

Contents

I Introduction	3	3. Convex Function	12
1. General Formulation of the Optimization Problem	3	3.1 Definition of The Convex Function and The Class \mathcal{F}^k	12
1.1 Natural Classification	3	Definition of The Convex Function	12
1.2 Classification based on the properties of feasible set	4	3.2 Properties of The Convex Function	12
1.3 Different Types of Solution	4	3.3 Equivalent Definitions	13
1.4 Example	4	3.4 Examples	13
Example 1	4	4. Smooth and Convex Function	13
Example 2	4	4.1 The Class $\mathcal{F}_L^{k,l}(\mathbb{R}^n)$	13
Example 3	4	4.2 Necessary and Sufficient Conditions for The Class $\mathcal{F}_L^{1,1}(\mathbb{R}^n)$	13
2. Interplay of Optimization and Machine Learning	4	4.3 Necessary and Sufficient Conditions for The Class $\mathcal{F}_L^{2,1}(\mathbb{R}^n)$	13
2.1 The Machine Learning (Supervised) Paradigm	4	5. Strongly Convex Function	13
2.2 Empirical Risk Minimization	4	5.1 Definition of The Strongly Convex and The Class $\mathcal{S}_\mu^1(\mathbb{R}^n)$	13
2.3 Classification: The 0 ~ 1 Loss	5	5.2 Property of Strongly Convex Function	13
The Surrogate of The 0 ~ 1 Loss	5	5.3 Equivalent Definitions	14
2.4 Maximum Likelihood Estimation View	5	5.4 Examples	14
2.5 Maximum A posteriori Estimation View	5	6. Smooth and Strongly Convex Function	14
3. Performance of Numerical Methods	5	6.1 The Class $\mathcal{S}_{\mu,L}^{1,1}(\mathbb{R}^n)$	14
3.1 Method and Class of Problem	5	6.2 Property of Smooth and Strongly Convex Function	14
3.2 General Iterative Scheme	6	7. Conclusion	14
3.3 Complexity	6	7.1 Upper Bounds on Functional Components	14
3.4 Oracle	6	7.2 Lower Bounds on Functional Components	14
4. Complexity Bounds for Global Optimization	6	7.3 Other Lower Bounds on Functional Components	14
4.1 N-dimensional Box Constraint Problem	6	III Descent Meethod	15
4.2 Uniform Grid Method	7	1. Gradient Descent	15
4.3 Upper Complexity Bound	7	1.1 Basic Scheme	15
4.4 Lower Complexity Bound	8	Gradient Descent Formulation	15
II Foundations of Smooth Optimization	9	Step Size	15
1. Relaxation and Approximation	9	1.2 Performance for $\mathcal{C}_L^{1,1}(\mathbb{R}^n)$	15
1.1 Concepts of Relaxation and Approximation	9	1.3 Performance for $\mathcal{F}_L^{1,1}(\mathbb{R}^n)$	16
1.2 First Order Approximation	9	1.4 Performance for $\mathcal{S}_{\mu,L}^{1,1}(\mathbb{R}^n)$	16
1.3 Second Order Approximation	11	2. General Descent Directions	16
2. Classes of differentiable function	11	2.1 Choosing the Direction	16
2.1 Class $\mathcal{C}_L^{k,p}(\mathbb{R}^n)$	11	3. Newton Method	16
2.2 Class $\mathcal{C}_L^{1,1}(\mathbb{R}^n)$	11	3.1 Basic Scheme	16
geometric interpretation	12	Historical Origins	16
2.3 Class $\mathcal{C}_M^{2,1}(\mathbb{R}^n)$	12		

Basic Scheme	16	4. Separation Theorem	20
Damped Newton Method	16	4.1 Projection	20
3.2 Local Convergence of The Newton Method	16	4.2 Main Theorems	21
Bound for	16	5. Subgradient	21
3.3 Convergence Analysis	16	5.1 Definition of Subgradient	21
4. Conjugate Gradient	16	5.2 Properties of Subgradient	21
4.1 Historical Origins	16	5.3 Rules for Computing	21
Krylov Subspace	16	6. General Lower Complexity Bounds	22
CG	17	6.1 Problem Class	22
4.2 Fundamental Theory	17	6.2 Resisting Oracle	22
4.3 CG Algorithm	17	6.3 Lower Bound	23
IV Acceleration Methods	17	7. Subgradient Method	23
1. Lower Complexity Bounds for $\mathcal{F}_L^{\infty,1}(\mathbb{R}^n)$	17	7.1 property of Subgradient	23
1.1 Problem Class	17	7.2 Main Lemma	23
1.2 Worst Function in $\mathcal{F}_L^{\infty,1}(\mathbb{R}^n)$	17	7.3 Scheme for Non-smooth Problem	23
1.3 Theoretical Analysis and Main Results	17	7.4 Main Theorem	23
2. Lower Complexity Bounds for $\mathcal{S}_L^{\infty,1}(\mathbb{R}^n)$	18	8. Frank-Wolfe Algorithm	23
2.1 Worst Function in $\mathcal{S}_L^{\infty,1}(\mathbb{R}^n)$	18	8.1 Problems	23
3. Basic Scheme	18	8.2 Examples	23
3.1 Difference bewteen Lower Bounds and Real Efficiency	18	8.3 Convergence Theory	23
Efficiency Estimation	18	VI Beyond The Black-box Model	24
Estimate Sequences	18	1. Proximal Gradient Method	24
3.2 Optimal Scheme	18	1.1 Proximal Operator	24
4. Theoretical Analysis and Variants	18	1.2 Properties of Proximal Operator	24
4.1 Analysis of Optimal Scheme	18	1.3 Analysis for Proximal Gradient Method	24
4.2 Variants of Optimal Scheme	18	1.4 Accelerated Proximal Gradient Method	24
V General Convex Problem	19	1.5 Special case: Proximal Point Method	24
1. Motivation and Definitions	19	2. Douglas-Rachford Splitting	24
1.1 Motivation	19	2.1 Different Setting for Convex Problem	24
1.2 Definition of Convex Function	19	2.2 Fixed Point for Nonsmooth Composition	24
1.3 Other Properties	19	2.3 Splitting Algorithm	24
2. Operation with Convex Function	20	3. Duality Principle	24
2.1 Invariant Operations	20	3.1 Duality Principle	24
3. Continuity and Differentiability	20	4. Lagrangian Duality and Algorithms	24
3.1 Continuity of Convex Function	20	4.1 Lagrangian Duality	24
3.2 Differentiability of Convex Function	20	4.2 KKT	25
		4.3 Algorithm using Lagrangian duality	25

5.	Fenchel conjugate and algorithm	25
5.1	Fenchel conjugate	25
5.2	Properties	25
5.3	Fenchel duality	25
6.	Smoothing Techniques	25
6.1	Introduction	25
6.2	Nesterov's Smoothing	25
	Proximity Function	25
6.3	Moreau-Yosida Regularization	25
7.	Generalized Distance: Mirror Descent	25
7.1	Motivation	25
7.2	Bregman Divergence	25
7.3	Mirror Descent	25
VII	Stochastic Optimization	26

I Introduction

1. General Formulation of the Optimization Problem

Let \mathbf{x} be an n -dimensional real vector

$$\mathbf{x} = \left(x^{(1)}, \dots, x^{(n)} \right)^T \in \mathbb{R}^n$$

and $f_0(x), \dots, f_m(x)$ be some **real-valued** function defined on a set $S \subseteq \mathbb{R}^n$. We consider different variants of the following general minimization problem:

$$\begin{aligned} \min f_0(\mathbf{x}) \\ \text{s. t. } f_j(\mathbf{x}) \leq 0, j = 1 \dots m \\ \mathbf{x} \in S \end{aligned} \quad (\text{I.1})$$

where the sign \leq can be \leq, \geq or $=$.

\leq 表示受约束. 可用 \leq 表示 \geq , 然后用 \leq, \geq 表示 $=$.

We call $f_0(x)$ the **objective function** (目标函数) of our problem, the vector function

$$\mathbf{f}(\mathbf{x}) = (f_1(x), \dots, f_m(x))^T$$

is called the vector of **functional constraints** (泛函约束), the set S is the **basic feasible set** (基本可行集), and the set

$$Q = \{ \mathbf{x} \in S \mid f_j(\mathbf{x}) \leq 0, j = 1, \dots, m \}$$

is the **(entire) feasible set** (完整可行集) of problem I.1. It's just a convention to consider minimization problems.

1.1 Natural Classification

There exists a natural classification of the types of minimization problems:

- Constrained problems (带约束问题): $Q \subset \mathbb{R}^n$
- Unconstrained problems (无约束问题): $Q \equiv \mathbb{R}^n$
- smooth problems (光滑问题): all $f_j(x)$ are differentiable
- nonsmooth problems: some $f_k(x)$ are non-differentiable
- linearly constrained problems: the functional constraints are affine (仿射) like

$$f_j(\mathbf{x}) = \sum_{i=1}^n a_j^{(i)} x^{(i)} + b_j \equiv \langle a_j, \mathbf{x} \rangle + b_j, j = 1, \dots, m$$

- linear optimization problem: $f_0(\cdot)$ is also affine
- quadratic optimization problem: $f_0(\cdot)$ is quadratic
- quadratically constrained problem: all the function $f_0(\cdot), \dots, f_m(\cdot)$ are quadratic

1.2 Classification based on the properties of feasible set

There is also a classification based on properties of the feasible set.

- Problem I.1 is called **feasible**, if $Q \neq \emptyset$
- Problem I.1 is called **strictly feasible**, if there exists an $\mathbf{x} \in \text{int } Q$ such that $f_j(\mathbf{x}) < 0$ (or > 0) for all inequality constraints and $f_j(\mathbf{x}) = 0$ for all equality constraints (Slater condition)

$\text{int } Q$ 表示不包含边界点的 Q .

1.3 Different Types of Solution

Finally, we distinguish different types of solution to I.1:

- \mathbf{x}^* is called the **global optimal solution** to I.1 if $f_0(\mathbf{x}^*) \leq f_0(\mathbf{x})$ for all $\mathbf{x} \in Q$. In this case, $f_0(\mathbf{x}^*)$ is called the (global) **optimal value** of the problem. ($\mathbf{x}^* = \arg \min \dots$, 是个集合)
- \mathbf{x}^* is called a **local optimal solution** to I.1 if $\exists \delta$, for all $\mathbf{x} \in NBR(\mathbf{x}^*, \delta) \cap Q$ (NBR means Neighborhood), we have

$$f_0(\mathbf{x}^*) \leq f_0(\mathbf{x})$$

1.4 Example

Example 1 Let $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$ be our **design variables**. Then we can fix some functional characteristics of our decision vector $\mathbf{x} : f_0(\mathbf{x}), \dots, f_m(\mathbf{x})$. For example, we can consider a price of the project, amount of required resources, reliability of the system, etc. We fix the most important characteristics, $f_0(\mathbf{x})$, as our objective. For all others, we impose some bounds: $a_j \leq f_j(\mathbf{x}) \leq b_j$. Thus, we come to the problem:

$$\begin{aligned} & \min f_0(\mathbf{x}) \\ \text{s. t. } & a_j \leq f_j(\mathbf{x}) \leq b_j, j = 1 \dots m \\ & \mathbf{x} \in S \end{aligned}$$

where S stands for the **structural** constraints like non-negativity, boundedness of some variables, etc.

Example 2 Let our initial problem be as follow:

$$\text{Find } \mathbf{x} \in \mathbb{R}^n \text{ such that } f_j(\mathbf{x}) = a_j, j = 1, \dots, m \quad (\text{I.2})$$

Then we consider the problem

$$\min_{\mathbf{x}} \sum_{j=1}^m (f_j(\mathbf{x}) - a_j)^2$$

perhaps even with some additional constraints on \mathbf{x} . If the optimal value of the latter problem is zero, we conclude that our initial problem I.2 has a solution.

Example 3 Sometimes our decision variables $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$ must be integers. This can be described by the following constraint:

$$\sin(\pi \mathbf{x}^{(i)}) = 0, i = 1, \dots, n$$

Thus, we can also treat integer optimization problems:

$$\begin{aligned} & \min f_0(\mathbf{x}) \\ \text{s. t. } & a_j \leq f_j(\mathbf{x}) \leq b_j, j = 1, \dots, m \\ & \mathbf{x} \in S \\ & \sin(\pi \mathbf{x}^{(i)}) = 0, i = 1, \dots, n \end{aligned}$$

2. Interplay of Optimization and Machine Learning

2.1 The Machine Learning (Supervised) Paradigm

Consider the following **spaces**: a space of *example* \mathcal{X} , a space of **labels** \mathcal{Y} , and a space of *hypothesis* \mathcal{H} that contains functions mapping \mathcal{X} to \mathcal{Y} . Also consider a **loss** function $L : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ and a **probability measure** P over the space $\mathcal{X} \times \mathcal{Y}$.

Definition I.1. For any hypothesis $h \in \mathcal{H}$, we define the *risk* of h with respect to L and P to be:

$$\mathcal{L}_P(h) = E_P(L(h(\mathbf{X}), Y)), \text{ with } (\mathbf{X}, Y) \sim P$$

The general problem of machine learning can be then cast as finding the hypothesis $h \in \mathcal{H}$ that solves the following optimization problems:

$$\min_{h \in \mathcal{H}} \mathcal{L}_P(h)$$

2.2 Empirical Risk Minimization

The practical way is to solve an *empirical risk minimization* problem. That is,

$$\min_{h \in \mathcal{H}} \sum_{i=1}^m L(h(\mathbf{x}_i), y_i) = \min_{\theta} \sum_{i=1}^m L(h_{\theta}(\mathbf{x}_i), y_i)$$

since we do not know the P exactly. That is

$$\theta_* = \arg \min_{\theta} \sum_{i=1}^m L(h_{\theta}(\mathbf{x}_i), y_i)$$

- **Classification**: the output variable takes **class labels**. Especially, for **binary classification**, $y_i \in \{1, -1\}$.

- **Regression**: the output variable takes **continuous values**. That is $y_i \in \mathbb{R}$.

2.3 Classification: The 0 ~ 1 Loss

By using the 0 ~ 1 loss, we can rewrite the objective function as:

$$\sum_{i=1}^m \mathbf{1}_{\text{sign}(h(\mathbf{x}_i)) \neq y_i} = \sum_{i=1}^m \mathbf{1}_{y_i h(\mathbf{x}_i) < 0}$$

where $\mathbf{1}_z$ is just the indicator function, that is

$$\mathbf{1}_z = \begin{cases} 1 & z \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

and $\text{sign}(x)$ is sign function,

$$\text{sign}(x) = \begin{cases} 1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases}$$

In fact, this problem is **NP-hard** so it's pointless to try to minimize this function.

The Surrogate of The 0 ~ 1 Loss To deal with this problem, we can introduce some *surrogate*.

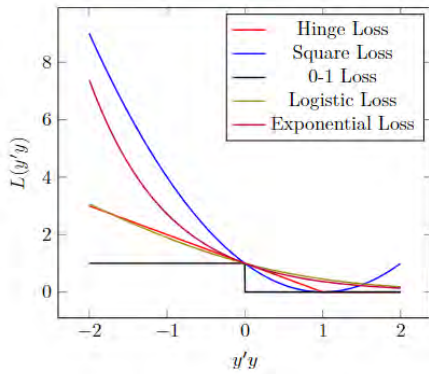


Figure I.1: The Surrogate of The 0 ~ 1 Loss

Some common loss function

- hinge loss or margin loss

$$\phi_h(y, y') = \max(0, 1 - y'y)$$

- Square loss:

$$\phi_s(y, y') = (1 - yy')^2$$

- Logistic loss:

$$\phi_l(y, y') = \frac{1}{\ln 2} \ln(1 + e^{-yy'})$$

- Exponential loss:

$$\phi_e(y, y') = e^{-yy'}$$

- Cross entropy loss:

$$\phi_c(y, y') = -t \ln(y') - (1 - t) \ln(1 - y')$$

where $t = \frac{(1+y)}{2}$

2.4 Maximum Likelihood Estimation View

Let's say, we have a likelihood function $P(\mathbf{X}, Y|\theta)$. Then the MLE for θ , the parameter we want to infer, is

$$\begin{aligned} \theta_{MLE} &= \underset{\theta}{\operatorname{argmax}} P(\mathbf{X}, Y|\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \log P(\mathbf{X}, Y|\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \log \prod_i P(x_i, y_i|\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log P(x_i, y_i|\theta) \end{aligned}$$

For case that $P(x_i, y_i|\theta) \sim \exp(-L(h_\theta(\mathbf{x}_i), y_i))$ (\sim 表示成比例), we have

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^n L(h_\theta(\mathbf{x}_i), y_i)$$

2.5 Maximum A posteriori Estimation View

If we replace the likelihood in the MLE formula above with the posterior, we get:

$$\begin{aligned} \theta_{MAP} &= \underset{\theta}{\operatorname{argmax}} P(\mathbf{X}, Y|\theta)P(\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \log\{P(\theta)P(\mathbf{X}, Y|\theta)\} \\ &= \underset{\theta}{\operatorname{argmax}} \log \left\{ P(\theta) \prod_i P(x_i, y_i|\theta) \right\} \\ &= \underset{\theta}{\operatorname{argmax}} \left\{ \log P(\theta) + \sum_{i=1}^n \log P(x_i, y_i|\theta) \right\} \end{aligned}$$

For case that $P(x_i, y_i|\theta) \sim \exp(-L(h_\theta(\mathbf{x}_i), y_i))$ and $P(\theta) \sim \exp(-\lambda \|\theta\|_2^2)$, we have

$$\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} \left\{ \sum_{i=1}^n L(h_\theta(\mathbf{x}_i), y_i) + \lambda \|\theta\|_2^2 \right\}$$

$P(\theta)$ 作为先验, 推导出正则项. 正则化可以让函数凸性更好, 更容易找到最小值.

3. Performance of Numerical Methods

3.1 Method and Class of Problem

考虑一种解决问题 1.1 的方法: 除了 $x^* = 0$ 之外, 什么都不做. 对于那些在原点恰好有最优解的问题, 这种方法的“性

能”是无法超越的。然而，这种方法对于其他任何问题都无法正常工作。因此，我们不能说某个特定问题 P 的最佳方法是什么，但我们可以针对一类问题 $\mathcal{F} \ni P$ 来讨论。

一个方法 \mathcal{M} 在整个类别 \mathcal{P} 上的性能可以自然地衡量其效率。考虑 \mathcal{M} 在类别 \mathcal{P} 上的性能。我们应该假设方法 \mathcal{M} 并没有关于特定问题 \mathcal{P} 的完整信息。

- **Model:** 已知 (对数值方案)“部分”问题 \mathcal{P} 被称为问题的模型。我们用 Σ 表示模型。通常，模型包括问题的表述、功能组件类别的描述等。

- **Oracle:** 我们通过预言的概念描述收集关于 \mathcal{P} 的特定信息的过程。预言 \mathcal{O} 只是一个回答方法连续问题的单位。

一般来说，每个问题都可以用不同的模型来描述。此外，对于每个问题，我们都可以开发出不同类型的预言。

Let us fix Σ and \mathcal{O} . In this case, it's natural to define the performance of \mathcal{M} on (Σ, \mathcal{O}) , as its performance on the worst \mathcal{P}_w from (Σ, \mathcal{O}) . Note that this \mathcal{P}_w can be bad only for \mathcal{M} .

Definition I.2. *The performance of \mathcal{M} on \mathcal{P} is the total amount of computational effort required by method \mathcal{M} to solve the problem \mathcal{P} .*

Solving the problem means finding an approximate solution on \mathcal{P} with some accuracy $\epsilon > 0$ (找逼近解就是解问题了)。

3.2 General Iterative Scheme

- Input: Starting point x_0 and accuracy $\epsilon > 0$
- Initialization: Set $k = 0, I_{-1} = \emptyset$

Main Loop:

- 1) Call oracle \mathcal{O} at point x_k
- 2) Update the information set: $I_k = I_{k-1} \cup (x_k, \mathcal{O}(x_k))$
- 3) Apply the rules of method \mathcal{M} to I_k and generate a new point x_{k+1}
- 4) Check criterion \mathcal{T}_ϵ . If **yes** then form an output \bar{x} . Otherwise set $k := k + 1$ and go to Step 1.

3.3 Complexity

In above scheme, we can see two potentially expensive steps. The first one is Step 1, where we call the oracle. The second one is Step 3, where we form the new test point. Thus, we can introduce two measures of complexity of problem \mathcal{P} for method \mathcal{M} :

- **Analytical Complexity:** The number of calls of the

oracle which is necessary to solve problem \mathcal{P} up to accuracy ϵ .

- **Arithmetical Complexity:** The total number of arithmetic operations (including the work of both oracle and method), which is necessary to solve problem \mathcal{P} up to accuracy ϵ .

For a particular method \mathcal{M} as applied to problem \mathcal{P} , arithmetical complexity can be easily obtained from the analytical complexity and complexity of the oracle.

3.4 Oracle

There is one standard assumption on the oracle which allows us to obtain the majority of results on analytical complexity for optimization schemes. This assumption, called the **Local Black Box Concept**, is listed in the following lines:

- 1) The only information available for the numerical scheme is the answer of the oracle.
- 2) The oracle is local: A small variation of the problem far enough from the test point x , which is compatible with the description of the problem class, does not change the answer at x

The standard formulation I.1 is called a functional model of optimization problems. Usually, for such models the standard assumptions are related to the level of smoothness of functional components. According to the degree of smoothness we can apply different types of oracle:

- zero-order oracle: return $f(x)$
- first-order oracle: return $f(x)$ and gradient $\nabla f(x)$
- second-order oracle: return $f(x)$, gradient $\nabla f(x)$ and Hessian $\nabla^2 f(x)$

4. Complexity Bounds for Global Optimization

4.1 N-dimensional Box Constraint Problem

Consider the following problem:

$$\min_{x \in \mathbb{B}_n} f(x) \quad (\text{I.3})$$

In our terminology, this is a constrained minimization problem with no functional constraints. The basic feasible set of this problem is \mathbb{B}_n , an n-dimensional box in \mathbb{R}^n :

$$\mathbb{B}_n = \{x \in \mathbb{R}^n \mid 0 \leq x^{(i)} \leq 1, i = 1, \dots, n\};$$

Let us measure distances in \mathbb{R}^n by the ℓ_∞ :

$$\|x\|_\infty = \max_{1 \leq i \leq n} |x^{(i)}|$$

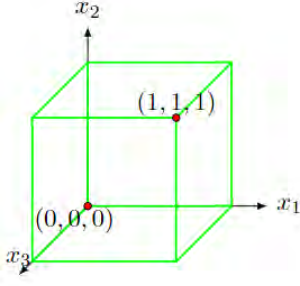


Figure I.2: N-dimensional Box Constraint Problem

Assume that, with respect to this norm, the objective function $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ is **Lipschitz continuous** on \mathbb{B}_n :

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq L \|\mathbf{x} - \mathbf{y}\|_\infty, \forall \mathbf{x}, \mathbf{y} \in \mathbb{B}_n \quad (\text{I.4})$$

with some constant L (**Lipschitz constant**).

4.2 Uniform Grid Method

Method $\mathcal{G}(p)$:

- 1) Form $(p+1)^n$ points

$$\mathbf{x}_{(i_1, \dots, i_n)} = \left(\frac{i_1}{p}, \frac{i_2}{p}, \dots, \frac{i_n}{p} \right)^T$$

where $(i_1, \dots, i_n) \in \{0, \dots, p\}^n$.

- 2) Among all points $\mathbf{x}_{(i_1, \dots, i_n)}$, find the point $\bar{\mathbf{x}}$ with the minimal value of the objective function.

- 3) The pair $(\bar{\mathbf{x}}, f(\bar{\mathbf{x}}))$ is the output of the method.

This method forms a **uniform grid** of the test points inside the box \mathbb{B}_n , computes the best value of the objective over this grid, and returns this value as an approximate solution to problem I.3.

- Zeroorder iterative method,
- Without any influence from the accumulated information on the sequence of test points.

Theorem I.3. Let f^* be a global optimal value of problem I.3. Then

$$f(\bar{\mathbf{x}}) - f^* \leq \frac{L}{2p}$$

Proof. For a multiindex (i_1, i_2, \dots, i_n) , let \mathbf{x}_* be a global minimum of our problem, Then there exist coordinates (i_1, i_2, \dots, i_n) such that

$$\mathbf{x} \equiv \mathbf{x}_{(i_1, i_2, \dots, i_n)} \leq \mathbf{x}_* \leq \mathbf{x}_{(i_1+1, i_2+1, \dots, i_n+1)} \equiv \mathbf{y}$$

Note that

$$\mathbf{y}^{(i)} - \mathbf{x}^{(i)} = \frac{1}{p}$$

for $i = 1, \dots, n$ and

$$\mathbf{x}_*^{(i)} \in [\mathbf{x}^{(i)}, \mathbf{y}^{(i)}], i = 1, \dots, n$$

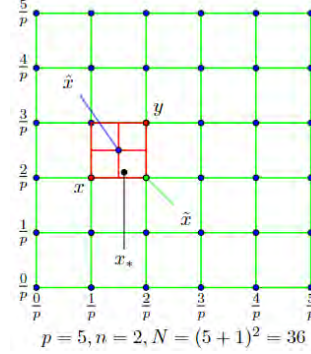


Figure I.3: Uniform Grid Method

Denote $\hat{\mathbf{x}} = (\mathbf{x} + \mathbf{y})/2$. Let's form a point $\tilde{\mathbf{x}}$ as follows:

$$\tilde{\mathbf{x}} = \begin{cases} \mathbf{y}^{(i)} & \text{if } \mathbf{x}_*^{(i)} \geq \hat{\mathbf{x}}^{(i)} \\ \mathbf{x}^{(i)} & \text{otherwise} \end{cases}$$

It's clear that $|\tilde{\mathbf{x}}^{(i)} - \mathbf{x}_*^{(i)}| \leq \frac{1}{2p}, i = 1, \dots, n$. Therefore

$$\|\tilde{\mathbf{x}} - \mathbf{x}_*\|_\infty = \max_{1 \leq i \leq n} |\tilde{\mathbf{x}}^{(i)} - \mathbf{x}_*^{(i)}| \leq \frac{1}{2p}$$

Since $\tilde{\mathbf{x}}$ belongs to our grid, we conclude that

$$f(\bar{\mathbf{x}}) - f(\mathbf{x}_*) \leq f(\tilde{\mathbf{x}}) - f(\mathbf{x}_*) \leq L \|\tilde{\mathbf{x}} - \mathbf{x}_*\|_\infty \leq \frac{L}{2p}$$

Q.E.D.

4.3 Upper Complexity Bound

Let us conclude with the definition of our problem class. We fix our goal as follows:

$$\text{Find } \bar{\mathbf{x}} \in \mathbb{B}_n : f(\bar{\mathbf{x}}) - f^* \leq \epsilon \quad (\text{I.5})$$

Then we immediately get the following result.

Corollary I.4. The analytical complexity of problem class I.3, I.4, I.5 for method \mathcal{G} is at most

$$\mathcal{A}(\mathcal{G}) = \left(\left\lceil \frac{L}{2\epsilon} \right\rceil + 2 \right)^n$$

Proof. Take $p = \lfloor \frac{L}{2\epsilon} \rfloor + 1$. Then $p \geq \frac{L}{2\epsilon}$. and , in view of Theorem I.3, we have

$$f(\bar{x}) - f^* \leq \frac{L}{2p} \leq \epsilon$$

Note that we construct $(p+1)^n$ points.

Q.E.D.

Remark:

- Each point has an access to oracle, so the number of accessing is equal to the number of points.
- Thus, $\mathcal{A}(\mathcal{G})$ justifies an upper complexity bound for our problem class.

4.4 Lower Complexity Bound

我们需要进一步探讨两种情况:

- 首先, 可能我们的证明过于粗糙, 方法 $\mathcal{G}(p)$ 的实际性能要好得多.
- 其次, 我们仍然不能确定 $\mathcal{G}(p)$ 是否是解决 I.3 的合理方法. 可能存在其他性能更高的方案.

为了回答这些问题, 我们需要为问题类别 I.3, I.4, I.5 推导出更低的复杂性界限. 这些界限的主要特点如下:

- 它们基于黑箱概念.
- 这些界限对所有合理的迭代方案都有效. 因此, 它们为我们提供了问题类别的分析复杂性的下限估计.
- 这样的界限往往采用抵抗预言 (resisting oracle) 的思想.

一个抵抗预言试图为每一种特定方法 (例如, $\mathcal{G}(p)$) 创建最糟糕的问题.

- 它从一个 “空” 函数开始, 并试图以最糟糕的方式回答方法的每个调用.
- 然而, 答案必须与先前的答案和问题类别的描述相符.

在方法终止后, 可以重构一个完全符合算法累积的最终信息集的问题. 此外, 如果我们在这个新生问题上运行该方法, 由于它将得到预言相同的答案序列, 因此它将重现相同的测试点序列.

Let us show how this works for problem I.3. Consider the class of problems \mathcal{C} defined as follows:

- Model: $\min_{x \in \mathbb{B}_n} f(x)$, where $f(x)$ is ℓ_∞ -Lipschitz continuous on \mathbb{B}_n .
- Oracle: Zero-order Local Black Box.
- Approximate solution: Find $\bar{x} \in \mathbb{B}_n : f(\bar{x}) - f^* \leq \epsilon$

Theorem I.5. For $\epsilon < \frac{1}{2}L$, the analytical complexity of problem class \mathcal{C} is at least $(\frac{L}{2\epsilon})^n$.

Proof. Denote $q = \lfloor \frac{L}{2\epsilon} \rfloor (\geq 1)$. Assume that there exists a method, which need $N < q^n$ calls of oracle to solve any problem from \mathcal{C} . Let us apply this method to the following resisting strategy:

Oracle returns $f(x) = 0$ at any test point x

Therefore this method can find only $\bar{x} \in \mathbb{B}_n$ with $f(\bar{x}) = 0$. However, note that there exists $\hat{x} \in \mathbb{B}_n$ such that

$$\hat{x} + \frac{1}{q}e \in \mathbb{B}_n, e = (1, \dots, 1)^T \in \mathbb{R}^n$$

and there were no test points inside the box $\mathbb{B} = \{x | \hat{x} \leq x \leq \hat{x} + \frac{1}{q}e\}$.

(“每个由 \hat{x} 构造的 \mathbb{B} 至少包含一个点” 这个表述是错误的, 因为和假设不符合, 因此没有测试点的 \mathbb{B} 一定存在)

Denotes $x_* = \hat{x} + \frac{1}{2q}e$. Consider the function

$$\bar{f}(x) = \min\{0, L\|x - x_*\|_\infty - \epsilon\}$$

Clearly,

- 1) this function is ℓ_∞ -Lipschitz continuous with the constant L and
- 2) its global optimal value is $-\epsilon$.

Moreover, $\bar{f}(x)$ differs from zero only inside the box $\mathbb{B}' = \{x | \|x - x_*\|_\infty \leq \frac{\epsilon}{L}\}$. Since $2q \leq L/\epsilon$, we conclude that

$$\mathbb{B}' \subseteq \mathbb{B} \equiv \{\|x - x_*\|_\infty \leq \frac{1}{2q}\}$$

Thus, $\bar{f}(x)$ is equal to zero at all test points of our methods. Since the accuracy of the rest of our method is ϵ , we come to the following conclusion: If the number of calls of the oracle is less than q^n , then the accuracy of the result cannot be better than ϵ .
Q.E.D.

Now we can say much more about the performance of the Uniform Grid Method. Let us compare its efficiency estimate with the lower bound:

$$\mathcal{G} : \left(\left\lfloor \frac{L}{2\epsilon} \right\rfloor + 2 \right)^n, \text{ Lower Bound: } \left(\left\lfloor \frac{L}{2\epsilon} \right\rfloor \right)^n$$

If $\epsilon \leq O(\frac{L}{n})$, then the lower and upper bounds coincide up to an absolute constant multiplicative factor. This means that, for such level of accuracy, $\mathcal{G}(p)$ is optimal for the problem class \mathcal{C} .

Meanwhile, the Theorem I.5 supports our original statement that general optimization problems are unsolvable.

Remark: If $\epsilon \leq O(\frac{L}{n})$, we have

$$\begin{aligned} \frac{(\lfloor \frac{L}{2\epsilon} \rfloor + 2)^n}{(\lfloor \frac{L}{2\epsilon} \rfloor)^n} &= \left(1 + \frac{1}{\frac{1}{2} \lfloor \frac{L}{2\epsilon} \rfloor}\right)^n \\ &\leq \left(1 + \frac{1}{cn}\right)^n = \left(\left(1 + \frac{1}{cn}\right)^{cn}\right)^{\frac{1}{c}} \stackrel{n \rightarrow \infty}{=} e^{\frac{1}{c}} \end{aligned}$$

where c is a constant.

II Foundations of Smooth Optimization

1. Relaxation and Approximation

1.1 Concepts of Relaxation and Approximation

We call the sequence $\{a_k\}_{k=0}^\infty$ a **relaxation sequence** if $\forall k \geq 0$,

$$a_{k+1} \leq a_k$$

Consider several methods for solving the following unconstrained minimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad (\text{II.1})$$

where $f(\mathbf{x})$ is a smooth function. In order to do so, we generate a relaxation sequence

$$\{f(\mathbf{x}_k)\}_{k=0}^\infty, \quad f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k), \quad k = 0, 1, \dots$$

This strategy has the following important advantages:

- 1) If $f(\mathbf{x})$ is bounded below on \mathbb{R}^n , then the sequence $\{f(\mathbf{x}_k)\}_{k=0}^\infty$ converges.
- 2) In any case we improve the initial value of the objective function.

In general, to **approximate** means to replace an initial complex object by a simplified one, which is close by its properties to the original.

In nonlinear optimization we usually apply local approximations based on derivatives(导数) of non-linear functions. These are the first order and the second order approximations (or, the linear and quadratic approximations).

1.2 First Order Approximation

Let $f(\mathbf{x})$ be differentiable at $\bar{\mathbf{x}}$. Then for $\mathbf{y} \in \mathbb{R}^n$, we have

$$f(\mathbf{y}) = f(\bar{\mathbf{x}}) + \langle \nabla f(\bar{\mathbf{x}}), \mathbf{y} - \bar{\mathbf{x}} \rangle + o(\|\mathbf{y} - \bar{\mathbf{x}}\|)$$

where $o(r)$ (peano, 皮亚诺余项) is some function of $r \geq 0$, such that

$$\lim_{r \rightarrow 0} \frac{1}{r} o(r) = 0, \quad o(0) = 0$$

In the sequel we fix the notation $\|\cdot\|$ for the standard Euclidean norm on \mathbb{R}^n :

$$\|\mathbf{x}\| = \left[\sum_{i=1}^n (x^{(i)})^2 \right]^{1/2}$$

The linear function $f(\bar{\mathbf{x}}) + \langle \nabla f(\bar{\mathbf{x}}), \mathbf{y} - \bar{\mathbf{x}} \rangle$ is called the linear approximation of f at $\bar{\mathbf{x}}$.

- The vector $\nabla f(\mathbf{x})$ is called the **gradient** of function f at \mathbf{x} .

Considering the points $\mathbf{y}_i = \bar{\mathbf{x}} + \epsilon \mathbf{e}_i$, where \mathbf{e}_i is the i -th coordinate vector in \mathbb{R}^n , and taking the limit in $\epsilon \rightarrow 0$, we obtain the following coordinate representation of the gradient

$$\nabla f(\mathbf{x}) = \left(\frac{\partial f(\mathbf{x})}{\partial x^{(1)}}, \dots, \frac{\partial f(\mathbf{x})}{\partial x^{(n)}} \right)^\top$$

- Denote by $\mathcal{L}_f(\alpha)$ the **level set**(等高线) of $f(\mathbf{x})$:

$$\mathcal{L}_f(\alpha) = \{\mathbf{x} \in \mathbb{R}^n | f(\mathbf{x}) \leq \alpha\}$$

- Consider the *set* of directions that are **tangent** to $\mathcal{L}_f(f(\bar{\mathbf{x}}))$ at $\bar{\mathbf{x}}$

$$S_f(\bar{\mathbf{x}}) = \left\{ \mathbf{s} \in \mathbb{R}^n \left| \lim_{\mathbf{y}_k \rightarrow \bar{\mathbf{x}}, f(\mathbf{y}_k)=f(\bar{\mathbf{x}})} \frac{\mathbf{y}_k - \bar{\mathbf{x}}}{\|\mathbf{y}_k - \bar{\mathbf{x}}\|} = \mathbf{s} \right. \right\}$$

Lemma II.1. If $\mathbf{s} \in S_f(\bar{\mathbf{x}})$, then $\langle \nabla f(\bar{\mathbf{x}}), \mathbf{s} \rangle = 0$

Proof. For $f(\mathbf{y}_k) = f(\bar{\mathbf{x}})$, we have,

$$\begin{aligned} f(\mathbf{y}_k) &= f(\bar{\mathbf{x}}) + \langle \nabla f(\bar{\mathbf{x}}), \mathbf{y}_k - \bar{\mathbf{x}} \rangle + o(\|\mathbf{y}_k - \bar{\mathbf{x}}\|) \\ &= f(\bar{\mathbf{x}}) \end{aligned}$$

Therefore, $\langle \nabla f(\bar{\mathbf{x}}), \mathbf{y}_k - \bar{\mathbf{x}} \rangle + o(\|\mathbf{y}_k - \bar{\mathbf{x}}\|) = 0$. Dividing this equation by $\|\mathbf{y}_k - \bar{\mathbf{x}}\|$, and taking the limit $\mathbf{y}_k \rightarrow \bar{\mathbf{x}}$, we obtain

$$\langle \nabla f(\bar{\mathbf{x}}), \lim_{\mathbf{y}_k \rightarrow \bar{\mathbf{x}}} \frac{\mathbf{y}_k - \bar{\mathbf{x}}}{\|\mathbf{y}_k - \bar{\mathbf{x}}\|} \rangle = \langle \nabla f(\bar{\mathbf{x}}), \mathbf{s} \rangle = 0$$

Q.E.D.

The directions $-\nabla f(\bar{\mathbf{x}})$ (the **antigradient**) is the directions of the **fastest local decrease** of $f(\mathbf{x})$ at point $\bar{\mathbf{x}}$. (证明: 梯度负方向下降最快)

Proof. Let \mathbf{s} be a direction in \mathbb{R}^n , $\|\mathbf{s}\| = 1$. Consider the local decrease of $f(\mathbf{x})$ along \mathbf{s} :

$$\Delta(\mathbf{s}) = \lim_{\alpha \rightarrow 0} \frac{1}{\alpha} [f(\bar{\mathbf{x}} + \alpha \mathbf{s}) - f(\bar{\mathbf{x}})]$$

Note that $f(\bar{\mathbf{x}} + \alpha \mathbf{s}) - f(\bar{\mathbf{x}}) = \alpha \langle \nabla f(\bar{\mathbf{x}}), \mathbf{s} \rangle + o(\alpha \|\mathbf{s}\|)$. Therefore, we have

$$\Delta(\mathbf{s}) = \langle \nabla f(\bar{\mathbf{x}}), \mathbf{s} \rangle$$

By leveraging the Cauchy-Schwartz inequality, that is $-\|x\| \cdot \|y\| \leq \langle x, y \rangle \leq \|x\| \cdot \|y\|$, we obtain,

$$\Delta(\mathbf{s}) = \langle \nabla f(\bar{\mathbf{x}}), \mathbf{s} \rangle \geq -\|\nabla f(\bar{\mathbf{x}})\|$$

For the lower bound $\bar{\mathbf{s}} = -\frac{\nabla f(\bar{\mathbf{x}})}{\|\nabla f(\bar{\mathbf{x}})\|}$ (取到下界), we have

$$\Delta(\bar{\mathbf{s}}) = -\frac{\langle \nabla f(\bar{\mathbf{x}}), \nabla f(\bar{\mathbf{x}}) \rangle}{\|\nabla f(\bar{\mathbf{x}})\|} = -\|\nabla f(\bar{\mathbf{x}})\|$$

Q.E.D.

Theorem II.2 (First-order Optimality Condition). Let \mathbf{x}^* be a local minimum of differentiable function $f(\mathbf{x})$. Then $\nabla f(\mathbf{x}^*) = 0$.

Proof. Since \mathbf{x}^* is a local minimum of $f(\mathbf{x})$, then there exists $r > 0$ such that $\forall \mathbf{y}$, $\|\mathbf{y} - \mathbf{x}^*\| \leq r$, we have $f(\mathbf{y}) \geq f(\mathbf{x}^*)$. Since f is differentiable, this implies that

$$f(\mathbf{y}) = f(\mathbf{x}^*) + \langle \nabla f(\mathbf{x}^*), \mathbf{y} - \mathbf{x}^* \rangle + o(\|\mathbf{y} - \mathbf{x}^*\|) \geq f(\mathbf{x}^*)$$

Thus, $\forall \mathbf{s}$, $\|\mathbf{s}\| = 1$, we have $\langle \nabla f(\mathbf{x}^*), \mathbf{s} \rangle \geq 0$. Consider the directions \mathbf{s} and $-\mathbf{s}$, we get

$$\langle \nabla f(\mathbf{x}^*), \mathbf{s} \rangle = 0, \forall \mathbf{s}, \|\mathbf{s}\| = 1$$

Finally, choosing $\mathbf{s} = \mathbf{e}_i, i = 1 \dots n$, where \mathbf{e}_i is the i -th coordinate vector in \mathbb{R}^n , we obtain $\nabla f(\mathbf{x}^*) = 0$. Q.E.D.

仅必要条件, 反过来是 stationary points (静态点).

e.g $f(x) = x^3, x \in \mathbb{R}^1$ at $x = 0$.

Corollary II.3. Let \mathbf{x}^* be a local minimum of a differentiable function $f(\mathbf{x})$ s.t. linear equality constraints

$$\mathbf{x} \in \mathcal{L} \equiv \{\mathbf{x} \in \mathbb{R}^n | A\mathbf{x} = \mathbf{b}\} \neq \emptyset$$

where A is an $m \times n$ matrix and $\mathbf{b} \in \mathbb{R}^m, m < n$. Then there exists a vector of multipliers λ^* such that

$$\nabla f(\mathbf{x}^*) = A^\top \lambda^* \quad (\text{II.2})$$

\mathcal{L} is the Null space of matrix A .

Proof. Consider some vector $\mathbf{u}_i, i = 1 \dots k$, that form a basis of the Null space (零空间) of matrix A . Then any $\mathbf{x} \in \mathcal{L}$ can be represented as follows:

$$\mathbf{x} = \mathbf{x}(\mathbf{y}) \equiv \mathbf{x}^* + \sum_{i=1}^k \mathbf{y}^{(i)} \mathbf{u}_i, \mathbf{y} \in \mathbb{R}^k$$

Moreover, the point $\mathbf{y} = 0$ is a local minimum of the function $\phi(\mathbf{y}) = f(\mathbf{x}(\mathbf{y}))$. In view of **Theorem II.2**, $\nabla \phi(0) = 0$. This means that

$$\frac{\partial \phi(0)}{\partial \mathbf{y}^{(i)}} = \frac{\partial \phi(0)}{\partial \mathbf{x}(\mathbf{y})} \cdot \frac{\partial \mathbf{x}(\mathbf{y})}{\partial \mathbf{y}^{(i)}} = \langle \nabla f(\mathbf{x}^*), \mathbf{u}_i \rangle = 0, i = 1 \dots k$$

and II.2 follows. (因为零空间与行空间正交)

Q.E.D.

1.3 Second Order Approximation

Let function $f(\mathbf{x})$ be twice differentiable at $\bar{\mathbf{x}}$. Then

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

The quadratic function

$$f(\bar{\mathbf{x}}) + \langle \nabla f(\bar{\mathbf{x}}), \mathbf{y} - \bar{\mathbf{x}} \rangle + \frac{1}{2} \langle \nabla^2 f(\bar{\mathbf{x}})(\mathbf{y} - \bar{\mathbf{x}}), \mathbf{y} - \bar{\mathbf{x}} \rangle$$

is called the **quadratic** (or **second order**) approximation of function f at $\bar{\mathbf{x}}$.

Recall that the $(n \times n)$ matrix $\nabla^2 f(\mathbf{x})$ has the following entries:

$$(\nabla^2 f(\mathbf{x}))^{(i,j)} = \frac{\partial^2 f(\mathbf{x})}{\partial x^{(i)} \partial x^{(j)}}$$

This matrix is called **Hessian** of function f at \mathbf{x} .

Note that the Hessian is a **symmetric** matrix:

$$\nabla^2 f(\mathbf{x}) = [\nabla^2 f(\mathbf{x})]^\top$$

Remark:

- Notation $A \succeq 0$ means that A is **positive semidefinite** (半正定):

$$\mathbf{x}^\top A \mathbf{x} \geq 0 \quad \forall \mathbf{x} \in \mathbb{R}^n$$

- Notation $A \succ 0$ means that A is **positive definite** (正定):

$$\mathbf{x}^\top A \mathbf{x} > 0 \quad \forall \mathbf{x} \in \mathbb{R}^n$$

Theorem II.4. Let \mathbf{x}^* be a local minimum of a twice differentiable function $f(\mathbf{x})$. Then

$$\nabla f(\mathbf{x}^*) = 0, \quad \nabla^2 f(\mathbf{x}^*) \succeq 0$$

Proof. Since \mathbf{x}^* is a local minimum of function $f(\mathbf{x})$, there exists $r > 0$ such that $\forall \mathbf{y}, \|\mathbf{y} - \mathbf{x}^*\| \leq r$, we have

$$f(\mathbf{y}) \geq f(\mathbf{x}^*)$$

In view of **Theorem II.2**, $\nabla f(\mathbf{x}^*) = 0$.

Therefore, for any such \mathbf{y} ,

$$\begin{aligned} f(\mathbf{y}) &= f(\mathbf{x}^*) + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}^*)(\mathbf{y} - \mathbf{x}^*), \mathbf{y} - \mathbf{x}^* \rangle + o(\|\mathbf{y} - \mathbf{x}^*\|^2) \\ &\geq f(\mathbf{x}^*) \end{aligned}$$

两项都除以 $\|\mathbf{y} - \mathbf{x}^*\|^2$, let $\mathbf{y} \rightarrow \mathbf{x}^*$, and let $\mathbf{s} = \frac{\mathbf{y} - \mathbf{x}^*}{\|\mathbf{y} - \mathbf{x}^*\|}$

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}^*} \left(\frac{1}{2} \frac{\langle \nabla^2 f(\mathbf{x}^*)(\mathbf{y} - \mathbf{x}^*), \mathbf{y} - \mathbf{x}^* \rangle}{\|\mathbf{y} - \mathbf{x}^*\|^2} + \frac{o(\|\mathbf{y} - \mathbf{x}^*\|^2)}{\|\mathbf{y} - \mathbf{x}^*\|^2} \right) \geq 0$$

$$\langle \nabla^2 f(\mathbf{x}^*) \mathbf{s}, \mathbf{s} \rangle \geq 0$$

$$\forall \mathbf{s}, \|\mathbf{s}\| = 1$$

Q.E.D. Let us give a sufficient condition for that inclusion.

Theorem II.5. Let function $f(\mathbf{x})$ be twice differentiable on \mathbb{R}^n and let \mathbf{x}^* satisfy the following conditions:

$$\nabla f(\mathbf{x}^*) = 0, \quad \nabla^2 f(\mathbf{x}^*) \succ 0$$

Then \mathbf{x}^* is a strict local minimum of $f(\mathbf{x})$.

Remark: A point $\bar{\mathbf{x}} \in \mathbb{R}^n$ is an unconstrained strict local minimum of a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ if $\exists \epsilon > 0$ such that $f(\bar{\mathbf{x}}) < f(\mathbf{x}) \quad \forall \mathbf{x} \in B(\bar{\mathbf{x}}, \epsilon), \mathbf{x} \neq \bar{\mathbf{x}}$, where $B(\bar{\mathbf{x}}, \epsilon) := \{\mathbf{x} \mid \|\mathbf{x} - \bar{\mathbf{x}}\| \leq \epsilon\}$. $\lambda_1(A)$ 表示 A 的最小特征值, $\lambda_{max}(A)$ 表示最大的特征值.

2. Classes of differentiable function

2.1 Class $C_L^{k,p}(\mathbb{R}^n)$

Consider a classes of **differentiable** functions which meet a **Lipschitz condition** for a derivative of certain order.

Let $Q \subseteq \mathbb{R}^n$. We denote by $C_L^{k,p}(Q)$ the class of functions with the following properties:

- Any $f \in C_L^{k,p}(Q)$ is k times continuously **differentiable** on Q .
- Its p -th derivative is **Lipschitz continuous** on Q with the constant L :

$$\|f^{(p)}(\mathbf{x}) - f^{(p)}(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$$

$\forall \mathbf{x}, \mathbf{y} \in Q$.

Clearly, we always have

- 1) $p \leq k$ 显然成立
- 2) If $q \geq k$, then $C_L^{q,p}(Q) \subseteq C_L^{k,p}$.
- 3) Note also that these classes possess the following property:

If $f_1 \in C_{L_1}^{k,p}(Q), f_2 \in C_{L_2}^{k,p}(Q)$ and $\alpha, \beta \in \mathbb{R}^1$, then for

$$L_3 = |\alpha|L_1 + |\beta|L_2$$

we have $\alpha f_1 + \beta f_2 \in C_{L_3}^{k,p}(Q)$.

Remark. We use notation $f \in C^k(Q)$ for a function f which is k times continuously differentiable on Q .

2.2 Class $C_L^{1,1}(\mathbb{R}^n)$

Consider $C_L^{1,1}(\mathbb{R}^n)$, the class of functions with **Lipschitz continuous gradient**. By definition, the inclusion $f \in C_L^{1,1}(\mathbb{R}^n)$ implies that $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$$

Lemma II.6. A function $f(\mathbf{x})$ belongs to $C_L^{2,1}(\mathbb{R}^n) \subset C_L^{1,1}(\mathbb{R}^n)$, if and only if

$$\|\nabla^2 f(\mathbf{x})\| \leq L, \forall \mathbf{x} \in \mathbb{R}^n$$

矩阵二范数

$$g(a) = \nabla f(x + a(y - x)), g(0) = \nabla f(x), g(1) = \nabla f(y)$$

$$g(1) = g(0) + \int_0^1 g'(\tau) d\tau$$

$$g'(\tau) = \frac{\partial \nabla f(x + \tau(y - x))}{\partial(x + \tau(y - x))} \frac{\partial(x + \tau(y - x))}{\partial \tau} = \nabla^2 f(x + \tau(y - x))(y - x)$$

e.g.

geometric interpretation The next statement is important for the geometric interpretation of function from $C_L^{1,1}(\mathbb{R}^n)$.

Lemma II.7. Let $f \in C_L^{1,1}(\mathbb{R}^n)$. Then $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, we have

$$|f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| \leq \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2$$

f is called smooth.

Remark.

$f(\mathbf{y})$ 和其一阶逼近 $g(\mathbf{y}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$ 的距离的上界.

强凸

2.3 Class $C_M^{2,1}(\mathbb{R}^n)$

Consider class $C_M^{2,2}(\mathbb{R}^n)$. That is $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, we have

$$\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\| \leq M \|\mathbf{x} - \mathbf{y}\|$$

Lemma II.8. Let $f \in C_M^{2,2}(\mathbb{R}^n)$. Then for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, we have

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x}) - \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x})\| \leq \frac{M}{2} \|\mathbf{y} - \mathbf{x}\|^2$$

Corollary II.9. Let $f \in C_M^{2,2}(\mathbb{R}^n)$ and $\|\mathbf{y} - \mathbf{x}\| = r$. Then

$$\nabla^2 f(\mathbf{x}) - MrI_n \succeq \nabla^2 f(\mathbf{y}) \succeq \nabla^2 f(\mathbf{x}) + MrI_n$$

where I_n is the unit matrix \mathbb{R}^n .

对于矩阵 A 和 B , $A \succeq B \iff A - B \succeq 0$.

3. Convex Function

3.1 Definition of The Convex Function and The Class \mathcal{F}^k

Consider the unconstrained minimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad (\text{II.3})$$

where the objective function $f(\mathbf{x})$ is smooth enough

定义 smooth 足够的问题, 但这个约束很弱, 在局部最小不能保证收敛, 不能保证上下界.

增加几个假设

Assumption II.10. For any $f \in \mathcal{F}$, first-order optimality condition is **sufficient** for a point to be a **global solution** to II.3

The possibility to verify the inclusion $f \in \mathcal{F}$ in a simple way.

Assumption II.11. If $f_1, f_2 \in \mathcal{F}$ and $\alpha, \beta \geq 0$, then $\alpha f_1 + \beta f_2 \in \mathcal{F}$.

Assumption II.12. Any linear function $f(\mathbf{x}) = \alpha + \langle a, \mathbf{x} \rangle$ belongs to \mathcal{F} .

Definition of The Convex Function Consider

$$f(\mathbf{y}) \geq f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{y} - \mathbf{x}_0 \rangle$$

函数大于线性逼近.

Definition II.13 (Convex Set). A set $\mathcal{Q} \subseteq \mathbb{R}^n$ is called **convex** if $\forall \mathbf{x}, \mathbf{y} \in \mathcal{Q}$, and $\alpha \in [0, 1]$, we have

$$\alpha \mathbf{x} + (1 - \alpha) \mathbf{y} \in \mathcal{Q}$$

(集合中任取两点, 连线都是属于集合的, \mathcal{F} 表示是凸的)

We denote by $\mathcal{F}^k(\mathcal{Q})$ the class we discussed above, and call it **class of the convex function**:

- Any $f \in \mathcal{F}^k(\mathcal{Q})$ is a **convex function**
- Any $f \in \mathcal{F}^k(\mathcal{Q})$ is k times continuously **differentiable** on \mathcal{Q}
- We assume $\mathcal{Q} = \mathbb{R}^n$ in this chapter

Definition II.14 (Convex Function). A continuously differentiable function $f(\cdot)$ is called **convex** on a convex set \mathcal{Q} (notation $f \in \mathcal{F}^1(\mathcal{Q}^n)$) if $\forall \mathbf{x}, \mathbf{y} \in \mathcal{Q}$ we have

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$$

If $-f(\mathbf{x})$ is convex, we call $f(\mathbf{x})$ **concave**.

3.2 Properties of The Convex Function

Theorem II.15 (Global Property). If $f \in \mathcal{F}^1(\mathbb{R}^n)$ and $\nabla f(\mathbf{x}^*) = 0$, then \mathbf{x}^* is the global minimum of $f(\mathbf{x})$ on \mathbb{R}^n

Lemma II.16 (Conic Combination). If f_1 and f_2 belong to $\mathcal{F}^1(\mathbb{R}^n)$, and $\alpha, \beta \geq 0$, then the function $f = \alpha f_1 + \beta f_2$ also belongs to $\mathcal{F}^1(\mathbb{R}^n)$.

Lemma II.17 (Affine Composition). If $f \in \mathcal{F}^1(\mathbb{R}^n)$, $\mathbf{b} \in \mathbb{R}^m$, and $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$, then

$$\phi(\mathbf{x}) = f(A\mathbf{x} + \mathbf{b}) \in \mathcal{F}^1(\mathbb{R}^m)$$

Lemma II.18 (Pointwise maximum and supremum).

If $f_i(\mathbf{x}), i \in I$, are convex, then

$$g(\mathbf{x}) = \max_{i \in I} f_i(\mathbf{x})$$

is also convex.

Lemma II.19 (Convex monotone composition).

1) If f is a convex function on \mathbb{R}^n and $F(\cdot)$ is a convex and non-decreasing function on \mathbb{R} , then $g(\mathbf{x}) = F(f(\mathbf{x}))$ is convex.

2) If $f_i, i = 1, \dots, m$ are convex functions on \mathbb{R}^n and $F(\mathbf{y}_1, \dots, \mathbf{y}_m)$ is convex and non-decreasing (component-wise) in each argument, then

$$g(\mathbf{x}) = F(f_1(\mathbf{x}), \dots, f_m(\mathbf{x}))$$

is convex.

Lemma II.20 (Partial minimization). If $f(\mathbf{x}, \mathbf{y})$ is convex in $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^n$ and Y is a convex set, then

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

is convex.

3.3 Equivalent Definitions

Theorem II.21. A continuously **differentiable** function f belongs to the class $\mathcal{F}^1(\mathbb{R}^n)$ if and only if $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and $\alpha \in [0, 1]$ we have

$$f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y})$$

(凸组合的函数 \leq 函数的凸组合)

Theorem II.22. A continuously differentiable function f belongs to the class $\mathcal{F}^1(\mathbb{R}^n)$ if and only if $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, we have

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq 0$$

Theorem II.23. A twice continuously differentiable function f belongs to class $\mathcal{F}^2(\mathbb{R}^n)$ if and only if $\forall \mathbf{x} \in \mathbb{R}^n$, we have

$$\nabla^2 f(\mathbf{x}) \succeq 0$$

3.4 Examples

1)

4. Smooth and Convex Function

4.1 The Class $\mathcal{F}_L^{k,l}(\mathbb{R}^n)$

- 1) Any function $f \in \mathcal{F}_L^{k,l}(\mathbb{R}^n)$ is convex
- 2) the meaning of the index is the same as $C_L^{k,l}(\mathbb{R}^n)$
- 3) The most important class of that type is $\mathcal{F}_L^{1,1}(\mathbb{R}^n)$, the class of convex functions with **Lipschitz continuous gradient**.

4.2 Necessary and Sufficient Conditions for The Class $\mathcal{F}_L^{1,1}(\mathbb{R}^n)$

Theorem II.24. All the conditions below, holding $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and $\alpha \in [0, 1]$, are equivalent to the inclusion $f \in \mathcal{F}_L^{1,1}(\mathbb{R}^n)$:

$$0 \leq f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \leq \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 \quad (\text{II.4})$$

$$f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2 \leq f(\mathbf{y}) \quad (\text{II.5})$$

$$\frac{1}{L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2 \leq \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \quad (\text{II.6})$$

$$0 \leq \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \leq L \|\mathbf{x} - \mathbf{y}\|^2 \quad (\text{II.7})$$

$$\begin{aligned} \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y}) &\geq f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) \\ &\quad + \frac{\alpha(1 - \alpha)}{2L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2 \end{aligned} \quad (\text{II.8})$$

$$\begin{aligned} 0 &\leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y}) - f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) \\ &\leq \alpha(1 - \alpha) \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 \end{aligned} \quad (\text{II.9})$$

4.3 Necessary and Sufficient Conditions for The Class $\mathcal{F}_L^{2,1}(\mathbb{R}^n)$

Theorem II.25. Twice continuously differentiable function $f \in \mathcal{F}_L^{2,1}(\mathbb{R}^n)$ if and only if $\forall \mathbf{x} \in \mathbb{R}^n$, we have

$$0 \preceq \nabla^2 f(\mathbf{x}) \preceq L I_n$$

5. Strongly Convex Function

5.1 Definition of The Strongly Convex and The Class $\mathcal{S}_\mu^1(\mathbb{R}^n)$

Definition II.26. A continuously differentiable function $f(\mathbf{x})$ is called **strongly convex** on \mathbb{R}^n (notation $f \in \mathcal{S}_\mu^1(\mathbb{R}^n)$) if there exists a constant $\mu > 0$ such that $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, we have

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2} \mu \|\mathbf{y} - \mathbf{x}\|^2$$

μ is called **convexity parameter** (凸参数) of f .

5.2 Property of Strongly Convex Function

Theorem II.27. If $f \in \mathcal{S}_\mu^1(\mathbb{R}^n)$ and $\nabla f(\mathbf{x}^*) = 0$, then

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) + \frac{1}{2} \mu \|\mathbf{x} - \mathbf{x}^*\|^2, \quad \forall \mathbf{x} \in \mathbb{R}^n$$

Lemma II.28. If $f_1 \in \mathcal{S}_{\mu_1}^1(\mathbb{R}^n), f_2 \in \mathcal{S}_{\mu_2}^1(\mathbb{R}^n)$ and $\alpha, \beta \geq 0$ then,

$$f = \alpha f_1 + \beta f_2 \in \mathcal{S}_{\alpha\mu_1 + \beta\mu_2}^1(\mathbb{R}^n)$$

加线性函数增加凸性, 让多解变为单解.

Theorem II.29. If $f \in \mathcal{S}_\mu^1(\mathbb{R}^n)$, then $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, we have

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2\mu} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2 \quad (\text{II.10})$$

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \leq \frac{1}{\mu} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2 \quad (\text{II.11})$$

$\phi(\nu)$

5.3 Equivalent Definitions

Theorem II.30. Let f be continuously differentiable. Both conditions below, holding $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and $\alpha \in [0, 1]$, are equivalent to inclusion $\mathcal{S}_\mu^1(\mathbb{R}^n)$

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \mu \|\mathbf{x} - \mathbf{y}\|^2 \quad (\text{II.12})$$

$$\alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}) \geq f(\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}) + \alpha(1 - \alpha)\frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|^2 \quad (\text{II.13})$$

$$\mathcal{S}_\mu^2(\mathbb{R}^n) \subseteq \mathcal{S}_\mu^1(\mathbb{R}^n)$$

Theorem II.31. Two times continuously differentiable functions f belongs to the class $\mathcal{S}_\mu^2(\mathbb{R}^n)$ if and only if $\mathbf{x} \in \mathbb{R}^n$

$$\nabla^2 f(\mathbf{x}) \succeq \mu I_n$$

$$\mathbf{s} \in \mathbb{R}^n$$

$$\langle I_n \mathbf{s}, \mathbf{s} \rangle$$

5.4 Examples

1)

6. Smooth and Strongly Convex Function

6.1 The Class $\mathcal{S}_{\mu,L}^{1,1}(\mathbb{R}^n)$

The value $Q_f = \frac{L}{\mu} \geq 1$ is called the **condition number**(条件数) of function f .

6.2 Property of Smooth and Strongly Convex Function

Theorem II.32. If $f \in \mathcal{S}_{\mu,L}^{1,1}(\mathbb{R}^n)$, then $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, we have

$$\begin{aligned} \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle &\geq \frac{\mu L}{\mu + L} \|\mathbf{x} - \mathbf{y}\|^2 \\ &\quad + \frac{1}{\mu + L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2 \end{aligned}$$

退化后和... 与... 相关

7. Conclusion

7.1 Upper Bounds on Functional Components

Lipschitz Continuity: $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,

- Zerothorder Condition:

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq L \|\mathbf{x} - \mathbf{y}\|$$

- First-order Condition:

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$$

- p -order Condition

$$\|\nabla^p f(\mathbf{x}) - \nabla^p f(\mathbf{y})\|_* \leq L \|\mathbf{x} - \mathbf{y}\|, \quad p \geq 2$$

7.2 Lower Bounds on Functional Components

Convexity: $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $\alpha \in [0, 1]$,

- Zerothorder Condition:

$$\alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}) - f(\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}) \geq 0$$

- First-order Condition:

$$\begin{aligned} D_f(\mathbf{x}, \mathbf{x}) &\triangleq f(\mathbf{y}) - \{f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle\} \geq 0 \\ \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle &\geq 0 \end{aligned}$$

- Second-order Condition:

$$\nabla^2 f(\mathbf{x}) \succeq 0$$

Strong Convexity: $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,

- Zerothorder Condition:

$$\begin{aligned} \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}) - f(\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}) \\ \geq \alpha(1 - \alpha)\frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|^2 \end{aligned}$$

- First-order Condition:

$$\begin{aligned} D_f(\mathbf{x}, \mathbf{y}) &\geq \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|^2 \\ \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle &\geq \mu \|\mathbf{x} - \mathbf{y}\|^2 \end{aligned}$$

- Second-order Condition

$$\nabla^2 f(\mathbf{x}) \succeq \mu I_n$$

7.3 Other Lower Bounds on Functional Components

- Weak Strong Convexity(WSC)

$$D_f(\mathbf{x}^*, \mathbf{y}) \geq \frac{\mu}{2} \|\mathbf{x}^* - \mathbf{x}\|^2, \quad \forall \mathbf{x} \in \mathbb{R}^n$$

- Restricted Secant Inequality(RSI)

$$\langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{x}^* \rangle \geq \mu \|\mathbf{x}^* - \mathbf{x}\|^2, \quad \forall \mathbf{x} \in \mathbb{R}^n$$

- Polyak-Łojaciewicz(PL)

$$\frac{1}{2} \|\nabla f(\mathbf{x})\|^2 \geq \mu(f(\mathbf{x}) - f(\mathbf{x}^*)), \quad \forall \mathbf{x} \in \mathbb{R}^n$$

- Error Bounds(EB)

$$\|\nabla f(\mathbf{x})\| \geq \mu \|\mathbf{x}^* - \mathbf{x}\|, \forall \mathbf{x} \in \mathbb{R}^n$$

- Quadratic Growth(QG)

$$f(\mathbf{x}) - f(\mathbf{x}^*) \geq \frac{\mu}{2} \|\mathbf{x}^* - \mathbf{x}\|^2, \forall \mathbf{x} \in \mathbb{R}^n$$

(SC)→(WSC)→(RSI)→(EB)≡(PL)→(QG)

If f is convex, (RSI)≡(EB)≡(PL)≡(QG)

Theorem II.33. μ -strongly convex function $\rightarrow \mu$ -PL

III Descent Meethod

1. Gradient Descent

1.1 Basic Scheme

Gradient Descent Formulation We will
Gradient Method

- Choose $\mathbf{x}_0 \in \mathbb{R}^n$
- Iterate $\mathbf{x}_{k+1} = \mathbf{x}_k - h_k \nabla f(\mathbf{x}_k), k = 0, 1, \dots$

Step Size

1.2 Performance for $C_L^{1,1}(\mathbb{R}^n)$

- 1) in advance
- 2) Full relaxation (贪心)
- 3) G-A (折中)

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}) + ..$$

$$f(\mathbf{x}_{k+1}) \geq f(\mathbf{x}) + ..$$

Geometric interpretation

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

where $f \in C_L^{1,1}(\mathbb{R}^n)$

ref (1.2.5)

$$f(\mathbf{y}) \leq f(\mathbf{x}) - h \left(1 - \frac{h}{2} L\right) \|\nabla f(\mathbf{x})\|^2$$

Remark: $\forall h \in (0, \frac{2}{L}), h(1 - \frac{h}{2} L) \|\nabla f(\mathbf{x})\|^2 \geq 0$.

step-size strategies:

- 1) $h_k = \frac{1}{L}$
- 2) $h_k = \frac{1}{L}$
- 3) $h_k \geq \frac{2}{L}(1 - \beta)$

步长选取影响 ω

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leq -\frac{\omega}{L} \|\nabla f(\mathbf{x}_k)\|^2$$

where ω is some positive constant.

$$\|\nabla f(\mathbf{x}_k)\| \rightarrow 0 \text{ as } k \rightarrow \infty$$

方法的收敛和梯度与步长相关. 每一步下降但最后是不一定收敛的 (需要一个下界).

$$g_T^* = \min g_k$$

$$g_T^* \leq \frac{1}{\sqrt{T+1}} \square$$

速度是 $\frac{1}{\sqrt{T+1}}$

可得到一个复杂度上界 $g_T^* \leq \epsilon$

步长收敛相当于主动构造序列, 保证了递减与有下界, 但收敛点的意义不能保证.

e.g.

1.3 Performance for $\mathcal{F}_L^{1,1}(\mathbb{R}^n)$

Theorem III.1.

$$f(\mathbf{x}_k) - f^* \leq \frac{2(f(\mathbf{x}_0) - f^*) \|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2\|\mathbf{x}_0 - \mathbf{x}^*\|^2 + kh(2 - Lh)(f(\mathbf{x}_0) - f^*)}$$

Corollary III.2. If

$$f(\mathbf{x}_k) - f^* \leq \frac{2L \|\mathbf{x}_0 - \mathbf{x}^*\|^2}{k + 4}$$

1.4 Performance for $\mathcal{S}_{\mu,L}^{1,1}(\mathbb{R}^n)$

Theorem III.3.

$$\|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left(1 - \frac{2h\mu L}{\mu + L}\right)^k \|\mathbf{x}_0 - \mathbf{x}^*\|^2$$

2. General Descent Directions

2.1 Choosing the Direction

3. Newton Method

3.1 Basic Scheme

Historical Origins 寻找单变量函数的根.

可扩展到解非线性系统. non degenerate

Basic Scheme 由非约束最小化问题到寻找非线性系统的根.

Remark:

- ss
- 是可能 diverge(发散) 的

Example:

Damped Newton Method 加阻尼防止发散

3.2 Local Convergence of The Newton Method

1)

考虑:

111

(III.1)

第一步 $\nabla f(\mathbf{x}_k)$ 进行 Taylor 零阶展开.

contracting mapping(收缩映射)

Bound for $[\nabla^2 f(\mathbf{x}_k)]^{-1}$

The rate of convergence of this type is called quadratic

$$\lim_{k \rightarrow \infty} \frac{a_{k+1}}{a_k^2} = r$$

Theorem III.4. $f(x)$

3.3 Convergence Analysis

- GD
- Lightweight Newton G_k
- Variable Metric(可变量方法) or Quasi-Newton (拟牛顿)
- New Inner Product (W.R.t)

Example: (α 应该是 \mathbf{a})

Variable metric method

- 1) Choose
- 2) iter
- a.

Quasi-Newton rule: 符合公式, 能迭代, 能收敛

Example:

- Rank-one correction
- (DFP)
- (BFGS)

Remark:

- 1) n iter
- 2) superlinear

$$\lim_{k \rightarrow \infty} \frac{a_{k+1}}{a_k} = 0$$

3) 对全局最优并不必 GD 好.

4) 拟牛顿最大的问题是有过多附加运算, 用以存储与计算 $n \times n$ 矩阵. 以此提出新的方法.

4. Conjugate Gradient

4.1 Historical Origins

迭代就是一个线性组合, 直接在子空间中寻找解.

Krylov Subspace To solve (这是一类方法, 例如)

Remark: 这些方法可以被看做为投影.

CG CG

$$\mathbf{x}_k = \operatorname{argmin}\{f(\mathbf{x}) | \mathbf{x} \in \mathbf{x}_0 + \mathcal{L}_k\}, \quad k \geq 1$$

4.2 Fundamental Theory

Lemma III.5. $\forall k \geq 1$

Lemma III.6. $\forall k, i \geq 0$

Corollary III.7. *finite*

Corollary III.8. $\forall p$

Lemma III.9. Let $\delta_i, \forall k \neq i$

4.3 CG Algorithm

Conjugate Gradient Method

1) 00

2) 11

a.

The specification of the coefficient β_k :

•

Recall:

• 11

• 做很多次 n 次最后可以收敛于 global

local 达到了 quadratic convergence

但全局并没有理论支持其更好.

IV Acceleration Methods

1. Lower Complexity Bounds for $\mathcal{F}_L^{\infty,1}(\mathbb{R}^n)$

1.1 Problem Class

Consider an optimization problem where objectives function comes from $\mathcal{F}_L^{\infty,1}(\mathbb{R}^n)$ (and $\mathcal{F}_L^{1,1}(\mathbb{R}^n)$). The problem class is as following.

Model	$\min_{\mathbf{x} \in \mathbb{B}^n} f(\mathbf{x}), \quad f \in \mathcal{F}_L^{1,1}(\mathbb{R}^n)$
Oracle	First-order Black Box
Approximate solution	$\bar{\mathbf{x}} \in \mathbb{R}^n : f(\bar{\mathbf{x}}) - f^* \leq \epsilon$

Assumption IV.1.

$$\mathbf{x}_k = \mathbf{x}_0 + \operatorname{Lin}\{\nabla f(\mathbf{x}_0), \dots, \nabla f(\mathbf{x}_{k-1})\}, \quad k \geq 1$$

不把 \mathbf{x}_0 写进 Lin , 表明对起点不关注.

1.2 Worst Function in $\mathcal{F}_L^{\infty,1}(\mathbb{R}^n)$

$$f_k(\mathbf{x}) = \frac{L}{4} \left\{ \frac{1}{2} \left[\left(\mathbf{x}^{(1)} \right)^2 + \sum_{i=1}^{k-1} \left(\mathbf{x}^{(i)} - \mathbf{x}^{(i+1)} \right)^2 + \left(\mathbf{x}^{(k)} \right)^2 \right] - \mathbf{x}^{(1)} \right\}$$

for $k = 1, \dots, n$

Remark: We have $f_k(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \left(\frac{L}{4} A_k \right) \mathbf{x} - \frac{L}{4} e_1 \mathbf{x}$.

$$\nabla^2 f(\mathbf{x}) = \frac{L}{4} A_k$$

$$A_k =$$

$$(A_k \preceq 4I, \text{ for } k \in [1, n])$$

As assumption $0 \preceq \nabla^2 f(\mathbf{x}) \preceq LI_n$

About the optimal value:

$$\bar{\mathbf{x}}^{(i)} = \begin{cases} 1 - \frac{i}{k+1} & i \in [1, k] \\ 0 & i \in [k+1, n] \end{cases}$$

Therefore, f_k

$$f_k =$$

Bounds of $\bar{\mathbf{x}}$

$\mathbb{R}^{k,n}$ sub-space

1.3 Theoretical Analysis and Main Results

Lemma IV.2.

$$\mathbf{x}_k \in \mathcal{L} =$$

Corollary IV.3.

$$f_p(\mathbf{x}_k) \geq f_k^*$$

Theorem IV.4. $1 \leq k \leq \frac{1}{2}(n-1)$

$$f(\mathbf{x}_k) - f^* \geq \frac{3L \|\mathbf{x}_0 - \mathbf{x}^*\|^2}{32(k+1)^2}$$

$$\|\mathbf{x}_k - \mathbf{x}^*\|^2 \geq \frac{1}{8} \|\mathbf{x}_0 - \mathbf{x}^*\|^2$$

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

- 1
- 有限维可能可以有更好的复杂度.

2. Lower Complexity Bounds for $\mathcal{S}_L^{\infty,1}(\mathbb{R}^n)$

2.1 Worst Function in $\mathcal{S}_L^{\infty,1}(\mathbb{R}^n)$

$$\|\mathbf{x}\|^2 = \sum_{i=1}^{\infty} \left(\mathbf{x}^{(i)}\right)^2 < \infty$$

收敛

$$f_{\mu, Q_f}(\mathbf{x}) = \frac{\mu(Q_f - 1)}{8} \left\{ \left(\mathbf{x}^{(1)}\right)^2 + \sum_{i=1}^{\infty} \left(\mathbf{x}^{(i)} - \mathbf{x}^{(i+1)}\right)^2 + -2\mathbf{x}^{(1)} \right\} + \frac{\mu}{2} \|\mathbf{x}\|^2$$

$$Q_{f_{\mu, Q_f}} = \frac{\mu Q_f}{\mu} = Q_f$$

Theorem IV.5.

$$\|\mathbf{x}_k - \mathbf{x}^*\|^2 \geq ()^{2k} \|\mathbf{x}_0 - \mathbf{x}^*\|^2$$

$$f(\mathbf{x}_k) - f^* \geq \frac{\mu}{2} ()^{2k} \|\mathbf{x}_0 - \mathbf{x}^*\|^2$$

3. Basic Scheme

3.1 Difference bewteen Lower Bounds and Real Efficiency

Efficiency Estimation upper vs. lower

Estimate Sequences

Definition IV.6. of function of $f(\mathbf{x})$

Lemma IV.7.

$$f(\mathbf{x}_k) \leq \phi_k^* \equiv \min_{\mathbf{x} \in \mathbb{R}^n} \phi_k(\mathbf{x})$$

Then $f(\mathbf{x}_k) - f^* \leq \lambda_k [\phi^*(\mathbf{x}^*) - f^*] \rightarrow 0$

However,

- 1) 不知道怎么估计
- 2) 不知道如何满足条件

Lemma IV.8. Assume that

- 1)

Lemma IV.9. Let

Constructing \mathbf{x}_k

3.2 Optimal Scheme

4. Theoretical Analysis and Variants

4.1 Analysis of Optimal Scheme

Lemma IV.10. *www*

Lemma IV.11.

$$\lambda_k \leq \min \{ \}$$

Lemma IV.12.

$$f(\mathbf{x}_k) - f^* \leq L \min \{ \}$$

4.2 Variants of Optimal Scheme

...

Theorem IV.13.

...

Other methods

- 1) 1
- 2) 2
- 3) 3

Lemma IV.14.

V General Convex Problem

1. Motivation and Definitions

1.1 Motivation

$$\begin{aligned} & \min f_0(\mathbf{x}) \\ \text{s. t. } & f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m \\ & \mathbf{x} \in Q \subseteq \mathbb{R}^n \end{aligned}$$

non-differentiable

Definition V.1 (Interior Point). An element $\mathbf{x} \in C \subset \mathbb{R}^n$ is called an interior point of C if there exists an $\epsilon > 0$ for which

$$\{\mathbf{y} \mid \|\mathbf{y} - \mathbf{x}\|_2 \leq \epsilon\} \subset C$$

以及另一个 related interior point

Definition V.2 (Open and Closed Set). A set C is open if $\text{int } C = C$, i.e. every point in C is an interior point. A set $C \subset \mathbb{R}^n$ is closed if its complement $\mathbb{R}^n / C = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{x} \notin C\}$ is open.

经常出现不光滑的优化问题.

- max-type functions
- 隐式给出问题

Denoted by

$$\text{dom } f = \{\mathbf{x} \in \mathbb{R}^n : |f(\mathbf{x})| < \infty\}$$

the domain of function f . We always assume that $\text{dom } f \neq \emptyset$

1.2 Definition of Convex Function

Definition V.3. A function $f(\mathbf{x})$ is called convex, if its domain is convex and $\forall \mathbf{x}, \mathbf{y} \in \text{dom } f$ and $\alpha \in [0, 1]$ the following inequality holds

$$f(\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y})$$

If this inequality is strict(严格小于), the function is called strictly convex. We call f concave if $-f$ is Convex.

之前学的优化都是 gradients, 对于不光滑的, 需要更多的技巧.

Lemma V.4 (Jensen's Inequality). For any $\mathbf{x}_1, \dots, \mathbf{x}_m \in \text{dom } f$ and positive coefficients $\alpha_1, \dots, \alpha_m$ such that

$$\sum_{i=1}^m \alpha_i = 1, \quad \alpha_i \geq 0, \quad i = 1, \dots, m \quad (\text{V.1})$$

we have

$$f\left(\sum_{i=1}^m \alpha_i \mathbf{x}_i\right) \leq \sum_{i=1}^m \alpha_i f(\mathbf{x}_i)$$

A point $\mathbf{x} = \sum_{i=1}^m \alpha_i \mathbf{x}_i$ with positive coefficients α_i satisfying the normalizing condition V.1 is called a convex combination(凸组合) of points $\{\mathbf{x}_i\}_{i=1}^m$.

Corollary V.5. Let \mathbf{x} be a convex combination of points $\mathbf{x}_1, \dots, \mathbf{x}_m$, then

$$f(\mathbf{x}) \leq \max_{1 \leq i \leq m} f(\mathbf{x}_i)$$

Corollary V.6. Let

$$\begin{aligned} \Delta &= \text{Conv}\{\mathbf{x}_1, \dots, \mathbf{x}_m\} \\ &= \left\{ \mathbf{x} = \sum_{i=1}^m \alpha_i \mathbf{x}_i \mid \alpha_i \geq 0, \sum_{i=1}^m \alpha_i = 1 \right\} \end{aligned}$$

Conv means convex combination.

$$\text{Then } \max_{\mathbf{x} \in \Delta} f(\mathbf{x}) = \max_{1 \leq i \leq m} f(\mathbf{x}_i)$$

Theorem V.7. A function f is convex iff $\forall \mathbf{x}, \mathbf{y} \in \text{dom } f$ and $\beta \geq 0$ such that $\mathbf{y} + \beta(\mathbf{y} - \mathbf{x}) \in \text{dom } f$, we have

$$f(\mathbf{y} + \beta(\mathbf{y} - \mathbf{x})) \geq f(\mathbf{y}) + \beta(f(\mathbf{y}) - f(\mathbf{x}))$$

Theorem V.8. A function f is convex iff its epigraph(函数线以上的区域)

$$\text{epi}(f) = \{(\mathbf{x}, t) \in \text{dom } f \times \mathbb{R} \mid t \geq f(\mathbf{x})\}$$

is a convex set.

Theorem V.9. If a function f is convex, then all level sets(层集)

$$\mathcal{L}_f(\beta) = \{\mathbf{x} \in \text{dom } f \mid f(\mathbf{x}) \leq \beta\}$$

are either convex or empty.

1.3 Other Properties

Definition V.10. A function f is called closed and convex if its epigraph is a closed set.

Theorem V.11. If convex function f is closed, then all its level sets are either empty or closed.

pathological(病态)

2. Operation with Convex Function

2.1 Invariant Operations

Definition V.12 (Lower Semicontinuity). A function f is lower-semicontinuity at a given vector $\bar{\mathbf{x}}$ if for every sequence $\{\mathbf{x}_k\}$ converging to $\bar{\mathbf{x}}$, we have

$$\liminf_{k \rightarrow \infty} f(\mathbf{x}_k) \geq f(\bar{\mathbf{x}})$$

We say that f is lower-semicontinuity over a set X if f is lower-semicontinuity at every $\mathbf{x} \in X$.

Note:

$$\liminf_{k \rightarrow \infty} \mathbf{x}_k \iff \lim_{k \rightarrow \infty} \left(\inf_{m \geq k} \mathbf{x}_m \right)$$

Theorem V.13. For a function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{-\infty, +\infty\}$ the following statements are equivalent

- 1) f is closed
- 2) Every level set of f is closed
- 3) f is lower-semicontinuity over \mathbb{R}^n

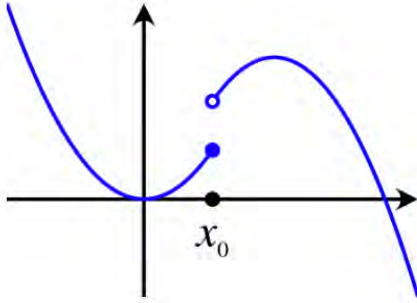


Figure V.1: f

Let us describe a set of invariant operations to create more complicated objects.

Theorem V.14. Let function f_1 and f_2 be closed and convex and let $\beta \geq 0$. Then all functions below are closed and convex:

- 1) $f(\mathbf{x}) = \beta f_1(\mathbf{x})$, $\text{dom } f = \text{dom } f_1$
- 2) $f(\mathbf{x}) = f_1(\mathbf{x}) + f_2(\mathbf{x})$, $\text{dom } f = (\text{dom } f_1) \cap (\text{dom } f_2)$
- 3) $f(\mathbf{x}) = \max\{f_1(\mathbf{x}), f_2(\mathbf{x})\}$, $\text{dom } f = (\text{dom } f_1) \cap (\text{dom } f_2)$

The following theorem demonstrates that convexity is an affine-invariant property.

Theorem V.15. Let function $\phi(\mathbf{y})$, $\mathbf{y} \in \mathbb{R}^n$ be convex and closed. Consider a linear operator

$$\mathcal{A}(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b} : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

Then $f(\mathbf{x}) = \phi(\mathcal{A}(\mathbf{x}))$ is a closed and convex function with domain

$$\text{dom } f = \{\mathbf{x} \in \mathbb{R}^n | \mathcal{A}(\mathbf{x}) \in \text{dom } \phi\}$$

The next theorem is one of the main suppliers of convex functions with implicit structure.

Theorem V.16. Let Δ be some set and

$$f(\mathbf{x}) = \sup_{\mathbf{y}} \{\phi(\mathbf{y}, \mathbf{x}) | \mathbf{y} \in \Delta\}$$

Suppose that for any fixed $\mathbf{y} \in \Delta$, the function $\phi(\mathbf{y}, \mathbf{x})$ is closed and convex in \mathbf{x} . Then $f(\mathbf{x})$ is a closed convex function with domain

$$\text{dom } f = \left\{ \mathbf{x} \in \bigcap_{\mathbf{y} \in \Delta} \text{dom } \phi(\mathbf{y}, \cdot) \mid \exists \gamma : \phi(\mathbf{y}, \mathbf{x}) \leq \gamma, \forall \mathbf{y} \in \Delta \right\}$$

3. Continuity and Differentiability

3.1 Continuity of Convex Function

Lemma V.17. Let f be convex and $\mathbf{x}_0 \in \text{int}(\text{dom } f)$. Then f is locally upper bounded at \mathbf{x}_0 .

Theorem V.18. Let f be convex and $\mathbf{x}_0 \in \text{int}(\text{dom } f)$. Then f is locally Lipschitz continuous at \mathbf{x}_0 .

Remark: 构造 \mathbf{z} 没有超出领域, \mathbf{y} 是 \mathbf{x}_0 与 \mathbf{z} 的凸组合.

3.2 Differentiability of Convex Function

Definition V.19. Let $\mathbf{x} \in \text{dom } f$. We call f differentiable in a direction \mathbf{p} at point \mathbf{x} if the following limit exists

$$f'(\mathbf{x}; \mathbf{p}) = \lim_{\alpha \downarrow 0} \frac{1}{\alpha} [f(\mathbf{x} + \alpha \mathbf{p}) - f(\mathbf{x})]$$

The value $f'(\mathbf{x}; \mathbf{p})$ is called the directional derivative of f at \mathbf{x} .

Theorem V.20. Convex f is differentiable in any direction at any interior point of its domain.

Lemma V.21. Let f be a convex function and $\mathbf{x} \in \text{int}(\text{dom } f)$. Then $f'(\mathbf{x}; \mathbf{p})$ is a convex function of \mathbf{p} , which is homogeneous of degree 1. For any $\mathbf{y} \in \text{dom } f$ we have

$$f(\mathbf{y}) \geq f(\mathbf{x}) + f'(\mathbf{x}; \mathbf{y} - \mathbf{x})$$

4. Separation Theorem

4.1 Projection

Definition V.22 (Hyperplane). Let Q be a convex set. We say that hyperplane

$$\mathcal{H}(\mathbf{g}, \gamma) = \{\mathbf{x} \in \mathbb{R}^n | \langle \mathbf{g}, \mathbf{x} \rangle = \gamma\}, \mathbf{g} \neq 0$$

is supporting to Q if any $\mathbf{x} \in Q$ satisfies inequality $\langle \mathbf{g}, \mathbf{x} \rangle \leq \gamma$.

We say that the hyperplane $\mathcal{H}(\mathbf{g}, \gamma)$ separates a point \mathbf{x}_0 from Q if

$$\langle \mathbf{g}, \mathbf{x} \rangle \leq \gamma \leq \langle \mathbf{g}, \mathbf{x}_0 \rangle \quad (\text{V.2})$$

$\forall \mathbf{x} \in Q$. If the second inequality in V.2 is strict, we call the separation strict.

Definition V.23 (Projection). Let Q be a closed set and $\mathbf{x}_0 \in \mathbb{R}^n$. Denote

$$\pi_Q(\mathbf{x}_0) = \operatorname{argmin} \{ \|\mathbf{x} - \mathbf{x}_0\| : \mathbf{x} \in Q \}$$

We call $\pi_Q(\mathbf{x}_0)$ the projection of point \mathbf{x}_0 onto the set Q .

Theorem V.24. If Q is a convex set, then there exists a unique projection $\pi_Q(\mathbf{x}_0)$

Remark: It's clear that $\pi_Q(\mathbf{x}_0) = \mathbf{x}_0 \iff \mathbf{x}_0 \in Q$.

Lemma V.25. Let Q be a closed and convex set and $\mathbf{x}_0 \notin Q$. Then $\forall \mathbf{x} \in Q$, we have

$$\langle \pi_Q(\mathbf{x}_0) - \mathbf{x}_0, \mathbf{x} - \pi_Q(\mathbf{x}_0) \rangle \geq 0$$

从 \mathbf{x} 到 $\pi_Q(\mathbf{x}_0)$ 的向量与从 $\pi_Q(\mathbf{x}_0)$ 到 \mathbf{x}_0 的向量的夹角是锐角.

Lemma V.26. $\forall \mathbf{x} \in Q$, we have

$$\|\mathbf{x} - \pi_Q(\mathbf{x}_0)\|^2 + \|\pi_Q(\mathbf{x}_0) - \mathbf{x}_0\|^2 \leq \|\mathbf{x} - \mathbf{x}_0\|^2$$

4.2 Main Theorems

Theorem V.27. Let Q be a closed set and $\mathbf{x}_0 \notin Q$. Then there exists a hyperplane $\mathcal{H}(\mathbf{g}, \gamma)$, which strictly separates \mathbf{x}_0 from Q . Namely, we can take

$$\begin{aligned} \mathbf{g} &= \mathbf{x}_0 - \pi_Q(\mathbf{x}_0) \neq 0 \\ \gamma &= \langle \mathbf{x}_0 - \pi_Q(\mathbf{x}_0), \pi_Q(\mathbf{x}_0) \rangle \end{aligned}$$

$$\psi(\mathbf{g}) = \sup \{ \langle \mathbf{g}, \mathbf{x} \rangle \mid \mathbf{x} \in Q \}$$

Corollary V.28. Let Q_1 and Q_2 be two closed convex sets.

- 1) If $\forall \mathbf{g} \in \operatorname{dom} \psi_{Q_2}$ we have $\psi_{Q_1}(\mathbf{g}) \leq \psi_{Q_2}(\mathbf{g})$, then $Q_1 \subseteq Q_2$.
- 2) Let $\operatorname{dom} \psi_{Q_1} = \operatorname{dom} \psi_{Q_2}$ and $\forall \mathbf{g} \in \operatorname{dom} \psi_{Q_1}$, we have $\psi_{Q_1}(\mathbf{g}) = \psi_{Q_2}(\mathbf{g})$. Then $Q_1 \equiv Q_2$.

Theorem V.29. Let Q be a closed and convex set, and \mathbf{x}_0 belong to boundary of set Q . Then there exists a hyperplane $\mathcal{H}(\mathbf{g}, \gamma)$, supporting to Q and passing through \mathbf{x}_0 .

Such a vector \mathbf{g} is called supporting to Q at \mathbf{x}_0 .

5. Subgradient

5.1 Definition of Subgradient

Definition V.30. Let f be a convex function. A vector \mathbf{g} is called a subgradient of function f at point $\mathbf{x}_0 \in \operatorname{dom} f$ if for any $\mathbf{x} \in \operatorname{dom} f$, we have

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \mathbf{g}, \mathbf{x} - \mathbf{x}_0 \rangle$$

The set of all subgradient of f at \mathbf{x}_0 , is called the subdifferential of function f at point \mathbf{x}_0 .

Subgradient and Convexity

Lemma V.31. Let for any $\mathbf{x} \in \operatorname{dom} f$ subdifferential $\partial f(\mathbf{x})$ be nonempty. Then f is a convex function.

Theorem V.32. Let f be closed and convex and $\mathbf{x}_0 \in \operatorname{int}(\operatorname{dom} f)$. Then $\partial f(\mathbf{x}_0)$ is a nonempty bounded set.

Theorem V.33. Let f be a closed convex function. For any $\mathbf{x}_0 \in \operatorname{int}(\operatorname{dom} f)$ and $\mathbf{p} \in \mathbb{R}^n$, we have

$$f'(\mathbf{x}_0; \mathbf{p}) = \max \{ \langle \mathbf{g}, \mathbf{p} \rangle \mid \mathbf{g} \in \partial f(\mathbf{x}_0) \}$$

5.2 Properties of Subgradient

Theorem V.34. We have

$$f(\mathbf{x}^*) = \min_{\mathbf{x} \in \operatorname{dom} f} f(\mathbf{x}) \iff 0 \in \partial f(\mathbf{x}^*)$$

Theorem V.35. For any $\mathbf{x}_0 \in \operatorname{dom} f$, all vector $\mathbf{g} \in \partial f(\mathbf{x}_0)$ are supporting to the level set $\mathcal{L}_f(f(\mathbf{x}_0))$:

$$\langle \mathbf{g}, \mathbf{x}_0 - \mathbf{x} \rangle \geq 0$$

$$\forall \mathbf{x} \in \mathcal{L}_f(f(\mathbf{x}_0)) \equiv \{ \mathbf{x} \in \operatorname{dom} f : f(\mathbf{x}) \leq f(\mathbf{x}_0) \}$$

Corollary V.36. Let $Q \subseteq \operatorname{dom} f$ be a closed convex set, $\mathbf{x}_0 \in Q$ and

$$\mathbf{x}\mathbf{x}^* = \operatorname{argmin} \{ f(\mathbf{x}) \mid \mathbf{x} \in Q \}$$

Then $\forall \mathbf{g} \in \partial f(\mathbf{x}_0)$, we have $\langle \mathbf{g}, \mathbf{x}_0 - \mathbf{x}^* \rangle \geq 0$.

5.3 Rules for Computing

Lemma V.37. Let f be closed and convex. Assume that it's differentiable on its domain. Then

$$\partial f(\mathbf{x}) = \{ \nabla f(\mathbf{x}) \}$$

$\forall \mathbf{x} \in \operatorname{int}(\operatorname{dom} f)$

Lemma V.38. Let $f(\mathbf{y})$ be closed and convex with $\text{dom } f \subseteq \mathbb{R}^m$. Consider a linear operator

$$\mathcal{A}(\mathbf{x}) = A\mathbf{x} + \mathbf{b} : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

Then $\phi(\mathbf{x}) = f(\mathcal{A}(\mathbf{x}))$ is closed convex function with domain $\text{dom } \phi = \{\mathbf{x} | \mathcal{A}(\mathbf{x}) \in \text{dom } f\}$. $\forall \mathbf{x} \in \text{int}(\text{dom } \phi)$, we have

$$\partial \phi(\mathbf{x}) = A^\top \partial f(\mathcal{A}(\mathbf{x}))$$

Lemma V.39. Let $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ be closed convex function and $\alpha_1, \alpha_2 \geq 0$. Then function $f(\mathbf{x}) = \alpha_1 f_1(\mathbf{x}) + \alpha_2 f_2(\mathbf{x})$ is closed and convex and

$$\partial f(\mathbf{x}) = \alpha_1 \partial f_1(\mathbf{x}) + \alpha_2 \partial f_2(\mathbf{x})$$

$$\forall \mathbf{x} \in \text{int}(\text{dom } f) = \text{int}(\text{dom } f_1) \cap \text{int}(\text{dom } f_2)$$

Lemma V.40. Let functions $f_i(\mathbf{x})$, $i = 1, \dots, m$ be closed and convex. Then function

$$f(\mathbf{x}) = \max_{1 \leq i \leq m} f_i(\mathbf{x})$$

is also closed and convex. $\forall \mathbf{x} \in \text{int}(\text{dom } f) = \bigcap_{i=1}^m \text{int}(\text{dom } f_i)$, we have

$$\partial f(\mathbf{x}) = \text{Conv}\{\partial f_i(\mathbf{x}) | i \in I(\mathbf{x})\}$$

where $I(\mathbf{x}) = \{i : f_i(\mathbf{x}) = f(\mathbf{x})\}$

Lemma V.41. Let Δ be a set and $f(\mathbf{x}) = \sup\{\phi(\mathbf{y}, \mathbf{x}) | \mathbf{y} \in \Delta\}$. Suppose that for any fixed $\mathbf{y} \in \Delta$, the function $\phi(\mathbf{y}, \mathbf{x})$ is closed and convex in \mathbf{x} . Then $f(\mathbf{x})$ is closed convex.

Moreover, for any \mathbf{x} from

$$\text{dom } f = \{\mathbf{x} \in \mathbb{R}^n | \exists \gamma : \phi(\mathbf{y}, \mathbf{x}) \leq \gamma, \forall \mathbf{y} \in \Delta\}$$

we have

$$\partial f(\mathbf{x}) \supseteq \text{Conv}\{\partial \phi_{\mathbf{x}}(\mathbf{y}, \mathbf{x}) | \mathbf{y} \in I(\mathbf{x})\}$$

where $I(\mathbf{x}) = \{\mathbf{y} | \phi(\mathbf{y}, \mathbf{x}) = f(\mathbf{x})\}$

Theorem V.42. Let $\|\cdot\|$ be a vector norm in \mathbb{R}^n , then

$$\partial \|\cdot\| = \left\{ V(\mathbf{x}) \triangleq \{\mathbf{v} \in \mathbb{R}^n | \langle \mathbf{v}, \mathbf{x} \rangle = \|\mathbf{x}\|, \|\mathbf{v}\| \leq 1\} \right\}$$

where $\|\mathbf{v}\|_*$ is the dual norm of $\|\cdot\|$, defined as

$$\|\mathbf{v}\|_* \triangleq \sup_{\|\mathbf{u}\| \leq 1} \langle \mathbf{x}, \mathbf{u} \rangle$$

6. General Lower Complexity Bounds

6.1 Problem Class

We have

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

Table V.1: Problem Class

Model
Oracle
Approximate solution
Method

6.2 Resisting Oracle

Special Function family

Let us fix some constant $\mu > 0$ and $\gamma > 0$,

$$f_k(\mathbf{x}) = \gamma \max_{1 \leq i \leq k} \mathbf{x}^{(i)} + \frac{\mu}{2} \|\mathbf{x}\|^2, \quad k = 1, \dots, n$$

$$\partial f_k(\mathbf{x}) = \mu \mathbf{x} + \gamma \text{Conv}\{e_i | i \in I(\mathbf{x})\}$$

$$I(\mathbf{x}) = \left\{ 1 \leq j \leq k, \mathbf{x}^{(j)} = \max_{1 \leq i \leq k} \mathbf{x}^{(i)} \right\}$$

1) Local Lipschitz Continuity of f_k

Therefore $\forall \mathbf{x}, \mathbf{y} \in B_2(0, \rho)$, $\rho > 0$, and $g_k(\mathbf{y}) \in \partial f_k(\mathbf{y})$, we have

$$\begin{aligned} f_k(\mathbf{y}) - f_k(\mathbf{x}) &\leq \langle g_k(\mathbf{y}), \mathbf{y} - \mathbf{x} \rangle \\ &\leq \|g_k(\mathbf{y})\| \cdot \|\mathbf{y} - \mathbf{x}\| \\ &\leq (\mu\rho + \gamma) \|\mathbf{y} - \mathbf{x}\| \end{aligned}$$

Thus f_k is Lipschitz continuous on $B_2(0, \rho)$ with Lipschitz constant

$$M = \mu\rho + \gamma$$

2) minimum of f_k

Further, consider the point \mathbf{x}_k^* with the coordinates

$$(\mathbf{x}_k^*)^{(i)} = \begin{cases} -\frac{\gamma}{\mu k} & 1 \leq i \leq k \\ 0 & K = 1 \leq i \leq n \end{cases}$$

It's easy to check that $0 \in \partial f_k(\mathbf{x}_k^*)$, and therefore \mathbf{x}_k^* is the minimum of $f_k(\mathbf{x})$. Note that

$$\begin{aligned} R_k &= \|\mathbf{x}_k^*\| = \frac{\gamma}{\mu\sqrt{k}} \\ f_k^* &= -\frac{\gamma^2}{\mu k} + \frac{\mu}{2} R_k^2 = -\frac{\gamma^2}{2\mu k} \end{aligned}$$

3) The subgradient of f_k

4) Characteristic of \mathbf{x}_{i+1} generated by resisting oracle
Let us choose starting point $\mathbf{x}_0 = 0$. Denote

$$\mathbb{R}^{p,n} = \left\{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{x}^{(i)} = 0, p+1 \leq i \leq n \right\}$$

6.3 Lower Bound

Theorem V.43. For any class $\mathcal{P}(\mathbf{x}_0, R, M)$ and any $k, 0 \leq k \leq n-1$, there exists a function $f \in \mathcal{P}(\mathbf{x}_0, R, M)$ such that

$$\mathbf{x}_k \in \mathbf{x}_0 + \text{Lin} \{g(\mathbf{x}_0), \dots, g(\mathbf{x}_{k-1})\}$$

for any optimization scheme, which generates a sequence $\{\mathbf{x}_k\}$ satisfying the condition

$$f(\mathbf{x}_k) - f^* \geq \frac{MR}{2(1 + \sqrt{k+1})}$$

7. Subgradient Method

7.1 property of Subgradient

$$\min \{f(\mathbf{x}) \mid \mathbf{x} \in Q\}$$

$$\langle g(\mathbf{x}), \mathbf{x} - \mathbf{x}^* \rangle \geq 0$$

- The distance between \mathbf{x} and \mathbf{x}^* is decreasing in the direction $-g(\mathbf{x})$

- In equality (4) cuts \mathbb{R}^n into two half-spaces. Only one of them contains \mathbf{x}^*

7.2 Main Lemma

Let us fix some $\bar{\mathbf{x}} \in \mathbb{R}^n$. For $\mathbf{x} \in \mathbb{R}^n$ with $g(\mathbf{x}) \neq 0$ define

$$v_f(\bar{\mathbf{x}}, \mathbf{x}) = \frac{1}{\|g(\mathbf{x})\|} \langle g(\mathbf{x}), \mathbf{x} - \bar{\mathbf{x}} \rangle$$

If $g(\mathbf{x}) = 0$, then define $v_f(\bar{\mathbf{x}}, \mathbf{x}) = 0$. Clearly, $v_f(\bar{\mathbf{x}}, \mathbf{x}) \leq \|\mathbf{x} - \bar{\mathbf{x}}\|$

Let us introduce a function that measures the variation of function f with respect to the point $\bar{\mathbf{x}}$. For $t \geq 0$, define

$$\omega_f(\bar{\mathbf{x}}; t) = \max \{f(\mathbf{x}) - f(\bar{\mathbf{x}}) \mid \|\mathbf{x} - \bar{\mathbf{x}}\| \leq t\}$$

If $t < 0$, we set $\omega_f(\bar{\mathbf{x}}; t) = 0$. Clearly, the function ω_f possesses the following Properties:

- 1) For all $t < 0$, $\omega_f(\bar{\mathbf{x}}; 0) = 0$
- 2) $\omega_f(\bar{\mathbf{x}}; t)$ is a nondecreasing function of t , $t \in \mathbb{R}^1$

$$3) f(\mathbf{x}) - f(\bar{\mathbf{x}}) \leq \omega_f(\bar{\mathbf{x}}; \|\mathbf{x} - \bar{\mathbf{x}}\|)$$

Lemma V.44. For any $\mathbf{x} \in \mathbb{R}^n$, we have

$$f(\mathbf{x}) - f(\bar{\mathbf{x}}) \leq \omega_f(\bar{\mathbf{x}}; v_f(\bar{\mathbf{x}}, \mathbf{x}))$$

If $f(\mathbf{x})$ is Lipschitz continuous on $B_2(\bar{\mathbf{x}}, R)$ with some M , then

$$f(\mathbf{x}) - f(\bar{\mathbf{x}}) \leq M(v_f(\bar{\mathbf{x}}, \mathbf{x}))_+$$

$\forall \mathbf{x} \in \mathbb{R}^n$ with $\omega_f(\bar{\mathbf{x}}; \mathbf{x}) \leq R$

Definition V.45. Let $\{\mathbf{x}_i\}_{i=0}^\infty$ be a sequence in Q . Define

$$S_k = \{\mathbf{x} \in Q \mid \langle g(\mathbf{x}_i), \mathbf{x}_i - \mathbf{x} \rangle \geq 0, i = 0, \dots, k\}$$

We call this set the location set of problem (3) generated by sequence $\{\mathbf{x}_i\}_{i=0}^\infty$.

Note that in view of inequality (4) $\forall k \geq 0$, we have $\mathbf{x}^* \in S_k$. Denote

$$v_i = v_f(\mathbf{x}^*; \mathbf{x}_i) \geq 0$$

$$v_k^* = \min_{0 \leq i \leq k} v_i$$

Thus

$$v_k^* = \max \{r \mid \langle g(\mathbf{x}_i), \mathbf{x}_i - \mathbf{x} \rangle \geq 0, i = 0, \dots, k, \forall \mathbf{x} \in B_2(\mathbf{x}^*, r)\}$$

Lemma V.46. Let $f_k^* = \min_{0 \leq i \leq k} f(\mathbf{x}_i)$. Then $f_k^* - f^* \leq \omega_f(\mathbf{x}^*; v_k^*)$

7.3 Scheme for Non-smooth Problem

7.4 Main Theorem

Theorem V.47. Let f be Lipschitz continuous on $B_2(\mathbf{x}^*, R)$ with constant M and $\mathbf{x}_0 \in B(\mathbf{x}^*, R)$. Then

$$f_k^* - f^* \leq M \frac{R^2 + \sum_{i=0}^k h_i^2}{2 \sum_{i=0}^k h_i}$$

8. Frank-Wolfe Algorithm

8.1 Problems

Algorithm 1 Frank-Wolfe Algorithm

8.2 Examples

8.3 Convergence Theory

Definition V.48. We will say that $\mathbf{x}^* \in \mathcal{D}$ is a stationary point if

$$\langle \nabla f(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle \geq 0$$

$\forall \mathbf{x} \in \mathcal{D}$

Definition V.49. We denote by g_k the Frank-Wolfe gap, defined as

$$g_k \equiv \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{s}_k \rangle$$

Lemma V.50.

VI Beyond The Black-box Model

1. Proximal Gradient Method

1.1 Proximal Operator

1.2 Properties of Proximal Operator

1.3 Analysis for Proximal Gradient Method

1.4 Accelerated Proximal Gradient Method

1.5 Special case: Proximal Point Method

2. Douglas-Rachford Splitting

2.1 Different Setting for Convex Problem

2.2 Fixed Point for Nonsmooth Composition

2.3 Splitting Algorithm

3. Duality Principle

3.1 Duality Principle

Example: Duality in Linear Programs

Primal LP:

$$\min_{\mathbf{x}} \mathbf{c}^\top \mathbf{x}$$

Duality Problem

$$\min_{\mathbf{x}} f(\mathbf{x}), \text{ s. t. } g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m$$

4. Lagrangian Duality and Algorithms

4.1 Lagrangian Duality

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s. t.} \quad & h(\mathbf{x}) \leq 0 \\ & l_i(\mathbf{x}) = 0 \end{aligned}$$

因为 min max 的互换, 所以有个 gap

4.2 KKT

4.3 Algorithm using Lagrangian duality

5. Fenchel conjugate and algorithm

5.1 Fenchel conjugate

5.2 Properties

5.3 Fenchel duality

6. Smoothing Techniques

6.1 Introduction

$$\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$$

where f is convex and nonsmooth. Approximate $f(\mathbf{x})$ by a smooth and convex $f_u(\mathbf{x})$.

$$\min_{\mathbf{x} \in \mathcal{X}} f_u(\mathbf{x})$$

where f_u is a L_u -Lipschitz continuous, smooth and convex.

e.g. (Motivation)

Huber function

$$f_u(\mathbf{x}) = \begin{cases} \frac{x^2}{2u} & |\mathbf{x}| \leq u \\ |\mathbf{x}| - \frac{u}{2} & |\mathbf{x}| > u \end{cases}$$

1) $f_u(\mathbf{x})$ is clearly continuous and differentiable everywhere.

$$2) f(\mathbf{x}) - \frac{u}{2} \leq f_u(\mathbf{x}) \leq f(\mathbf{x})$$

3) If $u \rightarrow 0$, then $f_u(\mathbf{x}) \rightarrow f(\mathbf{x})$

4) $|f_u''(\mathbf{x})| \leq \frac{1}{u}$. This implies that $f_u(\mathbf{x})$ is $\frac{1}{u}$ -Lipschitz continuous.

Major Techniques:

1) Nesterov's smoothing Techniques: proximity function(近邻函数)

2) Moreau-Yosida smoothing regularization: envelope(包络)

3) Ben-Tal-Teboulle smoothing based on recession function

4) Randomized smoothing

6.2 Nesterov's Smoothing

$$\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) \iff f_u(\mathbf{x}) = \max_{\mathbf{y} \in \text{dom } f^*} \{\mathbf{x}^\top \mathbf{y} - f^*(\mathbf{y}) - u d(\mathbf{y})\}$$

Assume that function f can be represented by

$$f(\mathbf{x}) = g(A\mathbf{x} + b) \triangleq \max_{\mathbf{y} \in \mathcal{Y}} \{\langle A\mathbf{x} + b, \mathbf{y} \rangle - \phi(\mathbf{y})\}$$

where $\phi(\mathbf{y})$ is a convex and continuous function and \mathcal{Y} is a convex and compact set.

Proximity Function The function $d(\mathbf{y})$

Proposition

Theorem VI.1. For any $u > 0$, let $D_{\mathcal{Y}}^2 = \max_{\mathbf{y} \in \mathcal{Y}} d(\mathbf{y})$, we have

$$f(\mathbf{x}) - u D_{\mathcal{Y}}^2 \leq f_u(\mathbf{x}) \leq f(\mathbf{x})$$

Analysis of Nesterov's smoothing

1)

e.g.

6.3 Moreau-Yosida Regularization

7. Generalized Distance: Mirror Descent

两大动机: general, bound 界更好

7.1 Motivation

7.2 Bregman Divergence

7.3 Mirror Descent

VII Stochastic Optimization