

Optimization Algorithm Notes

Depu Meng

Feb. 2019

1 Introduction to Optimization Algorithms

1.1 Goal of the Course

- Understand foundations of optimization
- Learn to analyze widely used optimization algorithms
- Be familiar with implementation of optimization algorithms

1.2 Topics Involved

- Unconstrained optimization
- Constrained optimization
- Convex optimization
- Sparse optimization
- Stochastic optimization
- Combinational optimization
- Global optimization

1.3 Basic Concepts

Problem Definition Find the value of the decision variable s.t. objective function is maximized/minimized under certain conditions.

$$\min f(x) \quad (1)$$

$$s.t. x \in \mathcal{S} \subset \mathbb{R}^n \quad (2)$$

Here, we call \mathcal{S} *feasible region*.

We often denote constrained optimization Problem as

$$\min f(x) \quad (3)$$

$$s.t. \quad g_i(x) \geq 0, i = 1, \dots, n \quad (4)$$

$$b_i(x) = 0, i = 1, \dots, m \quad (5)$$

Definition 1. *Global Optimality.* For global optimal value $x^* \in \mathcal{S}$,

$$f(x^*) \leq f(x), \forall x \in \mathcal{S} \quad (6)$$

Definition 2. *Local Optimality.* For local optimal value $x^* \in \mathcal{S}$, $\exists U(x^*)$, such that

$$f(x^*) \leq f(x), \forall x \in \mathcal{S} \cap U(x^*) \quad (7)$$

Definition 3. *Feasible direction.* Let $x \in \mathcal{S}$, $d \in \mathbb{R}^n$ is a non-zero vector. if $\exists \delta > 0$, such that

$$x + \lambda d \in \mathcal{S}, \forall \lambda \in (0, \delta) \quad (8)$$

Then d is a **feasible direction** at x . We denote $F(x, \mathcal{S})$ as the set of feasible directions at x .

Definition 4. *Descent direction.* $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$, $x \in \mathbb{R}^n$, d is a non-zero vector. If $\exists \delta > 0$, such that

$$f(x + \lambda d) < f(x), \forall \lambda \in (0, \delta) \quad (9)$$

Then d is a **descent direction** at x . We denote $D(x, f) = \{d \mid \nabla f(x)^T d < 0\}$ as the set of descent direction at x .

1.4 Optimal Conditions

Unconstrained Optimization

First-order necessary condition: $f(x)$ is differentiable at x ,

$$\nabla f(x) = 0 \quad (10)$$

Second-order necessary condition: $f(x)$ is second-order differentiable at x ,

$$\nabla f(x) = 0 \quad (11)$$

$$\nabla^2 f(x) \geq 0 \quad (12)$$

Constrained Optimization

Theorem 1. Fritz-John Condition

For constrained optimization problem

$$\min f(x) \quad (13)$$

$$\text{s.t. } g_i(x) \geq 0, i = 1, \dots, n \quad (14)$$

$$h_i(x) = 0, i = 1, \dots, m \quad (15)$$

Denote $I(x) = \{i \in \{1, \dots, n\} \mid g_i(x) = 0\}$. For $x \in \mathcal{S}$, f and $g_i, i \in I(x)$ is differentiable at x , $h_j(x)$ is continuously differentiable at x . If x is local optimal, then there exists non-trivial $\lambda_0, \lambda_i \geq 0, i \in I(x)$ and μ_j , such that

$$\lambda_0 \nabla f(x) - \sum_{i \in I(x)} \lambda_i \nabla g_i(x) - \sum_{j=1}^m \mu_j \nabla h_j(x) = 0 \quad (16)$$

Proof. (i) If $\{\nabla h_j(x)\}$ is linearly dependent, then there exists non-trivial μ_j , such that

$$\sum_{j=1}^m \nabla \mu_j h_j(x) = 0 \quad (17)$$

Let $\lambda_0, \lambda_i, i \in I(x) = 0$, then (13) holds.

(ii) If $\{\nabla h_j(x)\}$ is linearly independent, Denote

$$F_g = F(x, g) = \{d \mid \nabla g_i(x)^T d > 0, i \in I(x)\} \quad (18)$$

$$F_h = F(x, h) = \{d \mid \nabla h_j(x)^T d = 0, j = 1, \dots, m\} \quad (19)$$

If x is a optimal value, then apparently $F(x, \mathcal{S}) \cap D(x, f) = \emptyset$. Due to the independence of $\{\nabla h_j(x)\}$, we have $F_g \cap F_h \subset F(x, \mathcal{S})$, then

$$F_g \cap F_h \cap D(x, f) = \emptyset \quad (20)$$

that is

$$\begin{cases} \nabla f(x)^T d < 0 \\ \nabla g_i(x)^T d > 0, i \in I(x) \\ \nabla h_j(x)^T d = 0, j = 1, \dots, m \end{cases} \quad (21)$$

has no solution. Let

$$A = \{\nabla f(x)^T, -\nabla g_i(x)^T, i \in I(x)\} \quad (22)$$

$$B = \{-\nabla h_j(x)^T, j = 1, \dots, m\} \quad (23)$$

Then (21) is equivalent to

$$\begin{cases} A^T d < 0 \\ B^T d = 0 \end{cases} \quad (24)$$

has no solution.

Denote

$$S_1 = \left\{ \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \mid y_1 = A^T d, y_2 = B^T d, d \in \mathbb{R}^n \right\} \quad (25)$$

$$S_2 = \left\{ \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \mid y_1 < 0, y_2 = 0 \right\} \quad (26)$$

S_1, S_2 are non-trivial convex sets, and $S_1 \cap S_2 = \emptyset$. From *Hyperplane Separation Theorem*: $\exists \begin{pmatrix} p_1 \\ p_2 \end{pmatrix}$, such that

$$p_1^T A^T d + p_2^T B^T d \geq p_1^T y_1 + p_2^T y_2, \forall d \in \mathbb{R}^n, \forall \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \in CL(S_2) \quad (27)$$

Let $y_2 = 0, d = 0, y_1 < 0$, we have

$$p_1 \geq 0 \quad (28)$$

Let $\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \in CL(S_2)$ So that

$$(p_1^T A^T + p_2^T B^T)d \geq 0 \quad (29)$$

$$(Ap_1 + Bp_2)^T d \geq 0 \quad (30)$$

Let $d = -(Ap_1 + Bp_2)$, we have

$$Ap_1 + Bp_2 = 0 \quad (31)$$

From above, we have

$$\begin{cases} Ap_1 + Bp_2 = 0 \\ p_1 \geq 0 \end{cases} \quad (32)$$

Let $p_1 = \{\lambda_0, \dots, \lambda_{I(x)}\}$, $p_2 = \{\mu_1, \dots, \mu_m\}$, i.e.,

$$\begin{cases} \lambda_0 \nabla f(x) - \sum_{i \in I(x)} \lambda_i \nabla g_i(x) - \sum_{j=1}^m \mu_j \nabla h_j(x) = 0 \\ \lambda_i \geq 0 \end{cases} \quad (33)$$

Theorem 2. Kuhn-Tucker Condition

For constrained optimization problem

$$\min f(x) \quad (34)$$

$$s.t. \quad g_i(x) \geq 0, i = 1, \dots, n \quad (35)$$

$$h_i(x) = 0, i = 1, \dots, m \quad (36)$$

Denote $I(x) = \{i \in \{1, \dots, n\} | g_i(x) = 0\}$. For $x \in \mathcal{S}$, f and $g_i, i \in I(x)$ is differentiable at x , $h_j(x)$ is continuously differentiable at x . $\{\nabla g_i(x), i \in I(x); \nabla h_j(x), j = 1, \dots, m\}$ is linearly independent. If x is local optimal, then $\exists \lambda_i \geq 0$ and μ_j , such that

$$\nabla f(x) - \sum_{i \in I(x)} \lambda_i \nabla g_i(x) - \sum_{j=1}^m \mu_j \nabla h_j(x) = 0 \quad (37)$$

1.5 Descent function

Definition 5. *Descent function. Denote solution set $\Omega \in X$, \mathcal{A} is an algorithm on X , $\psi : X \rightarrow \mathbb{R}$. If*

$$\psi(y) < \psi(x), \quad \forall x \notin \Omega, y \in \mathcal{A}(x) \quad (38)$$

$$\psi(y) \leq \psi(x), \quad \forall x \in \Omega, y \in \mathcal{A}(x) \quad (39)$$

Then ψ is a **descent function** of (Ω, \mathcal{A}) .

1.6 Convergence of Algorithm

Theorem 3. \mathcal{A} is an algorithm on X , Ω is the solution set, $x^{(0)} \in X$. If $x^{(k)} \in \Omega$, then the iteration stops. Otherwise set $x^{(k+1)} = \mathcal{A}(x^{(k)})$, $k := k + 1$. If

- $\{x^{(k)}\}$ in a compact subset of X
- There exists a continuous function ψ , ψ is a descent function of (Ω, \mathcal{A})
- \mathcal{A} is closed on Ω^C

Then, any convergent subsequence of $\{x^{(k)}\}$ converges to $x, x \in \Omega$.

Proof.

1.7 Search Methods

Line Search

Generate $d^{(k)}$ from $x^{(k)}$,

$$x^{(k+1)} = x^{(k)} + \alpha_k d^{(k)} \quad (40)$$

. search α_k in 1-D space.

Trust Region

Generate local model $Q_k(s)$ of $x^{(k)}$,

$$s^{(k)} = \arg \min Q_k(s) \quad (41)$$

$$x^{(k+1)} = x^{(k)} + s^{(k)} \quad (42)$$

2 Unconstrained Optimization

2.1 Gradient Based Methods

Algorithm 1: Example of gradient based algorithm

Data: Solution set Ω , cost function f

$x^{(0)} \in \mathbb{R}^n, k := 0;$

while $x^{(k)} \notin \Omega$ **do**

$d^{(k)} = -H_k \nabla f(x^{(k)})$, (H_k is a positive definite symmetrical matrix);

 solve $\min_{\alpha_k \geq 0} f(x^{(k)} + \alpha_k d^{(k)})$;

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

end

2.2 Determine Search Direction

First-order gradient method

For unconstrained optimization problem

$$\min_{x \in \mathbb{R}^n} f(x) \quad (43)$$

We have

$$f(x) = f(x^{(k)}) + \nabla f(x^{(k)})^T (x - x^{(k)}) + O(\|x - x^{(k)}\|^2) \quad (44)$$

Set $d^{(k)} = -\nabla f(x^{(k)})$, when α_k is sufficiently small,

$$f(x^{(k)} + \alpha_k d^{(k)}) < f(x^{(k)}) \quad (45)$$

Second-order gradient method – Newton Direction

$$f(x) = f(x^{(k)}) + \nabla f(x^{(k)})^T (x - x^{(k)}) \quad (46)$$

$$+ \frac{1}{2} (x - x^{(k)})^T \nabla^2 f(x^{(k)}) (x - x^{(k)}) + O(\|x - x^{(k)}\|^3) \quad (47)$$

Set $d^{(k)} = -G_k^{-1} \nabla f(x^{(k)})$, where $G_k = \nabla^2 f(x^{(k)})$, i.e., Hesse matrix of f at $x^{(k)}$.

2.3 Determine Step Factor – Line Search