

MAT5010 Introduction to Nonlinear Optimization

Zhi-Quan (Tom) Luo

School of Science and Engineering

The Chinese University of Hong Kong, Shenzhen

luozq@cuhk.edu.cn

Organization

Instructor:

Dr. Tom Luo (luozq@cuhk.edu.cn), Daoyuan Building, Room 505

Textbook:

D. Bertsekas, *Nonlinear Programming*. Athena Scientific;
Course notes to be distributed in class

Prerequisites:

Calculus, Linear Algebra, Statistics, and Matlab (or C)

Lecture Time:

Monday, 10:00 - 13:00, Sept 5

Wednesdays 8:30 -11:30, Sept 21, 28, Oct 19, 26, Nov 23, 30, Dec 14

Thursdays 12:30 -15:30 , Sept 22, 29, Oct 20, 27 (midterm), Nov 24, Dec 1, 15

Location: Zhixin Building room 201A

Course Objectives and Grading

- to present the basic theory and important algorithms, concentrating on results that are useful in large scale computation and contemporary applications
- to give students the background required to use the methods in their own research or engineering work
- to give students a thorough understanding of how optimization problems are solved, and some experience in solving them

Intended Audience:

Anyone who uses or will use scientific computing or optimization in engineering or related work.

Course Requirement and Grading:

In-class midterm (35%), Final course project (40%),
Problem Sets and Programming Assignments (almost weekly) (25%)

Lecture Topics

- Introduction. Mathematical Review.
- Unconstrained Optimization - Optimality Conditions
- Gradient Methods: Convergence Analysis. Rate of Convergence
- Newton and Gauss–Newton Methods
- Optimization Over a Convex Set; Optimality Conditions
- Feasible Direction Methods
- Constrained Optimization; Lagrange Multipliers. Introduction to Duality
- Penalty Methods. Augmented Lagrangian Methods.
- Alternating Directions Method of Multipliers.
- A General Approximate Gradient Projection (AGP) Framework
- Rate of Convergence Analysis of AGP, Error Bounds
- AGP for Nonconvex, Nonsmooth Problems. Lasso, Group Lasso
- Proximity Operator, Bregman Distance, Proximal Gradient (PG) Method
- Accelerated PG Methods (Nesterov's Method), Incremental Gradient Methods
- Stochastic Approximation Methods

Additional Topics

If time permits, we will also go over

- Conic Programming (SDP, SOCP)
- Conic Duality
- Interior Point Methods for Conic Programming
- SDP Relaxation for Quadratic Optimization, Signal Processing Applications
- Approximation Quality of SDP Relaxation

What is Optimization?

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{subject to} & x \in X \end{array}$$

where

- x is the decision variable (discrete, continuous)
- f is the objective function or cost function (Differentiable, convex, linear,...)
- X is the feasible region (convex, nonempty, ...)

Key questions:

- When is the problem feasible? Does optimal solution exist?
- How do we determine if a candidate x is an optimal solution?
- How to find an optimal solution numerically? (should do better than exhaustive search)

The Importance and the Process of Optimization

Importance of Optimization:

- fundamental to engineering design and analysis, with ubiquitous applications in many disciplines
- a mathematical foundation for engineering and science, just like probability theory

The Process of Optimization:

- **pre-optimization analysis:** number of variables, number of constraints, continuous/discrete, convex/nonconvex, smooth/nonsmooth, practical engineering requirement on solution accuracy, solution time, current state of the art
- **choose a suitable optimization strategy:** decide on whether to code your own algorithm or use off-the-shelf implementation; decide on initial point and termination criterion
- **post-optimization analysis:** sensitivity, Lagrangian multipliers, feasibility/infeasibility, duality gap
- **troubleshoot and iterate:** repeat the above process

How Do You Use Optimization?

Questions:

- How do I know that the answer from my computer run is the global minimum?
- Which algorithm should I use?
- Is my problem convex?
- Why doesn't my algorithm terminate?
- My cost function is nondifferentiable, should I smooth it?
- Why does my algorithm run into numerical difficulty? It's so slow.
- How do I initialize my algorithm? Which stepsize should I use?

Answer: take this course.

Mathematical Review

- Notations: Sets, Inf, Sup, functions, derivatives, gradients
- Vectors, matrices
- Norms, sequences, limits, continuity
- Mean value theorems
- Implicit function theorem
- Contraction mappings

Notations

- **Sets:**

$$X, x \in X, X_1 \cap X_2, X_1 \cup X_2$$

The set of real (complex) numbers is denoted by \mathbb{R} (\mathbb{C}).

- **Inf and Sup:**

The supremum of a nonempty set $X \subset \mathbb{R}$ is the smallest scalar y such that

$$y \geq x, \text{ for all } x \in X.$$

Similarly, the infimum of a set $X \subset \mathbb{R}$ is the largest scalar y such that

$$y \leq x, \text{ for all } x \in X.$$

If $\sup X \in X$ ($\inf X \in X$), then we say $\sup X = \max X$ ($\inf X = \min X$).

$$\sup\{1/n : n \geq 1\} = ? \quad \inf\{\sin n : n \geq 1\} = ?$$

- **Function:**

$$f : X \mapsto Y, X \text{ is called the domain, } Y \text{ is called the range}$$

Monotonicity. Inverse function: f^{-1}

Vectors and Subspaces

- **Linear combination:** if $\mathbf{x} = (x_1, x_2, \dots, x_n)$ and $\mathbf{y} = (y_1, y_2, \dots, y_n)$, then

$$\alpha \mathbf{x} + \beta \mathbf{y} = (\alpha x_1 + \beta y_1, \alpha x_2 + \beta y_2, \dots, \alpha x_n + \beta y_n)$$

- **Subspaces and linear independence:**

$S \in \mathbb{R}^n$ is called a subspace if it is closed under linear combination.

A set of vectors $\{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^r\}$ are linearly independent if there does not exist a $(\alpha_1, \dots, \alpha_r) \neq 0$ s.t.

$$\alpha_1 \mathbf{x}^1 + \alpha_2 \mathbf{x}^2 + \dots + \alpha_r \mathbf{x}^r = \mathbf{0}.$$

Basis and dimension of a subspace;

Inner product: $\langle \mathbf{x}, \mathbf{y} \rangle = \sum_i x_i y_i$;

Orthogonality: $\mathbf{x} \perp \mathbf{y}$ iff $\langle \mathbf{x}, \mathbf{y} \rangle = 0$.

Orthogonal complement of a subspace S :

$$S^\perp := \{\mathbf{x} \mid \langle \mathbf{x}, \mathbf{y} \rangle = 0, \forall \mathbf{y} \in S\}.$$

- **Cauchy-Schwartz inequality:**

$$\langle \mathbf{x}, \mathbf{y} \rangle \leq \|\mathbf{x}\| \|\mathbf{y}\|$$

Matrices

- For any matrix \mathbf{A} , we use a_{ij} (or A_{ij}) to denote its (i, j) th entry.
- Matrix addition, multiplication, transpose, symmetric matrices $\mathbf{A} = \mathbf{A}'$.

$$[\mathbf{AB}]' = \mathbf{B}'\mathbf{A}', \mathbf{AB} \neq \mathbf{BA}.$$

- Let \mathbf{A} be a matrix of size $m \times n$. Range of \mathbf{A} : $R(\mathbf{A})$, null space of \mathbf{A} : $N(\mathbf{A})$.

$$R(\mathbf{A}) = N^\perp(\mathbf{A})$$

Rank of \mathbf{A} $\text{rank}(\mathbf{A})$. Full rank matrix \mathbf{A} : $\text{rank}(\mathbf{A}) = \min\{m, n\}$.

- Square matrix ($m = n$), identity matrix \mathbf{I} , determinant $\det(\mathbf{A})$, inverse \mathbf{A}^{-1} .

$$\mathbf{A}^{-1} \text{ exists iff } \det(\mathbf{A}) \neq 0$$

- Let $X \subset \mathbb{R}^n$, \mathbf{A} be a matrix of size $m \times n$, then the image of X under \mathbf{A} is

$$\mathbf{A}X = \{\mathbf{A}\mathbf{x} \mid \mathbf{x} \in X\}$$

- Inner product:

$$\langle \mathbf{A}, \mathbf{B} \rangle = \text{Tr}(\mathbf{AB}') = \sum_{i,j} A_{ij}B_{ij}$$

Square Matrices

- Useful identities: $\det(\mathbf{A}) = \det(\mathbf{A}')$, $(\mathbf{A}\mathbf{B})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$.
- Orthonormal matrices: $\mathbf{A}\mathbf{A}' = \mathbf{I}$.
- (Complex) Eigenvalue λ : $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ for some $\mathbf{x} \neq \mathbf{0}$.
- Spectral radius: $\rho(\mathbf{A}) = \max_i \{|\lambda_i| : \lambda_i \text{ is an eigenvalue of } \mathbf{A}\}$.
- Eigen-decomposition of a symmetric matrix:

$$\mathbf{A} = \mathbf{P}'\mathbf{\Lambda}\mathbf{P},$$

where \mathbf{P} is orthonormal, $\mathbf{\Lambda}$ is diagonal and real.

- Positive (semi-) definite matrix: $\mathbf{A} \succeq \mathbf{0}$.

$$\mathbf{A} \succeq \mathbf{0}, \mathbf{B} \succeq \mathbf{0} \Rightarrow \mathbf{A} + \mathbf{B} \succeq \mathbf{0}.$$

- Square root $\mathbf{A}^{1/2}$: $\mathbf{A}^{1/2} := \mathbf{P}\sqrt{\mathbf{\Lambda}}\mathbf{P}'$, where $\mathbf{A} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}'$ is the eigen-decomposition of $\mathbf{A} \succeq \mathbf{0}$.
- A useful property: $\mathbf{A} \succeq \mathbf{0} \iff \langle \mathbf{A}, \mathbf{B} \rangle \geq 0$, for all $\mathbf{B} \succeq \mathbf{0}$.

Matrices

- **Singular values** $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_n \geq 0$: σ_i^2 is an eigenvalue of $\mathbf{A}\mathbf{A}'$.
- **Condition number**: $\kappa(\mathbf{A}) = \sigma_1/\sigma_n$.
- **Singular values decomposition**: $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}$, with \mathbf{U} and \mathbf{V} orthonormal, $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_n) \succeq \mathbf{0}$ diagonal.

- **Norms:**

Frobenious norm: $\|\mathbf{A}\|_F = \left(\sum_{i,j} |A_{ij}|^2\right)^{1/2} = \left(\sum_i \sigma_i^2\right)^{1/2}$

Nuclear norm: $\|\mathbf{A}\|_* = \sum_i \sigma_i$

Matrix 2-norm: $\|\mathbf{A}\|_2 = \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} = \max_i \sigma_i$

and $\|\mathbf{A}\|^2 = \|\mathbf{A}\mathbf{A}'\| = \|\mathbf{A}'\mathbf{A}\|$.

- **Useful inequalities:**

$$\|\mathbf{A}\mathbf{x}\| \leq \|\mathbf{A}\|_2 \|\mathbf{x}\|, \quad \|\mathbf{A}\|_* \geq \|\mathbf{A}\|_F \geq \|\mathbf{A}\|_2 \geq \rho(\mathbf{A}).$$

- **Cauchy-Schwartz inequality:**

$$\langle \mathbf{A}, \mathbf{B} \rangle \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_F$$

Derivatives and Mean Value Theorems

Suppose $f : \mathbb{R}^n \mapsto \mathbb{R}$ is a continuously twice differentiable.

- **Derivative:**

$$\frac{\partial f(\mathbf{x})}{\partial x_i} := \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{e}_i) - f(\mathbf{x})}{t}$$

where \mathbf{e}_i is the i th unit vector of \mathbb{R}^n .

- **Gradient vector:**

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

- **Hessian matrix:**

$$\nabla^2 f = \left[\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j} \right]$$

- **Taylor expansion:**

$$f(\mathbf{y}) - f(\mathbf{x}) = \nabla f(\mathbf{x})'(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{x} - \mathbf{y})'\nabla^2 f(\mathbf{x})(\mathbf{x} - \mathbf{y}) + o(\|\mathbf{x} - \mathbf{y}\|^2)$$

- **Mean value theorem:** there exist ξ, η in the line segment connecting \mathbf{x} and \mathbf{y} such that

$$f(\mathbf{y}) - f(\mathbf{x}) = \nabla f(\xi)'(\mathbf{y} - \mathbf{x})$$

and

$$f(\mathbf{y}) - f(\mathbf{x}) = \nabla f(\mathbf{x})'(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{x} - \mathbf{y})'\nabla^2 f(\eta)(\mathbf{x} - \mathbf{y})$$

- **Jacobian matrix:** For a vector valued continuously differentiable function $f : \mathbb{R}^n \mapsto \mathbb{R}^m$, define the Jacobian matrix

$$\nabla f(\mathbf{x}) = [\nabla f_1(\mathbf{x}), \nabla f_2(\mathbf{x}), \dots, \nabla f_m(\mathbf{x})].$$

- **Chain rule:** for $f : \mathbb{R}^k \mapsto \mathbb{R}^m$ and $g : \mathbb{R}^m \mapsto \mathbb{R}^n$, let $h(\mathbf{x}) = g(f(\mathbf{x}))$. Then

$$\nabla h(\mathbf{x}) = \nabla f(\mathbf{x})\nabla g(f(\mathbf{x})).$$

For example, we have

$$\nabla(f(\mathbf{Ax})) = \mathbf{A}'\nabla f(\mathbf{Ax}), \quad \nabla^2 f(\mathbf{Ax}) = \mathbf{A}'\nabla^2 f(\mathbf{Ax})\mathbf{A}.$$

Implicit Function Theorem

- Let $G(x, y) : \mathbb{R}^2 \mapsto \mathbb{R}$ be a continuously differentiable function and $P = (x^0, y^0)$ is a point such that

$$\frac{\partial G}{\partial y}(x^0, y^0) \neq 0$$

If $G(x, y) = G(x^0, y^0)$, then y may be expressed as a function of x in a neighborhood containing P ; i.e., there exists a differentiable function $y = g(x)$ such that

$$y^0 = g(x^0), \quad \text{and} \quad G(x, g(x)) = 0 \text{ for all } x \text{ close to } x^0.$$

- Let $G(\mathbf{x}, \mathbf{y}) : \mathbb{R}^n \mapsto \mathbb{R}^m$ be a continuously differentiable mapping where $\mathbf{y} \in \mathbb{R}^m$ and

$$\nabla_{\mathbf{y}} G = \left[\frac{\partial G_j}{\partial y_i} \right]_{m \times m} \text{ is nonsingular at } P = (\mathbf{x}^0, \mathbf{y}^0)$$

If $G(\mathbf{x}, \mathbf{y}) = G(\mathbf{x}^0, \mathbf{y}^0)$, then \mathbf{y} may be expressed as a function of \mathbf{x} locally around P ; i.e., there exists a differentiable mapping over $\mathbf{y} = g(\mathbf{x})$ such that

$$\mathbf{y}^0 = g(\mathbf{x}^0), \quad \text{and} \quad G(\mathbf{x}, g(\mathbf{x})) = \mathbf{0} \text{ for all } \mathbf{x} \text{ close to } \mathbf{x}^0.$$

Contraction Mappings

- **Lipschitzian Property:** $f : \mathbb{R}^n \mapsto \mathbb{R}^m$ satisfies

$$\|f(\mathbf{x}) - f(\mathbf{y})\| \leq \gamma \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y}$$

γ is called the Lipschitz constant.

- If $\gamma \leq 1$, then f is called a *non-expansive mapping*.
- If $\gamma < 1$, then f is called a *contraction mapping*
- **Fixed point theorem:** If f is a contraction, then the iterated function sequence

$$\mathbf{x}, f(\mathbf{x}), f(f(\mathbf{x})), f(f(f(\mathbf{x}))), \dots$$

converges to a unique fixed point \mathbf{x}^* (independent of \mathbf{x}) satisfying

$$\mathbf{x}^* = f(\mathbf{x}^*).$$