样,计算出s<sup>(k)</sup>和y<sup>(k)</sup>后,可依上式估计在x<sup>(k+1)</sup>处的Hesse矩阵的逆。 我们有理由要求在迭代中构造出Hesse矩阵逆的近似 $H_{k+1}$ ,使其满足

$$H_{k+1}y^{(k)} = s^{(k)}.$$
 (22)

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

通常把式(22)称作正割条件,也称为拟牛顿条件。

#### 拟牛顿迭代算法的一般格式:

- (0) 选取初始点 $x^{(0)}$ , 令 $H_0 = I$ , k := 0.
- (1) 计算搜索方向d<sup>(k)</sup> =  $-H_k \nabla f(\mathbf{x}^{(k)})$ .
- (2) 采用一维搜索确定步长因子 $\alpha_k$ , 令 $x^{(k+1)} = x^{(k)} + \alpha_k d^{(k)}$ .
- (3) 基于x<sup>(k)</sup>到x<sup>(k+1)</sup>的梯度变化,更新Hesse矩阵逆的近似,即确定满 足正割条件的H<sub>k+1</sub>.置k := k + 1,返回第(1)步。

・ 何 ト ・ ヨ ト ・ ヨ ト …

# 下面我们就来讨论怎样构造及确定满足拟牛顿条件的Hesse矩阵逆的近 似*H<sub>k+1</sub>*.



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

设H<sub>k</sub>是第k次迭代的Hesse矩阵逆的近似,我们希望以H<sub>k</sub>来产生H<sub>k+1</sub>,即

$$H_{k+1}=H_k+E_k,$$

其中E<sub>k</sub>是一个低秩的矩阵。

为此,可采用对称秩一(SR1)校正

$$H_{k+1} = H_k + a u u^T$$
,  $(a \in \mathbb{R}, u \in \mathbb{R}^n)$ .

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ ― 臣 – ∽へ⊙

#### 由拟牛顿条件(22)知

$$H_{k+1}y^{(k)} = H_{k}y^{(k)} + (au^{T}y^{(k)})u = s^{(k)}$$
  
故u必与方向s^{(k)} - H\_{k}y^{(k)} - 致, 且假定s^{(k)} - H\_{k}y^{(k)} \neq 0.  
不妨取 $u = s^{(k)} - H_{k}y^{(k)},$ 此时 $a = \frac{1}{u^{T}y^{(k)}},$ 从而得到  

$$H_{k+1} = H_{k} + \frac{(s^{(k)} - H_{k}y^{(k)})(s^{(k)} - H_{k}y^{(k)})^{T}}{(s^{(k)} - H_{k}y^{(k)})^{T}y^{(k)}}.$$
(23)

上式称为对称秩一校正。

æ

イロト イヨト イヨト イヨト

#### 二次终止性

**定义:** 如果一种迭代法能在确知的有限步内找到二次函数的极小点,则称这种方法具有二次终止性。

YZW (USTC)

65 / 467

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

对称秩一校正的突出性质:

- **①** 针对二次函数具有遗传性,即  $H_k y^{(\ell)} = s^{(\ell)}, \ell = k 1, \dots, 1, 0.$
- ❷ 具有二次终止性,即对于二次函数不需要进行一维搜索而具有n步终止性质,且H<sub>n</sub> = [∇<sup>2</sup>f(x\*)]<sup>-1</sup>.

[思考题:请证明对称秩一校正拟牛顿法的上述性质。]

< 口 > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

对称秩一校正的缺点是,不能保持迭代矩阵 $H_{k+1}$ 的正定性。 仅当(s<sup>(k)</sup> –  $H_k y^{(k)}$ )<sup>T</sup> $y^{(k)} > 0$ 时,对称秩一校正才能保持正定性。而这个 条件往往很难保证,即使(s<sup>(k)</sup> –  $H_k y^{(k)}$ )<sup>T</sup> $y^{(k)} > 0$ 满足,它也可能很小 从而导致数值上的困难。

这些都使得对称秩一校正的拟牛顿法应用有较大局限性。

イロト (周) (三) (三) (三)

#### 采用对称秩二(SR2)校正

$$H_{k+1} = H_k + a u u^T + b v v^T,$$

并使得拟牛顿条件(22)成立,则有

$$H_{k+1}y^{(k)} = H_k y^{(k)} + (au^T y^{(k)})u + (bv^T y^{(k)})v = s^{(k)}.$$

这里u,v显然不是唯一确定的,但有一种明显的选择是:

$$\begin{cases} u = s^{(k)}, & au^{T}y^{(k)} = 1; \\ v = H_{k}y^{(k)}, & bv^{T}y^{(k)} = -1. \end{cases}$$

YZW (USTC)

68 / 467

イロト イ団ト イヨト イヨト 二日

#### 因此有

$$H_{k+1} = H_k + \frac{s^{(k)}s^{(k)}}{s^{(k)}y^{(k)}} - \frac{H_k y^{(k)}y^{(k)}H_k}{y^{(k)}H_k y^{(k)}}.$$
 (24)

上式称为 DFP(Davidon-Fletcher-Powell)校正公式,由Davidon(1959)提出,后经Fletcher & Powell(1963)修改而来。

3

DFP校正(24)是典型的拟牛顿校正公式,它有很多重要性质。

(一)对于二次函数(采用精确一维搜索)

② 二次终止性,即
$$H_n = [\nabla^2 f(\mathsf{x}^*)]^{-1}$$
.

头轭性,即当取H<sub>0</sub> = /时,迭代产生共轭方向。

(二) 对于一般非线性函数

- ❶ 校正保持正定性,因而d<sup>(k)</sup>总是下降方向。
- ④ 每次迭代需要3n<sup>2</sup> + O(n)次乘法运算。
- 方法具有超线性收敛速度。

#### 拟牛顿(正割)条件

$$H_{k+1}\mathsf{y}^{(k)}=\mathsf{s}^{(k)}$$

其中 $H_{k+1}$ 是Hesse矩阵逆的近似;

$$B_{k+1}\mathsf{s}^{(k)}=\mathsf{y}^{(k)}$$

其中 $B_{k+1}$ 是Hesse矩阵的近似。

YZW (USTC)

▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ 臣 のへで

#### 由对称秩二校正和拟牛顿条件 $H_{k+1}y^{(k)} = s^{(k)}$ 可得到 $H_k$ 的DFP校正公式

$$H_{k+1}^{(DFP)} = H_k + \frac{s^{(k)}s^{(k)}}{s^{(k)}{}^T y^{(k)}} - \frac{H_k y^{(k)} y^{(k)}{}^T H_k}{y^{(k)}{}^T H_k y^{(k)}}$$

YZW (USTC)

▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ 臣 のへで
BFGS (Broyden-Fletcher-Goldfarb-Shanno) 校正

类似地,我们可从拟牛顿条件 $B_{k+1}$ s<sup>(k)</sup> = y<sup>(k)</sup>得到关于 $B_k$ 的对称秩二校正公式

$$B_{k+1}^{(BFGS)} = B_k + \frac{y^{(k)}y^{(k)}}{y^{(k)}{}^{\mathsf{T}}{}_{\mathsf{S}}^{(k)}} - \frac{B_k s^{(k)}s^{(k)}{}^{\mathsf{T}}B_k}{s^{(k)}{}^{\mathsf{T}}B_k s^{(k)}}.$$
 (25)

把(25)式称为关于B<sub>k</sub>的BFGS校正。

如果我们对 $B_k$ 的BFGS校正"求逆",就可以得到关于 $H_k$ 的BFGS校正公式

$$H_{k+1}^{(BFGS)} = H_{k} + \left(1 + \frac{y^{(k)}{}^{T}H_{k}y^{(k)}}{s^{(k)}{}^{T}y^{(k)}}\right)\frac{s^{(k)}s^{(k)}{}^{T}}{s^{(k)}{}^{T}y^{(k)}} - \frac{H_{k}y^{(k)}s^{(k)}{}^{T} + s^{(k)}y^{(k)}{}^{T}H_{k}}{s^{(k)}{}^{T}y^{(k)}}.$$

[思考题:请给出 $H_{k+1}^{(BFGS)}$ 的对称秩二校正的特解,即a, u, b, v.]

크

# 进一步,若将(26)式中 { $H \leftrightarrow B$ , s $\leftrightarrow$ y}互换,便得到关于 $B_k$ 的DFP校正 公式

$$B_{k+1}^{(DFP)} = B_{k} + \left(1 + \frac{s^{(k)^{T}} B_{k} s^{(k)}}{y^{(k)^{T}} s^{(k)}}\right) \frac{y^{(k)} y^{(k)^{T}}}{y^{(k)^{T}} s^{(k)}} - \frac{B_{k} s^{(k)} y^{(k)^{T}} + y^{(k)} s^{(k)^{T}} B_{k}}{y^{(k)^{T}} s^{(k)}}.$$

YZW (USTC)

75 / 467

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

(27)

### 秩一校正的求逆公式

Sherman-Morrison定理:  $设A \in \mathbb{R}^{n \times n}$ 是非奇异阵,  $u, v \in \mathbb{R}^n$ 是任意向量。若 $1 + v^T A^{-1}u \neq 0$ , 则*A*的秩一校正 $A + uv^T$ 非奇异, 且其逆可以表示为

$$(A + uv^{T})^{-1} = A^{-1} - \frac{A^{-1}uv^{T}A^{-1}}{1 + v^{T}A^{-1}u}.$$
 (28)

[思考题:利用秩一校正的求逆公式,由 $H_{k+1}^{(DFP)}$ 推导 $B_{k+1}^{(DFP)}$ ]

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

### 进一步的参考资料

- R. Fletcher, Practical Methods of Optimization (2nd Edition). John Wiley & Sons, 1987.
- D. C. Liu and J. Nocedal, On the Limited Memory Method for Large Scale Optimization. Mathematical Programming B, 45(3), pp. 503-528, 1999.

...

### 共轭方向

**定义:** 设*G*是*n*×*n*正定阵, ℝ<sup>n</sup>中的任一组非零向量 { $d^{(0)}, d^{(1)}, \dots, d^{(k)}$ }, 如果 $d^{(i)^{T}}Gd^{(j)} = 0$ (*i* ≠ *j*), 则称 $d^{(0)}, d^{(1)}, \dots, d^{(k)}$ 是*G*-共轭的。

显然共轭是正交概念的推广,当取G = I时,共轭即为正交。

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

### 共轭方向法(类):

- (0) 给定正定阵*G*,选取初始点×<sup>(0)</sup>,计算g<sup>(0)</sup> =  $\nabla f(x^{(0)})$ 并构造d<sup>(0)</sup>使 得g<sup>(0) T</sup>d<sup>(0)</sup> < 0. 令*k* := 0.
- (1) 求精确的一维搜索步长 $\alpha_k$ , 即 $\alpha_k = \arg\min_{\alpha>0} f(\mathbf{x}^{(k)} + \alpha \mathsf{d}^{(k)})$ .
- (2) 更新迭代点x<sup>(k+1)</sup> = x<sup>(k)</sup> +  $\alpha_k d^{(k)}$ , 并构造d<sup>(k+1)</sup>使 得d<sup>(k+1)<sup>T</sup></sup>Gd<sup>(j)</sup> = 0, j = 0, 1, ..., k.
- (3) 置k := k + 1,返回第(1)步。

▲ 伺 ▶ ▲ ヨ ▶ ▲ ヨ ▶ …

共轭方向法是从研究二次函数的极小化问题中产生的,但它可以推广到 处理非二次函数的极小化问题。

共轭方向法的一个重要性质是,只要执行精确一维搜索,迭代算法就具 有二次终止性。

**共轭方向法基本定理:**严格凸二次函数 $f(x) = \frac{1}{2}x^T G x + c^T x$ , 共轭方向法执行精确一维搜索,则每步迭代点 $x^{(k+1)} \ge f(x)$ 在线性流形

$$\mathcal{V} = \{ \mathsf{x} \mid \mathsf{x} = \mathsf{x}^{(0)} + \sum_{j=0}^{k} \beta_j \mathsf{d}^{(j)}, \forall \beta_j \in \mathbb{R} \}$$

中的唯一极小点。

**证明:** 设共轭方向法产生的*G*-共轭方向为d<sup>(0)</sup>, d<sup>(1)</sup>,  $\cdots$ , d<sup>(k)</sup>. 由共轭方向的定义知,  $\{d^{(0)}, d^{(1)}, \cdots, d^{(k)}\}$ 线性无关。

下面只要证: 对所有k < n成立

$$g^{(k+1)} d^{(j)} = 0, \ j = 0, 1, \cdots, k.$$

即在点 $x^{(k+1)}$ 处的函数梯度 $g^{(k+1)} = \nabla f(x^{(k+1)})$ 与子空间 $span\{d^{(0)}, d^{(1)}, \dots, d^{(k)}\}$ 正交。

由此易得出定理的结论。

イロト (周) (三) (三) (三)

# **证明(续):** 直接由精确一维搜索知,对∀*j*成立 g<sup>(j+1)<sup>T</sup>d<sup>(j)</sup> = 0.</sup>

特别地, 当 j = k 时,  $g^{(k+1)} d^{(k)} = 0$ .

▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ 臣 のへで

### 证明(续): 事实上,由于

$$y^{(k)} = g^{(k+1)} - g^{(k)} = G(x^{(k+1)} - x^{(k)}) = Gs^{(k)} = \alpha_k Gd^{(k)}.$$

故当*j* < k时有

$$g^{(k+1)}{}^{T} d^{(j)} = g^{(j+1)}{}^{T} d^{(j)} + \sum_{i=j+1}^{k} y^{(i)}{}^{T} d^{(j)}$$
  
=  $g^{(j+1)}{}^{T} d^{(j)} + \sum_{i=j+1}^{k} \alpha_{i} d^{(i)}{}^{T} G d^{(j)}$   
=  $0 + 0$   
=  $0$ 

▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ 臣 のへで

### 证明(续): 综合上述,从而证明了

$${\sf g}^{(k+1)}{}^T{\sf d}^{(j)}=0,\; j=0,1,\cdots,k.$$

YZW (USTC)

▲ロ ▶ ▲圖 ▶ ▲ 圖 ▶ ▲ 圖 ▶ ● 圖 ● の Q @

### **推论:** 对于严格凸的二次函数,若沿着一组共轭方向搜索, 经有限步 迭代必达到极小点。

▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ 臣 のへで

由于共轭方向法具有二次终止性,人们希望能给出一个具体的算法(属于共轭方向法类)。通过修改最速下降法,使其搜索方向具有共轭性质, 这便是共轭梯度法。

下面我们先针对二次函数,给出共轭梯度法的具体描述。

设二次函数 $f(x) = \frac{1}{2}x^T G x + c^T x$ ,其中*G*是 $n \times n$ 正定阵,c是n维向量。 函数f的梯度向量为

$$g(x) = \nabla f(x) = Gx + c.$$

取 $d^{(0)} = -g^{(0)}$ ,因为 $x^{(1)} = x^{(0)} + \alpha_0 d^{(0)}$ 中步长 $\alpha_0$ 由精确一维搜索决定, 所以 $g^{(1)}^T d^{(0)} = 0$ .

现设 $d^{(1)} = -g^{(1)} + \beta_0^{(1)}d^{(0)}$ ,选择 $\beta_0^{(1)}$ 使 $d^{(1)}$ <sup>T</sup> $Gd^{(0)} = 0$ ,即得

$$\beta_0^{(1)} = \frac{g^{(1)} \tau_{g^{(1)}}}{g^{(0)} \tau_{g^{(0)}}}.$$

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ ― 臣 – ∽へ⊙

同理, 令d<sup>(2)</sup> =  $-g^{(2)} + \beta_0^{(2)}d^{(0)} + \beta_1^{(2)}d^{(1)}$ , 选择 $\beta_0^{(2)}, \beta_1^{(2)}$  使 得 $d^{(2)}^T G d^{(j)} = 0, j = 0, 1$ . 从而有

$$\beta_0^{(2)} = 0,$$
  
$$\beta_1^{(2)} = \frac{g^{(2)}{}^T g^{(2)}}{g^{(1)}{}^T g^{(1)}}$$

YZW (USTC)

▲ロ▶ ▲圖▶ ▲臣▶ ▲臣▶ ―臣 - 釣��

### 共轭梯度法

一般地,在第 k 次迭代中,令

$$d^{(k)} = -g^{(k)} + \sum_{j=0}^{k-1} \beta_j^{(k)} d^{(j)},$$

选择 $\beta_j^{(k)}$ 使得d<sup>(k) T</sup>Gd<sup>(j)</sup> = 0,  $j = 0, 1, \dots, k - 1$ , 则有

$$\beta_{j}^{(k)} = \frac{\mathsf{g}^{(k)}{}^{T} G \mathsf{d}^{(j)}}{\mathsf{d}^{(j)}{}^{T} G \mathsf{d}^{(j)}} = \frac{\mathsf{g}^{(k)}{}^{T} (\mathsf{g}^{(j+1)} - \mathsf{g}^{(j)})}{\mathsf{d}^{(j)}{}^{T} (\mathsf{g}^{(j+1)} - \mathsf{g}^{(j)})}.$$

又由于g<sup>(k) T</sup>g<sup>(j)</sup> = 0,  $j = 0, 1, \dots, k - 1$ , 故得

$$\beta_j^{(k)} = 0, \ j = 0, 1, \cdots, k - 2$$
$$\beta_{k-1}^{(k)} = \frac{g^{(k)}{}^T g^{(k)}}{g^{(k-1)}{}^T g^{(k-1)}}.$$

YZW (USTC)

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ ● ● ● ●

## 针对二次函数的共轭梯度算法 (Fletcher & Reeves, 1964) (0) 给定初始点x<sup>(0)</sup>, 计算g<sup>(0)</sup> = g(x<sup>(0)</sup>), 令d<sup>(0)</sup> = -g<sup>(0)</sup>, k := 0. (1) 迭代更新x<sup>(k+1)</sup> = x<sup>(k)</sup> + $\alpha_k d^{(k)}$ , 其中 $\alpha_k = \frac{g^{(k)^T}g^{(k)}}{d^{(k)^T}Gd^{(k)}}$ . (2) 计算g<sup>(k+1)</sup> = g(x<sup>(k+1)</sup>), 构造共轭梯度方 $\operatorname{pd}^{(k+1)} = -g^{(k+1)} + \beta_k d^{(k)}$ , 其中 $\beta_k = \frac{g^{(k+1)^T}g^{(k+1)}}{g^{(k)^T}g^{(k)}}$ . (3) 置k := k + 1, 返回第(1)步。

《曰》 《問》 《글》 《글》 []

## 共轭梯度法

YZW (USTC)



Ontin

æ

**共轭梯度法性质定理:** 设目标函数 $f(x) = \frac{1}{2}x^T G x + c^T x$ ,则采用精确一维搜索的共轭梯度法经 $m \le n$ 步迭代后终止,且对所有的 $1 \le k \le m$ 成立下列关系式:

[思考题:证明上述定理...]

< ロ > < 同 > < 三 > < 三 > 、

# 将共轭梯度法推广到非二次函数的极小化问题,其迭代为 $x^{(k+1)} = x^{(k)} + \alpha_k d^{(k)}.$

步长 $\alpha_k$ 由精确或者非精确一维搜索决定,而d<sup>(k+1)</sup>的构造如下: d<sup>(k+1)</sup> =  $-g^{(k+1)} + \beta_k d^{(k)}$ .

▲ロ▶ ▲圖▶ ▲臣▶ ▲臣▶ ―臣 - 釣��

## 共轭梯度法

## 其中

$$\beta_k := \frac{\mathsf{g}^{(k+1)^T} \mathsf{g}^{(k+1)}}{\mathsf{g}^{(k)^T} \mathsf{g}^{(k)}} \quad \text{(Fletcher - Reeves)}$$

$$\beta_k := \frac{\mathsf{g}^{(k+1)^T}(\mathsf{g}^{(k+1)} - \mathsf{g}^{(k)})}{\mathsf{d}^{(k)^T}(\mathsf{g}^{(k+1)} - \mathsf{g}^{(k)})} \quad (\text{Hestenes - Stiefel})$$

$$\beta_k := \frac{\mathsf{g}^{(k+1)^T}(\mathsf{g}^{(k+1)} - \mathsf{g}^{(k)})}{\mathsf{g}^{(k)^T}\mathsf{g}^{(k)}} \quad (\text{Polak} - \text{Ribiere} - \text{Polyak})$$

$$\beta_k := \frac{\mathsf{g}^{(k+1)^{\mathsf{T}}} \mathsf{g}^{(k+1)}}{-\mathsf{d}^{(k)^{\mathsf{T}}} \mathsf{g}^{(k)}} \quad \text{(Dixon)}$$

对于非二次函数,共轭梯度法迭代*n*步以后所产生的搜索方向  $d^{(k+1)} = -g^{(k+1)} + \beta_k d^{(k)}$ 可能不再是下降方向(由非精确一维搜索造成的)。

因此,*n*步以后我们应该周期性采用最速下降方向作为搜索方向,即  $令 d^{(\ell n)} = -g^{(\ell n)}, \ \ell = 1, 2, \dots$ 

这种策略称为重启动策略,这样的共轭梯度法也称作重启动共轭梯度 法。

《曰》 《問》 《글》 《글》 []

如上所述的共轭梯度法迭代对于一般的非线性函数的最小化也是照样适 用的,但迭代更新的步长因子无法显式表达,需要执行数值近似的非精 确一维搜索。

由于每*n*步迭代执行重启动策略,若记重新启动时得到的点列为{z<sup>(j)</sup>},则可证明这些相隔*n*次的迭代点列超线性收敛。受实际计算误差的影响, 在很多情形下仅能取得类似线性的收敛速率。

从实际计算效率及稳定性来看,共轭梯度法未必比拟牛顿法好。但是, 共轭梯度法中搜索方向的计算仅仅用到目标函数的梯度,而不必像拟牛 顿法那样在每次迭代中更新Hesse矩阵(或其逆)的近似阵并记忆之。 所以,当问题的规模大而且有稀疏结构时,共轭梯度法有高效执行计算 的好处。

#### The preconditioned conjugate gradient method

# 在大多数情况下,为确保共轭梯度法的快速收敛,预条件处理是必要的。



3

(日)

### 进一步的参考资料

- M. Hestenes and E. Stiefel, Methods of conjugate gradients for solving linear systems. Journal of Research of the National Bureau of Standards, 49 (6), 1952.
- K. Atkinson, An Introduction to Numerical Analysis (2nd Edition). John Wiley & Sons, 1988.
- M. Avriel, Nonlinear Programming: Analysis and Methods. Dover Publishing, 2003.
- G. Golub and C. Van Loan, Matrix Computations (3rd Edition). Johns Hopkins University Press.

...

### 为了保证迭代法的全局收敛性,之前我们采用了一维搜索策略。

一维搜索策略先确定一个搜索方向d<sup>(k)</sup>, 然后沿着这个方向选择适当的步长因子 $\alpha_k$ , 新的迭代点×<sup>(k+1)</sup> = ×<sup>(k)</sup> +  $\alpha_k$ d<sup>(k)</sup>.

现在,我们讨论另一种全局收敛策略 — 信赖域方法 (Trust-Region Method).

### 为了保证迭代法的全局收敛性,之前我们采用了一维搜索策略。

一维搜索策略先确定一个搜索方向d<sup>(k)</sup>, 然后沿着这个方向选择适当的步长因子 $\alpha_k$ , 新的迭代点x<sup>(k+1)</sup> = x<sup>(k)</sup> +  $\alpha_k$ d<sup>(k)</sup>.

现在,我们讨论另一种全局收敛策略 — 信赖域方法 (Trust-Region Method).

< ロ > < 同 > < 三 > < 三 > 、

### 为了保证迭代法的全局收敛性,之前我们采用了一维搜索策略。

一维搜索策略先确定一个搜索方向d<sup>(k)</sup>, 然后沿着这个方向选择适当的步长因子 $\alpha_k$ , 新的迭代点x<sup>(k+1)</sup> = x<sup>(k)</sup> +  $\alpha_k$ d<sup>(k)</sup>.

现在,我们讨论另一种全局收敛策略

— 信赖域方法 (Trust-Region Method).

< ロ > < 同 > < 三 > < 三 > 、

### 信赖域方法首先定义当前迭代点x<sup>(k)</sup>的邻域

$$\Omega_k = \{ \mathsf{x} \in \mathbb{R}^n \mid \|\mathsf{x} - \mathsf{x}^{(k)}\| \le e_k \},\$$

这里 $\Omega_k$ 称为信赖域,  $e_k$ 是信赖域半径。

假定在这个邻域里,二次模型 $q^{(k)}(s)$ 是目标函数f(x)的一个合适的近似,则在信赖域中极小化二次模型,得到近似极小点 $s^{(k)}$ ,并 $\operatorname{px}^{(k+1)} = x^{(k)} + s^{(k)}$ .

▲□▶ ▲冊▶ ▲臣▶ ▲臣▶ ―臣 \_ ����

信赖域方法利用二次模型在信赖域内求得方向步s<sup>(k)</sup>, 使得目标函数的下降比一维搜索更有效。

信赖域方法不仅具有全局收敛性,而且不要求目标函数的Hesse矩阵 (或其近似)是正定的。

A D N A (B) N A B N A B N B

### 信赖域子问题

min 
$$q^{(k)}(s) = f(x^{(k)}) + g^{(k)^T}s + \frac{1}{2}s^T B_k s$$
 (29)  
s.t.  $\|s\| \le e_k$ .

其中s = x - x<sup>(k)</sup>, g<sup>(k)</sup> =  $\nabla f(x^{(k)})$ , 对称阵 $B_k$ 是Hesse矩阵  $\nabla^2 f(x^{(k)})$ 或其近似,  $e_k > 0$ 为信赖域半径,  $\|\cdot\|$ 为某一范数。

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

如何选择信赖域半径e<sub>k</sub>?

我们将根据二次模型 $q^{(k)}(s)$ 对目标函数f(x)的拟合程度自适应地调整信赖域半径。



크

(日)
设子问题(29)的解s<sup>(k)</sup>, 令目标函数的下降量

$$Act_k = f(x^{(k)}) - f(x^{(k)} + s^{(k)})$$

为实际下降量,令二次模型函数的下降量

$$\operatorname{Pre}_k = q^{(k)}(0) - q^{(k)}(s^{(k)})$$

为预测下降量。定义比值

$$r_k = \frac{\operatorname{Act}_k}{\operatorname{Pre}_k} = \frac{f(x^{(k)}) - f(x^{(k)} + s^{(k)})}{q^{(k)}(0) - q^{(k)}(s^{(k)})}.$$

它衡量了二次模型与目标函数之间的一致程度。

イロト イ団ト イヨト イヨト 二日

当 $r_k$ 越接近1,表明二次模型函数 $q^{(k)}(s)$ 与目标函数f的一致性程度越好,此时可以增大半径 $e_k$ 以扩大信赖域。

如果 $r_k > 0$ 但不接近1,我们保持信赖域半径 $e_k$ 不变。

如果 $r_k$ 接近零或取负值,表明 $q^{(k)}(s)$ 与目标函数f的一致性程度不理想, 就减小半径 $e_k$ 以缩小信赖域。

< ロ > < 同 > < 三 > < 三 > 、

## 信赖域算法

(0) 给定初始点x<sup>(0)</sup>, 信赖域半径的上界ē, ε > 0,0 < γ<sub>1</sub> < γ<sub>2</sub> < 1,0 < η<sub>1</sub> < 1 < η<sub>2</sub>. 取e<sub>0</sub> ∈ (0,ē), 令k := 0.
(1) 如果||g<sup>(k)</sup>|| ≤ ε, 停止迭代。否则, 求解信赖域子问题(29)得到s<sup>(k)</sup>.
(2) 计算比值r<sub>k</sub>, 更新迭代点

$$\mathbf{x}^{(k+1)} = \begin{cases} \mathbf{x}^{(k)} + \mathbf{s}^{(k)} & \text{if } r_k > 0, \\ \mathbf{x}^{(k)} & \text{otherwise.} \end{cases}$$

· • @ • • = • • = • · ·

## 信赖域算法

(3) 调整信赖域半径, 令

$$e_{k+1} = \begin{cases} \eta_1 e_k & \text{if } r_k < \gamma_1, \\ e_k & \text{if } \gamma_1 \le r_k < \gamma_2, \\ \min(\eta_2 e_k, \bar{e}) & \text{if } r_k \ge \gamma_2. \end{cases}$$

(4) 置k := k + 1,返回第(1)步。

臣

## 信赖域方法的全局收敛性定理:

设水平集 $L(x^{(0)}) = \{x \mid f(x) \le f(x^{(0)})\}$ 有界,且f(x)在其上 $C^2$ 连续,则由 信赖域算法产生的迭代序列存在聚点 $x^{\infty}$ ,满足一阶和二阶必要条件,即

$$\mathbf{g}^{\infty} = 
abla f(\mathbf{x}^{\infty}) = \mathbf{0}, \quad G_{\infty} = 
abla^2 f(\mathbf{x}^{\infty}) \geq \mathbf{0}.$$

YZW (USTC)

< 口 > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

## 在信赖域算法中关键的一步是,解信赖域子问题(29).

这里我们介绍一种求解信赖域子问题的方法,即由 Powell (1970) 提出的折线法。

所谓折线法,是连接 Cauchy 点(由最速下降法产生的极小点)和牛顿 点(由牛顿法产生的极小点),其连线与信赖域边界的交点取为x<sup>(k+1)</sup>.

《曰》 《問》 《글》 《글》 \_ 글



## 折线法图示





113 / 467

æ

对于二次模型

$$q^{(k)}(-\alpha g^{(k)}) = f(x^{(k)}) - \alpha \|g^{(k)}\|^2 + \frac{1}{2} \alpha^2 g^{(k)T} B_k g^{(k)},$$

精确一维搜索的步长因子可表达为

$$\alpha_k = \frac{\|\mathbf{g}^{(k)}\|^2}{\mathbf{g}^{(k)}{}^T B_k \mathbf{g}^{(k)}}.$$

于是 Cauchy 步为

$$\mathbf{s}_{C}^{(k)} = -\alpha_{k}\mathbf{g}^{(k)} = -\frac{\|\mathbf{g}^{(k)}\|^{2}}{\mathbf{g}^{(k)}{}^{T}B_{k}\mathbf{g}^{(k)}}\mathbf{g}^{(k)}.$$

YZW (USTC)

æ

如果 $\|\mathbf{s}_{C}^{(k)}\| = \|\alpha_{k}\mathbf{g}^{(k)}\| \ge e_{k}$ , 取

$$\mathbf{s}^{(k)} = -\frac{e_k}{\|\mathbf{g}^{(k)}\|}\mathbf{g}^{(k)},$$

### 便得

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \frac{e_k}{\|\mathbf{g}^{(k)}\|} \mathbf{g}^{(k)}.$$

如果 $\|\mathbf{s}_{C}^{(k)}\| = \|\alpha_{k}\mathbf{g}^{(k)}\| < e_{k}$ , 再计算牛顿步

$$\mathsf{s}_N^{(k)} = -B_k^{-1}\mathsf{g}^{(k)}.$$

YZW (USTC)

▲□▶ ▲圖▶ ▲ 圖▶ ▲ 圖▶ ― 圖 … のへで

如果 $\|\mathbf{s}_N^{(k)}\| \le e_k$ ,取  $\mathbf{s}^{(k)} = \mathbf{s}_N^{(k)} = -B_k^{-1}\mathbf{g}^{(k)},$ 否则,取  $\mathbf{s}^{(k)} = \mathbf{s}_C^{(k)} + \lambda(\mathbf{s}_N^{(k)} - \mathbf{s}_C^{(k)}),$ 其中 $\lambda$ 使得  $\|\mathbf{s}_C^{(k)} + \lambda(\mathbf{s}_N^{(k)} - \mathbf{s}_C^{(k)})\| = e_k.$ 

YZW (USTC)

116 / 467

▲□▶ ▲圖▶ ▲ 圖▶ ▲ 圖▶ ― 圖 … のへで

## 综上所述,我们得到

$$\mathbf{x}^{(k+1)} = \begin{cases} \mathbf{x}^{(k)} - \frac{e_k}{\|\mathbf{g}^{(k)}\|} \mathbf{g}^{(k)} & \stackrel{\text{iff}}{=} \|\mathbf{s}_C^{(k)}\| \ge e_k, \\ \mathbf{x}^{(k)} - B_k^{-1} \mathbf{g}^{(k)} & \stackrel{\text{iff}}{=} \|\mathbf{s}_C^{(k)}\| < e_k \mathbf{E} \|\mathbf{s}_N^{(k)}\| \le e_k, \\ \mathbf{x}^{(k)} + \mathbf{s}_C^{(k)} + \lambda(\mathbf{s}_N^{(k)} - \mathbf{s}_C^{(k)}) & \stackrel{\text{iff}}{=} \|\mathbf{s}_C^{(k)}\| < e_k \mathbf{E} \|\mathbf{s}_N^{(k)}\| > e_k. \end{cases}$$

$$(30)$$

折线法满足下列性质:

- 1) 沿着 Cauchy  $dx_{C}^{(k+1)}$ 和牛顿 $dx_{N}^{(k+1)}$ 的连线, 到 $x^{(k)}$ 的距离单调增加;
- 2) 沿着 Cauchy  $A_C^{(k+1)}$ 和牛顿 $A_N^{(k+1)}$ 的连线,子问题模型函数值单 调减少。

[思考题:证明上述性质...]

イロト イ団ト イヨト イヨト 二日

- Exercise 1: 请写出上述基于Wolfe-Powell准则的非精确一维搜索算法中插值多项式p<sup>(1)</sup>(t), p<sup>(2)</sup>(t)的具体表达式。
- Exercise 2: 请证明基于Goldstein准则的非精确一维搜索算法的全局 收敛性。
- Exercise 3: 试将非线性方程组求根F(x) = 0的牛顿迭代,用于求解 无约束最优化问题 min<sub>x∈ℝ<sup>n</sup></sub> f(x). 请给出相应的迭代格式,并说明 理由。
- Exercise 4: 请证明对称秩一校正拟牛顿法具有二次终止性和遗传性。
- Exercise 5:利用秩一校正的求逆公式(Sherman-Morrison定理), 由H<sup>(DFP)</sup>推导B<sup>(DFP)</sup><sub>k+1</sub>.

- Exercise 6: 请证明共轭梯度法的性质定理。
- Exercise 7: 请证明折线法(信赖域方法)子问题模型的函数单调 性。
- Exercise 8: 在信赖域方法中,请给出一种与调整信赖域半径等效 的自适应模式算法。

3

## Outline I

## Unconstrained Optimization

## Constrained Optimization

- 二次规划
- 非线性约束最优化
- 3 Convex Optimization
  - Convex Set and Convex Function
  - Convex Optimization and Algorithms

## 4 Sparse Optimization

- Sparse Optimization Models
- Sparse Optimization Algorithms

## Optimization Methods for Machine Learning

YZW (USTC)

→ < ∃ →</p>

## Outline II

- Typical Form of Problems
- Stochastic Algorithms
- Other Popular Methods



æ

#### D Onconstrained Optimization

## 2 Constrained Optimization

- 二次规划
- 非线性约束最优化
- 3 Convex Optimization
  - Convex Set and Convex Function
  - Convex Optimization and Algorithms
- 4 Sparse Optimization
  - Sparse Optimization Models
  - Sparse Optimization Algorithms
- 5 Optimization Methods for Machine Learning
  - Typical Form of Problems
  - Stochastic Algorithms
  - Other Popular Methods yzw (ustc)

二次规划(Quadratic Programming)是指,在变量的线性等式和/或不等 式限制下求二次函数的极小点问题

min 
$$Q(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T G \mathbf{x} + \mathbf{c}^T \mathbf{x}$$
  
s.t.  $\mathbf{a}_i^T \mathbf{x} = b_i, i \in \mathcal{E} = \{1, \cdots, m_e\}$   
 $\mathbf{a}_i^T \mathbf{x} \ge b_i, i \in \mathcal{I} = \{m_e + 1, \cdots, m\}$  (31)

我们假定G为对称阵,  $a_i(i \in \mathcal{E})$ 是线性无关的。

3

二次规划

二次规划的约束可能不相容,也可能没有有限的最小值,这时称二次规 划问题无解。

如果矩阵G半正定,问题(31)是凸二次规划问题,它的任意局部解也是 整体解。

如果矩阵G正定,问题(31)是正定二次规划问题,只要存在解即是唯一的。

如果矩阵G不定,问题(31)是一般的二次规划问题,有可能出现非整体 解的局部解。

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ ― 臣 – ∽へで

## 等式约束二次规划问题

min 
$$Q(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T G \mathbf{x} + \mathbf{c}^T \mathbf{x}$$
  
s.t.  $A\mathbf{x} = \mathbf{b}$  (32)

这里A是 $m \times n$ 矩阵,且不失一般性可设 rank(A) = m.

YZW (USTC)

126 / 467

▲□▶ ▲圖▶ ▲ 圖▶ ▲ 圖▶ ― 圖 … のへで

二次规划

设有一种基分解x = 
$$\begin{pmatrix} x_B \\ x_N \end{pmatrix}$$
, 其中x $_B \in \mathbb{R}^m, x_N \in \mathbb{R}^{n-m}$ , 使得其约束矩  
阵的对应分块 $A = (A_B, A_N)$ 中 $A_B$ 可逆。于是,等式约束条件可写成

$$\mathsf{x}_B = A_B^{-1}(\mathsf{b} - A_N\mathsf{x}_N),$$

并将上式代入目标函数中得到无约束问题

$$\min_{\mathsf{x}_N \in \mathbb{R}^{n-m}} \frac{1}{2} \mathsf{x}_N^T \hat{G}_N \mathsf{x}_N + \hat{\mathsf{c}}_N^T \mathsf{x}_N.$$
(33)

YZW (USTC)

æ

二次规划

## 在上式中

$$\begin{split} \hat{G}_{N} &= G_{NN} - G_{NB} A_{B}^{-1} A_{N} - A_{N}^{T} A_{B}^{-T} G_{BN} + A_{N}^{T} A_{B}^{-T} G_{BB} A_{B}^{-1} A_{N}, \\ \hat{c}_{N} &= c_{N} - A_{N}^{T} A_{B}^{-T} c_{B} + G_{NB} A_{B}^{-1} b - A_{N}^{T} A_{B}^{-T} G_{BB} A_{B}^{-1} b, \\ & \textbf{以及对应分块形式} \end{split}$$

$$G = \left( \begin{array}{cc} G_{BB} & G_{BN} \\ G_{NB} & G_{NN} \end{array} \right), \quad c = \left( \begin{array}{c} c_B \\ c_N \end{array} \right).$$

YZW (USTC)

128 / 467

æ

二次规划

(1) 如果Ĝ<sub>N</sub>正定,则无约束问题的解可唯一地给出

$$\mathsf{x}_{\boldsymbol{N}}^{*} = -\hat{G}_{\boldsymbol{N}}^{-1}\hat{\mathsf{c}}_{\boldsymbol{N}},$$

进一步得原问题(32)的解为

$$\mathbf{x}^* = \begin{pmatrix} \mathbf{x}^*_B \\ \mathbf{x}^*_N \end{pmatrix} = \begin{pmatrix} A_B^{-1}\mathbf{b} \\ 0 \end{pmatrix} + \begin{pmatrix} A_B^{-1}A_N \\ -I \end{pmatrix} \hat{G}_N^{-1}\hat{\mathbf{c}}_N.$$

设x\*对应的Lagrange乘子向量为 $\lambda$ \*,则有

 $G\mathbf{x}^* + \mathbf{c} = A^T \lambda^* \Longrightarrow \lambda^* = A_B^{-T} (G_{BB} \mathbf{x}_B^* + G_{BN} \mathbf{x}_N^* + \mathbf{c}_B).$ 

YZW (USTC)

129 / 467

二次规划

(2) 如果 $\hat{G}_N$ 是半正定的,则在 $(I - \hat{G}_N \hat{G}_N^+)\hat{c}_N = 0$ 时,无约束问题有界, 且它的解可表示为

$$\mathsf{x}^*_{\mathcal{N}} = - \hat{G}^+_{\mathcal{N}} \hat{\mathsf{c}}_{\mathcal{N}} + (I - \hat{G}^+_{\mathcal{N}} \hat{G}_{\mathcal{N}}) \widetilde{\mathsf{y}},$$

其中 $\tilde{y} \in \mathbb{R}^{n-m}$ 为任意向量, $\hat{G}_N^+$ 表示 $\hat{G}_N$ 的广义逆矩阵。此时,原问题的解x\*和相应最优乘子 $\lambda$ \*可类似确定。

当( $I = \hat{G}_N \hat{G}_N^+$ ) $\hat{c}_N = 0$ 不成立时,则可推出无约束问题无下界,从而 原问题也无下界。



# (3) 如果 $\hat{G}_N$ 不定(即存在负的特征根),显然无约束问题无下界,故原问题不存在有限最优解。



▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ 臣 のへで



# 上述消去法的不足之处是,当A<sub>B</sub>接近奇异时,容易导致数值计算的不稳定。



▲ロ ▶ ▲圖 ▶ ▲ 圖 ▶ ▲ 圖 ▶ ● 圖 ● の Q @

二次规划

## 广义消去法

设 $Z = \{z_{m+1}, \dots, z_n\}$ 解空间Ker(A)的一组基, Y =  $\{y_1, \dots, y_m\}$ 是商空间ℝ<sup>n</sup>/Ker(A)的一组基, 则∀x ∈ ℝ<sup>n</sup>可作如下分解表达

$$\mathbf{x} = Y\mathbf{x}_Y + Z\mathbf{x}_Z.$$

## 从而有

$$A\mathbf{x} = \mathbf{b} \Longrightarrow AY\mathbf{x}_Y + AZ\mathbf{x}_Z = \mathbf{b} \Longrightarrow \mathbf{x}_Y = (AY)^{-1}\mathbf{b},$$

所以得

$$\mathsf{x} = \mathsf{Y}(\mathsf{A}\mathsf{Y})^{-1}\mathsf{b} + \mathsf{Z}\mathsf{x}_{\mathsf{Z}},$$

其中 $x_Z \in \mathbb{R}^{n-m}$ 是自由变量。

▲ロ▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ ▲ 国 ● のへで

二次规划

## 广义消去法

将上式代入目标函数中得无约束问题

$$\min_{\mathbf{x}_{Z} \in \mathbb{R}^{n-m}} \frac{1}{2} \mathbf{x}_{Z}^{T} (Z^{T} G Z) \mathbf{x}_{Z} + [Z^{T} G Y (A Y)^{-1} \mathbf{b} + Z^{T} \mathbf{c}]^{T} \mathbf{x}_{Z}.$$
(34)

假定 $Z^T GZ$ 正定,则有

$$\mathsf{x}_Z^* = -(Z^T G Z)^{-1} Z^T [G Y (A Y)^{-1} \mathsf{b} + \mathsf{c}].$$

YZW (USTC)

134 / 467

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

二次规划

## 广义消去法

## 从而得到原问题的最优解

$$x^* = Y(AY)^{-1} - Z(Z^T GZ)^{-1}Z^T[GY(AY)^{-1}b + c],$$

相应的Lagrange乘子为

$$\lambda^* = (AY)^{-T} Y^T (Gx^* + c).$$

YZW (USTC)

135 / 467

臣

二次规划

Lagrange方法是基于求解可行域内的(K-T)点,即Lagrange函数的稳定 点。

对于等式约束问题(32), 其Lagrange函数的稳定点就是如下线性方程组的解

$$\int G\mathbf{x} + \mathbf{c} = A^T \lambda,$$
$$A\mathbf{x} = \mathbf{b}.$$

写成矩阵形式得

$$\left(\begin{array}{cc} G & -A^{\mathsf{T}} \\ -A & 0 \end{array}\right) \left(\begin{array}{c} \mathsf{x} \\ \lambda \end{array}\right) = - \left(\begin{array}{c} \mathsf{c} \\ \mathsf{b} \end{array}\right).$$

э

(4日)

二次规划

设矩阵
$$\begin{pmatrix} G & -A^{T} \\ -A & 0 \end{pmatrix}$$
可逆,则存在矩  
阵 $U \in \mathbb{R}^{n \times n}, V \in \mathbb{R}^{m \times m}, W \in \mathbb{R}^{m \times n}$  使得 $\begin{pmatrix} G & -A^{T} \\ -A & 0 \end{pmatrix}^{-1} = \begin{pmatrix} U & W^{T} \\ W & V \end{pmatrix}$ 

从而可求得问题的唯一解

$$\left\{ \begin{array}{l} \mathsf{x}^* = - U \mathsf{c} - W^T \mathsf{b}, \\ \lambda^* = - W \mathsf{c} - V \mathsf{b}. \end{array} \right.$$

YZW (USTC)

137 / 467

Ξ.

イロン イ理 とく ヨン イ ヨン



# 上述Lagrange方法中的矩阵非奇异性并不一定要求 $G^{-1}$ 存在,可用不同的方法给出分块矩阵U, V, W的表达形式,从而导致不同的计算公式。

A D N A (B) N A B N A B N B

二次规划

当G可逆, A行满秩, 则( $AG^{-1}A^{T}$ )<sup>-1</sup>存在, 不难验证

$$\begin{cases} U = G^{-1} - G^{-1}A^{T}(AG^{-1}A^{T})^{-1}AG^{-1}, \\ V = -(AG^{-1}A^{T})^{-1}, \\ W = -(AG^{-1}A^{T})^{-1}AG^{-1}. \end{cases}$$

于是我们得到求解公式

$$\begin{cases} x^* = -G^{-1}c + G^{-1}A^T (AG^{-1}A^T)^{-1} (AG^{-1}c + b), \\ \lambda^* = (AG^{-1}A^T)^{-1} (AG^{-1}c + b). \end{cases}$$

YZW (USTC)

▲□▶ ▲圖▶ ▲ 圖▶ ▲ 圖▶ ― 圖 … のへで

二次规划

如果取*Y*, *Z*满足(*Y*, *Z*) = 
$$\begin{pmatrix} A \\ B \end{pmatrix}^{-1}$$
, 即*AY* = *I*<sub>m×m</sub>, *AZ* = 0.  
若另有*Z<sup>T</sup>GZ*可逆, 则知 $\begin{pmatrix} G & -A^T \\ -A & 0 \end{pmatrix}$ 可逆。此时  
 $\begin{cases} U = Z(Z^TGZ)^{-1}Z^T, \\ V = -Y^TGP^TY, \\ W = -Y^TP. \end{cases}$ 

其中 $P = I - GZ(Z^T GZ)^{-1}Z^T$ .

▲□▶ ▲圖▶ ▲国▶ ▲国▶ ▲国 ● のへで

二次规划

基于 $A^T$ 的QR分解,可给出(Y, Z)的一种特殊取法: 设

$$A^{T} = Q \left( egin{array}{c} R \\ 0 \end{array} 
ight) = (Q_{1}, Q_{2}) \left( egin{array}{c} R \\ 0 \end{array} 
ight),$$

1	
$\sim$	

$$A = (R^T, 0) \begin{pmatrix} Q_1^T \\ Q_2^T \end{pmatrix}.$$

其中Q为 $n \times n$ 正交阵, R为 $m \times m$ 上三角阵。于是 令 $Y = Q_1 R^{-T}, Z = Q_2, 则有$ 

 $AY = R^T Q_1^T Q_1 R^{-T} = I_{m \times m}, \quad AZ = R^T Q_1^T Q_2 = 0_{m \times (n-m)}.$ 

▲ロ▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ ▲ 国 ● のへで

二次规划

一般的二次规划

min 
$$Q(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T G \mathbf{x} + \mathbf{c}^T \mathbf{x}$$
  
s.t.  $\mathbf{a}_i^T \mathbf{x} = \mathbf{b}_i, i \in \mathcal{E} = \{1, \cdots, m_e\}$   
 $\mathbf{a}_i^T \mathbf{x} \ge \mathbf{b}_i, i \in \mathcal{I} = \{m_e + 1, \cdots, m\}$ 
(35)

直观上,不积极的不等式约束在解的附近不起作用,可去掉不予考虑; 而积极的不等式约束,由于它在解处等号成立,故我们可以用等式约束 来代替这些积极的不等式约束。
二次规划

**积极集基本定理:**设x\*是一般的二次规划问题(35)的局部极小点,则x\*也必是等式约束问题

(EQ) 
$$\begin{cases} \min \quad Q(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T G \mathbf{x} + \mathbf{c}^T \mathbf{x} \\ \text{s.t.} \quad \mathbf{a}_i^T \mathbf{x} = b_i, i \in \mathcal{E} \cup \mathcal{I}(\mathbf{x}^*) \end{cases}$$

的局部极小点。反之,如果x\*是一般问题(35)的可行点,同时 是(EQ)的K-T点,且相应的Lagrange乘子 $\lambda$ \*满足  $\lambda_i^* \ge 0, i \in \mathcal{I}(x^*),$ 则x\*必是原问题(35)的K-T点。

[习题6.1:证明上述定理...]

イロト イ理ト イヨト イヨト

二次规划

### 设x<sup>(k)</sup>为当前迭代点,且是问题(35)的可行点。

记 $\mathcal{E}_k = \mathcal{E} \cup \mathcal{I}(\mathbf{x}^{(k)}),$ 考虑等式约束问题

(EQ1) 
$$\begin{cases} \min & \frac{1}{2} \mathbf{s}^T G \mathbf{s} + (G \mathbf{x}^{(k)} + \mathbf{c})^T \mathbf{s} \\ \text{s.t.} & \mathbf{a}_i^T \mathbf{s} = 0, i \in \mathcal{E}_k \end{cases}$$

求得(EQ1)的解s<sup>(k)</sup>,及其相应的Lagrange乘子 $\lambda_i^{(k)}, i \in \mathcal{E}_k$ .

YZW (USTC)

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ ― 臣 – ∽へで

### 二次规划

(a) s<sup>(k)</sup> ≠ 0时, x<sup>(k)</sup>不可能是原问题的K-T点。
 (b) s<sup>(k)</sup> = 0时, x<sup>(k)</sup>是问题

(EQ2) 
$$\begin{cases} \min & \frac{1}{2} \mathbf{x}^T G \mathbf{x} + \mathbf{c}^T \mathbf{x} \\ \text{s.t.} & \mathbf{a}_i^T \mathbf{x} = b_i, i \in \mathcal{E}_k \end{cases}$$

的K-T点;如果 $\lambda_i^{(k)} \ge 0, i \in \mathcal{I}(x^{(k)}),$ 则 $x^{(k)}$ 也是原问题的K-T点。 (c)否则,由 $\lambda_{i_q}^{(k)} = \min_{i \in \mathcal{I}(x^{(k)})} \lambda_i^{(k)} < 0确定_{i_q}, 那么如下问题$ 

(EQ3) 
$$\begin{cases} \min & \frac{1}{2} \mathbf{s}^T G \mathbf{s} + (G \mathbf{x}^{(k)} + \mathbf{c})^T \mathbf{s} \\ \text{s.t.} & \mathbf{a}_i^T \mathbf{s} = 0, i \in \hat{\mathcal{E}} = \mathcal{E}_k \setminus \{i_q\}. \end{cases}$$

的解ŝ是原问题在当前点 $x^{(k)}$ 处的可行方向,即 $a_{i_a}^T$ ŝ  $\geq 0$ .

[思考题:证明上述(c)的结论...]

二次规划

### 积极集方法(Active Set Method)

(0) 给出可行点 $x^{(0)}$ , 令 $\mathcal{E}_0 = \mathcal{E} \cup \mathcal{I}(x^{(0)})$ , k := 0. (1) 求解等式约束问题(EQ1)得 $s^{(k)}$ , 若 $s^{(k)} \neq 0$ , 转第(3)步。 (2) 如果 $\lambda_i^{(k)} \ge 0, i \in \mathcal{I}(\mathbf{x}^{(k)})$ ,则停止;否则由 $\lambda_{i_q}^{(k)} = \min_{i \in \mathcal{I}(\mathbf{x}^{(k)})} \lambda_i^{(k)} < 0$ 确 定 $i_a$ 并令 $\mathcal{E}_k := \mathcal{E}_k \setminus \{i_a\}, x^{(k+1)} = x^{(k)}, 转第(4)步$ 。 (3) 由 $\alpha_k = \min\{1, \min_{\substack{i \notin \mathcal{E}_{k}, a_i^T \mathbf{s}^{(k)} < 0}} \frac{b_i - \mathbf{a}_i^T \mathbf{x}^{(k)}}{\mathbf{a}_i^T \mathbf{s}^{(k)}}\},$ 计算 $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{s}^{(k)}.$ 如 得 $a_p^T(\mathbf{x}^{(k)} + \alpha_k \mathbf{s}^{(k)}) = b_p$ , 并令 $\mathcal{E}_k := \mathcal{E}_k \cup \{p\}$ . (4)  $\mathcal{E}_{k+1} := \mathcal{E}_k, k := k+1$ , 返回第(1)步。

YZW (USTC)

146 / 467

### 本章作业(二次规划)

- Exercise 1: 请证明积极集基本定理。
- Exercise 2: 请证明前述(c)的结论。
- Exercise 3 (选做题): 试证明存在半正定的能量阵E<sub>T,m</sub>满足如下 模型解
   1

$$\frac{1}{n} \mathbf{y}^{\mathsf{T}} \mathsf{E}_{\mathsf{T},m} \mathbf{y} = \min_{\substack{\phi \in W_2^m(\Omega) \\ \phi(X_i) = y_i, i = 1, \cdots, n}} |\phi|_{\mathsf{T},m}^2$$

YZW (USTC)

先考虑等式约束问题  
min 
$$f(x)$$
  
s.t.  $c(x) = 0$ , (36)

其中c(x) =  $(c_1(x), \cdots, c_m(x))^T$ .



臣

イロン イ理 とく ヨン イ ヨン

$$\mathbf{i} \mathcal{A}(\mathsf{x}) = [\nabla \mathsf{c}(\mathsf{x})]^{\mathsf{T}} = (\nabla c_1(\mathsf{x}), \cdots, \nabla c_m(\mathsf{x}))^{\mathsf{T}}.$$

由最优性条件知: x是等式约束问题(36)的K-T点当且仅当存在乘  $\mathcal{F}_{\lambda} \in \mathbb{R}^{m}$ 使得

$$\nabla f(\mathbf{x}) - A(\mathbf{x})^T \lambda = \mathbf{0},$$

且×是一可行点,即c(×) = 0.

イロト イポト イヨト イヨト 一日

### 于是得到联立方程组

$$\begin{cases} \nabla f(\mathbf{x}) - A(\mathbf{x})^T \lambda = 0, \\ -\mathbf{c}(\mathbf{x}) = 0. \end{cases}$$

### 我们可用Newton-Raphson迭代法求解上述联立方程组。

YZW (USTC)

3

イロト 不得 トイヨト イヨト

记x和 $\lambda$ 的计算增量分别为 $\delta_x, \delta_\lambda$ , Newton-Raphson迭代满足:

$$\begin{pmatrix} W(\mathbf{x},\lambda) & -A(\mathbf{x})^{T} \\ -A(\mathbf{x}) & 0 \end{pmatrix} \begin{pmatrix} \delta_{\mathbf{x}} \\ \delta_{\lambda} \end{pmatrix} = -\begin{pmatrix} \nabla f(\mathbf{x}) - A(\mathbf{x})^{T} \lambda \\ -c(\mathbf{x}) \end{pmatrix}, \quad (37)$$

其中  $W(\mathbf{x}, \lambda) = \nabla^2 f(\mathbf{x}) - \sum_{i=1}^m \lambda_i \nabla^2 c_i(\mathbf{x}).$ 

イロト イポト イヨト イヨト 一日

### 上述方法称为Lagrange-Newton法,最早由Wilson(1963)提出的。

其实质上是用Newton-Raphson迭代求问题(36)的Lagrange函数  $L(x, \lambda)$ 的稳定点。



3

イロト 不得 トイヨト イヨト

### 在此,我们定义价值函数

$$\psi(\mathbf{x},\lambda) = \|\nabla f(\mathbf{x}) - A(\mathbf{x})^T \lambda\|^2 + \|\mathbf{c}(\mathbf{x})\|^2.$$
(38)

显然,  $\psi(x, \lambda)$ 是关于Lagrange-Newton法的下降函数, 即满足

$$abla \psi(\mathsf{x},\lambda)^{\mathcal{T}} \left( egin{array}{c} \delta_{\mathsf{x}} \ \delta_{\lambda} \end{array} 
ight) = -2\psi(\mathsf{x},\lambda) \leq \mathbf{0}.$$

YZW (USTC)

153 / 467

イロト イ団ト イヨト イヨト 二日

### Lagrange-Newton法:

- (0) 给定 $\mathbf{x}^{(0)}$ ,  $\lambda \in \mathbb{R}^m$ ,  $\beta \in (0, 1)$ ,  $\varepsilon \ge 0$ , 令k := 0.
- (1) 计算价值函数 $\psi(x^{(k)}, \lambda^{(k)})$ , 如果 $\psi(x^{(k)}, \lambda^{(k)}) \le \varepsilon$ , 则停止; 否则 在 $(x^{(k)}, \lambda^{(k)})$ 处求解(37)得到 ( $\delta_{x^{(k)}}, \delta_{\lambda^{(k)}}$ ), 并令 $\alpha_k = 1$ .
- (2) 若 $\psi(\mathbf{x}^{(k)} + \alpha_k \delta_{\mathbf{x}^{(k)}}, \lambda^{(k)} + \alpha_k \delta_{\lambda^{(k)}}) \leq (1 \beta \alpha_k) \psi(\mathbf{x}^{(k)}, \lambda^{(k)}),$ 第(3)步; 否则令 $\alpha_k = \frac{1}{4} \alpha_k$ , 返回第(2)步。
- (3) 置x<sup>(k+1)</sup> = x<sup>(k)</sup> +  $\alpha_k \delta_{x^{(k)}}, \lambda^{(k+1)} = \lambda^{(k)} + \alpha_k \delta_{\lambda^{(k)}}, k := k + 1, 返回 第(1)步。$

- ロ ト - 4 同 ト - 4 回 ト - - - 三

Lagrange-Newton法的收敛性结果

**定理:** 设Lagrange-Newton法产生的迭代点列{ $(x^{(k)}, \lambda^{(k)})$ }有界,如 果f(x)和 $c_i(x)$ 都是二次连续可微,且逆矩阵

$$\left( egin{array}{cc} W({\sf x},\lambda) & -{\cal A}({\sf x})^{{\cal T}} \ -{\cal A}({\sf x}) & 0 \end{array} 
ight)^{-1}$$

一致有界,则{ $(x^{(k)}, \lambda^{(k)})$ }的任何聚点都是方程 $\psi(x, \lambda) = 0$ 的根,从 而{ $x^{(k)}$ }的聚点是问题(36)的K-T点。

注: 在一定条件下, 还可进一步证明Lagrange-Newton法具有二阶收敛 速度。

ヘロト 人間 とくほ とくほ とう



Lagrange-Newton法的一大重要贡献是,在其基础上发展出了逐步二次 规划方法(Sequential Quadratic Programming Methods)。而后者已成为 求解一般非线性约束最优化问题的一类十分重要的方法。



我们可将式(37)写成如下形式:

$$\begin{cases} W(\mathbf{x},\lambda)\delta_{\mathbf{x}} + \nabla f(\mathbf{x}) = A(\mathbf{x})^{T}(\lambda + \delta_{\lambda}), \\ c(\mathbf{x}) + A(\mathbf{x})\delta_{\mathbf{x}} = 0. \end{cases}$$

由最优性条件知, $\delta_{x^{(k)}}$ 即为下列二次规划问题

$$\min_{\substack{k=0 \\ \text{s.t.}}} \frac{1}{2} d^T W(\mathbf{x}^{(k)}, \lambda^{(k)}) d + \nabla f(\mathbf{x}^{(k)})^T d$$

$$\text{s.t.} \quad \mathbf{c}(\mathbf{x}^{(k)}) + A(\mathbf{x}^{(k)}) d = 0$$

$$(39)$$

的K-T点。

3

イロト イヨト イヨト イヨト



Lagrange-Newton法可以理解为逐步求解上述等式约束二次规划的方法。

设d<sup>(k)</sup>是二次规划问题(39)的最优解,那么可迭代更新

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)},$$

其中 $\alpha_k$ 为第k次迭代的步长。

设 $\overline{\lambda}^{(k)}$ 是(39)对应的Lagrange乘子向量,那么对 $k \ge 1$ 有 $\lambda^{(k+1)} = \lambda^{(k)} + \alpha_k (\overline{\lambda}^{(k)} - \lambda^{(k)}).$ 

イロト (個) イヨト イヨト ヨー のくべ



现考虑一般的非线性约束最优化问题

$$\begin{array}{ll} \min & f({\sf x}) \\ {\rm s.t.} & c_i({\sf x}) = 0, i \in \mathcal{E} = \{1, \cdots, m_e\}, \\ & c_i({\sf x}) \ge 0, i \in \mathcal{I} = \{m_e + 1, \cdots, m\}. \end{array}$$

类似地,在第k次迭代里求解子问题

$$\begin{array}{ll} \min & \frac{1}{2} \mathbf{d}^{T} W_{k} \mathbf{d} + \mathbf{g}^{(k)}{}^{T} \mathbf{d} \\ \text{s.t.} & c_{i}(\mathbf{x}^{(k)}) + \mathbf{a}_{i}(\mathbf{x}^{(k)}){}^{T} \mathbf{d} = 0, i \in \mathcal{E}, \\ & c_{i}(\mathbf{x}^{(k)}) + \mathbf{a}_{i}(\mathbf{x}^{(k)}){}^{T} \mathbf{d} \ge 0, i \in \mathcal{I}. \end{array}$$

$$\tag{41}$$

YZW (USTC)

159 / 467

イロト イ団ト イヨト イヨト 二日



在这里,  $W_k$ 是原问题Lagrange函数的Hesse阵或其近似,  $g^{(k)} = \nabla f(x^{(k)}), A(x^{(k)}) = (a_1(x^{(k)}), \dots, a_m(x^{(k)}))^T = [\nabla c(x^{(k)})]^T.$ 

记子问题(41)的解为d<sup>(k)</sup>,相应Lagrange乘子向量为 $\bar{\lambda}^{(k)}$ ,故有

$$\begin{cases} W_k d^{(k)} + g^{(k)} = A(x^{(k)})^T \bar{\lambda}^{(k)} \\ \bar{\lambda}_i^{(k)} \ge 0, i \in \mathcal{I}, \\ c(x^{(k)}) + A(x^{(k)}) d^{(k)} = 0. \end{cases}$$

YZW (USTC)

160 / 467

ヘロト 人間 とくほ とくほ とう



逐步二次规划法的迭代以d<sup>(k)</sup>作为搜索方向。该搜索方向有很好的性质。 它是许多罚函数的下降方向,例如L<sub>1</sub>罚函数

$$P(\mathbf{x},\sigma) = f(\mathbf{x}) + \sigma \left( \sum_{i=1}^{m_e} |c_i(\mathbf{x})| + \sum_{i=m_e+1}^m |c_i(\mathbf{x})| \right)$$

其中c(x)-定义如下:

$$\begin{cases} c_i(\mathbf{x})_- = c_i(\mathbf{x}), & i \in \mathcal{E}, \\ c_i(\mathbf{x})_- = \min\{0, c_i(\mathbf{x})\}, & i \in \mathcal{I}. \end{cases}$$

YZW (USTC)

161 / 467

3



### 下面的算法是Han(1977)提出的逐步二次规划方法:

- (0) 给定 $x^{(0)}, W_0 \in \mathbb{R}^{n \times n}, \sigma > 0, \rho \in (0, 1), \varepsilon \ge 0$ , 令k := 0.
- (1) 求解子问题(41)给出d<sup>(k)</sup>, 如果 $\|d^{(k)}\| \le \varepsilon$ , 则停止; 否则 求 $\alpha_k \in [0, \rho]$ 使得

$$P(\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}, \sigma) \le \min_{0 \le \alpha \le \rho} P(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}, \sigma) + \epsilon_k.$$

(2) 置 $x^{(k+1)} = x^{(k)} + \alpha_k \delta_{x^{(k)}}$ , 计算 $W_{k+1}$ , 令k := k+1, 返回第(1)步。

ヘロト 人間 とくほ とくほ とう



可证明前述逐步二次规划法的收敛性结果如下:

**定理**: 假定f(x)和 $c_i(x)$ 连续可微,且存在常数 $M_1, M_2 > 0$ 使得

 $M_1 \|\mathbf{d}\|^2 \leq \mathbf{d}^T W_k \mathbf{d} \leq M_2 \|\mathbf{d}\|^2, \forall k \in \mathbb{N}, \forall \mathbf{d} \in \mathbb{R}^n,$ 

如果 $\|\lambda^{(k)}\|_{\infty} \leq \sigma$ 均成立,则Han(1977)算法产生的点列 $\{x^{(k)}\}$ 的任何聚 点都是问题(40)的K-T点。

《曰》《問》《曰》《王》 []

#### 对于非线性约束最优化问题

$$\begin{array}{ll} \min & f(\mathbf{x}) \\ \text{s.t.} & c_i(\mathbf{x}) = 0, i \in \mathcal{E} = \{1, \cdots, m_e\} \\ & c_i(\mathbf{x}) \geq 0, i \in \mathcal{I} = \{m_e + 1, \cdots, m\} \end{array}$$

的罚函数,是指利用目标函数*f*(x)和约束方程c(x)所构造的具有"罚性 质"的函数

$$P(\mathsf{x}) = P(f(\mathsf{x}), \mathsf{c}(\mathsf{x})).$$

3

イロト 不得 トイヨト イヨト



# 所谓"罚性质",即要求对问题的可行点 $x \in S$ 均有P(x) = f(x),而当约 束条件破坏时有P(x) > f(x).



165 / 467

▲ロ▶ ▲圖▶ ▲臣▶ ▲臣▶ ―臣 - 釣��

### 为了描述约束条件被破坏的程度,我们定义c(x)\_如下:

$$\begin{cases} c_i(\mathbf{x})_- = c_i(\mathbf{x}), & i \in \mathcal{E}, \\ c_i(\mathbf{x})_- = \min\{0, c_i(\mathbf{x})\}, & i \in \mathcal{I}. \end{cases}$$

YZW (USTC)

臣

イロト イヨト イヨト イヨト

# 罚函数一般可取为目标函数与"罚项"之和,即 $P(x) = f(x) + \phi(c(x)_{-}).$

罚项 $\phi(c(x)_{-})$ 是定义在 $\mathbb{R}^{m}$ 上的函数,它满足

$$\phi(\mathsf{0}) = \mathsf{0}, \quad \lim_{\|\mathsf{c}\| \to \infty} \phi(\mathsf{c}) = +\infty.$$

YZW (USTC)

167 / 467

イロト イ団ト イヨト イヨト 二日

### 如Courant罚函数:

$$P_{\sigma}(\mathbf{x}) = f(\mathbf{x}) + \sigma \|\mathbf{c}(\mathbf{x})_{-}\|_{2}^{2},$$

其中 $\sigma > 0$ 是罚因子。



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

考虑简单罚函数

$$P_{\sigma}(\mathsf{x}) = f(\mathsf{x}) + \sigma \|\mathsf{c}(\mathsf{x})_{-}\|^{2}.$$

# $ilx(\sigma)$ 是无约束问题 $\min_{x \in \mathbb{R}^n} P_{\sigma}(x)$ 的最优解,我们有如下引理。

**引理:** 若 $x(\sigma)$ 同时是非线性约束最优化问题(42)的可行点,则 $x(\sigma)$ 也是 原问题的最优解。

3

上述引理表明,只要选取充分大的罚因子 $\sigma > 0$ ,则通过求解无约束最优化问题应可找到相应约束最优化问题的最优解。

然而在实际计算中,确定大小合适的 $\sigma$ 往往比较困难,故通常是选取一个单调增的罚因子序列 $\{\sigma_k\}$ .

通过求解一系列无约束问题来获得约束最优化问题的解,这称为序贯无 约束极小化技术(SUMT)。

至此,我们可以给出罚函数法的迭代步骤:

- (0) 任选初始点 $x^{(0)}$ , 给定初始罚因子 $\sigma_0 > 0$ 及 $\beta > 1, \varepsilon > 0$ . 令k := 0.
- (1) 以x<sup>(k)</sup>作为初始迭代点求解无约束问题的极小点,即

$$\mathbf{x}(\sigma_k) = \arg\min_{\mathbf{x}\in\mathbb{R}^n} P_{\sigma_k}(\mathbf{x}).$$

 (2) 若 ||c(x(σ<sub>k</sub>))<sub>-</sub>|| < ε, 则停止迭代并取x(σ<sub>k</sub>)为原约束问题的近似最优 解; 否则, 置x<sup>(k+1)</sup> = x(σ<sub>k</sub>), σ<sub>k+1</sub> = βσ<sub>k</sub>, 令<sub>k</sub> := k + 1返回 第(1)步。

易得如下三个引理

引理1: 设 $\sigma_{k+1} > \sigma_k > 0$ , 则有  $P_{\sigma_k}(\mathbf{x}(\sigma_k)) \le P_{\sigma_{k+1}}(\mathbf{x}(\sigma_{k+1}))$ ,

 $\|\mathsf{c}(\mathsf{x}(\sigma_k))_-\| \geq \|\mathsf{c}(\mathsf{x}(\sigma_{k+1}))_-\|, \quad f(\mathsf{x}(\sigma_k)) \leq f(\mathsf{x}(\sigma_{k+1})).$ 

**引理2:** 设令求是原问题(42)的最优解,则对任意的  $\sigma_k > 0$  成立 $f(\bar{x}) \ge P_{\sigma_k}(x(\sigma_k)) \ge f(x(\sigma_k)).$ 

**引理3:** 令 $\delta = \|c(x(\sigma))_{-}\|, Mx(\sigma)$ 也是约束问题

$$\begin{array}{ll} \min & f(\mathbf{x}) \\ \text{s.t.} & \|\mathbf{c}(\mathbf{x})_{-}\| \leq \delta \end{array}$$

的最优解。

YZW (USTC)

### [思考题:证明上述引理1...]

[思考题:证明上述引理3...]





< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

这里只给出引理2的证明:

由题设易得

$$P_{\sigma_k}(\mathbf{x}(\sigma_k)) \geq f(\mathbf{x}(\sigma_k)).$$

因为求是原问题的最优解,自然为可行点,于是  $\sigma_k \| c(\bar{x})_- \|^2 = 0$ .

又因为 $x(\sigma_k) = \arg\min_{x \in \mathbb{R}^n} P_{\sigma_k}(x),$ 则有

$$f(\bar{\mathsf{x}}) = P_{\sigma_k}(\bar{\mathsf{x}}) \ge P_{\sigma_k}(\mathsf{x}(\sigma_k)).$$

YZW (USTC)

A D N A (B) N A B N A B N B

关于罚函数法的收敛性,我们有如下结果

**定理1:** 设罚函数法中的ε满足

 $\varepsilon > \min_{\mathbf{x} \in \mathbb{R}^n} \| \mathbf{c}(\mathbf{x})_- \|,$ 

则算法必有限终止。

[思考题:证明上述定理...]

イロト イポト イヨト イヨト 一日

# 该定理表明,如果原约束问题存在可行点,则对任意给定的 $\varepsilon > 0$ ,算法都将有限终止于问题

 $\begin{array}{ll} \min & f(\mathsf{x}) \\ \text{s.t.} & \|\mathsf{c}(\mathsf{x})_{-}\| \leq \delta \end{array}$ 

的解,且 $\delta \leq \varepsilon$ .

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ ― 臣 – ∽へで

**定理2:** 如果算法不有限终止,则必有  $\min_{x \in \mathbb{R}^n} ||c(x)_-|| \ge \varepsilon$ , 且  $\lim_{k \to \infty} ||c(x(\sigma_k))_-|| = \min_{x \in \mathbb{R}^n} ||c(x)_-||$ .此时, { $x(\sigma_k)$ }的任何聚点 $x^*$ 都是问题 min f(x)

i.t. 
$$\|c(x)_{-}\| = \min_{y \in \mathbb{R}^{n}} \|c(y)_{-}\|$$

 $\mathbf{S}$ 

的解。

э

(日) (同) (三) (三)

# 非线性约束最优化 乘子罚函数

为了叙述简单,仅考虑等式约束问题

$$\begin{array}{l} \min \quad f(\mathbf{x}) \\ \mathrm{s.t.} \quad \mathbf{c}(\mathbf{x}) = \mathbf{0} \end{array} \tag{43}$$

其中 $c(x) = (c_1(x), \cdots, c_{m_e}(x))^T$ .

设x\*是上述问题的最优解且λ\*是相应的Lagrange乘子,由Kuhn-Tucher定 理知, x\*必是Lagrange函数

$$L(\mathbf{x}, \lambda^*) = f(\mathbf{x}) - (\lambda^*)^T \mathbf{c}(\mathbf{x})$$

的稳定点。但一般而言,x\*并不是Lagrange函数的极小点。

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >
## 非线性约束最优化 乘子罚函数

我们考虑乘子罚函数(也称增广Lagrange函数)

$$P(\mathsf{x},\lambda,\sigma) = L(\mathsf{x},\lambda) + rac{\sigma}{2} \|\mathsf{c}(\mathsf{x})\|_2^2.$$

由于增广Lagrange函数的性态,只要取足够大的罚因子 $\sigma$ 而不必趋向无穷大,就可通过极小化 $P(x, \lambda, \sigma)$ 求得原问题的最优解。

イロト イ理ト イヨト イヨト

我们事先并不知道最优乘子向量 $\lambda^*$ ,因此用乘子 $\lambda$ 代替,得到增广Lagrange罚函数:

$$\mathsf{P}(\mathsf{x},\lambda,\sigma) = f(\mathsf{x}) - \lambda^{\mathsf{T}}\mathsf{c}(\mathsf{x}) + rac{\sigma}{2} \|\mathsf{c}(\mathsf{x})\|_2^2.$$

一般的策略是,先给定充分大的 $\sigma$ 和乘子向量的初始估计 $\lambda$ ,然后在迭代 过程中修正乘子 $\lambda$ 力图使之趋向最优乘子 $\lambda^*$ .

如何修正?

3

イロト イ理ト イヨト イヨト

# 非线性约束最优化 乘子罚函数

基于增广Lagrange函数的迭代算法:

- (0) 给定初始点 $x^{(0)}$ 和乘子向量初始估计 $\lambda^{(0)}$ , 给定罚因子 $\sigma_0 > 0$ , 常数 $0 < \alpha < 1, \beta > 1$ 及容许误差 $\varepsilon > 0$ . 令k := 0.
- (1) 以x<sup>(k)</sup>为初点求解无约束问题的极小点,即

$$\mathbf{x}^{(k+1)} = \arg\min_{\mathbf{x}\in\mathbb{R}^n} P(\mathbf{x}, \lambda^{(k)}, \sigma).$$

(2) 若||c(x<sup>(k+1)</sup>)|| < ε, 则停止迭代并取x<sup>(k+1)</sup>作为原问题的近似最优解;
 否则,更新乘子向量

$$\lambda^{(k+1)} = \lambda^{(k)} - \sigma \mathsf{c}(\mathsf{x}^{(k+1)}).$$

(3) 如果 $\frac{\|\mathbf{c}(\mathbf{x}^{(k+1)})\|}{\|\mathbf{c}(\mathbf{x}^{(k)})\|} \ge \alpha$ ,则置 $\sigma := \beta \sigma$ . 令k := k + 1返回第(1)步。

・ 同 ト ・ ヨ ト ・ ヨ ト …



#### $i\partial A(x) = [\nabla c(x)]^T$ , 由于 $c(x^*) = 0$ , 我们易得

$$\nabla_{\mathbf{x}} P(\mathbf{x}^*, \lambda^*, \sigma) = \nabla_{\mathbf{x}} L(\mathbf{x}^*, \lambda^*) = 0,$$

$$\nabla_{xx}^2 P(\mathbf{x}^*, \lambda^*, \sigma) = \nabla_{xx}^2 L(\mathbf{x}^*, \lambda^*) + \sigma A(\mathbf{x}^*) A(\mathbf{x}^*)^T.$$

YZW (USTC)

182 / 467

臣

イロト イヨト イヨト イヨト

## 非线性约束最优化 乘子罚函数

# 设在x\*处满足二阶充分条件,即对 $\forall$ d使得 $A(x^*)^T$ d = 0的非零向量,均有 d<sup>T</sup> $\nabla^2_{xx}L(x^*, \lambda^*)$ d > 0.

#### 因此,在二阶充分条件的假定下,对于充分大的 $\sigma$ ,可证 $\nabla^2_{xx} P(x^*, \lambda^*, \sigma)$ 是正定阵。

A D N A (B) N A B N A B N B



# **定理:** 设x\*和 $\lambda$ \*满足等式约束问题(43)局部最优解的二阶充分条件,则存在 $\sigma_0$ 使得当 $\sigma > \sigma_0$ 时, x\*是函数 $P(x, \lambda^*, \sigma)$ 的严格局部极小点。

[思考题:证明上述定理...]



イロト (周) (三) (三) (三)



**证明:** 由题设知(满足局部最优解的二阶充分条件), x\*必为问题(43)的(K-T)点,所以

$$\nabla_{\mathsf{x}} P(\mathsf{x}^*, \lambda^*, \sigma) = \nabla_{\mathsf{x}} L(\mathsf{x}^*, \lambda^*) = 0.$$

下面证明,在x\*处的Hessian矩阵 $\nabla_{xx}^2 P(x^*, \lambda^*, \sigma)$ 是正定的。

《曰》《御》《曰》《曰》 []

## 非线性约束最优化 乘子罚函数

#### 证明:

$$\begin{aligned} \nabla_{xx}^2 P(\mathbf{x}^*, \lambda^*, \sigma) &= \nabla_{xx}^2 L(\mathbf{x}^*, \lambda^*) + \sigma A(\mathbf{x}^*) A(\mathbf{x}^*)^T \\ &= \bar{Q} + \sigma \bar{A} \bar{A}^T \\ \mathbf{\xi} \mathbf{\psi} \bar{Q} = \nabla_{xx}^2 L(\mathbf{x}^*, \lambda^*), \ \bar{A} = A(\mathbf{x}^*). \end{aligned}$$

æ

<ロト <問 > < 回 > < 回 > .



# **定理:** 若x是等式约束问题(43)的可行解,且对于某个 $\bar{\lambda}$ , x满 足 $P(x, \bar{\lambda}, \sigma)$ 的极小点二阶充分条件,则x是问题(43)的严格局部最优解。

[思考题:证明上述定理...]





< 口 > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

基本思想 在迭代中总是从内点出发,并通过引入障碍函数使之保持在 可行域内部进行搜索。因此,这种方法适用于不等式约束的非线性最优 化问题。

min 
$$f(\mathbf{x})$$
  
s.t.  $g_i(\mathbf{x}) \ge 0, i = 1, \cdots, m.$  (44)

现将可行域内部记作 int S, 其中 $S = \{x \mid g_i(x) \ge 0, i = 1, \dots, m\}$ . 保持 迭代点含于可行域内部的方法是定义如下障碍函数:

$$B(\mathbf{x},\theta) = f(\mathbf{x}) + \theta \psi(\mathbf{x})$$

其中障碍因子 $\theta$ 是很小的正数, $\psi(x)$ 是连续函数,当x趋于可行域边界时, $\psi(x) \rightarrow +\infty$ .

两种重要的障碍形式是:

$$\psi(\mathsf{x}) = \sum_{i=1}^{m} \frac{1}{g_i(\mathsf{x})}$$
 and  $\psi(\mathsf{x}) = -\sum_{i=1}^{m} \log g_i(\mathsf{x})$ 

这样,当x趋向可行域边界时,函数 $B(x, \theta) \rightarrow +\infty$ .否则,由于 $\theta$ 很小,函数 $B(x, \theta)$ 的取值近似于f(x).

イロト イポト イヨト イヨト 二日

因此,我们可通过求解下列问题得到原问题(44)的近似解:

$$\begin{array}{ll} \min & B(\mathbf{x}, \theta) \\ \text{s.t.} & \mathbf{x} \in \operatorname{int} S \end{array} \tag{45}$$

由于 $\psi(x)$ 的存在,在可行域边界形成一道"围墙",因此上述障碍问题(45)的解 $\bar{x}(\theta)$ 必含于可行域的内部。

需要解释的是,障碍问题(45)表面上看起来仍是带约束的最优化问题, 且它的约束条件比原来的约束还要复杂。但是,由于函数 $\psi(x)的障碍阻$ 挡作用是自动实现的,因此从计算观点看,求解(45) 完全可当作无约束问题来处理。

< 口 > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

因此,我们可通过求解下列问题得到原问题(44)的近似解:

$$\begin{array}{ll} \min & B(\mathbf{x}, \theta) \\ \text{s.t.} & \mathbf{x} \in \operatorname{int} S \end{array} \tag{45}$$

由于 $\psi(x)$ 的存在,在可行域边界形成一道"围墙",因此上述障碍问题(45)的解 $\bar{x}(\theta)$ 必含于可行域的内部。

需要解释的是,障碍问题(45)表面上看起来仍是带约束的最优化问题, 且它的约束条件比原来的约束还要复杂。但是,由于函数 $\psi(x)$ 的障碍阻 挡作用是自动实现的,因此从计算观点看,求解(45) 完全可当作无约束 问题来处理。

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ・ クへぐ

- 于是,我们可以给出障碍函数法的计算步骤如下:
- (0) 给定初始点 $x^{(0)} \in int S$ ,初始障碍因子 $\theta_0 > 0, \beta \in (0,1), \varepsilon > 0.$ 令k := 0.
- (1) 以×<sup>(k)</sup>作为初始迭代点求解下列问题:

 $\begin{array}{ll} \min & f(\mathbf{x}) + \theta_k \psi(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in \operatorname{int} \mathcal{S} \end{array}$ 

记求得的极小点为 $x(\theta_k)$ .

(2) 若θ<sub>k</sub>ψ(x(θ<sub>k</sub>)) < ε, 则停止计算并取x(θ<sub>k</sub>)为原问题的近似最优解; 否则, 置x<sup>(k+1)</sup> = x(θ<sub>k</sub>), θ<sub>k+1</sub> = βθ<sub>k</sub>, 令k := k + 1返回第(1)步。

定理: 设 $\theta_k > \theta_{k+1} > 0$ , 记 $x(\theta) = \arg\min_x B(x, \theta)$ , 则有  $B(x(\theta_k), \theta_k) \ge B(x(\theta_{k+1}), \theta_{k+1}),$   $\psi(x(\theta_k)) \le \psi(x(\theta_{k+1})),$   $f(x(\theta_k)) \ge f(x(\theta_{k+1})).$ 

[思考题:证明上述定理...]

3

イロト 不得 トイヨト イヨト



$$\begin{array}{ll} \min & f(\mathbf{x}) \\ \mathrm{s.t.} & \mathsf{c}_E(\mathbf{x}) = 0 \\ & \mathsf{c}_I(\mathbf{x}) \geq 0 \end{array} \tag{46}$$

$$\begin{array}{ll} \min & f(\mathbf{x}) \\ \text{s.t.} & c_E(\mathbf{x}) = 0 \\ & c_I(\mathbf{x}) - \mathbf{s} = 0 \\ & \mathbf{s} \ge 0 \end{array} \tag{47}$$

YZW (USTC)

193 / 467

◆□▶ ◆□▶ ◆ヨ▶ ◆ヨ▶ ─ヨ ─ の々ぐ

The Karush-Kuhn-Tucker (KKT) conditions for the nonlinear program (47) can be written as

$$\nabla f(x) - A_E(x)^T y - A_I(x)^T z = 0$$
  

$$Sz - \mu 1 = 0$$
  

$$c_E(x) = 0$$
  

$$c_I(x) - s = 0$$
(48)

with  $\mu = 0$ , together with s  $\geq 0, z \geq 0$ .

Here  $A_E(x)$  and  $A_I(x)$  are the Jacobian matrices of the functions  $c_E(x)$  and  $c_I(x)$ , respectively, and y and z are their Lagrange multipliers. We define S and Z to be the diagonal matrices whose diagonal entries are given by the vectors s and z, respectively, and let  $1 = (1, \dots, 1)^T$ .

Applying Newton's method to the KKT system (48), in the variables x, s, y, z, we obtain

$$\begin{pmatrix} \nabla_{xx}^{2}L & 0 & -A_{E}(x)^{T} & -A_{I}(x)^{T} \\ 0 & Z & 0 & S \\ -A_{E}(x) & 0 & 0 & 0 \\ -A_{I}(x) & I & 0 & 0 \end{pmatrix} \begin{pmatrix} p_{x} \\ p_{s} \\ p_{y} \\ p_{z} \end{pmatrix}$$

$$= -\begin{pmatrix} \nabla f(x) - A_{E}(x)^{T}y - A_{I}(x)^{T}z \\ Sz - \mu 1 \\ -c_{E}(x) \\ -c_{I}(x) + s \end{pmatrix}$$

$$(49)$$

where L(x, s, y, z) denotes the Lagrange function

$$L(\mathbf{x}, \mathbf{s}, \mathbf{y}, \mathbf{z}) = f(\mathbf{x}) - \mathbf{y}^{\mathsf{T}} \mathbf{c}_{\mathsf{E}}(\mathbf{x}) - \mathbf{z}^{\mathsf{T}} (\mathbf{c}_{\mathsf{I}}(\mathbf{x}) - \mathbf{s}).$$

The system (49) is called the primal-dual system. After the step  $p = (p_x, p_s, p_y, p_z)$  has been determined, we compute the new iterate  $(x^+, s^+, y^+, z^+)$  as

$$\begin{aligned} \mathbf{x}^+ &= \mathbf{x} + \alpha_s^{\max} \mathbf{p}_s, \quad \mathbf{s}^+ &= \mathbf{s} + \alpha_s^{\max} \mathbf{p}_s, \\ \mathbf{y}^+ &= \mathbf{y} + \alpha_z^{\max} \mathbf{p}_y, \quad \mathbf{z}^+ &= \mathbf{z} + \alpha_z^{\max} \mathbf{p}_z, \end{aligned}$$

where

$$\alpha_s^{\max} = \max\{\alpha \in (0, 1] : s + \alpha p_s \ge (1 - \tau)s\}, \alpha_z^{\max} = \max\{\alpha \in (0, 1] : z + \alpha p_z \ge (1 - \tau)z\},$$
(50)

with  $\tau \in (0, 1)$  (A typical value of  $\tau$  is 0.995). The condition (50), called the fraction to the boundary rule, prevents the variables s and z from approaching their lower bounds of 0 too quickly.

- Exercise 4: 证明(38)中定义的ψ(x, λ)是关于Lagrange-Newton法的 下降函数。
- Exercise 5: 证明罚函数法求解带误差界近似问题的算法有限终止性。
- Exercise 6: 给出约束最优化问题的二阶充分最优性条件,并用于 说明增广Lagrange函数的极小点与原问题最优解的等价性。

# Outline I

- Unconstrained Optimization
- 2 Constrained Optimization
  - 二次规划
  - 非线性约束最优化
- 3 Convex Optimization
  - Convex Set and Convex Function
  - Convex Optimization and Algorithms
  - Sparse Optimization
    - Sparse Optimization Models
    - Sparse Optimization Algorithms

#### Optimization Methods for Machine Learning

YZW (USTC)

→ < Ξ →</p>

## Outline II

- Typical Form of Problems
- Stochastic Algorithms
- Other Popular Methods



э

- 2 Constrained Optimization
  - 二次规划
  - 非线性约束最优化
- 3 Convex Optimization
  - Convex Set and Convex Function
  - Convex Optimization and Algorithms
  - 4 Sparse Optimization
    - Sparse Optimization Models
    - Sparse Optimization Algorithms
  - 5 Optimization Methods for Machine Learning
    - Typical Form of Problems
    - Stochastic Algorithms
    - Other Popular Methods yzw (ustc)

**Convex optimization** is a subfield of mathematical optimization that studies the problem of minimizing convex functions over convex sets. Whereas many classes of convex optimization problems admit polynomial-time algorithms, mathematical optimization is in general NP-hard.

We introduce the main definitions and results of convex optimization needed for the analysis of algorithms presented in the section.

#### 定义 (affine set)

A set  $C \subseteq \mathbb{R}^n$  is *affine* if  $\forall x_1, x_2 \in C$  and  $\theta \in \mathbb{R}$ , we have

$$\theta x_1 + (1 - \theta) x_2 \in C$$

i.e., if it contains the line through any two distinct points in it.

It can be generalized to more than two points: If C is an affine set,  $x_1, \ldots, x_k \in C$  and  $\theta_1 + \ldots + \theta_k = 1$ , then  $\theta_1 x_1 + \ldots + \theta_k x_k \in C$ .

We refer to a point of the form  $\theta_1 x_1 + \ldots + \theta_k x_k$  where  $\theta_1 + \ldots + \theta_k = 1$ , as an *affine combination* of the points  $x_1, \ldots, x_k$ .

・ロト ・四ト ・ヨト

If C is an affine set and  $x_0 \in C$ , then the set

$$V = C - x_0 = \{x - x_0 | x \in C\}$$

is a (linear) subspace. We can express C as

$$C = V + x_0 = \{ v + x_0 | v \in V \}.$$

The *dimesion* of an affine set C is the dimesion of the subspace  $V = C - x_0$ .

#### 例 (Solution set of linear equations)

For  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^m$ , the set  $C = \{x | Ax = b\}$  is affine. Let  $V = \{v | Av = 0\}$  be a subspace and  $Ax_0 = b$ , then  $C = V + x_0$ .

A D F A B F A B F A B F

#### 定义 (affine hull)

The set of all affine combinations of points in some set  $C \subseteq \mathbb{R}^n$  is called the *affine hull* of *C*, denoted **aff***C*:

$$\mathsf{aff} C = \{\theta_1 x_1 + \ldots + \theta_k x_k | x_1, \ldots, x_k \in C, \theta_1 + \ldots + \theta_k = 1\}.$$

The affine hull is the smallest affine set that contains C.

Image: Image:

定义 (convex set)

A set C is *convex* if  $\forall x_1, x_2 \in C$  and  $0 \leq \theta \leq 1$ , we have

$$\theta x_1 + (1-\theta)x_2 \in C$$

i.e., if it contains the line segment between any two points in it.

Generalization to more than two points: for any  $k \ge 1, x_1, \ldots, x_k \in C$  and  $\theta_1 + \ldots + \theta_k = 1$  where  $\theta_i \ge 0, i = 1, \ldots, k$ , we have

$$\theta_1 x_1 + \ldots + \theta_k x_k \in C.$$

The form  $\theta_1 x_1 + \ldots + \theta_k x_k$  is called the *convex combination* of the points  $x_1, \ldots, x_k$ , where  $\theta_1, \ldots, \theta_k \ge 0$  and  $\sum_{i=1}^k \theta_i = 1$ .

YZW (USTC)

#### 定义 (convex hull)

The *convex hull* of a set C, denoted **conv**C, is the set of all convex combinations of points in C:

 $\operatorname{conv} C = \{\theta_1 x_1 + \ldots + \theta_k x_k | x_i \in C, \theta_i \ge 0, i = 1, \ldots, k, \theta_1 + \ldots + \theta_k = 1\}.$ 

The convex hull is the smallest convex set that contains C.

イロト 不得 トイヨト イヨト

#### Convex set and convex hull



Figure: (a) A convex set (polyhydron). (b) A non-convex set. (c) The convex hull of (b).

∃ ▶ ∢ ∃ ▶

#### 定义 (cone)

A set C is called a *cone*, if  $\forall x \in C$  and  $\theta \ge 0$  we have  $\theta x$  in C. A set C is a convex cone if it's convex and a cone, i.e.,  $\forall x_1, x_2 \in C$  and  $\theta_1, \theta_2 \ge 0$ , we have

 $\theta_1 x_1 + \theta_2 x_2 \in C.$ 

A point of the form  $\theta_1 x_1 + \ldots + \theta_k x_k$  with  $\theta_1, \ldots, \theta_k \ge 0$  is called a *conic combination* of  $x_1, \ldots, x_k$ .

#### 定义 (conic hull)

The conic hull of a set C is the set of all conic combinations of points in C, i.e.,

$$\{\theta_1 x_1 + \ldots + \theta_k x_k | x_i \in C, \theta_i \ge 0, i = 1, \ldots, k\}.$$

#### Conic hull



Figure: Left. The shaded set is the conic hull of a set of fifteen points (not including the origin). Right. The shaded set is the conic hull of the non-convex kidney-shaped set that is surrounded by a curve.

• Hyperplane: A hyperplane is a set of the form

$$\{x \mid a^\top x = b\}.$$

It's also affine.

• Halfspace: A (closed) halfspace is a set of the form

 $\{x \mid a^{\top}x \leq b\}.$ 

A hyperplane divides  $\mathbb{R}^n$  into two halfspaces.

• *Polyhedra:* A polyhedron is defined as the solution set of a finite number of linear equalities and inequalities:

$$\mathcal{P} = \{x \mid a_j^\top x \leq b_j, j = 1, \ldots, m, c_k^\top x = d_k, k = 1, \ldots, p\}.$$

• Ball: A (Euclidean) ball in  $\mathbb{R}^n$  has the form

$$B(x_c,r) = \{x \mid ||x - x_c||_2 \leqslant r\}$$

where r > 0 and  $||u||_2 = (u^{\top}u)^{1/2}$  denotes the Euclidean norm.

• Norm balls and norm cones:

Suppose  $\|\cdot\|$  is any norm on  $\mathbb{R}^n$ , a norm ball of radius *r* and center  $x_c$  is given by

$$\{x\|\|x-x_c\|\leqslant r\}.$$

The norm cone associated with the norm  $\|\cdot\|$  is the set

$$C = \{(x,t) | \|x\| \leq t\} \subseteq \mathbb{R}^{n+1}$$

It's a convex cone.

 The positive semidefinite cone: The set of symmetric n × n matrices S<sup>n</sup>:

$$\mathsf{S}^n = \{ X \in \mathbb{R}^{n \times n} | X = X^\top \},\$$

the set of symmetric positive semidefinite matrices  $S_{+}^{n}$ :

$$\mathsf{S}^n_+ = \{ X \in \mathsf{S}^n | X \succeq \mathsf{0} \},\$$

and the set of symmetric positive definite matrices  $S_{++}^n$ :

$$\mathsf{S}_{++}^n = \{X \in \mathsf{S}^n | X \succ \mathsf{0}\}$$

are all convex.
## Proper cones and generalized inequalities

A cone  $K \subseteq \mathbb{R}^n$  is called a *proper cone* if it satisfies the following:

- K is convex.
- K is closed.
- K is solid, which means it has nonempty interior.
- K is pointed, which means that it contains no line, i.e.,

$$x \in K$$
 and  $-x \in K \Rightarrow x = 0$ .

A proper cone K can be used to define a generalized inequality:

$$x \preceq_{\mathcal{K}} y \iff y - x \in \mathcal{K},$$

which is a partial ordering on  $\mathbb{R}^n$ . Similarly, we define an associated strict partial ordering by

$$x \prec_K y \iff y - x \in \operatorname{int} K$$

YZW (USTC)

## Properties of generalized inequalities

- If  $x \preceq_{\mathcal{K}} y$  and  $u \preceq_{\mathcal{K}} v$ , then  $x + u \preceq_{\mathcal{K}} y + v$ .
- If  $x \preceq_{\mathcal{K}} y$  and  $y \preceq_{\mathcal{K}} z$  then  $x \preceq_{\mathcal{K}} z$ .
- If  $x \preceq_{\mathcal{K}} y$  and  $\alpha \ge 0$  then  $\alpha x \preceq_{\mathcal{K}} \alpha y$ .
- $x \preceq_{K} x$ .
- If  $x \preceq_{K} y$  and  $y \preceq_{K} x$  then x = y.
- If  $x_i \leq_K y_i$  for  $i = 1, 2, ..., x_i \rightarrow x$  and  $y_i \rightarrow y$  as  $i \rightarrow \infty$ , then  $x \leq_K y$ .

### Minimum and minimal elements

x ∈ S is the minimum element of S (with respect to the generalized inequality ≤<sub>K</sub>) if for every y ∈ S we have x ≤<sub>K</sub> y, i.e.,

$$S \subseteq x + K$$
,

where  $x + K = \{x + z | z \in K\}$ .

x ∈ S is a minimal element of S (with respect to the generalized inequality ⊥<sub>K</sub>) if y ∈ S, y ⊥<sub>K</sub> x only if y = x, i.e.,

$$(x-K)\cap S=\{x\},$$

where  $x - K = \{x - z | z \in K\}.$ 

Maximum element and maximal element are defined in a similar way.

YZW (USTC)

### Minimum and minimal elements



Figure: Let  $K = \{(u, v) | u, v \ge 0\}$ . (a)  $x_1$  is the minimum element of  $S_1$ . (b)  $x_2$  is a minimal element of  $S_2$ .

If x is the minimum element of S, then x must be a minimal element of S (with respect to the generalized inequality  $\leq_{\mathcal{K}}$ ).

Brief proof: Suppose  $S \subseteq x + K$ , and  $y \in (x - K) \cap S$ , i.e.,  $\exists z \in K$  such that y = x - z. By  $y \in S \subseteq x + K$ , there exists  $w \in K$  such that y = x + w. Then we have w = -z, which leads to  $-w = z \in K$  and  $w \in K$ . Since K is a proper cone, w = 0 and y = x.

But the reverse proposition doesn't hold.

Simple example: Let  $K = \{(u, v) | u, v \ge 0\}$  and  $L = \{(x, y) | x = -y\}$ . Then every point of L is a minimal element, but none of them is the minimum element of L.

YZW (USTC)

### 定义 (convex function)

A function  $f : \mathbb{R}^n \to \mathbb{R}$  is *convex* if **dom** *f* is a convex set and if  $\forall x, y \in \mathbf{dom} f$  and  $\theta$  with  $0 \leq \theta \leq 1$ , we have

$$f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y).$$
(51)

A function is *strictly convex* if strict inequality holds in (51) whenever  $x \neq y$  and  $0 < \theta < 1$ .

We say f is concave if -f is convex, and strictly concave if -f is strictly convex.

Geometrically, Eq.(51) means that the line segment between (x, f(x)) and (y, f(y)) lies above the graph of f (as shown in Fig.4).



Figure: Graph of a convex function.

Suppose f is differentiable, i.e., its gradient  $\nabla f$  exists at each point in **dom** f.

Function f is convex if and only if **dom** f is convex and for  $\forall x, y \in \mathbf{dom} f$ , the following holds:

$$f(y) \ge f(x) + \nabla f(x)^{\top} (y - x).$$

**Remark.** As a simple result, if  $\nabla f(x^*) = 0$ , then for all  $y \in \text{dom} f$ ,  $f(y) \ge f(x^*)$ , *i.e.*,  $x^*$  is a global minimizer of the function f.

### First-order conditions



Figure: The tangent to a convex function.

YZW (USTC)

222 / 467

Function f is strictly convex if and only if **dom** f is convex and for  $\forall x, y \in \mathbf{dom} f, x \neq y$ , we have

$$f(y) > f(x) + \nabla f(x)^{\top} (y - x).$$

Correspondingly, f is concave if and only if **dom** f is convex and for  $\forall x, y \in \mathbf{dom} f$ , we have

$$f(y) \leqslant f(x) + \nabla f(x)^{\top} (y - x).$$

Assume that f is twice differentiable.

Function f is convex if and only if **dom** f is convex and for  $\forall x \in \mathbf{dom} f$ ,

$$\nabla^2 f(x) \succeq 0.$$

Similarly, f is concave if and only if **dom** f is convex and  $\nabla^2 f(x) \preceq 0$  for  $\forall x \in \mathbf{dom} f$ .

Strict convexity can be partially characterized by second-order conditions.

If  $\nabla^2 f(x) \succ 0$  for  $\forall x \in \mathbf{dom} f$ , then f is strictly convex.

However, the converse is not true. For example,  $f : \mathbb{R} \to \mathbb{R}$  given by  $f(x) = x^4$  is strictly convex but has zero second derivative at x = 0.

イロト 不得 トイヨト イヨト

- Exponential:  $e^{ax}$  is convex on  $\mathbb{R}$ , for any  $a \in \mathbb{R}$ .
- Powers:

 $x^a$  is convex on  $\mathbb{R}_{++}$  when  $a \ge 1$  or  $a \le 0$ , and concave for  $0 \le a \le 1$ .

- Powers of absolute value: |x|<sup>p</sup>, for p ≥ 1, is convex on ℝ.
- Logarithm: log x is concave on ℝ<sub>++</sub>.
- Negative entropy:
   x log x is convex on ℝ<sub>+</sub>, where 0 log 0 defined to be 0.

# Examples

• Norms:

Every norm on  $\mathbb{R}^n$  is convex.

- Max function:
   f(x) = max{x<sub>1</sub>,...,x<sub>n</sub>} is convex on ℝ<sup>n</sup>.
- Log-sum-exp:

Then function  $f(x) = \log(e^{x_1} + \ldots + e^{x_n})$  is convex on  $\mathbb{R}^n$ . This function can be interpreted as a differentiable approximation of the max function, since for all x,

$$\max\{x_1,\ldots,x_n\} \leqslant f(x) \leqslant \max\{x_1,\ldots,x_n\} + \log n.$$

- Geometric mean:  $f(x) = (\prod_{i=1}^{n} x_i)^{1/n}$  is concave on **dom** $f = \mathbb{R}^n_{++}$ .
- Log-determinant:  $f(X) = \log \det X$  is concave on  $\operatorname{dom} f = S^n_{++1}$ ,  $\operatorname{det} f = \operatorname{det} X$

YZW (USTC)

The inequality (51), *i.e.*,  $f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y)$ , is sometimes called *Jensen's inequality*.

It is easily extended to convex combinations of more than two points:

If f is convex,  $x_1, \ldots, x_k \in \mathbf{dom} f$ , and  $\theta_1, \ldots, \theta_k \ge 0$  with  $\theta_1 + \ldots + \theta_k = 1$ , then

$$f(\theta_1 x_1 + \ldots + \theta_k x_k) \leq \theta_1 f(x_1) + \ldots + \theta_k f(x_k).$$

YZW (USTC)

・ロト ・四ト ・ヨト

#### • Nonnegative weighted sums:

If  $f_1,\ldots,f_m$  are convex and  $w_1,\ldots,w_m \geqslant 0$ , then

$$f = w_1 f_1 + \ldots + w_m f_m$$

is convex.

• These properties extend to infinite sums and integrals:

If f(x, y) is convex in x for each  $y \in A$ , and  $w(y) \ge 0$  for each  $y \in A$ , then the function

$$g(x) = \int_{\mathcal{A}} w(y) f(x, y) dy$$

is convex in x (provided the integral exists).

### • Composition with an affine mapping: Suppose $f : \mathbb{R}^n \to \mathbb{R}$ , $A \in \mathbb{R}^{n \times m}$ , and $b \in \mathbb{R}^n$ . Define $g : \mathbb{R}^m \to \mathbb{R}$ by

$$g(x)=f(Ax+b),$$

with  $\mathbf{dom}g = \{x | Ax + b \in \mathbf{dom}f\}$ . Then if f is convex, so is g; if f is concave, so is g.

• Pointwise maximum:

If  $f_1$  and  $f_2$  are convex functions, then

 $f(x) = \max\{f_1(x), f_2(x)\},\$ 

with  $dom f = dom f_1 \cap dom f_2$ , is also convex.

• Extension to the pointwise supremum:

If for each  $y \in A$ , f(x, y) is convex in x, then

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

is convex in x, where

$$\mathbf{dom}g = \{x | (x, y) \in \mathbf{dom}f \text{ for all } y \in \mathcal{A}, \sup_{y \in \mathcal{A}} f(x, y) < \infty\}.$$

• Quasi-convex function: A function  $f : \mathbb{R}^n \to \mathbb{R}$  such that its domain and all its sublevel sets

$$S_{\alpha} = \{ x \in \operatorname{dom} f | f(x) \leq \alpha \}, \quad \alpha \in \mathbb{R}$$

are convex.

• Log-concave function: A function  $f : \mathbb{R}^n \to \mathbb{R}$  such that  $f(x) > 0, \forall x \in \text{dom} f$  and  $\log f$  is concave.

## Basic terminology

min 
$$f_0(x)$$
  
s.t.  $f_i(x) \le 0, \quad i = 1, ..., m$   
 $h_j(x) = 0, \quad j = 1, ..., p$  (52)

 $\begin{array}{ll} x \in \mathbb{R}^n & \text{the optimization variable} \\ f_0 : \mathbb{R}^n \to \mathbb{R} & \text{the objective function or cost function} \\ f_i(x) \leqslant 0 & \text{the inequality constraints} \\ f_i : \mathbb{R}^n \to \mathbb{R} & \text{the inequality constraint functions} \\ h_j(x) = 0 & \text{the equality constraints} \\ h_i : \mathbb{R}^n \to \mathbb{R} & \text{the equality constraint functions} \end{array}$ 

If there are no constraints (*i.e.*, m = p = 0) we say the problem is unconstrained.

< 日 > < 同 > < 回 > < 回 > .

• The domain of the optimization problem (52) is given as

$$\mathcal{D}=\bigcap_{i=0}^m \mathbf{dom} f_i\cap \bigcap_{j=1}^p \mathbf{dom} h_j.$$

- A point  $x \in D$  is *feasible* if  $f_i(x) \leq 0, i = 1, ..., m$ , and  $h_j(x) = 0, j = 1, ..., p$ .
- The problem (52) is said to be feasible if there exists at least one feasible point, and *infeasible* otherwise.

The optimal value  $v^*$  of the problem (52) is defined as

 $v^* = \inf\{f_0(x) | f_i(x) \leq 0, i = 1, \dots, m, h_j(x) = 0, j = 1, \dots, p\}$ 

If the problem is infeasible, we have  $v^* = \infty$ .

- We say  $x^*$  is an *optimal point*, or solves the problem (52), if  $x^*$  is feasible and  $f_0(x^*) = v^*$ .
- We say a feasible points  $\bar{x}$  is *locally optimal* if there is a constant  $\delta > 0$  such that

$$f_0(\bar{x}) = \inf\{f_0(z) | f_i(z) \leq 0, i = 1, \dots, m, \\ h_j(z) = 0, j = 1, \dots, p, ||z - \bar{x}||_2 \leq \delta\}.$$

イロト 不得 トイヨト イヨト

#### A convex optimization problem is one of the form

min 
$$f_0(x)$$
  
s.t.  $f_i(x) \leq 0, \quad i = 1, \dots, m$   
 $a_j^\top x = b_j, \quad j = 1, \dots, p$  (53)

where  $f_0, f_1, \ldots, f_m$  are convex functions.

Any locally optimal point of a convex optimization problem is also globally optimal.

イロト 不得 トイヨト イヨト 二日

Suppose that the objective  $f_0$  in a convex optimization problem is differentiable. Let X denote the feasible set, *i.e.*,

$$X = \{x | f_i(x) \leq 0, i = 1, \dots, m, h_j(x) = 0, j = 1, \dots, p\}.$$

Then x is optimal if and only if  $x \in X$  and

$$\nabla f_0(x)^\top (y-x) \ge 0, \ \forall y \in X.$$
(54)



#### For an unconstrained problem, the condition (54) reduces to

$$\nabla f_0(x) = 0 \tag{55}$$

for x to be optimal.



For a convex problem with equality constraints only, *i.e.*,

min  $f_0(x)$ s.t. Ax = b

We assume that the feasible set is nonempty. The optimality condition can be expressed as:

$$\nabla f_0(x)^\top u \ge 0$$
 for all  $u \in \mathcal{N}(A)$ .

In other words,

 $\nabla f_0(x) \perp \mathcal{N}(A).$ 

YZW (USTC)

### A general linear program (LP) has the form

min 
$$q^{\top}x + r$$
  
s.t.  $Gx \le h$  (56)  
 $Ax = b$ 

イロト 不得 トイヨト イヨト

where  $G \in \mathbb{R}^{m \times n}$  and  $A \in \mathbb{R}^{p \times n}$ . It is common to omit the constant r in the objective function.

A convex optimization problem is called *quadratic program* (QP) if it has the form

$$\min \frac{1}{2} x^{\top} P x + q^{\top} x + r$$
  
s.t.  $G x \le h$   
 $A x = b$  (57)

where  $P \in S^n_+$ ,  $G \in \mathbb{R}^{m \times n}$ , and  $A \in \mathbb{R}^{p \times n}$ .

QPs include LPs as a special case by taking P = 0.

If the objective in (53) as well as the inequality constraint functions are (convex) quadratic, as in

min 
$$\frac{1}{2}x^{\top}P_0x + q_0^{\top}x + r_0$$
  
s.t.  $\frac{1}{2}x^{\top}P_ix + q_i^{\top}x + r_i \leq 0, \quad i = 1, \dots, m$   
 $Ax = b$  (58)

where  $P_i \in S_+^n$ , i = 0, 1, ..., m, and the problem is called a *quadratically* constrained quadratic program (QCQP).

QCQPs include QPs as a special case by taking  $P_i = 0$  for i = 1, ..., m.

イロト イヨト イヨト --

A problem that is closely related to quadratic programming is the *second-order cone program* (SOCP):

min 
$$f^{\top}x$$
  
s.t.  $\|L_i x + g_i\|_2 \leq c_i^{\top}x + d_i, \quad i = 1, \dots, m$  (59)  
 $Ax = b$ 

where  $x \in \mathbb{R}^n$  is the optimization variable,  $L_i \in \mathbb{R}^{n_i \times n}$ , and  $A \in \mathbb{R}^{p \times n}$ .

When  $c_i = 0, i = 1, ..., m$ , the SOCP is equivalent to a QCQP. However, second-order cone programs are more general than QCQPs (and of course, LPs).

## Transform a QCQP into an SOCP

For a QCQP problem (58), let y be an auxiliary variable with constraint:

$$\frac{1}{2}x^{\top}P_0x + q_0^{\top}x + r_0 \leqslant y,$$

then (58) becomes

$$\begin{array}{l} \min_{x,y} \ y \\ \text{s.t.} \ \frac{1}{2} x^\top P_i x + q_i^\top x + r_i \leqslant 0, \quad i = 1, \dots, m \\ \frac{1}{2} x^\top P_0 x + q_0^\top x - y + r_0 \leqslant 0 \\ Ax = b \end{array}$$

whose objective is linear. To transform it into an SOCP, we need only translate quadratic constraints into second-order conic ones.

YZW (USTC)

244 / 467

< ロト < 同ト < 三ト < 三ト

## Transform a QCQP into an SOCP

For a quadratic constraint

$$\frac{1}{2}x^{\top}Px + q^{\top}x + r \leqslant 0$$

with  $P \in S_+^n$ , let  $L \in S_+^n$  be the square root of P, *i.e.*, LL = P. Let

$$\widetilde{L} = \begin{bmatrix} L \\ q^{\top} \end{bmatrix}, \quad \widetilde{g} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ r + \frac{1}{2} \end{bmatrix} \in \mathbb{R}^{n+1},$$

then the constraint is equivalent to

$$\|\tilde{L}x+\tilde{g}\|_2 \leq -(q^{\top}x+r-\frac{1}{2}).$$

Consider an optimization problem in the standard form (52):

min 
$$f_0(x)$$
  
s.t.  $f_i(x) \le 0, \quad i = 1, ..., m$   
 $h_j(x) = 0, \quad j = 1, ..., p.$  (60)

We assume its domain  $\mathcal{D} = \bigcap_{i=0}^{m} \operatorname{dom} f_i \cap \bigcap_{j=1}^{p} \operatorname{dom} h_j$  is nonempty, and denote the optimal value of (60) by  $v^*$ , but do not assume the problem (60) is convex.

The basic idea of Lagrangian duality is to take the constraints in (60) into account by augmenting the objective function with a weighted sum of the constraint functions.

A D F A B F A B F A B F

# The Lagrangian

We define the Lagrangian  $L : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^p \to \mathbb{R}$  associated with the problem (60) as

$$L(x,\lambda,\nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{j=1}^p \nu_j h_j(x)$$

with **dom** $L = \mathcal{D} \times \mathbb{R}^m \times \mathbb{R}^p$ .

- Refer to λ<sub>i</sub> as the Lagrange multiplier associated with the *i*th inequality constraint f<sub>i</sub>(x) ≤ 0.
- Refer to ν<sub>j</sub> as the Lagrange multiplier associated with the *j*th equality constraint h<sub>j</sub>(x) = 0.
- The vectors λ and ν are called Lagrange multiplier vectors or the dual variables associated with the problem (60).

We define the Lagrange dual function (or just dual function)  $g: \mathbb{R}^m \times \mathbb{R}^p \to \mathbb{R}$  as

$$g(\lambda,\nu) = \inf_{x\in\mathcal{D}} L(x,\lambda,\nu) = \inf_{x\in\mathcal{D}} \left( f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{j=1}^p \nu_j h_j(x) \right).$$

Since the dual function is the pointwise infimum of a family of affine functions of  $(\lambda, \nu)$ , it is concave, even when the problem (60) is not convex.

(I)
# Lower bounds on optimal value

Let  $v^*$  be the optimal value of the primal problem (60). For any  $\lambda \ge 0$  and any  $\nu$  we have

$$g(\lambda,\nu) \leqslant v^*.$$
 (61)

#### Proof.

Suppose  $\tilde{x}$  is a feasible point for (60), then we have

$$\sum_{i=1}^m \lambda_i f_i(\tilde{x}) + \sum_{j=1}^p \nu_j h_j(\tilde{x}) \leqslant 0.$$

Hence

$$g(\lambda,\nu) = \inf_{x\in\mathcal{D}} L(x,\lambda,\nu) \leq L(\tilde{x},\lambda,\nu) \leq f_0(\tilde{x}).$$

Since  $g(\lambda, \nu) \leq f_0(\tilde{x})$  holds for every feasible point  $\tilde{x}$ , the inequality (61) follows.

The dual function gives a nontrivial lower bound on  $v^*$  only when  $\lambda \ge 0$ and  $(\lambda, \nu) \in \operatorname{dom} g$ , *i.e.*,  $g(\lambda, \nu) > -\infty$ .

We refer to a pair  $(\lambda, \nu)$  with  $\lambda \ge 0$  and  $(\lambda, \nu) \in \mathbf{dom}g$  as *dual feasible*.

イロト イヨト イヨト 一旦

Let  $I_- : \mathbb{R} \to \mathbb{R} \cup \{\infty\}$  and  $I_0 : \mathbb{R} \to \mathbb{R} \cup \{\infty\}$  to be the indicator function for the nonpositive reals and  $\{0\}$  respectively:

$$I_{-}(u) = \left\{ egin{array}{ccc} 0 & u \leqslant 0 \ \infty & u > 0 \end{array} 
ight., \quad I_{0}(u) = \left\{ egin{array}{ccc} 0 & u = 0 \ \infty & u 
eq 0 \end{array} 
ight.$$

Then the primal problem (60) can be reformulated as an unconstrained problem:

min 
$$f_0(x) + \sum_{i=1}^m I_-(f_i(x)) + \sum_{j=1}^p I_0(h_j(x)).$$
 (62)

YZW (USTC)

# Linear approximation interpretation

We replace the function  $I_{-}(u)$  with the linear function  $\lambda_{i}u$ , where  $\lambda_{i} \ge 0$ , and the function  $I_{0}(u)$  with  $\nu_{j}u$ . The objective becomes the Lagrangian function, i.e.,

min 
$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{j=1}^p \nu_j h_j(x).$$

In this formulation, we use a linear or "soft" displeasure function in place of  ${\it I}_-$  and  ${\it I}_0.$ 

Linear function is an *underestimator* of the indicator function. Since  $\lambda_i u \leq I_-(u)$  and  $\nu_j u \leq I_0(u)$  for all u, we see immediately that the dual function yields a lower bound on the optimal value of the primal problem.

YZW (USTC)

A D F A B F A B F A B F

To attain the *best* lower bound that can be obtained from the Lagrange dual function leads to the optimization problem

$$\begin{array}{ll} \max & g(\lambda,\nu) \\ \text{s.t.} & \lambda \geq 0 \end{array} \tag{63}$$

This problem is called the *Lagrange dual problem* associated with the problem (60). Correspondingly, the problem (60) is called the *primal problem*.

The term *dual feasible*, to describe a pair  $(\lambda, \nu)$  with  $\lambda \ge 0$  and  $g(\lambda, \nu) > -\infty$ , now makes sense.

We refer to  $(\lambda^*, \nu^*)$  as dual optimal or optimal Lagrange multipliers if they are optimal for the Lagrange dual problem (63).

The Lagrange dual problem (63) is a convex optimization problem no matter the primal problem is convex or not, since the objective to be maximized is concave and the constraint is convex.

#### For the optimal value of the Lagrange dual problem $g^*$ , we have

$$g^* \leqslant v^*. \tag{64}$$

(日)

#### This property is called weak duality.

 $v^* - g^*$  is the optimal duality gap of the primal problem.

YZW (USTC)

If the equality

$$g^* = v^* \tag{65}$$

holds, then we say that stong duality holds.

- Strong duality does not, in general, hold.
- For a convex primal problem, there are many additional conditions on the primal problem, under which strong duality holds.

One simple condition is Slater's condition:

There exists an  $x \in \operatorname{relint} \mathcal{D}$  such that

$$f_i(x) < 0, \quad i = 1, \dots, m, \qquad Ax = b,$$
 (66)

where  $\operatorname{relint} \mathcal{D} = \{x \in \mathcal{D} | B(x, r) \cap \operatorname{aff} \mathcal{D} \subseteq \mathcal{D} \text{ for some } r > 0\}$ . Such a point is called *relative feasible interior point*.

Slater's theorem states that strong duality holds if Slater's condition holds (and the problem is convex).

# Optimality conditions

Dual feasible points allow us to bound how suboptimal a given feasible point is, without knowing the exact value of  $v^*$ .

If x is primal feasible and  $(\lambda, \nu)$  is dual feasible, then

$$f_0(x) - v^* \leqslant f_0(x) - g(\lambda, 
u)$$

and

$$\mathbf{v}^* \in [g(\lambda, \nu), f_0(x)], \quad g^* \in [g(\lambda, \nu), f_0(x)].$$

It leads to

$$g(\lambda,\nu) = f_0(x) \Longrightarrow v^* = f_0(x) = g(\lambda,\nu) = g^*.$$

We refer to  $f_0(x) - g(\lambda, \nu)$  as the *duality gap* associated with the primal feasible point x and dual feasible point  $(\lambda, \nu)$ .

Suppose that the primal and dual optimal values are attained and equal, let  $x^*$  be a primal optimal and  $(\lambda^*, \nu^*)$  be a dual optimal points, then

$$\begin{aligned} f_0(x^*) &= g(\lambda^*, \nu^*) \\ &= \inf_x \left( f_0(x) + \sum_{i=1}^m \lambda_i^* f_i(x) + \sum_{j=1}^p \nu_j^* h_j(x) \right) \\ &\leq f_0(x^*) + \sum_{i=1}^m \lambda_i^* f_i(x^*) + \sum_{j=1}^p \nu_j^* h_j(x^*) \\ &\leq f_0(x^*) \end{aligned}$$

YZW (USTC)

By 
$$\lambda_i^* \ge 0, f_i(x^*) \le 0, i = 1, \dots, m$$
, we have

$$\lambda_i^* f_i(x^*) = 0, \quad i = 1, \dots, m.$$
 (67)

This condition is known as *complementary slackness*.

We can express it as

$$\lambda_i^* > 0 \Longrightarrow f_i(x^*) = 0, \ f_i(x^*) < 0 \Longrightarrow \lambda_i^* = 0.$$

YZW (USTC)

261 / 467

э

We now assume that the functions  $f_0, \ldots, f_m, h_1, \ldots, h_p$  are differentiable. As above, let  $x^*$  and  $(\lambda^*, \nu^*)$  be any primal and dual optimal points with zero duality gap.

Since  $x^*$  minimizes  $L(x, \lambda^*, \nu^*)$  over x, it follows

$$\nabla f_0(x^*) + \sum_{i=1}^m \lambda_i^* \nabla f_i(x^*) + \sum_{j=1}^p \nu_j^* \nabla h_j(x^*) = 0.$$

Together with constraints and complementary slackness, we have

$$\begin{cases} f_{i}(x^{*}) \leq 0, & i = 1, \dots, m \\ h_{j}(x^{*}) = 0, & j = 1, \dots, p \\ \lambda_{i}^{*} \geq 0, & i = 1, \dots, m \\ \lambda_{i}^{*} f_{i}(x^{*}) = 0, & i = 1, \dots, m \\ \nabla f_{0}(x^{*}) + \sum_{i=1}^{m} \lambda_{i}^{*} \nabla f_{i}(x^{*}) + \sum_{j=1}^{p} \nu_{j}^{*} \nabla h_{j}(x^{*}) = 0 \end{cases}$$

$$(68)$$

which are called the Karush-Kuhn-Tucker (KKT) conditions.

A D F A B F A B F A B F

For *any* optimization problem with differentiable objective and constraint functions for which strong duality obtains, any pair of primal and dual optimal points must satisfy the KKT conditions.

When the primal problem is convex, the KKT conditions are also sufficient for the points to be primal and dual optimal.

There is *no* analytical formula for the solution of convex optimization problems, not to mention general nonlinear optimization problems.

Thus we describe numerical methods for solving convex optimization problems in the section.

## To solve an unconstrained optimization problem

# $\min f(x)$

where f(x) is differentiable and convex, we usually employ descent methods.



イロト 不得 トイヨト イヨト

Given a starting point  $x^{(0)}$ , a descent method produces a sequence  $x^{(k)}, k = 1, \ldots$ , where

$$x^{(k+1)} = x^{(k)} + \alpha_k \delta_x^{(k)}, \quad f(x^{(k+1)}) < f(x^{(k)}).$$
(69)

We usually drop the superscripts and use the notation  $x := x + \alpha \delta_x$  to focus on one iteration of an algorithm.  $\alpha > 0$  is called step size and  $\delta_x$  called search direction. Different methods differ from choices of  $\alpha$  or/and  $\delta_x$ .

Given a descent direction  $\delta_{\rm X},$  we usually use line search to determine step size  $\alpha.$ 

Different search directions:

• Negative gradient:

$$\delta_x = -\nabla f(x).$$

• Normalized steepest descent direction (with respect to the norm  $\|\cdot\|$ ):

$$\delta_{x_{\mathsf{nsd}}} = \arg\min\{\nabla f(x)^\top v \,|\, \|v\| = 1\}.$$

• Newton step:

$$\delta_{x_{\rm nt}} = -\nabla^2 f(x)^{-1} \nabla f(x).$$

YZW (USTC)

### A convex optimization problem with equality constraints has the form

$$\begin{array}{l} \min \quad f(x) \\ \text{s.t.} \quad Ax = b, \end{array} \tag{70}$$

where  $f : \mathbb{R}^n \to \mathbb{R}$  is convex and twice continuously differentiable, and  $A \in \mathbb{R}^{p \times n}$  with  $\operatorname{rank} A = p < n$ . We assume that an optimal solution  $x^*$  exists and  $v^* = f(x^*)$ .

Recall the KKT conditions for (70): a point  $x^* \in \mathbf{dom} f$  is optimal if and only if there is a multiplier  $\nu^* \in \mathbb{R}^p$  such that

$$Ax^* = b, \quad \nabla f(x^*) + A^{\top} \nu^* = 0.$$
 (71)

The first set of equations,  $Ax^* = b$ , are called the *primal feasibility* equations.

The second set of equations,  $\nabla f(x^*) + A^{\top}\nu^* = 0$ , are called the *dual feasibility equations*.

(I)

*Newton's method with equality constraints* is almost the same as Newton's method without constraints, except for two differences:

- The initial point must be feasible (*i.e.*,  $x \in \text{dom} f$  and Ax = b).
- The definition of Newton step  $\delta_{\rm x_{nt}}$  is modified to take the equality constraints into account.

To derive the Newton step  $\delta_{x_{nt}}$  for problem (70) at the feasible point x, we replace the objective with its second-order Taylor approximation near x

min 
$$\hat{f}(x+s) = f(x) + \nabla f(x)^{\top}s + \frac{1}{2}s^{\top}\nabla^2 f(x)s$$
  
s.t.  $A(x+s) = b$  (72)

with variable s. Suppose  $\delta_{x_{nt}}$  is optimal for (72). By KKT conditions, there exists associated optimal dual variable  $w \in \mathbb{R}^p$  such that

$$\begin{bmatrix} \nabla^2 f(x) & A^{\top} \\ A & 0 \end{bmatrix} \begin{bmatrix} \delta_{x_{nt}} \\ w \end{bmatrix} = \begin{bmatrix} -\nabla f(x) \\ 0 \end{bmatrix}.$$
 (73)

We can also derive the Newton Step  $\delta_{x_{nt}}$  by simply replacing  $x^*$  and  $\nu^*$  in the KKT conditions for problem (70):

$$Ax^* = b$$
,  $\nabla f(x^*) + A^\top \nu^* = 0$ 

with  $x + \delta_{x_{nt}}$  and w, respectively, and replace the gradient term in the second equation by its linearized approximation near x, to obtain the equations

$$\begin{aligned} & \mathcal{A}(x+\delta_{x_{\mathsf{n}\mathsf{t}}})=b,\\ & \nabla f(x+\delta_{x_{\mathsf{n}\mathsf{t}}})+\mathcal{A}^\top w\approx \nabla f(x)+\nabla^2 f(x)\delta_{x_{\mathsf{n}\mathsf{t}}}+\mathcal{A}^\top w=0. \end{aligned}$$

YZW (USTC)

Using Ax = b, these become

$$A\delta_{x_{nt}} = 0, \quad \nabla^2 f(x)\delta_{x_{nt}} + A^\top w = -\nabla f(x),$$

which are precisely the equations (73).

YZW (USTC)

э

イロト 不得 トイヨト イヨト

#### The Newton decrement is defined as

$$\kappa(\mathbf{x}) = (\delta_{\mathbf{x}_{\mathsf{nt}}}^\top \nabla^2 f(\mathbf{x}) \delta_{\mathbf{x}_{\mathsf{nt}}})^{1/2}.$$

Since

$$\frac{\mathrm{d}}{\mathrm{d}\alpha}f(x+\alpha\delta_{x_{\mathsf{nt}}})\Big|_{\alpha=0}=\nabla f(x)^{\top}\delta_{x_{\mathsf{nt}}}=-\kappa(x)^{2},$$

the algorithm should terminate when  $\kappa(x)$  is small.

э

Image: Image:

Algorithm. Newton's method for equality constrained minimization.

**given** starting point  $x \in \text{dom} f$  with Ax = b, tolerance  $\epsilon > 0$ . repeat

- **(**) Compute the Newton step  $\delta_{x_{nt}}$  and the decrement  $\kappa(x)$ .
- 2 Stopping criterion. quit if  $\kappa^2/2 \leq \epsilon$ .
- **3** Line search Choose step size  $\alpha$  by backtracking line search.

• Update. 
$$x := x + \alpha \delta_{x_{nt}}$$
.

Newton's method described above is a feasible descent method. Now we describe a generalization of Newton's method that works with initial points and iterates that are *not* feasible.

Let x denote the current point, which we do not assume to be feasible, but we do assume satisfies  $x \in \text{dom} f$ .

Our goal is to find a step  $\delta_x$  so that  $x + \delta_x$  satisfies the optimality conditions (71), *i.e.*,  $x + \delta_x \approx x^*$ .

Similarly, we substitute  $x+\delta_x$  for  $x^*$  and  $\mu$  for  $\nu^*$  in

$$Ax^* = b$$
,  $\nabla f(x^*) + A^\top \nu^* = 0$ 

and use the first-order approximation for the gradient to obtain

$$A(x+\delta_x)=b,$$

$$abla f(x + \delta_x) + A^\top \mu \approx \nabla f(x) + \nabla^2 f(x) \delta_x + A^\top \mu = 0.$$

This is a set of linear equations for  $\delta_x$  and  $\mu$ ,

$$\begin{bmatrix} \nabla^2 f(x) & A^{\top} \\ A & 0 \end{bmatrix} \begin{bmatrix} \delta_x \\ \mu \end{bmatrix} = -\begin{bmatrix} \nabla f(x) \\ Ax - b \end{bmatrix}.$$
 (74)

We express the optimality conditions (71) as  $r(x^*, \nu^*) = 0$ , where  $r : \mathbb{R}^n \times \mathbb{R}^p \mapsto \mathbb{R}^n \times \mathbb{R}^p$  is defined as

$$r(x,\nu) = (r_{\mathsf{dual}}(x,\nu), r_{\mathsf{pri}}(x,\nu)).$$

Here

$$r_{\text{dual}}(x, \nu) = \nabla f(x) + A^{\top} \nu, \quad r_{\text{pri}}(x, \nu) = Ax - b$$

are the dual residual and primal residual, respectively.

YZW (USTC)

(I)

The first-order Taylor approximation of r, near our current point  $y = (x, \nu)$ , is

$$r(y + \delta_y) \approx \hat{r}(y + \delta_y) = r(y) + J[r(y)]\delta_y,$$

where  $J[r(y)] \in \mathbb{R}^{(n+p) \times (n+p)}$  is the derivative (Jacobian) of r, evaluated at y.

(1) マン・ション・ (1) マン・

We define  $\delta_{y_{pd}}$  as the primal-dual Newton step for which  $\hat{r}(y + \delta_y) = 0$ , *i.e.*,  $J[r(y)]\delta_{y_{pd}} = -r(y)$ . (75)

Note that  $\delta_{y_{pd}} = (\delta_{x_{pd}}, \delta_{\nu_{pd}})$  gives both a primal and a dual step.

・ロト ・四ト ・ヨト ・ヨト ・ヨ

Equations (75) can be expressed as

$$\begin{bmatrix} \nabla^2 f(x) & A^{\top} \\ A & 0 \end{bmatrix} \begin{bmatrix} \delta_{x_{pd}} \\ \delta_{\nu_{pd}} \end{bmatrix} = -\begin{bmatrix} r_{dual} \\ r_{pri} \end{bmatrix} = -\begin{bmatrix} \nabla f(x) + A^{\top}\nu \\ Ax - b \end{bmatrix}.$$
 (76)

Writing  $\nu + \delta_{\nu_{\rm pd}}$  as  $\mu$ , we find it coincide with (74)

$$\left[ egin{array}{cc} 
abla^2 f(x) & A^{ op} \\ 
A & 0 \end{array} 
ight] \left[ egin{array}{cc} \delta_x \\ 
\mu \end{array} 
ight] = - \left[ egin{array}{cc} 
abla f(x) \\ 
Ax-b \end{array} 
ight].$$

YZW (USTC)

イロト 不得 トイヨト イヨト

The Newton direction at an infeasible point is not necessarily a descent direction for f.

The primal-dual interpretation, however, shows that the norm of the residual decreases in the Newton direction. By calculation we have

$$\frac{\mathrm{d}}{\mathrm{d}\alpha} \| r(y + \alpha \delta_{y_{\mathsf{pd}}}) \|_2 \Big|_{\alpha = 0} = - \| r(y) \|_2.$$

This allows us to use  $||r||_2$  to measure the progress of the infeasible start Newton method.

#### Algorithm. Infeasible start Newton method.

given starting point  $x \in \text{dom} f$ , tolerance  $\epsilon > 0$ ,  $\tau \in (0, 1/2), \gamma \in (0, 1)$ . repeat

**1** Compute primal and dual Newton steps  $\delta_{x_{nt}}, \delta_{\nu_{nt}}$ .

until Ax = b and  $||r(x, \nu)||_2 \leq \epsilon$ .
#### The convex optimization problems that include inequality constraints:

min 
$$f_0(x)$$
  
s.t.  $f_i(x) \leq 0, \quad i = 1, \dots, m$  (77)  
 $Ax = b$ 

where  $f_0, \ldots, f_m : \mathbb{R}^n \to \mathbb{R}$  are convex and twice continuously differentiable, and  $A \in \mathbb{R}^{p \times n}$  with rank A = p < n.

We assume that an optimal  $x^*$  exists and denote the optimal value  $v^* = f_0(x^*)$ .

We also assume that the problem is strictly feasible, *i.e.*,  $\exists x \in D$  satisfying Ax = b and  $f_i(x) < 0$  for i = 1, ..., m.

This means that Slater's constraint qualification holds, and therefore strong duality holds, so there exists dual optimal  $\lambda^* \in \mathbb{R}^m, \nu^* \in \mathbb{R}^p$ , which together with  $x^*$  satisfy the KKT conditions:

$$Ax^{*} = b, \quad f_{i}(x^{*}) \leq 0, \quad i = 1, \dots, m$$

$$\lambda^{*} \geq 0$$

$$\nabla f_{0}(x^{*}) + \sum_{i=1}^{m} \lambda_{i}^{*} \nabla f_{i}(x^{*}) + A^{\top} \nu^{*} = 0$$

$$\lambda_{i}^{*} f_{i}(x^{*}) = 0, \quad i = 1, \dots, m.$$
(78)

Interior-point methods solve the problem (77) by applying Newton's method to a sequence of equality constrained problems, or to a sequence of modified versions of the KKT conditions.

We will introduce two particular interior-point algorithms:

- The barrier method
- The primal-dual interior-point method

Rewrite the problem (77) and make the inequality constraints implicit in the objective:

min 
$$f_0(x) + \sum_{i=1}^m I_-(f_i(x))$$
  
s.t.  $Ax = b$ , (79)

(日)

where

$$I_{-}(u) = \begin{cases} 0 & u \leq 0 \\ \infty & u > 0. \end{cases}$$

The basic idea of the barrier method is to approximate the indicator function  $I_{-}$  by the function

$$\hat{I}_{-}(u) = -(1/t)\log(-u), \quad \mathbf{dom}\hat{I}_{-} = -\mathbb{R}_{++}$$

where t is a parameter that sets the accuracy of the approximation.

Obviously,  $\hat{I}_{-}$  is convex, nondecreasing and differentiable.

## Logarithmic barrier function



Figure: The dashed lines show the function  $I_{-}(u)$ , and the solid curves show  $\hat{I}_{-}(u) = -(1/t)\log(-u)$ , for t = 0.5, 1, 2. The curve for t = 2 gives the best approximation.

YZW (USTC)

291 / 467

Substituting  $\hat{I}_{-}$  for  $I_{-}$  in (79) gives the approximation

min 
$$f_0(x) + \sum_{i=1}^{m} -(1/t) \log(-f_i(x))$$
  
s.t.  $Ax = b$ . (80)

The function

$$\phi(x) = -\sum_{i=1}^{m} \log(-f_i(x)),$$
(81)

is called the *logarithmic barrier* for the problem (77). Its domain is the set of points that satisfy the inequality constraints of (77) strictly:

$$\mathbf{dom}\phi = \{x \in \mathbb{R}^n | f_i(x) < 0, i = 1, \dots, m\}.$$

The gradient and Hessian of  $\phi$  are given by

$$\nabla \phi(x) = \sum_{i=1}^m \frac{1}{-f_i(x)} \nabla f_i(x),$$

$$\nabla^2 \phi(x) = \sum_{i=1}^m \frac{1}{f_i(x)^2} \nabla f_i(x) \nabla f_i(x)^\top + \sum_{i=1}^m \frac{1}{-f_i(x)} \nabla^2 f_i(x).$$

э

(日)

We multiply the objective of (80) by t, and consider the equivalent problem

nin 
$$tf_0(x) + \phi(x)$$
  
s.t.  $Ax = b$ . (82)

We assume problem (82) can be solved via *Newton's method*, and, that it has a unique solution for each t > 0.

For t > 0 we define  $x^*(t) = \arg \min_x \{tf_0(x) + \phi(x) \text{ s.t. } Ax = b\}$  as the solution of (82).

The *central path* associated with problem (77) is defined as the set of points  $\{x^*(t) \mid t > 0\}$ , which we call the *central points*.

(日)

Points on the central path are characterized by the following necessary and sufficient conditions:  $x^*(t)$  is strictly feasible, *i.e.*, satisfies

$$Ax^*(t) = b, \quad f_i(x^*(t)) < 0, \ i = 1, \dots, m$$

and  $\exists \hat{\nu} \in \mathbb{R}^{p}$  such that

$$0 = t \nabla f_0(x^*(t)) + \nabla \phi(x^*(t)) + A^{\top} \hat{\nu}$$
  
=  $t \nabla f_0(x^*(t)) + \sum_{i=1}^m \frac{1}{-f_i(x^*(t))} \nabla f_i(x^*(t)) + A^{\top} \hat{\nu}$  (83)

holds.

- E - - E -

Every central point yields a dual feasible point.

Define

$$\lambda_i^*(t) = -\frac{1}{tf_i(x^*(t))}, \ i = 1, \dots, m, \quad \nu^*(t) = \frac{\hat{\nu}}{t}.$$
 (84)

Because  $f_i(x^*(t)) < 0, i = 1, \dots, m$ , it's clear that  $\lambda^*(t) > 0$ .

YZW (USTC)

By expressing (83) as

$$abla f_0(x^*(t)) + \sum_{i=1}^m \lambda_i^*(t) \nabla f_i(x^*(t)) + A^\top \nu^*(t) = 0,$$

we see that  $x^*(t)$  minimizes the Lagrangian

$$L(x,\lambda,\nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \nu^\top (Ax - b)$$

for  $\lambda = \lambda^*(t)$  and  $\nu = \nu^*(t)$ . Thus  $(\lambda^*(t), \nu^*(t))$  is a dual feasible pair.

YZW (USTC)

Therefore the dual function  $g(\lambda^*(t), \nu^*(t)) = \min_x L(x, \lambda^*(t), \nu^*(t))$  is finite and

$$g(\lambda^*(t), \nu^*(t)) = f_0(x^*(t)) + \sum_{i=1}^m \lambda^*_i(t) f_i(x^*(t)) + \nu^*(t)^\top (Ax^*(t) - b)$$
  
=  $f_0(x^*(t)) - m/t.$ 

• As an important consequence, we have

$$f_0(x^*(t))-v^*\leqslant m/t.$$

• This confirms that  $x^*(t)$  converge to an optimal point as  $t \to \infty$ .

(日)

Since we have assumed that  $x^*(t)$  is the unique solution to problem (82) for each t > 0, a point is equal to  $x^*(t)$  if and only if  $\exists \lambda, \nu$  such that

$$Ax = b, f_i(x) \leq 0, \quad i = 1, \dots, m$$
  

$$\lambda \geq 0$$
  

$$\nabla f_0(x) + \sum_{i=1}^m \lambda_i \nabla f_i(x) + A^\top \nu = 0$$
  

$$-\lambda_i f_i(x) = 1/t, \quad i = 1, \dots, m.$$
(85)

The only difference between (85) and the KKT condition (78) is that the complementarity condition  $-\lambda_i f_i(x) = 0$  is replaced by the condition  $-\lambda_i f_i(x) = 1/t$ .

In particular, for large t,  $x^*(t)$  and  $\lambda^*(t)$ ,  $\nu^*(t)$  'almost' satisfy the KKT optimality conditions for the problem (77).

・ロト ・四ト ・ヨト ・ヨト ・ヨ

### Algorithm. Barrier method

given strictly feasible  $x, t := t^{(0)} > 0, \gamma > 1$ , tolerance  $\epsilon > 0$ . repeat

- Centering step. Starting at x, compute  $x^*(t)$  by minimizing  $tf_0(x) + \phi(x)$ , subject to Ax = b.
- **2** Update.  $x := x^*(t)$
- **3** Stopping criterion. quit if  $m/t < \epsilon$ .
- Increase t. Let  $t := \gamma t$ .

An execution of step 1 is called an *outer iteration*. We assume that Newton's method is used in step 1, and we refer to the Newton iterations or steps executed during the centering step as *inner iterations*.

- Computing  $x^*(t)$  exactly is not necessary.
- Choice of t<sup>(0)</sup>: If t<sup>(0)</sup> is chosen too large, the first outer iteration will require too many iterations. If t<sup>(0)</sup> is chosen too small, the algorithm will require extra outer iterations.
- The choice of the parameter γ involves a trade-off: If γ is small (*i.e.*, near 1) then centering step will be easy since the previous iterate x is a very good starting point but of course there will be a large number of outer iterations. On the other hand, a large γ resulting in fewer outer iterations but more inner iterations.

イロト 不得 ト イヨト イヨト

In the step 1 of the barrier method, the Newton step  $\delta_{x_{\rm nt}}$  and associated dual variable are given by the linear equations

$$\begin{bmatrix} t\nabla^2 f_0(x) + \nabla^2 \phi(x) & A^{\top} \\ A & 0 \end{bmatrix} \begin{bmatrix} \delta_{x_{\rm nt}} \\ \nu_{\rm nt} \end{bmatrix} = -\begin{bmatrix} t\nabla f_0(x) + \nabla \phi(x) \\ 0 \end{bmatrix}.$$
 (86)

These Newton steps for the centering problem can be interpreted as Newton steps for directly solving the modified KKT equations

$$\nabla f_0(x) + \sum_{i=1}^m \lambda_i \nabla f_i(x) + A^\top \nu = 0$$
  
- $\lambda_i f_i(x) = 1/t, \quad i = 1, \dots, m$   
$$Ax = b.$$
 (87)

YZW (USTC)

## Newton step for modified KKT equations

Let  $\lambda_i = 1/(-tf_i(x))$ . This transforms (87) into

$$\nabla f_0(x) + \sum_{i=1}^m \frac{1}{-tf_i(x)} \nabla f_i(x) + A^{\top} \nu = 0, \quad Ax = b.$$
 (88)

For small  $\delta_x$ ,

$$\begin{aligned} \nabla f_0(x+\delta_x) + \sum_{i=1}^m \frac{1}{-tf_i(x+\delta_x)} \nabla f_i(x+\delta_x) \\ \approx \nabla f_0(x) + \nabla^2 f_0(x) \delta_x + \sum_{i=1}^m \frac{1}{-tf_i(x)} \nabla f_i(x) + \sum_{i=1}^m \frac{1}{-tf_i(x)} \nabla^2 f_i(x) \delta_x \\ + \sum_{i=1}^m \frac{1}{tf_i(x)^2} \nabla f_i(x) \nabla f_i(x)^\top \delta_x. \end{aligned}$$

YZW (USTC)

(日)

# Newton step for modified KKT equations

Let

$$g = \nabla f_0(x) + \sum_{i=1}^m \frac{1}{-tf_i(x)} \nabla f_i(x),$$
  

$$H = \nabla^2 f_0(x) + \sum_{i=1}^m \frac{1}{-tf_i(x)} \nabla^2 f_i(x) + \sum_{i=1}^m \frac{1}{tf_i(x)^2} \nabla f_i(x) \nabla f_i(x)^\top.$$

Observe that

$$g = 
abla f_0(x) + (1/t) 
abla \phi(x), \quad H = 
abla f_0(x) + (1/t) 
abla^2 \phi(x).$$

The Newton step for (88) is

$$H\delta_x + A^{\top}\nu = -g, \quad A\delta_x = 0.$$

Comparing this with (86) shows that

$$\delta_x = \delta_{x_{\text{nt}}}, \quad \nu = \frac{\nu_{\text{nt}}}{t}$$

YZW (USTC)

- The barrier method requires a strictly feasible starting point  $x^{(0)}$ .
- When such a point is not known, the barrier method is preceded by a preliminary stage, called *phase I*, in which a strictly feasible point is computed and used as the starting point for the barrier method.

To find a strictly feasible solution of inequalities and equalities

$$f_i(x) < 0, \quad i = 1, \dots, m, \quad Ax = b,$$
 (89)

we form and solve the following optimization problem

min s  
s.t. 
$$f_i(x) \leq s$$
,  $i = 1, ..., m$  (90)  
 $Ax = b$ 

in the variable  $x \in \mathbb{R}^n$ ,  $s \in \mathbb{R}$ . It's always strictly feasible, and called the *phase I optimization problem* associated with the inequality and equality system (89).

Let  $\bar{v}^*$  be the optimal value of (90).

- If v
  <sup>\*</sup> < 0, then (89) has a strictly feasible solution. In fact, we can terminate solving the problem (90) when s < 0.</li>
- If v
  <sup>\*</sup> > 0, then (89) is infeasible. In fact, we can terminate when a central point give a positive lower bound of v
  <sup>\*</sup> > 0.
- If v
   <sup>\*</sup> = 0 and the minimum is attained at x<sup>\*</sup> and s<sup>\*</sup> = 0, then the set
   of inequalities is feasible but not strictly feasible. If v
   <sup>\*</sup> = 0 and the
   minimum is not attained, then the inequalities are infeasible.

The modified KKT conditions (87) can be expressed as  $r_t(x, \lambda, \nu) = 0$ , where t > 0 and

$$Y_t(x,\lambda,\nu) = \begin{bmatrix} \nabla f_0(x) + J[f(x)]^\top \lambda + A^\top \nu \\ -\operatorname{diag}(\lambda)f(x) - (1/t)\mathbf{1} \\ Ax - b \end{bmatrix}.$$
 (91)

Here  $f : \mathbb{R}^n \to \mathbb{R}^m$  and J[f] are given by

$$f(x) = \begin{bmatrix} f_1(x) \\ \vdots \\ f_m(x) \end{bmatrix}, \quad J[f(x)] = \begin{bmatrix} \nabla f_1(x)^\top \\ \vdots \\ \nabla f_m(x)^\top \end{bmatrix}$$

YZW (USTC)

- 4 個 ト 4 ヨ ト 4 ヨ ト -

## Primal-dual search direction

If  $x, \lambda, \nu$  satisfy  $r_t(x, \lambda, \nu) = 0$  (and  $f_i(x) < 0$ ), then  $x = x^*(t)$ ,  $\lambda = \lambda^*(t)$ and  $\nu = \nu^*(t)$ .

• The first block component of  $r_t$ ,

$$r_{\text{dual}} = \nabla f_0(x) + J[f(x)]^\top \lambda + A^\top \nu$$

is called the *dual residual*.

- The last block component,  $r_{pri} = Ax b$ , is called the *primal residual*.
- The middle block

$$r_{\text{cent}} = -\operatorname{diag}(\lambda)f(x) - (1/t)\mathbf{1},$$

is the *centrality residual*, *i.e.*, the residual for the modified complementarity condition.

Let  $y = (x, \lambda, \nu)$  denote the current point and  $\delta_y = (\delta_x, \delta_\lambda, \delta_\nu)$  denote the Newton step for solving the equation  $r_t(x, \lambda, \nu) = 0$ , for fixed t where  $f(x) < 0, \lambda > 0$ .

The Newton step is characterized by

$$r_t(y+\delta_y)\approx r_t(y)+J[r_t(y)]\delta_y=0.$$

In terms of  $x, \lambda, \nu$ , we have

$$\begin{array}{ccc} \nabla^2 f_0(x) + \sum_{i=1}^m \lambda_i \nabla^2 f_i(x) & J[f(x)]^\top & A^\top \\ -\mathbf{diag}(\lambda) J[f(x)] & -\mathbf{diag}(f(x)) & 0 \\ A & 0 & 0 \end{array} \right] \begin{bmatrix} \delta_x \\ \delta_\lambda \\ \delta_\nu \end{bmatrix} = - \begin{bmatrix} r_{\mathsf{dual}} \\ r_{\mathsf{cent}} \\ r_{\mathsf{pri}} \end{bmatrix}$$
(92)

The primal-dual search direction  $\delta_{y_{pd}} = (\delta_{x_{pd}}, \delta_{\lambda_{pd}}, \delta_{\nu_{pd}})$  is defined as the solution of (92).

In the primal-dual interior-point method the iterates  $x^{(k)}$ ,  $\lambda^{(k)}$  and  $\nu^{(k)}$  are not necessarily feasible. We cannot easily evaluate a duality gap as we do in the barrier method.

Instead, we define the *surrogate duality gap*, for any x that satisfies f(x) < 0 and  $\lambda \ge 0$ , as

$$\hat{\eta}(\mathbf{x},\lambda) = -f(\mathbf{x})^{\top}\lambda.$$

**Remark:** The surrogate gap  $\hat{\eta}$  would be the duality gap, if x were primal feasible and  $\lambda, \nu$  were dual feasible. Note that the value of the parameter t corresponding to the surrogate duality gap  $\hat{\eta}$  is  $m/\hat{\eta}$ .

### Algorithm. Primal-dual interior-point method.

given x that satisfies  $f_1(x) < 0, \ldots, f_m(x) < 0, \lambda > 0, \gamma > 1, \epsilon_{feas} > 0, \epsilon > 0.$  repeat

• Determine t. Set  $t := \gamma(m/\hat{\eta})$ .

2 Compute primal-dual search direction  $\delta_{y_{pd}}$ .

3 Line search and update.

Determine step length  $\alpha > 0$  and set  $y := y + \alpha \delta_{y_{pd}}$ .

**until**  $||r_{pri}||_2 \leq \epsilon_{feas}, ||r_{dual}||_2 \leq \epsilon_{feas}, \text{ and } \hat{\eta} \leq \epsilon.$ 

## Line search in primal-dual interior-point method

The line search in step 3 is a standard backtracking line search.

For a step size  $\alpha$ , let

$$y^{+} = \begin{bmatrix} x^{+} \\ \lambda^{+} \\ \nu^{+} \end{bmatrix} = \begin{bmatrix} x \\ \lambda \\ \nu \end{bmatrix} + \alpha \begin{bmatrix} \delta_{x_{pd}} \\ \delta_{\lambda_{pd}} \\ \delta_{\nu_{pd}} \end{bmatrix}$$

Let

$$\alpha^{\max} = \sup\{\alpha \in [0,1] \mid \lambda + \alpha \delta_{\lambda} \ge 0\} = \min\left\{1, \min\{\frac{-\lambda_{i}}{\delta_{\lambda_{i}}} \mid \delta_{\lambda_{i}} < 0\}\right\}$$

to be the largest positive step length that gives  $\lambda^+ \ge 0$ .

We start backtracking with  $\alpha = 0.99 \alpha^{\text{max}}$ , and multiply  $\alpha$  by  $\beta \in (0, 1)$  until we have  $f(x^+) < 0$ . We continue multiplying  $\alpha$  by  $\beta$  until we have

$$\|r_t(x^+,\lambda^+,\nu^+)\|_2 \leq (1-\tau\alpha)\|r_t(x,\lambda,\nu)\|_2.$$

Here  $\tau$  is typically chosen in the range [0.01, 0.1].

Ex 1. Let  $C \subseteq \mathbb{R}^n$  be the solution set of a quadratic inequality,

$$C = \{ x \in \mathbb{R}^n | x^\top A x + b^\top x + c \leq 0 \},\$$

with  $A \in S^n$ ,  $b \in \mathbb{R}^n$ , and  $c \in \mathbb{R}$ .

- (a) Show that C is convex if  $A \succeq 0$ .
- (b) Show that the intersection of C and the hyperplane defined by  $g^{\top}x + h = 0$  (where  $g \neq 0$ ) is convex if  $A + \lambda gg^{\top} \succeq 0$  for some  $\lambda \in \mathbb{R}$ .
- Ex 2. Let  $\lambda_1(X) \ge \lambda_2(X) \ge \ldots \ge \lambda_n(X)$  denote the eigenvalues of a matrix  $X \in S^n$ . Prove that the maximum eigenvalue  $\lambda_1(X)$  is convex. Moreover, Show that  $\sum_{i=1}^k \lambda_i(X)$  is convex on  $S^n$ . Hint. Use the variational characterization

$$\sum_{i=1}^k \lambda_i(X) = \sup\{\operatorname{tr}(V^\top X V) | V \in \mathbb{R}^{n \times k}, V^\top V = I\}.$$

YZW (USTC)

Ex 3. Find the dual function of the LP

$$\begin{array}{ll} \min & c^T x \\ \text{s.t.} & Gx \leq h \\ & Ax = b. \end{array}$$

Give the dual problem, and make the implicit equality constraints explicit.

Ex 4. Consider the equality constrained least-squares problem

$$\begin{array}{ll} \min & \|Ax - b\|_2^2 \\ \text{s.t.} & Gx = h \end{array}$$

where  $A \in \mathbb{R}^{m \times n}$  with  $\operatorname{rank} A = n$ , and  $G \in \mathbb{R}^{p \times n}$  with  $\operatorname{rank} G = p$ . Give the KKT conditions, and derive expressions for the primal solution  $x^*$  and the dual solution  $\nu^*$ .

YZW (USTC)

Ex 5. Suppose  $Q \succeq 0$ . The problem

min 
$$f(x) + (Ax - b)^{\top}Q(Ax - b)$$
  
s.t.  $Ax = b$ 

is equivalent to the primal equality constrained optimization problem (70). What is the Newton step for this problem? Is it the same as that for the primal problem?

- E - - E -

- Ex 6. Suppose we use the infeasible start Newton method to minimize f(x) subject to  $a_i^{\top} x = b_i$ , i = 1, ..., p.
  - (a) Suppose the initial point x<sup>(0)</sup> satisfies the linear equality a<sub>i</sub><sup>⊤</sup>x<sup>(0)</sup> = b<sub>i</sub>. Show that the linear equality will remain satisfied for future iterates, *i.e.*, a<sub>i</sub><sup>⊤</sup>x<sup>(k)</sup> = b<sub>i</sub> for all k.
  - (b) Suppose that one of the equality constraints becomes satisfied at iteration k, i.e., we have a<sup>T</sup><sub>i</sub>x<sup>(k-1)</sup> ≠ b<sub>i</sub>, a<sup>T</sup><sub>i</sub>x<sup>(k)</sup> = b<sub>i</sub>. Show that at iteration k, all the equality constraints are satisfied.

ヘロト 人間 とくほとく ほとう

Ex 7. Suppose we add the constraint  $x^{\top}x \leq R^2$  to the problem (77):

min 
$$f_0(x)$$
  
s.t.  $f_i(x) \le 0$ ,  $i = 1, ..., m$   
 $Ax = b$   
 $x^T x \le R^2$ 

Let  $\tilde{\phi}$  denote the logarithmic barrier function for this modified problem. Find a > 0 for which  $\nabla^2(tf_0(x) + \phi(x)) \succeq aI$  holds, for all feasible x.
Ex 8. Consider the problem (77), with central path  $x^*(t)$  for t > 0, defined as the solution of (82). For  $u > p^*$ , let  $z^*(u)$  denote the solution of

min 
$$-\log(u - f_0(x)) - \sum_{i=1}^m \log(-f_i(x))$$
  
s.t.  $Ax = b$ 

Show that the curve define by  $z^*(u)$ , for  $u > p^*$ , is the central path. (In other words, for each  $u > p^*$ , there is a t > 0 for which  $x^*(t) = z^*(u)$ , and conversely, for each t > 0, there is a  $u > p^*$  for which  $z^*(u) = x^*(t)$ ).

# Outline I

- Unconstrained Optimization
- 2 Constrained Optimization
  - 二次规划
  - 非线性约束最优化
- 3 Convex Optimization
  - Convex Set and Convex Function
  - Convex Optimization and Algorithms

#### 4 Sparse Optimization

- Sparse Optimization Models
- Sparse Optimization Algorithms

#### Optimization Methods for Machine Learning

YZW (USTC)

→ < Ξ →</p>

# Outline II

- Typical Form of Problems
- Stochastic Algorithms
- Other Popular Methods



æ

Many problems of recent interest in statistics and related areas can be posed in the framework of sparse optimization. Due to the explosion in size and complexity of modern data analysis (BigData), it is increasingly important to be able to solve problems with a very large number of features, training examples, or both.

$$(P_0) \quad \min_{x} \|x\|_0 \quad \text{s.t.} \quad Ax = b.$$
 (93)

# $(P_0^{\epsilon}) \quad \min_{\mathbf{x}} \|\mathbf{x}\|_0 \quad \text{s.t.} \quad \|\mathbf{b} - A\mathbf{x}\| \leq \epsilon.$ (94)

YZW (USTC)

326 / 467

臣

イロト イヨト イヨト イヨト

### Greedy algorithms

Greedy strategies are usually adopted in solving the 0-norm problems. The following algorithm is known in the literature of signal processing by the name *Orthogonal Matching Pursuit* (OMP).

**Task:** Approximate the solution of  $(P_0)$ :  $\min_{\mathbf{x}} \|\mathbf{x}\|_0$  subject to  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .

**Parameters:** We are given the matrix A, the vector b, and the error threshold  $\epsilon_0$ .

**Initialization:** Initialize k = 0, and set

- The initial solution  $\mathbf{x}^0 = 0$ .
- The initial residual  $\mathbf{r}^0 = \mathbf{b} \mathbf{A}\mathbf{x}^0 = \mathbf{b}$ .
- The initial solution support  $S^0 = Support\{\mathbf{x}^0\} = \emptyset$ .

**Main Iteration:** Increment k by I and perform the following steps:

- Sweep: Compute the errors  $\epsilon(j) = \min_{z_j} \|\mathbf{a}_j z_j \mathbf{r}^{k-1}\|_2^2$  for all j using the optimal choice  $z_j^* = \mathbf{a}_j^T \mathbf{r}^{k-1} / \|\mathbf{a}_j\|_2^2$ .
- Update Support: Find a minimizer  $j_0$  of  $\epsilon(j)$ :  $\forall j \notin S^{k-1}$ ,  $\epsilon(j_0) \leq \epsilon(j)$ , and update  $S^k = S^{k-1} \cup \{j_0\}$ .
- Update Provisional Solution: Compute  $\mathbf{x}^k$ , the minimizer of  $\|\mathbf{A}\mathbf{x} \mathbf{b}\|_2^2$  subject to  $Support\{\mathbf{x}\} = S^k$ .
- Update Residual: Compute  $r^k = b Ax^k$ .
- Stopping Rule: If  $\|\mathbf{r}^k\|_2 < \epsilon_0$ , stop. Otherwise, apply another iteration.

**Output:** The proposed solution is  $\mathbf{x}^k$  obtained after k iterations.

The optimization model of dictionary learning for sparse and redundant representations:

$$\min_{D,X} \|Y - DX\|_{Frob} \quad \text{s.t.} \quad \|x_j\|_0 \le k_0, \ j = 1, \cdots, N$$
(95)

where

$$Y = (y_1, \cdots, y_N) \in \mathbb{R}^{n \times N},$$

$$D = (\mathsf{d}_1, \cdots, \mathsf{d}_m) \in \mathbb{R}^{n \times m},$$

$$X = (\mathsf{x}_1, \cdots, \mathsf{x}_N) \in \mathbb{R}^{m \times N}.$$

YZW (USTC)

イロト 不得下 イヨト イヨト

There are two training mechanisms, the first named Method of Optimal Directions (MOD) by Engan et al., and the second named K-SVD, by Aharon et al..

MOD

K-SVD



YZW (USTC)

Convex relaxation technique is a way to render 0-norm more tractable.

Convexifying with the  $\ell_1$  norm, we come to the new optimization problem

$$(P_1) \quad \min_{x} \|Wx\|_1 \quad \text{s.t.} \quad Ax = b$$
 (96)

where W is a diagonal positive-definite matrix that introduces the precompensating weights.

An error-tolerant version of  $(P_1)$  is defined by

$$(P_1^{\epsilon}) \quad \min_{\mathbf{x}} \|W\mathbf{x}\|_1 \quad \text{s.t.} \quad \|\mathbf{b} - A\mathbf{x}\| \leq \epsilon.$$
(97)

It was named *Basis Pursuit* (BP) when all the columns of A are normalized (and thus W = I).

$$(BP) \quad \min_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{s.t.} \quad A\mathbf{x} = \mathbf{b}.$$



YZW (USTC)

#### **Optimization Algorithms**

331 / 467

æ

$$\begin{array}{ll} (BP_{\tau}) & \min_{\mathbf{x}} \|A\mathbf{x} - \mathbf{b}\|_{2}^{2} & \text{s.t.} & \|\mathbf{x}\|_{1} \leqslant \tau, \\ (BP_{\mu}) & \min_{\mathbf{x}} \|\mathbf{x}\|_{1} + \frac{\mu}{2} \|A\mathbf{x} - \mathbf{b}\|_{2}^{2}, \\ (BP_{\delta}) & \min_{\mathbf{x}} \|\mathbf{x}\|_{1} & \text{s.t.} & \|A\mathbf{x} - \mathbf{b}\|_{2} \leqslant \delta. \end{array}$$

Questions:

- Are they equivalent? and in what sense?
- How to choose parameters?

E + 4 E +

#### Sparse under basis $\Psi$

$$\min_{s}\{\|s\|_{1}:A\Psi s=b\}$$



If  $\Psi$  is orthogonal, the problem is equivalent to

$$\min_{\mathbf{x}}\{\|\Psi^*\mathbf{x}\|_1: A\mathbf{x}=\mathsf{b}\}.$$

∃▶ ∢ ∃▶

$$\min_{\mathbf{x}}\{\|\mathcal{L}\mathbf{x}\|_1: A\mathbf{x} = \mathbf{b}\}$$

Examples of  $\mathcal{L}$ :

- DCT, wavelets, curvelets, ridgelets, ...
- tight frames, Gabor, ...
- total (generalized) variation

**Ref**: E. J. Cands, Y. Eldar, D. Needell and P. Randall. Compressed sensing with coherent and redundant dictionaries. Applied and Computational Harmonic Analysis, 31(1): 59-73, 2011.

Decompose  $\{1, 2, \dots, n\} = \mathcal{G}_1 \bigcup \mathcal{G}_2 \bigcup \dots \bigcup \mathcal{G}_S$ , and  $\mathcal{G}_i \bigcap \mathcal{G}_j = \emptyset, i \neq j$ .

Joint/group sparse recovery model:

$$\min_{\mathbf{x}}\{\|\mathbf{x}\|_{\mathcal{G},2,1}:A\mathbf{x}=\mathbf{b}\}$$

where

$$\|\mathbf{x}\|_{\mathcal{G},2,1} = \sum_{s=1}^{S} w_s \|\mathbf{x}_{\mathcal{G}_s}\|_2.$$

YZW (USTC)

▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ 臣 のへで

- Nonnegativity:  $x \ge 0$
- Box constraints:  $lb \le x \le ub$
- Linear inequalities:  $Qx \leq c$

They generate "corners" and can be very effective in practice.

- Shrinkage is popular in sparse optimization algorithms
- In optimization, non-smooth functions like  $\ell_1$  has difficulty using general smooth optimization methods.
- But,  $\ell_1$  is component-wise separable, so it does get along well with separable (smooth or non-smooth) functions.
- For example,

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1 + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{z}\|_2^2$$

is equivalent to solving  $\min_{x_i} |x_i| + \frac{1}{2\tau} |x_i - z_i|^2$  over each *i*.

(1) マン・ション・ (1) マン・

#### Soft-thresholding shrinkage

The problem is separable and has an explicit solution

$$(\operatorname{shrink}(\mathsf{z}, au))_i = \left\{ egin{array}{ccc} z_i - au & z_i > au, \ 0 & - au \leq z_i \leq au, \ z_i + au & z_i < - au. \end{array} 
ight.$$



The shrinkage operator can be written in Matlab code as:  $\mathbf{x} = \max(abs(\mathbf{z})-tau, 0).*sign(\mathbf{z}).$ 

YZW (USTC)

Optimization Algorithms

338 / 467

• The following problem is called Moreau-Yosida regularization

$$\min_{\mathbf{x}} r(\mathbf{x}) + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{z}\|_2^2.$$

• For example  $r(x) = ||x||_2$ , the solution to

$$\min_{\mathbf{x}} \|\mathbf{x}\|_2 + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{z}\|_2^2$$

is, if we treat 0/0 = 0,

$$x_{opt} = \max\{\|z\|_2 - \tau, 0\} \cdot (z/\|z\|_2).$$

• Used in joint/group-sparse recovery algorithms.

#### Soft-thresholding shrinkage

• Consider the following nuclear norm optimization

$$\min_{\mathsf{X}} \|\mathsf{X}\|_* + \frac{1}{2\tau} \|\mathsf{X} - \mathsf{Z}\|_F^2.$$

Let  $Z = U\Sigma V^T$  be the singular value decomposition of Z.

• Let  $\hat{\Sigma}$  be the diagonal matrix with diagonal entries

$$\operatorname{diag}(\hat{\Sigma}) = \operatorname{shrink}(\operatorname{diag}(\Sigma), \tau)),$$

then

$$X_{opt} = U\hat{\Sigma}V^{T}.$$

YZW (USTC)

#### Prox-linear algorithm

Consider the general form

$$\min_{\mathbf{x}} r(\mathbf{x}) + f(\mathbf{x}).$$

where r is the regularization function and f is the data fidelity function.

The prox-linear algorithm is:

$$\mathsf{x}^{k+1} = \arg\min_{\mathsf{x}} r(\mathsf{x}) + f(\mathsf{x}^k) + < \nabla f(\mathsf{x}^k), \mathsf{x} - \mathsf{x}^k > + \frac{1}{2\delta_k} \|\mathsf{x} - \mathsf{x}^k\|_2^2.$$

The last term keeps  $x^{k+1}$  close to  $x^k$ , and the parameter  $\delta_k$  determines the step size. It is equivalent to

$$\mathbf{x}^{k+1} = \arg\min_{\mathbf{x}} r(\mathbf{x}) + \frac{1}{2\delta_k} \|\mathbf{x} - (\mathbf{x}^k - \delta_k \nabla f(\mathbf{x}^k))\|_2^2.$$

YZW (USTC)

The Alternating Direction Method of Multipliers (ADMM) was developed in the 1970s, with roots in the 1950s, and is equivalent or closely related to many other algorithms, such as dual decomposition, the method of multipliers, Douglas-Rachford splitting, Spingarns method of partial inverses, Dykstras alternating projections, Bregman iterative algorithms for 1-norm problems, proximal methods, and others. The ADMM can be applied to a wide variety of statistical and machine learning problems of recent interest, including the lasso, sparse logistic regression, basis pursuit, covariance selection, support vector machines, and many others.

$$\min_{\mathsf{X}\in C^{n\times T}} \mu \|\mathsf{X}\|_p + \|\mathsf{A}\mathsf{X} - \mathsf{B}\|_q \tag{98}$$

Let  $p := \{2, 1\}, q := \{1, 1\}$  which denote joint convex norm, we have

$$\min_{\mathsf{X} \in C^{n \times T}} \mu \|\mathsf{X}\|_{2,1} + \|\mathsf{A}\mathsf{X} - \mathsf{B}\|_{1,1}$$

where  $\|X\|_{2,1} = \sum_{i=1}^{n} \sqrt{\sum_{j=1}^{T} x_{ij}^2}$ ,  $\|X\|_{1,1} = \sum_{i=1}^{n} \sum_{j=1}^{T} |x_{ij}|$ .

For example T = 1,

$$\min_{\mathbf{x}\in C^n} \mu \|\mathbf{x}\|_p + \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_q.$$

YZW (USTC)

$$\min_{\mathsf{X}\in C^{n\times T}} \mu \|\mathsf{X}\|_{\rho} + \|\mathsf{A}\mathsf{X} - \mathsf{B}\|_{q}$$
(98)

Let  $p := \{2,1\}, q := \{1,1\}$  which denote joint convex norm, we have

$$\min_{\mathsf{X}\in C^{n\times T}} \mu \|\mathsf{X}\|_{2,1} + \|\mathsf{A}\mathsf{X} - \mathsf{B}\|_{1,1}$$

where 
$$\|X\|_{2,1} = \sum_{i=1}^n \sqrt{\sum_{j=1}^T x_{ij}^2}$$
,  $\|X\|_{1,1} = \sum_{i=1}^n \sum_{j=1}^T |x_{ij}|$ .

For example T = 1,

$$\min_{\mathbf{x}\in C^n} \mu \|\mathbf{x}\|_p + \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_q.$$

YZW (USTC)

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

$$\min_{\mathsf{X}\in C^{n\times T}} \mu \|\mathsf{X}\|_{\rho} + \|\mathsf{A}\mathsf{X} - \mathsf{B}\|_{q}$$
(98)

Let  $p := \{2,1\}, q := \{1,1\}$  which denote joint convex norm, we have

$$\min_{\mathsf{X}\in C^{n\times T}} \mu \|\mathsf{X}\|_{2,1} + \|\mathsf{A}\mathsf{X} - \mathsf{B}\|_{1,1}$$

where 
$$\|X\|_{2,1} = \sum_{i=1}^{n} \sqrt{\sum_{j=1}^{T} x_{ij}^2}$$
,  $\|X\|_{1,1} = \sum_{i=1}^{n} \sum_{j=1}^{T} |x_{ij}|$ .

For example T = 1,

$$\min_{\mathbf{x}\in C^n} \mu \|\mathbf{x}\|_p + \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_q.$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

min 
$$\mu \|\mathbf{z}\|_{p} + \|\mathbf{y}\|_{q}$$
  
s.t.  $\mathbf{x} - \mathbf{z} = \mathbf{0}$  (99)  
 $A\mathbf{x} - \mathbf{y} = \mathbf{b}$ 

イロト イヨト イヨト イヨト

$$L(\mathbf{x}, \mathbf{y}, \mathbf{z}, \lambda_y, \lambda_z, \rho) = \mu \|\mathbf{z}\|_{\rho} + \|\mathbf{y}\|_{q} + \operatorname{Re}(\lambda_z^{T}(\mathbf{x} - \mathbf{z}) + \lambda_y^{T}(\mathbf{A}\mathbf{x} - \mathbf{y} - \mathbf{b})) + \frac{\rho}{2}(\|\mathbf{x} - \mathbf{z}\|_{2}^{2} + \|\mathbf{A}\mathbf{x} - \mathbf{y} - \mathbf{b}\|_{2}^{2})$$
(100)  
where  $\lambda_y \in C^n, \lambda_z \in C^m$  are the Lagrangian multipliers and  $\rho > 0$  is a

where  $\lambda_y \in C^n, \lambda_z \in C^m$  are the Lagrangian multipliers and  $\rho > 0$  is a penalty parameter.

æ

$$\min \mu \|\mathbf{z}\|_{p} + \|\mathbf{y}\|_{q}$$
s.t.  $\mathbf{x} - \mathbf{z} = \mathbf{0}$ 

$$A\mathbf{x} - \mathbf{y} = \mathbf{b}$$

$$(99)$$

$$L(x, y, z, \lambda_{y}, \lambda_{z}, \rho) = \mu ||z||_{\rho} + ||y||_{q} + \operatorname{Re}(\lambda_{z}^{T}(x - z) + \lambda_{y}^{T}(Ax - y - b)) + \frac{\rho}{2}(||x - z||_{2}^{2} + ||Ax - y - b||_{2}^{2})$$
(100)

where  $\lambda_y \in C^n, \lambda_z \in C^m$  are the Lagrangian multipliers and  $\rho > 0$  is a penalty parameter.

▲□▶ ▲圖▶ ▲ 圖▶ ▲ 圖▶ ― 圖 … のへで

$$\begin{cases} x^{k+1} := \arg\min\frac{1}{2}(\|\mathbf{x} - \mathbf{z}^{k} + \mathbf{u}_{z}^{k}\|_{2}^{2} + \|\mathbf{A}\mathbf{x} - \mathbf{y}^{k} - \mathbf{b} + \mathbf{u}_{y}^{k}\|_{2}^{2}) \\ y^{k+1} := \arg\min\|\mathbf{y}\|_{q} + \frac{\rho}{2}\|\mathbf{y} - (\mathbf{A}\mathbf{x}^{k+1} - \mathbf{b}) - \mathbf{u}_{y}^{k}\|_{2}^{2} \\ z^{k+1} := \arg\min\mu\|\mathbf{z}\|_{\rho} + \frac{\rho}{2}\|\mathbf{z} - \mathbf{x}^{k+1} - \mathbf{u}_{z}^{k}\|_{2}^{2} \end{cases}$$
(101)

After solving three subproblems, we update the Lagrangian multipliers as follows:

$$\begin{cases} u_{z}^{k+1} = u_{z}^{k} + \gamma(x^{k+1} - z^{k+1}) \\ u_{y}^{k+1} = u_{y}^{k} + \gamma(Ax^{k+1} - y^{k+1} - b) \end{cases}$$
(102)

where  $u_y = \frac{1}{\rho} \lambda_y$ ,  $u_z = \frac{1}{\rho} \lambda_z$ ,  $\gamma > 0$  is the step size.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

# Outline I

- Unconstrained Optimization
- 2 Constrained Optimization
  - 二次规划
  - 非线性约束最优化
- 3 Convex Optimization
  - Convex Set and Convex Function
  - Convex Optimization and Algorithms
- Sparse Optimization
  - Sparse Optimization Models
  - Sparse Optimization Algorithms

#### Optimization Methods for Machine Learning

YZW (USTC)

• • = •

# Outline II

- Typical Form of Problems
- Stochastic Algorithms
- Other Popular Methods



æ

Mathematical optimization is one of the pillars of machine learning. Large-scale machine learning, where the amount of both the training data and the parameters is large, represents a distinctive setting in which traditional nonlinear optimization techniques typically falter.

We will briefly introduce some typical optimization problems arising from machine learning and then turn to stochastic algorithms—the main content of this section—and other popular methods together with specific models applicable to them.

For simplicity, we focus on the problems that arise in the context of *supervised classification*; i.e., we focus on the optimization of prediction functions for labeling unseen data based on information contained in a set of labeled training data.

Such a focus is reasonable as many unsupervised and other learning techniques reduce to optimization problems of comparable form.

For example:

- Regression. Although the methodology of dealing with regression is quite different from that of classification, regression does share a model similar to supervised classification. Supervised classification and regression are collectively called supervised learning.
- Deep reinforcement learning. In deep Q-learning network (DQN), the samples are attained by interacting with environment, and to train the agent is to solve the Bellman equation in a regression fashion.
- Generative adversarial network. The GAN is composed of a generator and a discriminator, which are usually trained alternately. The training process of each part could be treated as a supervised classification, where the label means whether the sample comes from the data distribution or not.

A 目 > A 目 > A 目 >

- Goal determine a prediction function  $h : \mathcal{X} \to \mathcal{Y}$  from an input space  $\mathcal{X}$  to an output space  $\mathcal{Y}$ .
- Request *h* should avoid rote memorization and instead generalizes the concepts that can be learned from a given set of examples.
- Scheme choose h by attempting to minimize a risk measure over an adequately selected family of prediction functions, call it H.

Suppose that samples are sampled from a joint probability distribution function Pr(x, y).

*h* to be sought should yield a small *expected risk* of misclassification *over all possible inputs*, i.e., minimize

$$R(h) = \Pr[h(x) \neq y] = E\left[1_{[h(x)\neq y]}\right].$$
(103)

Such a framework is *variational* since we are optimizing over a set of functions, and is *stochastic* since the objective function involves an expectation.

In practice, the expectation is taken on samples  $\{(x_i, y_i)\}_{i=1}^n$  and h should minimize the *empirical risk* of misclassification

$$R_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{[h(x_i) \neq y_i]}, \text{ where } \mathbb{1}_{[A]} = \begin{cases} 1 & \text{if } A \text{ is true} \\ 0 & \text{otherwise} \end{cases}.$$
(104)

YZW (USTC)

The family of function  ${\cal H}$  should be determined with three potentially competing goals in mind.

1. Adequate capacity:  $\mathcal{H}$  should contain prediction functions that are able to achieve a low empirical risk over the training set, so as to avoid underfitting the data.

This can be achieved by selecting a rich family of functions or by using a *priori* knowledge to select a well-targeted family.
### Choice of Prediction Function Family

2. Low generalization error: The gap between expected risk and empirical risk  $R(h) - R_n(h)$  should be small over all  $h \in \mathcal{H}$ .

Generally this gap decreases when one uses more training examples but it increases when one uses richer families of functions, due to potential overfitting.

 Efficient training: H should be selected so that one can efficiently solve the corresponding optimization problem, the difficulty of which may increase when one employs a richer family of functions and/or a larger training set. By certain laws of large numbers, the Hoeffding inequality guarantees that, with probability  $1-\eta,$ 

$$|R(h) - R_n(h)| \leqslant \sqrt{rac{1}{2n}\log\left(rac{2}{\eta}
ight)}$$
 for a given  $h \in \mathcal{H}$ 

This bound offers the intuitive explanation that the gap decreases as one uses more training examples.

For a uniform generalization error bound, one often turns to *uniform laws* of large numbers and the concept of the Vapnik-Chervonenkis (VC) dimension of  $\mathcal{H}$ , a measure of the capacity of such a family of functions.

イロト 不得 ト イヨト イヨト

# Generalization Error

Roughly speaking, the VC dimension of a family of functions is the minimal size of samples on which all the functions in the family fail.

For the intuition behind this concept, consider, e.g., a binary classification scheme in  $\mathbb{R}^2$  where one assigns a label of 1 for points above a polynomial and -1 for points below. Then the set of linear polynomials has a low capacity with VC dimension of 3.

With  $d_{\mathcal{H}}$  defined as the VC dimension of  $\mathcal{H}$ , one has with probability at least  $1 - \eta$  that

$$\sup_{h \in \mathcal{H}} |R(h) - R_n(h)| \leq \mathcal{O}\left(\sqrt{\frac{1}{2n}\log\left(\frac{2}{\eta}\right) + \frac{d_{\mathcal{H}}}{n}\log\left(\frac{n}{d_{\mathcal{H}}}\right)}\right).$$
(105)

(105) is one of the most important results in learning theory.

・ロト ・雪 ト ・ヨ ト

Rather than choose a generic family of prediction functions (difficult to optimize and estimate the generalization error) one chooses a *structure*, i.e., a collection of nested function families.

For instance, such a structure can be formed as a collection of subsets of a given family  $\mathcal{H}$  in the following manner: given a preference function  $\Omega$ , choose various values of a *hyperparameter C*, according to each of which one obtains the subset  $\mathcal{H}_C := \{h \in \mathcal{H} \mid \Omega(h) \leq C\}$ .

Given a fixed number of examples, increasing C reduces the empirical risk, but after some point it typically increases the gap between expected and empirical risks, as illustrated in Fig 7.

Other ways to introduce structures are to consider a regularized empirical risk  $R_n(h) + \lambda \Omega(h)$ .

## Structural Risk Minimization



#### Figure: Illustration of structural risk minimization

< □ > < 凸

359 / 467

∃ ▶ ∢ ∃ ▶

One can avoid estimating the gap between empirical and expected risk by splitting the available data into three subsets: a *training set*, a *validation set* and a *testing set*.

Specifically, over the training set one minimizes an empirical risk measure  $R_n$  over  $\mathcal{H}_C$  for various values of C. This results in a handful of candidate functions.

The validation set is then used to estimate the expected risk corresponding to each candidate solution, after which one chooses the function yielding the lowest estimated risk value.

The testing set is used to estimate the expected risk for the candidate that is ultimately chosen.

A 回 > A 回 > A 回 >

Now we assume that the prediction function h has a fixed form and is parameterized by a real vector  $w \in \mathbb{R}^d$  over which the optimization is to be performed.

Formally, for some given  $h(\cdot, \cdot) : \mathbb{R}^{d_x} \times \mathbb{R}^d \to \mathbb{R}^{d_y}$ , we consider the family of prediction functions

$$\mathcal{H} := \left\{ h(\cdot, w) \mid w \in \mathbb{R}^d \right\}.$$

To measure the losses incurred from inaccurate predictions, we assume a given loss function  $\ell : \mathbb{R}^{d_y} \times \mathbb{R}^{d_y} \to \mathbb{R}$ . An input-output pair (x, y) yields the predicted output h(x, w) and the loss  $\ell(h(x, w), y)$ .

### More Practical Statements

We have the expected risk

$$R(w) = E_{(x,y)\sim Pr(x,y)} \left[ \ell(h(x,w), y) \right],$$
(106)

and the empirical risk

$$R_n(w) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i, w), y_i).$$
(107)

To simplify the notation, let  $\xi$  be a sample (x, y) and  $f(w, \xi) = \ell(h(x, w), y)$ , then the expected risk is

$$R(w) = E_{\xi} [f(w, \xi)].$$
(108)

For a set of samples  $\{\xi_i\}_{i=1}^n$ , let us define  $f_i(w)$  to be  $f(w, \xi_i)$  and then the empirical risk is

$$R_n(w) = \frac{1}{n} \sum_{i=1}^n f_i(w).$$
 (109)

YZW (USTC)

362 / 467

# A Brief Introduction

Recall the batch (ordinary) gradient descent method. To minimize the empirical risk (as (109)), w is updated by

$$w_{k+1} \leftarrow w_k - \alpha_k \nabla R_n(w_k) = w_k - \frac{\alpha_k}{n} \sum_{i=1}^n \nabla f_i(w_k)$$
(110)

where  $\alpha_k > 0$  is a stepsize. Computing the step  $-\alpha_k \nabla R_n(w_k)$  is expensive since it needs accessing all the samples.

*Stochastic gradient* (SG) meanwhile uses only one sample at each iteration:

$$w_{k+1} \leftarrow w_k - \alpha_k \nabla f_{i_k}(w_k) \tag{111}$$

where  $i_k$  is chosen randomly from  $\{1, \ldots, n\}$ . While  $-\nabla f_{i_k}(w_k)$  might not be one of descent from  $w_k$ , if it is a descent direction *in expectation*, then the sequence  $\{w_k\}$  can be guided toward a minimizer of  $R_n$ .

To generalize SG method, we consider two ways:

- reduce the noise (variance) of each iteration by generating a batch of samples instead of a single sample.
- make use of second-order information and compute a stochastic Newton or quasi-Newton direction rather than a gradient direction.

# A Brief Introduction



Figure: Schematic of a two dimensional spectrum of optimization methods for machine learning.

YZW (USTC)

365 / 467

э

(日)

Here we give a general framework of stochastic gradient methods by introducing a general  $\xi_k$  and a general direction  $g(w_k, \xi_k)$ :

#### Algorithm 1 Stochastic Gradient

- 1: Choose an initial iterate  $w_1$ .
- 2: for k = 1, 2, ... do
- 3: Generate a realization of the random variable  $\xi_k$ .
- 4: Compute a direction  $g(w_k, \xi_k)$ .
- 5: Choose a stepsize  $\alpha_k > 0$ .
- 6: Set the new iterate as  $w_{k+1} \leftarrow w_k \alpha_k g(w_k, \xi_k)$ .
- 7: end for

 $\xi_k$  could be either one sample or a set of samples, and our analysis cover the following choices of  $g(w_k, \xi_k)$ :

$$g(w_{k},\xi_{k}) = \begin{cases} \nabla_{w}f(w_{k},\xi_{k}) \\ \frac{1}{n_{k}}\sum_{i=1}^{n_{k}}\nabla_{w}f(w_{k},\xi_{k,i}) \\ H_{k}\frac{1}{n_{k}}\sum_{i=1}^{n_{k}}\nabla_{w}f(w_{k},\xi_{k,i}) \end{cases}$$
(112)

where  $H_k$  is a symmetric positive definite scaling matrix and the third choice represents a stochastic Newton or quasi-Newton direction.

Before establishing the convergence guarantees for SG, we need to make an assumption of smoothness of the objective function:

#### Assumption (Lipschitz-continuous objective gradients)

The objective function  $F : \mathbb{R}^d \to \mathbb{R}$  is continuously differentiable and the gradient  $\nabla F : \mathbb{R}^d \to \mathbb{R}^d$ , is Lipschitz continuous with Lipschitz constant L > 0, i.e.,

$$\|\nabla F(w) - \nabla F(\bar{w})\|_2 \leq L \|w - \bar{w}\|_2, \quad \forall w, \bar{w} \in \mathbb{R}^d.$$
(113)

Under the above assumption, we obtain the following lemma.

### 引理

Under Assumption (113), the iterates of SG (Algorithm 1) satisfy the following inequality for all  $k \in \mathbb{N}$ :

$$E_{\xi_{k}}[F(w_{k+1})] - F(w_{k}) \leq -\alpha_{k} \nabla F(w_{k})^{\top} E_{\xi_{k}}[g(w_{k},\xi_{k})] \\ + \frac{\alpha_{k}^{2}L}{2} E_{\xi_{k}}\left[\|g(w_{k},\xi_{k})\|_{2}^{2}\right].$$

$$(114)$$

Noting that  $w_{k+1}$  but not  $w_k$  depends on  $\xi_k$ , we can derive this equation immediately by simply applying the second-order expansion of  $F(w_{k+1}) - F(w_k)$  and the assumption (113) then taking expectations.

# Two Fundamental Lemmas

To get further, we need another assumption about the first and second moments of the stochastic vectors  $\{g(w_k, \xi_k)\}$ .

#### Assumption (First and second moment limits)

The objective function and SG satisfy the following:

- The sequence of iterates {w<sub>k</sub>} is contained in an open set over which F is bounded below by a scalar F<sub>inf</sub>.
- **)** There exist scalars  $\mu_{\mathsf{G}} \ge \mu > 0$  such that, for all  $k \in \mathbb{N}$ ,

$$\nabla F(w_k)^{\top} E_{\xi_k} \left[ g(w_k, \xi_k) \right] \ge \mu \left\| \nabla F(w_k) \right\|_2^2 \text{ and } (115a) \left\| E_{\xi_k} \left[ g(w_k, \xi_k) \right] \right\|_2 \le \mu_G \left\| \nabla F(w_k) \right\|_2.$$
 (115b)

There exist scalars  $M \ge 0$  and  $M_V \ge 0$  such that, for all  $k \in \mathbb{N}$ ,

$$Var_{\xi_k} [\|g(w_k,\xi_k)\|] \leq M + M_V \|\nabla F(w_k)\|_2^2.$$
 (116)

YZW (USTC)

## Two Fundamental Lemmas

By the definition of variance, it requires that the second moment of  $g(w_k, \xi_k)$  satisfies

 $E_{\xi_k} \left[ \|g(w_k, \xi_k)\|_2^2 \right] \leqslant M + M_G \|\nabla F(w_k)\|_2^2 \text{ with } M_G := M_V + \mu_G^2 \geqslant \mu^2 > 0.$ (117)

### 引理

Under the above two assumptions, the iterates of SG satisfy the following inequalities for all  $k \in \mathbb{N}$ :

$$E_{\xi_{k}}[F(w_{k+1})] - F(w_{k}) \leq -\mu\alpha_{k} \|\nabla F(w_{k})\|_{2}^{2} + \frac{\alpha_{k}^{2}L}{2} E_{\xi_{k}}[\|g(w_{k},\xi_{k})\|_{2}^{2}]$$
(118a)

$$\leq -\left(\mu\alpha_{k}-\frac{\alpha_{k}^{2}L}{2}M_{G}\right)\|\nabla F(w_{k})\|_{2}^{2}+\frac{\alpha_{k}^{2}L}{2}M.$$
(118b)

The most benign setting for analyzing the SG method is in the context of minimizing a strongly convex objective function. We formalize a strong convexity assumption as the following.

#### Assumption (Strong convexity)

The objective function  $F : \mathbb{R}^d \to \mathbb{R}$  is strongly convex in that there exists a constant c > 0 such that

$$F(\bar{w}) \geq F(w) + \nabla F(w)^{\top}(\bar{w} - w) + \frac{c}{2} \|\bar{w} - w\|_{2}^{2}, \ \forall (\bar{w}, w) \in \mathbb{R}^{d} \times \mathbb{R}^{d}$$
(119)

Hence, F has a unique minimizer, denoted as  $w^* \in \mathbb{R}^d$  with  $F_* := F(w^*)$ .

A useful fact is that, under the above assumption, we can bound the optimality gap at a given point:

$$F(w) - F_* \leqslant \frac{1}{2c} \|\nabla F(w)\|_2^2, \quad \forall w \in \mathbb{R}^d.$$
(120)

Noting that

$$\begin{array}{rcl} F(w) - F_* & \leqslant & -\nabla F(w)^\top (w^* - w) - \frac{c}{2} \|w^* - w\|_2^2 \\ & = & -\|\sqrt{\frac{1}{2c}} \nabla F(w) + \sqrt{\frac{c}{2}} (w^* - w)\|_2^2 + \frac{1}{2c} \|\nabla F(w)\|_2^2 \\ & \leqslant & \frac{1}{2c} \|\nabla F(w)\|_2^2 \,. \end{array}$$

We now state our first convergence theorem for SG.

- We use *E*[·] to denote an expected value taken with respect to the joint distribution of all random variables.
- For example, since w<sub>k</sub> is determined by {ξ<sub>1</sub>, ξ<sub>2</sub>,...,ξ<sub>k-1</sub>}, the total expectation of F(w<sub>k</sub>) for any k ∈ N can be taken as

$$E[F(w_k)] = E_{\xi_1}E_{\xi_2}\dots E_{\xi_{k-1}}[F(w_k)]$$

### 定理 (Strongly Convex Objective, Fixed Stepsize)

Under the above three assumptions (with  $F_{inf} = F_*$ ), suppose that the SG method is run with a fixed stepsize,  $\alpha_k = \bar{\alpha}$  for all  $k \in \mathbb{N}$ , satisfying

$$0 < \bar{\alpha} \leqslant \frac{\mu}{LM_G}.$$
 (121)

Then the expected optimality gap satisfies the following inequality for all  $k \in \mathbb{N}$ :

$$E[F(w_k) - F_*] \leq \frac{\bar{\alpha}LM}{2c\mu} + (1 - \bar{\alpha}c\mu)^k \left(F(w_1) - F_* - \frac{\bar{\alpha}LM}{2c\mu}\right)$$

$$\xrightarrow{k \to \infty} \frac{\bar{\alpha}LM}{2c\mu}.$$
(122)

This theorem illustrates the interplay between the stepsizes and bound on the variance of the stochastic directions.

- If the variance of  $g(w_k, \xi_k)$  is 0 or if noise is to decay with  $\|\nabla F(w_k)\|_2^2$ , then we can obtain linear convergence to the optimal value.
- On the other hand, when the gradient computation is noisy, a fixed and small enough stepsize can assure the expected objective values will converge linearly to a neighborhood of the optimal value, but the noise in the gradient estimates prevent further progress.

It's natural to ask if diminishing stepsizes will bring a better result.

### 定理 (Strongly Convex Objective, Diminishing Stepsizes)

Under the assumptions of Lipschitz-continuous objective gradients, first and second moment limits and strong convexity, suppose that SG method is run with a step size sequence such that, for all  $k \in \mathbb{N}$ ,

$$\alpha_k = \frac{\beta}{\gamma + k} \text{ for some } \beta > \frac{1}{c\mu} \text{ and } \gamma > 0 \text{ such that } \alpha_1 \leqslant \frac{\mu}{LM_G}.$$
(123)

Then, for all  $k \in \mathbb{N}$ , the expected optimality gap satisfies

$$E\left[F(w_k)-F_*\right] \leqslant \frac{\nu}{\gamma+k},\tag{124}$$

where

$$\nu := \max\left\{\frac{\beta^2 LM}{2(\beta c\mu - 1)}, (\gamma + 1)(F(w_1) - F_*)\right\}.$$
 (125)

Many important machine learning models lead to nonconvex optimization problems. Analyzing the SG method when minimizing nonconvex objectives is more challenging since such functions may possess multiple local minima and other stationary points.

Still, one can provide meaningful guarantees for the SG method in nonconvex settings.

While one cannot bound the expected optimality gap as in the convex case, we can bound the average norm of the gradient of the objective function observed on  $\{w_k\}$  visited during the first K iterations.

# SG for General Objectives

#### 定理 (Nonconvex Objective, Fixed Stepsize)

Under the assumptions of Lipschitz-continuous objective gradients (113) and first and second moment limits (115)-(116), suppose the SG method is run with a fixed stepsize  $\alpha_k = \bar{\alpha}$  satisfying

$$0 < \bar{\alpha} \leqslant \frac{\mu}{LM_G}.$$
 (126)

Then the expected sum-of-squares and average-squared gradients of F corresponding to the SG iterates satisfy the following inequalities for all  $K \in \mathbb{N}$ :

$$E\left[\sum_{k=1}^{K} \|\nabla F(w_{k})\|_{2}^{2}\right] \leq \frac{K\bar{\alpha}LM}{\mu} + \frac{2(F(w_{1}) - F_{inf})}{\mu\bar{\alpha}} \quad (127a)$$
  
and therefore  $E\left[\frac{1}{K}\sum_{k=1}^{K} \|\nabla F(w_{k})\|_{2}^{2}\right] \leq \frac{\bar{\alpha}LM}{\mu} + \frac{2(F(w_{1}) - F_{inf})}{K\mu\bar{\alpha}} \quad (127b)$   
 $\xrightarrow{K \to \infty} \frac{\bar{\alpha}LM}{\mu}.$ 

YZW (USTC)

#### Proof.

Taking the total expectation of (118b) and from (126),

$$\mathbb{E}\left[F(w_{k+1})\right] - \mathbb{E}\left[F(w_k)\right] \leqslant -\left(\mu - \frac{\bar{\alpha}LM_G}{2}\right)\bar{\alpha}\mathbb{E}\left[\left\|\nabla F(x_k)\right\|_2^2\right] + \frac{\bar{\alpha}^2LM}{2} \\ \leqslant -\frac{\mu\bar{\alpha}}{2}\mathbb{E}\left[\left\|\nabla F(w_k)\right\|_2^2\right] + \frac{\bar{\alpha}^2LM}{2}.$$

Summing both sides of this inequality for  $k \in \{1, ..., K\}$  and recalling (a) of the assumption on first and second moment limits gives

$$F_{inf} - F(w_1) \leqslant E\left[F(w_{K+1})\right] - F(w_1) \leqslant -\frac{\mu\bar{\alpha}}{2}\sum_{k=1}^{K} E\left[\left\|\nabla F(w_k)\right\|_2^2\right] + \frac{K\bar{\alpha}^2 LM}{2}.$$

Rearranging yields (127a), and dividing further by K yields (127b).

< ロト < 同ト < 三ト < 三ト

### 定理 (Nonconvex Objective, Diminishing Stepsize)

Under the assumptions of Lipschitz-continuous objective gradients (113) and first and second moment limits (115)-(116), suppose that the SG method is run with a stepsize sequence satisfying

$$\sum_{k=1}^{\infty} \alpha_k = \infty, \quad \sum_{k=1}^{\infty} \alpha_k^2 < \infty.$$
(128)

Then

$$\liminf_{k\to\infty} \left( E\left[ \|\nabla F(w_k)\|_2^2 \right] \right) = 0.$$
(129)

While not the strongest result in this context, this theorem is perhaps the easiest to interpret and remember. The proof of this theorem follows based on the stronger results given in the next theorem.

### 定理 (Nonconvex Objective, Diminishing Stepsize)

Under the assumptions of Lipschitz-continuous objective gradients (113) and first and second moment limits (115)-(116), suppose that the SG method is run with a stepsize sequence satisfying (128). Then, with  $A_{K} := \sum_{k=1}^{K} \alpha_{k}$ ,

$$E\left[\sum_{k=1}^{K} \alpha_{k} \|\nabla F(w_{k})\|_{2}^{2}\right] < \infty$$
(130a)  
and therefore  $E\left[\frac{1}{A_{K}}\sum_{k=1}^{K} \alpha_{k} \|\nabla F(w_{k})\|_{2}^{2}\right] \xrightarrow{K \to \infty} 0.$  (130b)

YZW (USTC)

### 推论

Suppose the conditions of the last theorem hold. For any  $K \in \mathbb{N}$ , let  $k(K) \in \{1, ..., K\}$  represent a random index chosen with probabilities proportional to  $\{\alpha_k\}_{k=1}^K$ . Then  $\|\nabla F(w_{k(K)})\|_2 \xrightarrow{K \to \infty} 0$  in probability.

### 推论

Under the conditions of the last theorem, if we further assume that the objective function F is twice differentiable, and that the mapping  $w \mapsto \|\nabla F(w)\|_2^2$  has Lipschitz-continuous derivatives, then

$$\lim_{k\to\infty} E\left[\|\nabla F(w_k)\|_2^2\right] = 0.$$

ヘロト ヘ団ト ヘヨト

SG suffers from the adverse effect of noisy gradient estimates. To address this limitation, methods endowed with *noise reduction* capabilities have been developed.



Figure: View of the schematic with a focus on noise reduction methods.

YZW (USTC)

< □ > < □ > < □ > < □ > < □ > < □ >

The first two classes of methods achieve noise reduction in a manner that allows them to possess a linear rate of convergence to the optimal value using a fixed stepsize. The third class of methods employing a stepsize sequence of order  $O(1/\sqrt{k})$  rather than O(1/k).

- **Dynamic sampling methods** achieve noise reduction by gradually increasing the mini-batch size used in the gradient computation.
- **Gradient aggregation methods** improve the quality of the search directions by storing gradient estimates in previous iterations, updating one (or some) of these estimates in each iteration, and defining the search direction as a weighted average of these estimates.
- Iterate averaging methods accomplish noise reduction by maintaining an average of iterates computed during the optimization process.

Recall the first lemma in this section

$$\begin{split} E_{\xi_k}\left[F(w_{k+1})\right] - F(w_k) &\leqslant -\alpha_k \nabla F(w_k)^\top E_{\xi_k}\left[g(w_k,\xi_k)\right] \\ &+ \frac{\alpha_k^2 L}{2} E_{\xi_k}\left[\left\|g(w_k,\xi_k)\right\|_2^2\right]. \end{split}$$

If we are able to decrease  $E_{\xi_k}[||g(w_k, \xi_k)||_2^2]$  fast enough, then the noise will not prevent the convergence.

We'll show that the sequence of expected optimality gaps vanishes at a linear rate as long as the variance of the stochastic vectors, denoted by  $Var_{\xi_k}[g(w_k, \xi_k)]$ , decreases geometrically.

#### 定理 (Strongly Convex Objective, Noise Reduction)

Under the assumptions of Lipschitz-continuous objective gradients and first and second moment limits and strong convexity, but with (116) refined to the existence of constants  $M \ge 0$  and  $\zeta \in (0, 1)$  such that

$$Var_{\xi_k}[g(w_k,\xi_k)] \leqslant M\zeta^{k-1}, \ \forall k \in \mathbb{N}.$$
 (131)

In addition, suppose that the SG method is run with a fixed stepsize,  $\alpha_k = \bar{\alpha}$  satisfying

$$0 < \bar{\alpha} \leqslant \min\left\{\frac{\mu}{L\mu_{G}^{2}}, \frac{1}{c\mu}\right\}.$$
(132)

# Dynamic Sampling Methods

定理 (Strongly Convex Objective, Noise Reduction)

Then the expected optimality gap satisfies

$$E\left[F(w_k) - F_*\right] \leqslant \omega \rho^{k-1},\tag{133}$$

where

$$\omega := \max\{\frac{\bar{\alpha}LM}{c\mu}, F(w_1) - F_*\} \text{ and } \rho := \max\{1 - \frac{\bar{\alpha}c\mu}{2}, \zeta\} < 1.$$

The restriction on the stepsize  $\bar{\alpha}$  is not unrealistic in practical situations, considering the typical magnitudes of the constants  $\mu$ , L,  $\mu_G$  and c.

Now a natural question is how to design efficient optimization methods for attaining the critical bound (131) on the variance of the stochastic directions.

Consider the iteration

$$w_{k+1} \leftarrow w_k - \bar{\alpha}g(w_k, \xi_k), \qquad (134)$$

where the stochastic directions are computed for some au>1 as

$$g(w_k,\xi_k) := \frac{1}{n_k} \sum_{i \in \mathcal{S}_k} \nabla f(w_k;\xi_{k,i}) \text{ with } n_k := |\mathcal{S}_k| = \lceil \tau^{k-1} \rceil.$$
(135)

That is, a mini-batch SG iteration with a fixed stepsize in which the mini-batch size increases geometrically as a function of the iteration counter k.

## Dynamic Sampling Methods

If we assume that each stochastic gradient  $\nabla f(w_k; \xi_{k,i})$  has an expectation equal to the true gradient  $\nabla F(w_k)$ , then (115) holds with  $\mu_G = \mu = 1$ . If, in addition, the variance of each such stochastic gradient is equal and is bounded by  $M \ge 0$ , then for arbitrary  $i \in S_k$  we have

$$Var_{\xi_k}\left[g(w_k,\xi_k)\right] \leqslant \frac{Var_{\xi_k}\left[\nabla f(w_k;\xi_{k,i})\right]}{n_k} \leqslant \frac{M}{n_k}.$$
 (136)

This bound combined with the rate of increase in  $n_k$  given in (135) yields (131). We state these formally as the following corollary.

### 推论

Let  $\{w_k\}$  be the iterates generated by (134)-(135) with  $E_{\xi_{k,i}}[\nabla f(w_k;\xi_{k,i})] = \nabla F(w_k), \forall k \in \mathbb{N}, i \in S_k$ . Then, the variance condition (131) is satisfied and if all other assumptions of the theorem of noise reduction for strongly convex objective holds, then the expected optimality gap vanishes linearly in the sense of (133).

YZW (USTC)
But, comparing to classical SG approach, is it meaningful to describe a method as linearly convergent if the per-iteration cost increases without bound?

To address this question, let's estimate the number of evaluations of the individual gradients  $\nabla f(w_k, \xi_{k,i})$  required to compute an  $\epsilon$ -optimal solution, i.e., to achieve

$$E\left[F(w_k) - F_*\right] \leqslant \epsilon. \tag{137}$$

(日)

As previously mentioned, the number of stochastic gradient evaluations required by the SG method to guarantee (137) is  $\mathcal{O}(\epsilon^{-1})$ .

#### 定理

Suppose the dynamic sampling SG method (134)-(135) is run with a stepsize  $\bar{\alpha}$  satisfying (132) and some  $\tau \in (1, (1 - \frac{\bar{\alpha}c\mu}{2})^{-1}]$ . In addition, suppose that the three assumptions hold. Then the total number of evaluations of a stochastic gradient of the form  $\nabla f(w_k, \xi_{k,i})$  required to obtain (137) is  $\mathcal{O}(\epsilon^{-1})$ .

Rather than compute increasingly more *new* stochastic gradient information in each iteration, *gradient aggregation methods* achieve a lower variance by *reusing* and/or *revising* previously computed information.

If the current iterate has not been displaced too far from previous iterates, then stochastic gradient information from previous iterates may still be useful.

The first method we consider is composed of outer and inner iterations.

At each step of outer iteration, an iterate  $w_k$  is available at which the algorithm computes a batch gradient  $\nabla R_n(w_k) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(w_k)$ .

Then, after initializing  $\tilde{w}_1 \leftarrow w_k$ , *m* inner iterations indexed by *j* are performed:

$$\tilde{w}_{j+1} \leftarrow \tilde{w}_j - \alpha \tilde{g}_j$$

where

$$\tilde{g}_j \leftarrow \nabla f_{i_j}(\tilde{w}_j) - \left(\nabla f_{i_j}(w_k) - \nabla R_n(w_k)\right)$$
(138)

and  $i_j \in \{1, \ldots, n\}$  is chosen at random.

Interpretation:

Since  $E_{i_j} [\nabla f_{i_j}(w_k)] = \nabla R_n(w_k)$ , we can view  $\nabla f_{i_j}(w_k) - \nabla R_n(w_k)$  as the bias in the gradient estimate  $\nabla f_{i_j}(w_k)$ . Thus the stochastic gradient  $\nabla f_{i_j}(\tilde{w}_j)$  evaluated at the current inner iterate  $\tilde{w}_j$  is corrected based on a perceived bias.

Overall,  $\tilde{g}_j$  represents an unbiased estimator of  $\nabla R_n(\tilde{w}_j)$ , with a smaller variance than simply choosing  $\nabla f_{i_j}(\tilde{w}_j)$  (as in simple SG). This is the reason why the method is referred to as the *stochastic variance* reduced gradient (SVRG) method.

# Gradient Aggregation Methods

Algorithm 2 SVRG Methods for Minimizing an Empirical Risk  $R_n$ 

- 1: Choose an initial iterate  $w_1 \in \mathbb{R}^d$ , stepsize  $\alpha > 0$ , positive integer m.
- 2: for k = 1, 2, ... do
- 3: Compute the batch gradient  $\nabla R_n(w_k)$ .
- 4: Initialize  $\tilde{w}_1 \leftarrow w_k$ .
- 5: **for** j = 1, ..., m **do**
- 6: Choose  $i_j$  uniformly from  $\{1, \ldots, n\}$ .

7: 
$$\tilde{g}_j \leftarrow \nabla f_{i_j}(\tilde{w}_j) - (\nabla f_{i_j}(w_k) - \nabla R_n(w_k)).$$

8: Set 
$$\tilde{w}_{j+1} \leftarrow \tilde{w}_j - \alpha \tilde{g}_j$$
.

9: end for

10: Option (a): Set 
$$w_{k+1} = \tilde{w}_{m+1}$$

- 11: Option (b): Set  $w_{k+1} = \frac{1}{m} \sum_{j=1}^{m} \tilde{w}_{j+1}$
- 12: Option (c): Choose j uniformly from  $\{1, \ldots, m\}$  and set  $w_{k+1} = \tilde{w}_{j+1}$ .
- 13: end for

For both options (b) and (c), it can achieve a linear rate of convergence when  $R_n$  is strongly convex.

More precisely, if the stepsize  $\alpha$  and the length of the inner cycle m are chosen so that

$$\rho := \frac{1}{1 - 2\alpha L} \left( \frac{1}{mc\alpha} + 2L\alpha \right) < 1,$$

then, given that the algorithm has reached  $w_k$ , one obtains

$$E_{i_j}\left[R_n(w_{k+1})-R_n(w_*)\right] \leqslant \rho E_{i_j}\left[R_n(w_{\kappa})-R_n(w_*)\right].$$

Each step (of outer iteration) requires 2m + n evaluations of component gradients, which is much more expensive than one in SG, and in fact is comparable to a full gradient iteration.

< 日 > < 同 > < 回 > < 回 > .

The second method does not include inner loop nor does it compute batch gradients (except possibly at the initial point).

Instead, in each iteration, it computes a stochastic vector  $g_k$  as the average of stochastic gradients evaluated at previous iterates.

Specifically, in iteration k, the method will have stored  $\nabla f_i(w_{[i]})$  for all  $i \in \{1, ..., n\}$  where  $w_{[i]}$  represents the latest iterate at which  $\nabla f_i$  was evaluated. An integer  $j \in \{1, ..., n\}$  is then chosen at random and the stochastic vector is set by

$$g_k \leftarrow \nabla f_j(w_k) - \nabla f_j(w_{[j]}) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(w_{[i]}).$$
(139)

Taking the expectation of  $g_k$  w.r.t. all choices of  $j \in \{1, ..., n\}$ , we have  $E[g_k] = \nabla R_n(w_k)$ . Thus the gradient estimates is unbiased with variances that are expected to be less than the stochastic gradients in a basic SG.

Algorithm 3 SAGA Methods for Minimizing an Empirical Risk  $R_n$ 

- 1: Choose an initial iterate  $w_1 \in \mathbb{R}^d$  and stepsize  $\alpha > 0$ .
- 2: for k = 1, 2, ... do
- 3: Compute  $\nabla f_i(w_1)$ .
- 4: Store  $\nabla f_i(w_{[i]}) \leftarrow \nabla f_i(w_1)$ .
- 5: end for
- 6: for k = 1, 2, ... do
- 7: Choose j uniformly in  $\{1, \ldots, n\}$ .
- 8: Compute  $\nabla f_j(w_k)$ .
- 9: Set  $g_k \leftarrow \nabla f_j(w_k) \nabla f_j(w_{[j]}) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(w_{[i]})$ .
- 10: Store  $\nabla f_j(w_{[j]}) \leftarrow \nabla f_j(w_k)$ .
- 11: Set  $w_{k+1} \leftarrow w_k \alpha g_k$ .

12: end for

(4回) (4 回) (4 回)

Beyond its initialization phase, the per-iteration cost of it is the same as in a basic SG method. However, it can achieve a linear rate of convergence when  $R_n$  is strongly convex. With  $\alpha = 1/(2(cn + L))$ , we have

$$E\left[\|w_k - w_*\|_2^2\right] \leqslant \left(1 - \frac{c}{2(cn+L)}\right)^k \left(\|w_1 - w_*\|_2^2 + \frac{nD}{cn+L}\right)$$

where  $D := R_n(w_1) - R_n(w_*) - \nabla R_n(w_*)^\top (w_1 - w_*).$ 

Alternative initialization techniques could be used in practice. For example, one could perform one epoch of simple SG, or assimilate iterates one-by-one and compute  $g_k$  only using the gradients available up to that point.

One important drawback of Algorithm 3 is the need to store *n* stochastic gradient vectors. Note, however, that if the component functions are of the form  $f_i(w_k) = \hat{f}(x_i^\top w_k)$ , then

$$abla f_i(w_k) = \hat{f}'(x_i^\top w_k) x_i.$$

That is, when the feature vectors  $\{x_i\}$  are already available in storage, one need only store the scalar  $\hat{f}'(x_i^{\top}w_k)$  to construct  $\nabla f_i(w_k)$  at a later iteration. This occurs in logistic and least squares regression.

Although the gradient aggregation methods above enjoy a faster rate of convergence than SG, they should not be regarded as clearly superior to SG.

Following similar analysis as before, the computing time for SG can be shown to be  $\mathcal{T}(n,\epsilon) \sim \kappa^2/\epsilon$  with  $\kappa := L/c$ . On the other hand, the computing times for SVRG and SAGA are  $\mathcal{T}(n,\epsilon) \sim (n+\kappa)\log(1/\epsilon)$ .

For very large n, gradient aggregation methods are comparable to batch algorithms and therefore cannot beat SG in this regime.

SG generates noisy iterate sequences that tend to oscillate around minimizers. Hence, a natural idea is to compute a corresponding sequence of *iterate averages* that would automatically possess less noisy behavior.

Specifically, for minimizing a continuously differentiable F with unbiased gradient estimates, it employs the iteration

$$w_{k+1} \leftarrow w_k - lpha_k g(w_k, \xi_k)$$
  
and  $ilde w_{k+1} \leftarrow rac{1}{k+1} \sum_{j=1}^{k+1} w_j.$  (140)

However, convergence properties better than SG of this method is elusive when using classical stepsize sequences that diminish with a rate of  $\mathcal{O}(1/k)$ .

An idea is to employ the iteration (140) but with stepsizes diminishing at a slower rate of  $\mathcal{O}(1/(k^a))$  for some  $a \in (\frac{1}{2}, 1)$ . When minimizing strongly convex objectives, it follows from this choice that

$$E[\|w_k - w_*\|_2^2] = \mathcal{O}(1/(k^a))$$
 while  $E[\|\tilde{w}_k - w_*\|_2^2] = \mathcal{O}(1/k)$ .

Besides reducing the noise in the stochastic directions, another manner to move beyond classical SG is to address the adverse effects of high nonlinearity and ill-conditioning of the objective function through the use of second-order information.

Deterministic methods are known to benefit from the use of second-order information, e.g., Newton's method achieves a locally quadratic convergence.

We start by considering a *Hessian-free Newton* method that employs exact second-order information in a judicious manner that exploits the stochastic nature of the objective function.

Then we describe methods that attempt to mimic the behavior of a Newton algorithm through first-order information computed over sequences of iterates, including *quasi-Newton*, *Gauss-Newton* and related algorithms that employ only diagonal re-scalings.

Finally we will sketch the natural gradient method.

# Second-Order Methods



We use double-sided arrows for the methods that can be effective throughout the spectrum between the stochastic and batch regimes. Single-sides arrows are used for those methods that are effective only with at least a moderate batch size in the stochastic gradient estimates.

A D F A B F A B F A B F

When minimizing a twice-continuously differentiable F, a Newton iteration is

$$w_{k+1} \leftarrow w_k + \alpha_k s_k \tag{141a}$$

where 
$$\nabla^2 F(w_k) s_k = -\nabla F(w_k)$$
. (141b)

This iteration demands much in terms of computation and storage. However, we can instead only solve (141b) inexactly through an iterative approach such as the conjugate gradient (CG) method.

By ensuring that the linear solves are accurate enough, such an *inexact Newton-CG* method can enjoy a superlinear convergence.

CG applied to (141b) does not require access to the Hessian itself, only Hessian-vector products. Such a method may be called *Hessian-free*.

## Subsampled Hessian-Free Newton Methods

In inexact Newton methods, the Hessian matrix need not be as accurate as the gradient to yield an effective iteration. It means that the iteration is more tolerant to noise in the Hessian estimate than it is to noise in the gradient estimate.

We employ a smaller sample for defining the Hessian than for the stochastic gradient estimate. Let the stochastic gradient estimate be

$$abla f_{\mathcal{S}_k}(w_k,\xi_k) = rac{1}{|\mathcal{S}_k|} \sum_{i \in \mathcal{S}_k} 
abla f(w_k,\xi_{k,i})$$

and let the stochastic Hessian estimate be

$$\nabla^2 f_{\mathcal{S}_k^H}(w_k, \xi_k^H) = \frac{1}{|\mathcal{S}_k^H|} \sum_{i \in \mathcal{S}_k^H} \nabla^2 f(w_k, \xi_{k,i})$$
(142)

where  $\mathcal{S}_k^H \subseteq \mathcal{S}_k$ .

YZW (USTC)

- If the subsample size  $|S_k^H|$  small enough, then the cost of each product involving the Hessian approximation can be reduced significantly, thus reducing the cost of each CG iteration.
- On the other hand, one should choose  $|S_k^H|$  large enough so that the curvature information captured through the Hessian-vector products is productive.

### Subsampled Hessian-Free Newton Methods

Algorithm 4 Subsampled Hessian-Free Inexact Newton Method

- 1: Choose an initial iterate  $w_1$ .
- 2: Choose constants  $\rho \in (0,1), \eta \in (0,1)$ , and  $\max_{cg} \in \mathbb{N}$ .
- 3: for k = 1, 2, ... do
- 4: Generate a realizations of  $\xi_k$  and  $\xi_k^H$  corresponding to  $\mathcal{S}_k^H \subseteq \mathcal{S}_k$ .
- 5: Compute  $s_k$  by applying Hessian-free CG to solve

$$\nabla^2 f_{\mathcal{S}_k^H}(w_k, \xi_k^H) s = -\nabla f_{\mathcal{S}_k}(w_k, \xi_k)$$
(143)

until  $\max_{cg}$  iterations have been performed or a trial solution yields

$$\|r_k\|_2 := \left\|\nabla^2 f_{\mathcal{S}_k}^H(w_k,\xi_k^H)s + \nabla f_{\mathcal{S}_k}(w_k,\xi_k)\right\|_2 \leq \rho \left\|\nabla f_{\mathcal{S}_k}(w_k,\xi_k)\right\|_2.$$

6: Set  $w_{k+1} \leftarrow w_k + \alpha_k s_k$ , where  $\alpha_k \in \{\gamma^0, \gamma^1, \gamma^2, \ldots\}$  is the largest element with

$$f_{\mathcal{S}_{k}}(w_{k+1},\xi_{k}) \leqslant f_{\mathcal{S}_{k}}(w_{k},\xi_{k}) + \eta \alpha_{k} \nabla f_{\mathcal{S}_{k}}(w_{k},\xi_{k})^{\top} s_{k}.$$
(144)

7: end for

(a)

If the algorithm were to operate in the stochastic regime of SG where  $|S_k|$  is small and gradients are very noisy, then it may be necessary to choose  $|S_k^H| > |S_k|$  so that Hessian approximations do not corrupt the step.

Therefore, the subsampled Hessian-free Newton method outlined here is only recommended when  $S_k$  is large.

When full gradients are always used, it's easy to establish the convergence of Algorithm 4 for minimizing a strongly convex empirical risk measure  $F = R_n$  with  $S_k^H = S_k = \{1, ..., n\}$ .

When the Hessians are subsampled, it has not been shown that the rate of convergence is faster than linear.

<ロ> <四> <四> <日> <日> <日> <日</p>

When Hessian-free Newton methods are applied for the solution of nonconvex problems, it's common to employ a *trust region* instead of a line search and to add an additional condition in Step 5 of Algorithm 4: terminate CG if a candidate solution  $s_k$  is a direction of negative curvature, i.e.,  $s_k^{\top} \nabla^2 f_{\mathcal{S}_k^H}(w_k; \xi_k^H) s_k < 0$ .

Instead of coping with indefiniteness, one can focus on strategies for ensuring positive (semi)definite Hessian approximations. One of the most attractive ways of doing this in the context of machine learning is to employ a (subsampled) Gauss-Newton approximation to the Hessian, which we will explain later.

The quasi-Newton iteration for minimizing a twice continuously differentiable function F has the form

$$w_{k+1} \leftarrow w_k - \alpha_k H_k \nabla F(w_k),$$
 (145)

where  $H_k$  is a approximation of  $(\nabla^2 F(w_k))^{-1}$ . The most popular quasi-Newton scheme is BFGS.

In BFGS, the sequence  $\{H_k\}$  is updated dynamically, without the need for second-order derivative computations nor any linear system solves. It enjoys a local superlinear convergence with only first-order information.

But  $H_k$  is often a dense matrix, even when the exact Hessian is sparse, restricting its use to small and midsize problems. A common solution for this is to employ a *limited memory scheme*, leading to a method such as L-BFGS. In this case,  $H_k$  need not be formed explicitly.

・ロト ・雪 ト ・ヨ ト ・

Now we consider the iterations taking the form

$$w_{k+1} \leftarrow w_k - \alpha_k H_k g(w_k, \xi_k). \tag{146}$$

Since we are interested in large-scale problems, we assume that (146) implements an L-BFGS scheme. A number of questions arise when considering (146), and we list them now with some proposed solutions:

**Theoretical Limitations** The convergence rate of a stochastic iteration such as (146) cannot be faster than sublinear. Since SG also has a sublinear rate of convergence, what benefit could come from incorporating  $H_k$  in (146)?

Benefit: The constant that appears in the sublinear rate.

For SG, the constant depends on the conditioning of  $\{\nabla^2 F(w_k)\}$ . This is typical of first-order methods. In contrast, if the sequence of Hessian approximations in (146) satisfies  $\{H_k\} \rightarrow \nabla^2 F(w_*)^{-1}$ , then the constant is independent of the conditioning of the Hessian.

・ロト ・四ト ・ヨト ・ヨト ・ヨ

**Additional Per-Iteration Costs** The product  $H_kg(w_k, \xi_k)$  requires 4md operations where m is the memory in the L-BFGS updating scheme. Assuming the cost of evaluating  $g(w_k, \xi_k)$  is exactly d operations (using only one sample) and m is set to the typical value of 5, then the stochastic quasi-Newton method is 20 times more expensive than SG. Can we offset this additional per-iteration cost?

When employing mini-batch gradient estimates, the additional cost of the iteration (146) is only marginal. The use of mini-batches may be considered essential. Mini-batch should not be less than, say, 20 or 50, and mini-batches of size 256 are common in practice.

< 日 > < 同 > < 回 > < 回 > .

**Conditioning of the Scaling Matrices** Updating  $H_k$  involves differences in gradient estimates computed in consecutive iterations.  $\{g(w_k, \xi_k)\}$  are noisy estimates of  $\{\nabla F(w_k)\}$ , which can cause the updating process to yield poor curvature estimates. How could such effects be avoided in the stochastic regime?

One possibility is to employ the same sample when computing gradient differences. An alternative approach is to *decouple* the step computation and the Hessian update.

Replacing deterministic gradients with stochastic gradients, we have

$$s_k := w_{k+1} - w_k$$
 and  $v_k := 
abla f_{\mathcal{S}_k}(w_{k+1}, \xi_k) - 
abla f_{\mathcal{S}_k}(w_k, \xi_k).$  (147)

and  $H_k$  is defined recursively by

$$H_{k+1} \leftarrow \left(I - \frac{v_k s_k^\top}{s_k^\top v_k}\right)^\top H_k \left(I - \frac{v_k s_k^\top}{s_k^\top v_k}\right) + \frac{s_k s_k^\top}{s_k^\top v_k}.$$

Note that the use of the same realization  $\xi_k$  in the two gradient estimates, in order to address the issues related to noise mentioned above.

A worrisome feature is that updating the inverse Hessian approximation with *every* step may not be warranted and could easily represent a poor approximation of the action of the true Hessian of F.

Here's an alternative strategy for this issue. Since  $\nabla f(w_{k+1}) - \nabla F(w_k) \approx \nabla^2 F(w_k)(w_{k+1} - w_k)$ , we can define

$$\mathbf{v}_k := \nabla^2 f_{\mathcal{S}_k^H}(\mathbf{w}_k, \boldsymbol{\xi}_k^H) \mathbf{s}_k, \tag{148}$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

where  $\nabla^2 f_{\mathcal{S}_k^H}(w_k; \xi_k^H)$  is a subsampled Hessian and  $|\mathcal{S}_k^H|$  is large enough to provide useful curvature information.

When  $|S_k^H|$  is much larger than  $|S_k|$ , the computation of  $v_k$  can be performed only after a sequence of iterations, to amortize the cost of quasi-Newton updating.

YZW (USTC)

This leads to the idea of decoupling the step computation from the quasi-Newton update. This approach, which we refer to as SQN, performs a sequence of iterations of (146) with  $H_k$  fixed, then computes a new displacement pair  $(s_k, v_k)$  with  $s_k$  defined as in (147) and  $v_k$  set using one of the strategies outlined above.

To formalize all of these alternatives, we state the general stochastic quasi-Newton method presented as Algorithm 5.

# Stochastic Quasi-Newton Methods

#### Algorithm 5 Stochastic Quasi-Newton Framework

- 1: Choose an initial iterate  $w_1$  and initialize  $\mathcal{P} \leftarrow \emptyset$ .
- 2: Choose a constant  $m \in \mathbb{N}$ .
- 3: Choose a stepsize sequence  $\{\alpha_k\} \subset \mathbb{R}_{++}$ .
- 4: for k = 1, 2, ..., do
- 5: Generate realizations of  $\xi_k$  and  $\xi_k^H$  corresponding to  $\mathcal{S}_k^H \subseteq \mathcal{S}_k$
- 6: Compute  $\hat{s}_k = H_k g(w_k, \xi_k)$  using the two-loop recursion based on the set  $\mathcal{P}$ .
- 7: Set  $s_k \leftarrow -\alpha_k \hat{s}_k$ .
- 8: Set  $w_{k+1} \leftarrow w_k + s_k$ .
- 9: if update pairs then
- 10: Compute  $s_k$  and  $v_k$  (based on the sample  $\mathcal{S}_k^H$ ).
- 11: Add the new displacement pair  $(s_k, v_k)$  to  $\mathcal{P}$ .
- 12: If  $|\mathcal{P}| > m$ , then remove eldest pair from  $\mathcal{P}$ .
- 13: end if
- 14: end for

The Gauss-Newton method is a classical approach for nonlinear least squares. It constructs an approximation to the Hessian using only first-order information, and this approximation is guaranteed to be positive semidefinite, even when the full Hessian itself may be indefinite.

#### Gauss-Newton Methods

Given an input-output pair  $(x_{\xi}, y_{\xi})$ , the loss incurred by a parameter vector w is measured via a squared norm discrepancy between  $h(x_{\xi}, w) \in \mathbb{R}^d$  and  $y \in \mathbb{R}^d$ :

$$f(w,\xi) = \ell(h(x_{\xi},w),y_{\xi}) = \frac{1}{2} \|h(x_{\xi},w) - y_{\xi}\|_{2}^{2}.$$

Let  $J_h(\cdot,\xi)$  represent the Jacobian of  $h(x_{\xi},\cdot)$  with respect to w. The affine approximation of  $h(x_{\xi}, w)$  is

$$h(x_{\xi},w) \approx h(x_{\xi},w_k) + J_h(w_k,\xi)(w-w_k),$$

which leads to

$$\begin{split} f(w,\xi) \approx &\frac{1}{2} \|h(x_{\xi},w_{k}) + J_{h}(w_{k},\xi)(w-w_{k}) - y_{\xi}\|_{2}^{2} \\ = &\frac{1}{2} \|h(x_{\xi},w_{k}) - y_{\xi}\|_{2}^{2} + (h(x_{\xi},w_{k}) - y_{\xi})^{\top} J_{h}(w_{k},\xi)(w-w_{k}) \\ &+ &\frac{1}{2} (w-w_{k})^{\top} J_{h}(w_{k},\xi)^{\top} J_{h}(w_{k},\xi)(w-w_{k}). \end{split}$$

YZW (USTC)

< □ > < □ > < □ > < □ > < □ > < □ >

It is similar to a second-oder Taylor series model, except that the terms involving the second derivatives of h with respect to w have been dropped, and the remaining second-order terms are resulting from the positive curvature of the quadratic loss  $\ell$ .

This leads to replacing the subsample Hessian matrix by the Gauss-Newton matrix

$$G_{\mathcal{S}_{k}^{H}}(w_{k},\xi_{k}^{H}) = \frac{1}{|\mathcal{S}_{k}^{H}|} \sum_{i \in \mathcal{S}_{k}^{H}} J_{h}(w_{k},\xi_{k,i})^{\top} J_{h}(w_{k},\xi_{k,i}).$$
(149)
A challenge of Gauss-Newton method is that Gauss-Newton matrix is often singular or nearly singular. In practice, this is handled by regularizing it by adding to it a positive multiple of the identity matrix.

The computational cost of the Gauss-Newton method depends on the dimensionality of the prediction function. It should be remarked that in machine learning, computing the stochastic gradient vector  $\nabla f(w,\xi)$  does not usually require the explicit computation of all rows of the Jacobian matrix. And there are some new ways to solve a Gauss-Newton iterate at a low cost.

Consider a slightly more general situation in which loss between a prediction function h and output y is measured by an arbitrary convex loss function  $\ell(h, y)$ . Combining the affine approximation of the prediction function  $h(x_{\xi}, w)$  with a second order Taylor expansion of the loss function  $\ell$  leads to the generalized Gauss-Newton matrix

$$G_{\mathcal{S}_{k}^{H}}(w_{k},\xi_{k}^{H}) = \frac{1}{|\mathcal{S}_{k}^{H}|} \sum_{i \in \mathcal{S}_{k}^{H}} J_{h}(w_{k},\xi_{k,i})^{\top} H_{\ell}(w_{k},\xi_{k,i}) J_{h}(w_{k},\xi_{k,i})$$
(150)

where  $H_{\ell}(w_k,\xi) = \frac{\partial^2 \ell}{\partial h^2}(h(x_{\xi}, w_k), y_{\xi})$  captures the curvature of the loss function  $\ell$ .

(日)

We have seen that the added per-iteration costs of second-order methods (such as L-BFGS) can be as little as 4*md* operations. A strategy to further reduce this multiplicative factor is to restrict attention to *diagonal* or *block-diagonal* scaling matrices.

The incorporation of a diagonal scaling matrix will only scale the individual search direction components. This can be efficiently achieved by multiplying the individual search direction components.

A D F A B F A B F A B F

A first family of algorithms directly computes the diagonal terms of the Hessian or Gauss-Newton matrix, then divides each coefficient of the stochastic gradient vector  $g(w_k, \xi_k)$  by the corresponding diagonal term.

For instance, each iteration of the proposed algorithm picks a training example, computes  $g(w_k, \xi_k)$ , updates a running estimate of the diagonal coefficients of the Gauss-Newton matrix by

$$[G_k]_i = (1 - \lambda)[G_{k-1}]_i + \lambda \left[J_h(w_k, \xi_k)^\top J_h(w_k, \xi_k)\right]_{ii} \text{ for some } 0 < \lambda < 1,$$

then performs the scaled stochastic weight update

$$[w_{k+1}]_i = [w_k]_i - \left(\frac{\alpha}{[G_k]_i + \mu}\right) [g(w_k, \xi_k)]_i.$$

The small regularization constant  $\mu > 0$  is introduced to deal with a singular or nearly singular Gauss-Newton matrix.

It's more enlightening to view such an algorithm as a scheme to periodically retune a first-order SG approach rather than as a complete second-order method.

Instead of explicitly computing the diagonal terms of the curvature matrix, one can follow the template of quasi-Newton method and directly estimate the diagonal  $[H_k]_i$  of the inverse Hessian using displacement pairs  $\{(s_k, v_k)\}$ .

For instance,  $[H_k]_i$  can be computed with the running average

$$[H_{k+1}]_i = (1-\lambda)[H_k]_i + \lambda \operatorname{Proj}\left(rac{[s_k]_i}{[v_k]_i}
ight),$$

where  $Proj(\cdot)$  represents a projection onto a predefined positive interval. But a direct application of (147) after a parameter update introduces a correlated noise that ruins the curvature estimate, which is hard to correct because of the chaotic behavior of the rescaling factors  $[H_k]_i$ .

YZW (USTC)

These problems can be addressed with a combination of two ideas.

First, estimate the diagonal of the Hessian instead of its inverse.

Second, ensure the effective stepsizes are monotonically decreasing by replacing the running average by the sum

$$[G_{k+1}]_i = [G_k]_i + \operatorname{Proj}\left(rac{[v_k]_i}{[s_k]_i}
ight).$$

Keeping the curvature estimates in a fixed positive interval ensures the effective stepsizes decrease at the rate  $\mathcal{O}(\frac{1}{k})$ .

YZW (USTC)

## Natural Gradient Method

The essential idea of natural gradient method consists of formulating the gradient descent algorithm in the space of prediction functions rather than specific parameters. The actual computation of course takes place with respect to the parameters, but the algorithm will move the parameters more quickly along directions that have a small impact on the decision function.

The space  $\mathcal{H}$  of prediction functions is a family of densities  $h_w(x)$  parametrized by  $w \in \mathcal{W}$  and satisfying the normalization condition

$$\int h_w(x)dx = 1, \quad \forall w \in \mathcal{W}.$$

And we assume sufficient regularity, i.e.,

$$\forall t > 0, \quad \int \frac{\partial^t h_w(x)}{\partial w^t} dx = \frac{\partial^t}{\partial w^t} \int h_w(x) dx = \frac{\partial^t 1}{\partial w^t} = 0. \tag{151}$$

YZW (USTC)

To quantify how the density  $h_w$  changes when adding a small quantity  $\delta w$  to its paramter, we use the Kullback-Leibler (KL) divergence

$$D_{KL}(h_w \| h_{w+\delta w}) = E_{h_W} \left[ \log \left( \frac{h_w(x)}{h_{w+\delta w}(x)} \right) \right].$$
(152)

Approximating the divergence with a second-order Taylor expansion, we have

$$D_{\mathcal{K}L}(h_w \| h_{w+\delta w}) = E_{h_w} [\log(h_w(x)) - \log(h_{w+\delta w}(x))] \\\approx -\delta w^\top E_{h_w} \left[ \frac{\partial \log(h_w(x))}{\partial w} \right] - \frac{1}{2} \delta w^\top E_{h_w} \left[ \frac{\partial^2 \log(h_w(x))}{\partial w^2} \right] \delta w.$$

By (151),

$$D_{KL}(h_w \| h_{w+\delta w}) \approx \frac{1}{2} \delta w^\top G(w) \delta w.$$
(153)

Natural gradient method minimizes a functional  $F: h_w \in \mathcal{H} \mapsto F(h_w) = F(w) \in \mathbb{R}$ . A greedy strategy is

$$h_{w_{k+1}} = \underset{h \in \mathcal{H}}{\arg\min} F(h) \quad \text{s.t.} \quad D_{\mathcal{KL}}(h_{w_k} \| h) \leqslant \eta_k^2.$$
(154)

Use (153) we can reformulate it in terms of the parameters:

$$w_{k+1} = \underset{w \in \mathcal{W}}{\operatorname{arg\,min}} F(w) \quad \text{s.t.} \quad \frac{1}{2} (w - w_k)^\top G(w_k) (w - w_k) \leqslant \eta_k^2. \quad (155)$$

Lagrangian formulation is customarily used to handle this problem. Assuming  $\eta_k$  small, we can replace F(w) with  $F(w_k) + \nabla F(w_k)^\top (w - w_k)$ . These two choices lead to

$$w_{k+1} = \operatorname*{arg\,min}_{w\in\mathcal{W}} 
abla F(w_k)^{ op} (w-w_k) + rac{1}{2lpha_k} (w-w_k)^{ op} G(w_k) (w-w_k),$$

and the optimization of the right-hand side leads to the natural gradient iteration

$$w_{k+1} = w_k - \alpha_k G^{-1}(w_k) \nabla F(w_k).$$
 (156)

G(w) is a called Fisher information matrix, with expression

$$G(w) := -E_{h_w} \left[ \frac{\partial^2 \log(h_w(x))}{\partial w^2} \right]$$
  
=  $E_{h_w} \left[ \left( \frac{\partial \log(h_w(x))}{\partial w} \right) \left( \frac{\partial \log(h_w(x))}{\partial w} \right)^\top \right],$  (157)

where the latter equality follows from (151).

A sampled version of  $G(w_k)$  is

$$ilde{G}(w_k) = rac{1}{|\mathcal{S}_k|} \sum_{i \in \mathcal{S}_k} \left( rac{\partial \log(h_w(x_i))}{\partial w} \Big|_{w_k} 
ight) \left( rac{\partial \log(h_w(x_i))}{\partial w} \Big|_{w_k} 
ight)^ op$$

YZW (USTC)

With an initial point  $w_1 = w_0$ , scalar sequences  $\{\alpha_k\}$  and  $\{\beta_k\}$ , the iteration of gradient methods with momentum is

$$w_{k+1} \leftarrow w_k - \alpha_k \nabla F(w_k) + \beta_k (w_k - w_{k-1}).$$
(158)

The latter is referred to as the *momentum* term. It is named after the fact that it represents a discretization of a certain second-order ordinary differential equation with friction.

When  $\beta_k = 0$  for all  $k \in \mathbb{N}$ , it reduces to the steepest descent method.

When  $\alpha_k = \alpha$  and  $\beta_k = \beta$  for some constants  $\alpha > 0$  and  $\beta > 0$ , it is referred to as the heavy ball method, which yields a linear convergence with a superior rate compared to steepest descent with a fixed stepsize for certain functions.

< 日 > < 同 > < 回 > < 回 > .

Additional connection with (158) can be made when F is a strictly convex quadratic. If  $(\alpha_k, \beta_k)$  is chosen optimally in the sense that

$$(\alpha_k, \beta_k) = \underset{(\alpha, \beta)}{\arg\min} F(w_k - \alpha \nabla F(w_k) + \beta(w_k - w_{k-1})), \quad (159)$$

then (158) is exactly the linear conjugate gradient (CG) algorithm.

An alternative view of the heavy ball method is obtained by expanding (158) as:

$$w_{k+1} \leftarrow w_k - \alpha \sum_{j=1}^k \beta^{k-j} \nabla F(w_k);$$

thus, each step can be viewed as an exponentially average of past gradients.

イロト 不得 ト イヨト イヨト

Nesterov accelerated gradient method is similar to (158) but with its own unique properties. It involves the updates

$$\widetilde{w}_k \leftarrow w_k + \beta_k (w_k - w_{k-1})$$
  
and  $w_{k+1} \leftarrow \widetilde{w}_k - \alpha_k \nabla F(\widetilde{w}_k),$  (160)

which leads to the condensed form

$$w_{k+1} \leftarrow w_k - \alpha_k \nabla F(w_k + \beta_k(w_k - w_{k-1})) + \beta_k(w_k - w_{k-1}). \quad (161)$$

Compared with gradient method with momentum, it applies the momentum term first, then takes a steepest descent step at  $\tilde{w}_k$ .

When *F* is convex and continuously differentiable with a Lipschitz continuous gradient, with appropriately chosen  $\alpha_k = \alpha > 0$  for all  $k \in \mathbb{N}$  and  $\{\beta_k\} \nearrow 1$  leads to an *optimal* iteration complexity.

While the convergence rate of steepest descent method is  $\mathcal{O}(\frac{1}{k})$ , the iteration (161) converges with a rate  $\mathcal{O}(\frac{1}{k^2})$ , which is provably the best rate that can be achieved by a gradient method.

Unfortunately, no intuitive explanation as to how Nesterov's method achieves this optimal rate has been widely accepted.

ヘロト 人間 とくほとく ほとう

Coordinate descent (CD) methods operate to a single variable while all others are kept fixed, then other variables are updated similarly.

The CD method for minimizing  $F : \mathbb{R}^d \to \mathbb{R}$  is given by the iteration

$$w_{k+1} \leftarrow w_k - \alpha_k \nabla_{i_k} F(w_k) e_{i_k}, \qquad (162)$$

where  $\nabla_{i_k} F(w_k) := \frac{\partial F}{\partial w^{i_k}}(w_k)$ ,  $w^{i_k}$  represents the  $i_k$ -th element of the parameter vector, and  $e_{i_k}$  represents the  $i_k$ -th coordinate vector for some  $i_k \in \{1, \ldots, d\}$ .

(日)

Specific versions of the CD method are defined by the manner in which the sequences  $\{\alpha_k\}$  and  $\{i_k\}$  are chosen.

 $\{\alpha_k\}$ :

- Choose  $\alpha_k$  as the global minimizer of F from  $w_k$  along the  $i_k$ -th coordinate.
- Choose  $\alpha_k$  yielding a sufficient reduction in F from  $w_k$ .
- Compute  $\alpha_k$  as the minimizer of a quadratic model of F along the  $i_k$ -th coordinate direction. (so-called second-order CD methods)

 $\{i_k\}$ :

- Cycle through  $\{1, \ldots, d\}$ .
- Cycle through a random reordering of  $\{1, \ldots, d\}$ , with the indexes reordered after each set of *d* steps.
- Simply choose an index randomly with replacement in each iteration.

The latter two strategies for  $\{i_k\}$  have superior theoretical properties than the first strategy.

(I)

A CD method is not guaranteed to converge when applied to minimize any given continuously differentiable function. This is in contrast with the full gradient method, which guarantees convergence to stationarity even when the objective is nonconvex.

However, if the objective F is strongly convex, the CD method will not fail. The analysis is very simple when using a constant stepsize. Assume that  $\nabla F$  is coordinate-wise Lipschitz continuous in the sense that for all  $w \in \mathbb{R}^d, i \in \{1, \ldots, d\}$ , and  $\Delta w^i \in \mathbb{R}$ , there exists a constant  $L_i > 0$  such that

$$|\nabla_i F(w + \Delta w^i e_i) - \nabla_i F(w)| \leq L_i |\Delta w^i|.$$
(163)

And we define  $\hat{L} := \max_{i \in \{1,...,d\}} L_i$ .

#### 定理

Suppose that the objective function  $F : \mathbb{R}^d \to \mathbb{R}$  is continuously differentiable, strongly convex with constant c > 0, and has a gradient that is coordinate-wise Lipshcitz continuous with constants  $\{L_1, \ldots, L_d\}$ . In addition, suppose that  $\alpha_k = 1/\hat{L}$  and  $i_k$  is chosen independently and uniformly from  $\{1, \ldots, d\}$  for all  $k \in \mathbb{N}$ . Then for all  $k \in \mathbb{N}$ , the iteration (162) yields

$$E\left[F(w_{k+1})\right] - F_* \leqslant \left(1 - \frac{c}{d\hat{L}}\right)^k \left(F(w_1) - F_*\right). \tag{164}$$

A simple randomized CD method is linearly convergent with constant dependent on the parameter dimension d. If d coordinate updates can be performed at a cost similar to the evaluation of one full gradient, the method is competitive with a full gradient method both theoretically and in practice.

This kind of problems include those in which the objective function is

$$F(w) = \frac{1}{n} \sum_{j=1}^{n} \tilde{F}_{j}(x_{j}^{\top}w) + \sum_{i=1}^{d} \hat{F}_{i}(w^{i}), \qquad (165)$$

where  $\forall j \in \{1, \ldots, n\}$ ,  $\tilde{F}_j$  is continuously differentiable and dependent on the *sparse* data vector  $x_j$ , and  $\forall i \in \{1, \ldots, d\}$ ,  $\hat{F}_i$  is a regularization function (potentially nonsmooth).

For example, consider an objective function of the form

$$f(w) = \frac{1}{2} \|Xw - y\|_2^2 + \sum_{i=1}^d \hat{F}_i(w^i) \text{ with } X = [x_1 \dots x_n].$$

In this setting,

$$abla_{i_k} f(w_{k+1}) = x_{i_k}^\top r_{k+1} + \hat{F}'_{i_k}(w_{k+1}^{i_k}) \text{ with } r_{k+1} := Aw_{k+1} - b,$$

where, with  $w_{k+1} = w_k + \beta_k e_{i_k}$ , we have  $r_{k+1} = r_k + \beta_k x_{i_k}$ .

Since the residuals  $\{r_k\}$  can be updated with cost proportional to the number of nonzeros in  $x_{i_k}$ , call it  $nnz(x_{i_k})$ , the overall cost of computing the search direction in iteration k + 1 is also  $\mathcal{O}(nnz(x_{i_k}))$ . On the other hand, an evaluation of the entire gradient requires a cost of  $\mathcal{O}(\sum_{j=1}^n nnz(x_j))$ .

(日)

## Stochastic Dual Coordinate Ascent

Consider minimizing a convex objective function of the form (165) by maximizing its dual.

Defining the convex conjugate of  $\tilde{F}_j$  as  $\tilde{F}_j^*(u) := \max_w (w^\top u - \tilde{F}_j(w))$ when  $\hat{F}_i(\cdot) = \frac{\lambda}{2}(\cdot)^2$  for all  $i \in \{1, \ldots, d\}$  is given by

$$F_{dual}(v) = \frac{1}{n} \sum_{j=1}^{n} \left[ -\tilde{F}_{j}^{*}(-v_{j}) \right] - \frac{\lambda}{2} \left\| \frac{1}{\lambda n} \sum_{j=1}^{n} v_{j} x_{j} \right\|_{2}^{2}.$$

The stochastic dual coordinate ascent (SDCA) method applied to a function of this form has an iteration similar to (162), except that negative gradient steps are replaced by gradient steps.

When the algorithm terminates, the corresponding primal solution can be obtained as  $w \leftarrow \frac{1}{\lambda n} \sum_{j=1}^{n} v_j x_j$ .

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Consider a multicore system in which the parameter vector w is stored in shared memory.

Each core can then execute a CD iteration independently and in an asynchronous manner, where if d is large compared to the number of cores, then it is unlikely that two cores are attempting to update the same variable at the same time.

Each update is being made based on slightly stale information. However, convergence of the method can be proved, and improves when one can bound the degree of staleness of each update.

The discussion of structural risk minimization highlighted the key role played by regularization functions.

The optimization methods we have presented in this section are all applicable for objectives involving smooth regularizers, such as the squared  $\ell_2$ -norm. And we expand our investigation by considering optimization methods that handle the regularization as a distinct entity, in particular when the function is *nonsmooth*, for example,  $\ell_1$ -norm, which induces sparsity in the optimal solution vector.

For machine learning, sparsity can be seen as a form of *feature selection*.

We focuses on the nonsmooth optimization problem

$$\min_{w \in \mathbb{R}^d} \Phi(w) := F(w) + \lambda \Omega(w), \tag{166}$$

where  $F : \mathbb{R}^d \to \mathbb{R}$  includes the composition of a loss and prediction function,  $\lambda > 0$  is a regularization parameter, and  $\Omega : \mathbb{R}^d \to \mathbb{R}$  is a convex, nonsmooth regularization function.

Specifically, we pay special attention to methods for solving the problem

$$\min_{w \in \mathbb{R}^d} \phi(w) := F(w) + \lambda \|w\|_1.$$
(167)

イロト 不得 ト イヨト イヨト

For solving problem (166), the *proximal gradient* method represents a fundamental approach.

Given an iterate  $w_k$ , a generic proximal gradient iteration with  $\alpha_k > 0$  is given by

$$w_{k+1} \leftarrow \operatorname*{arg\,min}_{w \in \mathbb{R}^d} \left( F(w_k) + \nabla F(w_k)^\top (w - w_k) + \frac{1}{2\alpha_k} \|w - w_k\|_2^2 + \lambda \Omega(w) \right)$$
(168)

The term *proximal* refers to the presence of the third term in the minimization problem on the right-hand side, which encourage the new iterate to be close to  $w_k$ . If the last term were not present, then (168) exactly recovers the gradient method update  $w_{k+1} \leftarrow w_k - \alpha_k \nabla F(w_k)$ ; hence we refer to  $\alpha_k$  as the stepsize parameter.

#### 定理

Suppose that  $F : \mathbb{R}^d \to \mathbb{R}$  is continuously differentiable, strongly convex with constant c > 0, and has a gradient that is Lipschitz continuous with constant L > 0. In addition, suppose that  $\alpha_k = \alpha \in (0, 1/L)$  for all  $k \in \mathbb{N}$ . Then, for all  $k \in \mathbb{N}$ , the iteration (168) yields

$$\Phi(w_{k+1}) - \Phi(w_*) \leqslant (1 - \alpha c)^k (\Phi(w_1) - \Phi(w_*)),$$

where  $w_* \in \mathbb{R}^d$  is the unique global minimizer of  $\Phi$  in (166).

The proximal gradient iteration (168) is practical only when the proximal mapping

$$prox_{\lambda\Omega,lpha_k}( ilde{w}) := rgmin_{w\in\mathbb{R}^n} \left(\lambda\Omega(w) + rac{1}{2lpha_k}\|w - ilde{w}\|_2^2
ight)$$

can be computed efficiently. Situations when the proximal mapping is inexpensive to compute include when  $\Omega$  is the indicator function for a simple set, when it is the  $\ell_1$ -norm, or when it is separable.

A stochastic version of the proximal gradient method can be obtained by replacing  $\nabla F(w_k)$  in (168) by a stochastic approximation  $g(w_k, \xi_k)$ . The resulting method attains similar behavior as a stochastic gradient method.

For solving the  $\ell_1\text{-norm}$  regularized problem (167), the proximal gradient method is

$$w_{k+1} \leftarrow \operatorname*{arg\,min}_{w \in \mathbb{R}^d} \left( F(w_k) + \nabla F(w_k)^\top (w - w_k) + \frac{1}{2\alpha_k} \|w - w_k\|_2^2 + \lambda \|w\|_1 \right)$$
(169)

The solution can be written component-wise in closed form, with  $(\cdot)_+:=\max\{\cdot,0\},$  as

$$w_{k+1} \leftarrow \mathcal{T}_{\alpha_k \lambda}(w_k - \alpha_k \nabla F(w_k)), \text{ where } [\mathcal{T}_{\alpha_k \lambda}]_i = (|\tilde{w}_i| - \alpha_k \lambda)_+ \operatorname{sgn}(\tilde{w}_i).$$
(170)

 $\mathcal{T}_{\alpha_k\lambda}$  is referred to as the soft-thresholding operator, which leads to the name *iterative soft-thresholding algorithm* (ISTA). It is clear from (170) that the ISTA iteration induces sparsity in the iterates.

< ロ > < 同 > < 三 > < 三 > < 三 > <

# Bound-constrained Methods for $\ell_1$ -norm Regularized Problems

An equivalent *smooth* reformulation of problem (167) is easily derived, by writing w = u - v where u and v play the *positive part* and *negative part* of w respectively:

$$\min_{u,v)\in\mathbb{R}^d\times\mathbb{R}^d}\tilde{\phi}(u,v) \quad \text{s.t.} \quad (u,v) \ge 0 \tag{171}$$

where 
$$\tilde{\phi}(u, v) = F(u - v) + \lambda \sum_{i=1}^{d} (u_i + v_i)$$
.

The fundamental iteration for solving bound-constrained optimization problems is the *gradient projection* method. In the context of (171), the iteration reduces to

$$\begin{bmatrix} u_{k+1} \\ v_{k+1} \end{bmatrix} \leftarrow P_+ \left( \begin{bmatrix} u_k \\ v_k \end{bmatrix} - \alpha_k \begin{bmatrix} \nabla_u \tilde{\phi}(u_k, v_k) \\ \nabla_v \tilde{\phi}(u_k, v_k) \end{bmatrix} \right) = P_+ \left( \begin{bmatrix} u_k - \alpha_k \nabla F(u_k - v_k) - \alpha_k \lambda e \\ v_k + \alpha_k \nabla F(u_k - v_k) - \alpha_k \lambda e \end{bmatrix} \right)$$
(172)

where  $P_+$  projects onto the nonnegative orthant and  $e \in \mathbb{R}^d$  is a vector of ones.

YZW (USTC)

The iteration (172) is expected to inherit the property of being globally linearly convergent when F satisfies the assumptions of the last theorem. However, since the variables in (171) have been split into positive and negative parts, this property is maintained *only if* the iteration maintains complementarity of each iterate pair, i.e., if  $[u_k]_i [v_k]_i = 0, \forall k \in \mathbb{N}, i \in \{1, \ldots, d\}.$ 

A stochastic projected gradient method, with  $\nabla F(w_k)$  replaced by  $g(w_k, \xi_k)$ , has similar convergence properties as a standard SG method.

ヘロト 人間 とくほとく ほとう

For solving problem (167), a *proximal Newton* method is one that constructs, at each  $k \in \mathbb{N}$ , a model

$$q_{k}(w) = F(w_{k}) + \nabla F(w_{k})^{\top}(w - w_{k}) + \frac{1}{2}(w - w_{k})^{\top}H_{k}(w - w_{k}) + \lambda ||w||_{1},$$
(173)

where  $H_k$  represents  $\nabla^2 F(w_k)$  or a quasi-Newton approximation of it.

A proximal Newton method would involve (approximately) minimizing this model to compute a trial iterate  $\tilde{w}_k$ , then a step size  $\alpha_k > 0$  would be taken from a predetermined sequence or chosen by a line search to ensure that the new iterate  $w_{k+1} \leftarrow w_k + \alpha_k(\tilde{w}_k - w_k)$  yields  $\phi(w_{k+1}) < \phi(w_k)$ .

Proximal Newton methods are more challenging to design, analyze and implement than proximal gradient methods. Assuming  $H_k$  has been chosen to be positive definite, here are three essential ingredients in proximal Newton method:

**Choice of Subproblem Solver**  $q_k$  is nonsmooth and is challenging to minimize. One choice is coordinate descent, since the global minimizer of  $q_k$  along a coordinate descent direction can be computed analytically.

**Inaccurate Subproblem Solves** It's impractical to minimize  $q_k$  accurately for all  $k \in \mathbb{N}$ . Thus we need a practical and theoretically sufficient termination criteria.

A D F A B F A B F A B F

Interestingly, the norm of an ISTA step is an appropriate measure. Let  $ista_k(w)$  represent the result of an ISTA step applied to  $q_k$  from w. A trial point  $\tilde{w}_k$  represents a sufficiently accurate minimizer of  $q_k$  if, for some  $\eta \in [0, 1)$ , one finds

 $\|\operatorname{ista}_k(\widetilde{w}_k) - \widetilde{w}_k\|_2 \leqslant \eta \|\operatorname{ista}_k(w_k) - w_k\|_2 \text{ and } q_k(\widetilde{w}_k) < q_k(w_k).$ 

**Elimination of Variables** Due to the structure created by the  $\ell_1$ -norm regularizer, it can be effective in some applications to first identify a set of *active* variables then compute an approximate minimizer of  $q_k$  over the remaining *free* variables.

Our second class of second-order methods is based on the observation that  $\ell_1$ -norm regularized objective  $\phi$  in problem (167) is smooth in any orthant in  $\mathbb{R}^d$ .

In every iteration, orthant-based methods construct a smooth quadratic model of the objective, then produce a search direction by minimizing this model.

After performing a line search designed to reduce the objective function, a new orthant is selected and the process is repeated.
#### Orthant-based Methods

With the minimum norm subgradient of  $\phi$  at  $w \in \mathbb{R}^d$ , which is given component-wise for all  $i \in \{1, \ldots, d\}$  by

$$\hat{g}_i(w) = \begin{cases} [\nabla F(w)]_i + \lambda & \text{if } w_i > 0 \text{ or } \{w_i = 0 \text{ and } [\nabla F(w)]_i + \lambda < 0\} \\ [\nabla F(w)]_i - \lambda & \text{if } w_i < 0 \text{ or } \{w_i = 0 \text{ and } [\nabla F(w)]_i - \lambda > 0\} \\ 0 & \text{otherwise,} \end{cases}$$
(174)

the active orthant for an iterate  $w_k$  is characterized by the sign vector

$$\zeta_{k,i} = \begin{cases} sgn([w_k]_i) & \text{if } [w_k]_i \neq 0\\ sgn(-[\hat{g}(w_k)]_i) & \text{if } [w_k]_i = 0. \end{cases}$$
(175)

463 / 467

Along these lines, define the subsets of  $\{1,\ldots,d\}$  given by

$$\mathcal{A}_{k} = \{i : [w_{k}]_{i} = 0 \text{ and } |[\nabla F(w_{k})]_{i}| \leq \lambda\}$$
(176)  
and  $\mathcal{F}_{k} = \{i : [w_{k}]_{i} \neq 0\} \cup \{i : [w_{k}]_{i} = 0 \text{ and } |[\nabla F(w_{k})]_{i}| > \lambda\},$ (177)

where  $A_k$  represents the indices of variables that are active and kept at zero while  $\mathcal{F}_k$  represents those that are free to move  $\mathcal{P} \times \mathcal{P} \times \mathcal{P} \times \mathcal{P} \times \mathcal{P}$ YZW (USTC) Optimization Algorithms

#### Orthant-based Methods

Given these quantities, an orthant-based method proceeds as follows. First, compute the (approximate) solution  $d_k$  of the (smooth) quadratic problem

$$\min_{d \in \mathbb{R}^n} \hat{g}(w_k)^\top d + \frac{1}{2} d^\top H_k d$$
  
s.t.  $d_i = 0, \ i \in \mathcal{A}_k,$ 

where  $H_k$  represents  $\nabla^2 F(w_k)$  or an approximation of it.

Then the algorithm performs a line search—over a path contained in the current orthant—to compute the next iterate.

One option is a projected backtracking line search along  $d_k$ , computing the largest  $\alpha_k$  in a decreasing geometric sequence so

$$F(P_k(w_k + \alpha_k d_k)) < F(w_k),$$

where  $P_k(w)$  projects  $w \in \mathbb{R}^d$  onto the orthant defined by  $\zeta_k$ .

イロト イヨト イヨト イヨト 二日

## Outline I

- Unconstrained Optimization
- 2 Constrained Optimization
  - 二次规划
  - 非线性约束最优化
- 3 Convex Optimization
  - Convex Set and Convex Function
  - Convex Optimization and Algorithms
- 4 Sparse Optimization
  - Sparse Optimization Models
  - Sparse Optimization Algorithms

#### Optimization Methods for Machine Learning

YZW (USTC)

Image: A image: A

### Outline II

- Typical Form of Problems
- Stochastic Algorithms
- Other Popular Methods



æ

(日)

# Thanks for your attention!



æ

A D > A B > A B > A B >