

$$F \text{ norm : } \|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$$



- 日期:
- ★ ℓ_1 norm: $\|x\|_1 := \sum_{i=1}^n |x_i|.$
 - ★ ℓ_2 norm: $\|x\|_2 := \sqrt{\sum_{i=1}^n |x_i|^2}.$
 - ★ ℓ_∞ norm: $\|x\|_\infty := \max_{1 \leq i \leq n} |x_i|.$

Theorem 1.1: Let $\|\cdot\|$ be a norm in \mathbb{R}^n . Then there exist positive numbers C_1 and C_2 so that for all $x \in \mathbb{R}^n$,

$$C_1 \sum_{i=1}^n |x_i| \leq \|x\| \leq C_2 \sum_{i=1}^n |x_i|$$

Theorem 1.2: Let $\|\cdot\|$ be a norm. Then the following function defines a matrix norm

$$\|A\| := \max_{\|x\|=1} \|Ax\|.$$

$$\|AB\| \leq \|A\| \cdot \|B\| , \quad \|A+B\| \leq \|A\| + \|B\|.$$

- $\|A\|_1 = \max_j \sum_{i=1}^n |a_{ij}|$ (maximum of the ℓ_1 norms of columns).
Moreover, this is an induced matrix norm:

$$\|A\|_1 = \max_{\|x\|_1=1} \|Ax\|_1.$$

- $\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)}$. Moreover, this is an induced matrix norm:
 $\|A^T\|_2 = \|A\|_2$

$$\|A\|_2 = \max_{\|x\|=1} \|Ax\|_2.$$

- $\|A\|_\infty = \max_i \sum_{j=1}^n |a_{ij}|$ (maximum of the ℓ_1 norms of rows).
Moreover, this is an induced matrix norm:

$$\|A\|_\infty = \max_{\|x\|_\infty=1} \|Ax\|_\infty.$$

- $\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$. This is known as the Frobenius norm.

$$\|\text{vec}(A)\|_2 \quad \|AB\|_F \leq \|A\|_F \|B\|_F$$

Theorem 1.9. (Taylor's theorem in \mathbb{R}^n with remainder term)

- Let $f \in C^1(\mathbb{R}^n)$, x and $y \in \mathbb{R}^n$. Then there exists $\xi \in \{(1-s)x + sy : s \in (0, 1)\}$ such that

$$f(y) = f(x) + [\nabla f(\xi)]^T (y - x).$$
- Let $f \in C^2(\mathbb{R}^n)$, x and $y \in \mathbb{R}^n$. Then there exists $\xi \in \{(1-s)x + sy : s \in (0, 1)\}$ such that

$$f(y) = f(x) + [\nabla f(x)]^T (y - x) + \frac{1}{2} (y - x)^T \nabla^2 f(\xi) (y - x).$$

Definition: We say that x^* is a stationary point of f if $\nabla f(x^*) = 0$.

Theorem 2.2. Let $f \in C^2(\mathbb{R}^n)$.

1. If x^* is a local minimizer of f , then $\nabla^2 f(x^*) \succeq 0$.
2. If x^* is a stationary point of f and $\nabla^2 f(x^*) \succ 0$, then x^* is a local minimizer.

日期

Newton's methodLet $x^0 \in \mathbb{R}^n$. For $k = 0, 1, 2, \dots$, update

$$x^{k+1} = x^k - [\nabla^2 f(x^k)]^{-1} \nabla f(x^k).$$

Solve for d

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

Newton direction

Solve

$$f'(x_0) = 0 \iff g(x_0) = f'(x_0) = 0$$

$$x_{k+1} = x_k - \frac{g(x_k)}{g'(x_k)}$$

Theorem 2.3. (Quadratic convergence of Newton's method)

Let $g \in C^2(\mathbb{R})$ and x_* satisfies $g(x_*) = 0$ and $g'(x_*) \neq 0$. Then there exists $\epsilon > 0$ so that if $|x_0 - x_*| < \epsilon$ then the Newton's iterate $x_{k+1} = x_k - g(x_k)/g'(x_k)$ is well defined and there exists $M > 0$ so that

$$|x_{k+1} - x_*| \leq M |x_k - x_*|^2.$$

Note:

0.1
0.2

$$x_k \rightarrow x_*$$

convergent

Steepest descent cont.

- $-\nabla f(x)$ is called the steepest descent direction.
- A natural greedy algorithm is

Steepest descent with exact line searchStart at $x^0 \in \mathbb{R}^n$. For each $k = 0, 1, 2, \dots$,★ Set $d^k = -\nabla f(x^k)$.★ Pick α_k so that

the set of minimizers

$$\alpha_k \in \underbrace{\text{Arg min}}_{\alpha \geq 0} \{f(x^k + \alpha d^k) : \alpha \geq 0\}. \quad (1)$$

★ Set $x^{k+1} = x^k + \alpha_k d^k$.Note: The update $x^{k+1} = x^k + \alpha_k d^k$ is prototypical in optimization.

- d^k is called the search direction. In the above algorithm, $d^k = -\nabla f(x^k)$.
- α_k is called the step size. In the above algorithm, it is chosen according to the exact line search criterion (1).

$$f(x_0 + d) \approx f(x_0) + (\nabla f(x_0))^T d$$

$$f(x_0 + 2 \cdot \nabla f(x_0)) \approx f(x_0) + 2 \cdot \nabla f(x_0)^T \cdot \nabla f(x_0) < f(x_0).$$

$$d^k = -\nabla f(x^k).$$

Armijo rule

日期:

In contrast to exact line search, usually inexact line search strategy is performed. One commonly used rule is:

Armijo rule:

Let $\sigma \in (0, 1)$, $x \in \mathbb{R}^n$ and $d \in \mathbb{R}^n$. Find $\alpha > 0$ so that

$$f(x + \alpha d) \leq f(x) + \alpha \sigma [\nabla f(x)]^T d.$$

Definition. Let $f \in C^1(\mathbb{R}^n)$ and $x \in \mathbb{R}^n$. A $d \in \mathbb{R}^n$ is said to be a descent direction of f at x if

$$[\nabla f(x)]^T d < 0.$$

Examples: At an x that is not stationary,

- $d = -\nabla f(x)$ is a descent direction;
- More generally, if $D \succ 0$, then $d = -D \nabla f(x)$ is a descent direction.

Is the Newton direction $-[\nabla^2 f(x)]^{-1} \nabla f(x)$ a descent direction?

Newton direction is in general not a descent direction.

Armijo rule cont. $\{\bar{\alpha}, \bar{\alpha}\beta, \bar{\alpha}\beta^2, \dots\}$

How to execute Armijo rule in practice?

Armijo line search by backtracking:

Fix $\sigma \in (0, 1)$ and $\beta \in (0, 1)$. Given $x \in \mathbb{R}^n$, $d \in \mathbb{R}^n$ and $\bar{\alpha} > 0$. Find the smallest nonnegative integer $j = j_0$ so that

$$f(x + \bar{\alpha}\beta^j d) \leq f(x) + \bar{\alpha}\beta^j \sigma [\nabla f(x)]^T d. \quad (2)$$

The stepsize generated is then $\bar{\alpha}\beta^{j_0}$.

$$\alpha^k = \bar{\alpha} \cdot \beta^{j_k}, \quad \beta < 1, \quad j_k = 0, 1, 2, \dots$$

$f(x + \bar{\alpha}\beta^j d) \leq f(x) + \bar{\alpha}\beta^j \sigma [\nabla f(x)]^T d$, find j satisfies it.

Convergence under Armijo rule

Theorem 2.5:

Let $f \in C^1(\mathbb{R}^n)$ with $\inf_{x \in \mathbb{R}^n} f > -\infty$. Let $\{\bar{\alpha}_k\} \subset \mathbb{R}$ satisfy $0 < \inf_k \bar{\alpha}_k \leq \sup_k \bar{\alpha}_k < \infty$, and fix $\sigma \in (0, 1)$ and $\beta \in (0, 1)$. Suppose $\{x^k\}$ is generated as

$$x^{k+1} = x^k + \bar{\alpha}_k d^k, \quad (3)$$

where

- (2) $d^k := -D_k \nabla f(x^k)$ here $\{D_k\}$ is a bounded sequence of positive definite matrices with $D_k - \delta I \succeq 0$ for some $\delta > 0$;
 • α_k is generated via the Armijo line search by backtracking with $x = x^k$, $d = d^k$ and $\bar{\alpha} = \bar{\alpha}_k$, and σ and β defined above.

Then any accumulation point of $\{x^k\}$ is a stationary point of f .

(3) $\lambda_{\max}(A), \lambda_{\min}(A)$ is bounded.

$$\|A\|_2 = \sqrt{\lambda_{\max}(A)}.$$

Any convergent subsequence (if exists) will converge

to a stationary ptm of f .

日期:



Special case

$$\nabla f(x)$$

Find to

Corollary 2.1: (Steepest descent with constant stepsize)

Let $f \in C^2(\mathbb{R}^n)$ with $\inf f > -\infty$. Suppose that there exists $L > 0$ so that

$$L \geq \|\nabla^2 f(x)\|_2 \text{ for all } x.$$

Fix any $\gamma \in (0, 2)$ and consider the sequence generated as

$$x^{k+1} = x^k - \frac{\gamma}{L} \nabla f(x^k).$$

Then any accumulation point of $\{x^k\}$ is a stationary point of f .

③

$$\text{step size } \in (0, \frac{2}{L})$$

$$\text{step size } < \frac{2}{L}$$

A chain rule

Let $h \in C^2(\mathbb{R}^m)$ and let $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$.

$$\nabla f(x) = A^T \nabla h(Ax - b)$$

Define $f(x) := h(Ax - b)$. Then $f \in C^2(\mathbb{R}^n)$ and

$$\nabla f(x) = A^T \nabla h(Ax - b) \text{ and } \nabla^2 f(x) = A^T \nabla^2 h(Ax - b) A.$$

Quasi-Newton method: Basic version

Given $f \in C^1(\mathbb{R}^n)$.

Quasi-Newton based on B_k

Initialize at $x^0 \in \mathbb{R}^n$ and $B_0 \succ 0$.

For $k = 0, 1, 2, \dots$

1. Find d^k via $B_k d^k = -\nabla f(x^k)$ ↪ Solve for d^k

$$\rightarrow d^k \rightarrow x^{k+1} \rightarrow y^k \rightarrow s^k \rightarrow B_{k+1}$$

2. Update $x^{k+1} = x^k + d^k$. Or, more generally, $x^{k+1} = x^k + \alpha_k d^k$ for some $\alpha_k > 0$.

3. Set $y^k = \nabla f(x^{k+1}) - \nabla f(x^k)$ and $s^k = x^{k+1} - x^k$. Compute B_{k+1} .

Quasi-Newton based on H_k

Initialize at $x^0 \in \mathbb{R}^n$ and $H_0 \succ 0$.

For $k = 0, 1, 2, \dots$

1. Find d^k via $d^k = -H_k \nabla f(x^k)$ ↪ Not solving system of equations

2. Update $x^{k+1} = x^k + d^k$. Or, more generally, $x^{k+1} = x^k + \alpha_k d^k$ for some $\alpha_k > 0$.

3. Set $y^k = \nabla f(x^{k+1}) - \nabla f(x^k)$ and $s^k = x^{k+1} - x^k$. Compute H_{k+1} .

$$H_k \succ 0?$$

Proposition 3.1

Let $H_k \succ 0$ and $y^k T s^k > 0$. Let H_{k+1} be given by BFGS update.

Then $H_{k+1} \succ 0$. The same conclusion holds if H_k and H_{k+1} are replaced by B_k and B_{k+1} , respectively.



Popular update formulae

Initialize B_0 (or H_0) at a **positive definite matrix**.

Method	$B_{k+1} =$	$H_{k+1} =$
DFP	$\left(I - \frac{y^k s^{kT}}{y^{kT} s^k}\right) B_k \left(I - \frac{s^k y^{kT}}{y^{kT} s^k}\right) + \frac{y^k y^{kT}}{y^{kT} s^k}$	$H_k + \frac{s^k s^{kT}}{y^{kT} s^k} - \frac{H_k y^{kT} y^k H_k}{y^{kT} H_k y^k}$ <i>(symmetric rank 1 matrix)</i>
BFGS	$B_k + \frac{y^k y^{kT}}{y^{kT} s^k} - \frac{B_k s^k s^{kT} B_k}{s^{kT} B_k s^k}$	$\left(I - \frac{s^k y^{kT}}{y^{kT} s^k}\right) H_k \left(I - \frac{y^k s^{kT}}{y^{kT} s^k}\right) + \frac{s^k s^{kT}}{y^{kT} s^k}$
SR1	$B_k + \frac{(y^k - B_k s^k)(y^k - B_k s^k)^T}{(y^k - B_k s^k)^T s^k}$	$H_k + \frac{(s^k - H_k y^k)(s^k - H_k y^k)^T}{(s^k - H_k y^k)^T y^k}$

Remark: **symmetric rank 1**.

Wolfe conditions

In view of **Proposition 3.1**, it suffices to guarantee that $H_0 \succ 0$ and make sure that $y^k T s^k > 0$ for each $k \geq 0$.

The latter can be guaranteed if line search is performed to guarantee the **Wolfe conditions**.

Wolfe conditions:

Let $0 < c_1 < c_2 < 1$, $x \in \mathbb{R}^n$ and $d \in \mathbb{R}^n$. Find $\alpha > 0$ so that

$$f(x + \alpha d) \leq f(x) + \alpha c_1 [\nabla f(x)]^T d, \quad \text{← Armijo rule}$$

$$- [\nabla f(x + \alpha d)]^T d \leq -c_2 [\nabla f(x)]^T d.$$

Remark:

- The first inequality in Wolfe conditions is the Armijo rule.
- The second relation is called **curvature condition**.

Wolfe conditions cont.

Theorem 3.3 (Wolfe conditions are not void)

Let $f \in C^1(\mathbb{R}^n)$ with $\inf f > -\infty$, $x \in \mathbb{R}^n$, and $d \in \mathbb{R}^n$ be a descent direction at x . Let $0 < c_1 < c_2 < 1$. Then there exists $\alpha > 0$ with

$$f(x + \alpha d) \leq f(x) + \underline{\alpha c_1 [\nabla f(x)]^T d}, \quad (3)$$

$$- [\nabla f(x + \underline{\alpha d})]^T d \leq -c_2 [\nabla f(x)]^T d.$$

Quasi-Newton method: Wolfe line search

日期:

Quasi-Newton using H_k in BFGS for $f \in C^1(\mathbb{R}^n)$ with $\inf f > -\infty$:

Pick $0 < c_1 < c_2 < 1$, $x^0 \in \mathbb{R}^n$, $H_0 = \eta I$ for some $\eta > 0$.

For $k = 0, 1, 2, \dots$

1. Find d^k via $d^k = -H_k \nabla f(x^k)$.
2. Compute α_k that satisfies the Wolfe conditions.
3. Update $x^{k+1} = x^k + \alpha_k d^k$.
4. Set $y^k = \nabla f(x^{k+1}) - \nabla f(x^k)$, $s^k = x^{k+1} - x^k$ and compute H_{k+1} as in BFGS.

difference



Convergence under Wolfe conditions

conditions:

Theorem 3.4: (Zoutendijk theorem)

Let $f \in C^1(\mathbb{R}^n)$ with $\inf f > -\infty$, $x^0 \in \mathbb{R}^n$ and $\exists \ell > 0$ so that

$$\textcircled{1} \quad \textcircled{3} \quad \|\nabla f(x) - \nabla f(y)\|_2 \leq \ell \|x - y\|_2 \quad \text{Lipchitz's condition}$$

whenever $\max\{f(x), f(y)\} \leq f(x^0)$. Let $\{x^k\}$ be a sequence of non-stationary points generated as

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

with d^k being a descent direction and α_k satisfying the Wolfe conditions. Then it holds that

conclusion:

$$\text{as } \theta_k \text{ is bounded} \quad \sum_{k=0}^{\infty} \cos^2 \theta_k \|\nabla f(x^k)\|_2^2 < \infty,$$

where $\cos \theta_k := \frac{-\langle \nabla f(x^k), d^k \rangle}{\|\nabla f(x^k)\|_2 \|d^k\|_2}$. TM

(2)

Then we prove sequence $\{x^k\}$, $\|\nabla f(x^k)\|_2 \rightarrow 0$ ($k \rightarrow \infty$).

$$\|x^k\| \text{ bounded} \Rightarrow x^k \rightarrow x^* \quad (j \rightarrow \infty)$$

Thm: $f \in C^2(\mathbb{R}^n)$, $\|\nabla f(x) - \nabla f(y)\|_2 \leq \max_i \|\nabla^2 f(x)\|_2 \cdot \|x - y\|_2$



A closedness result

Theorem 4.4:

Let $A \in \mathbb{R}^{m \times n}$. Then the set

$$S := \{Ay : y \geq 0\}$$

is closed and convex, where $y \geq 0$ means $y_i \geq 0$ for each i .

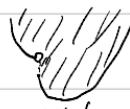
$$A(\lambda y_1 + (1-\lambda)y_2) \in S$$

Def:

$$f: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$$

$$\textcircled{1} \quad \text{dom } f := \{x | f(x) < \infty\}$$

$$\textcircled{2} \quad \text{epi } f := \{(x, r) \in \mathbb{R}^{n+1} | r \geq f(x)\}$$



epi f convex $\Leftrightarrow f$ convex

日期:

Theorem 4.7:

Suppose that $f \in C^2(\mathbb{R}^n)$. Then f is convex if and only if $\nabla^2 f(x) \succeq 0$ for all $x \in \mathbb{R}^n$.

Thm: f convex $\Leftrightarrow \nabla^2 f(x) \succeq 0$

x^* is local minimizer $\Rightarrow \nabla^2 f(x^*) \succeq 0$

$\nabla^2 f(x^*) \succ 0, \nabla f(x^*) = 0 \Rightarrow x^*$ is local minimizer.

Calculus of convex functions

Proposition 4.1:

Let $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$, $g : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ both be convex, $A \in \mathbb{R}^{n \times p}$, $b \in \mathbb{R}^p$, $H(x) := Ax - b$, and $\alpha > 0$. Then the following functions are convex: ①

• $f + g$

• $f \circ H$; ② ↴

• αf ; ③ ↴

• $\max\{f, g\}$. ④ ↴

$f, \Rightarrow f(H)$

f, g convex

Example: prove $\sum_{i=1}^{n-1} |x_{i+1} - x_i|$ convex,

$$\sum_{i=1}^{n-1} |x_{i+1} - x_i| = \|Dx\|_1 = \|\cdot\|_1 \circ D(x)$$

$\|\cdot\|_1$ is convex, D linear $\Rightarrow \|Dx\|_1$ convex

Thm: if h is convex

g is convex and non-decreasing

$\Rightarrow f = g \circ h$ is convex.

日期:

Strong duality I (necessary)

Theorem 5.1 (Strong duality for LP, version I)

Let $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$ and $c \in \mathbb{R}^n$. Consider

$$v_p := \sup_x \{c^T x : Ax = b, x \geq 0\}, \quad v_d := \inf_y \{b^T y : c \leq A^T y\}.$$

implies $v_p \geq v_d$

Suppose that there exists $\hat{x} \geq 0$ with $A\hat{x} = b$. Then $v_p = v_d$.

primal:
$$\begin{array}{ll} \max & c^T x \\ \text{s.t.} & Ax = b \\ & x \geq 0 \end{array}$$

dual:
$$\begin{array}{ll} \min & b^T y \\ \text{s.t.} & c \leq A^T y \end{array}$$

$\exists \hat{x} \geq 0$, and $A\hat{x} = b \Rightarrow v_p = v_d$.
(constraints)

Strong duality II

Theorem 5.2 (Strong duality theorem for LP)

Let $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$ and $c \in \mathbb{R}^n$. Consider

$$v_p := \sup_x \{c^T x : Ax = b, x \geq 0\}, \quad v_d := \inf_y \{b^T y : c \leq A^T y\}.$$

Suppose that either

- there exists $\hat{x} \geq 0$ with $A\hat{x} = b$; or $f(\hat{x}) = v_p$, since \hat{x} is unique feasible element.
- there exists \hat{y} with $c \leq A^T \hat{y}$.

Then $v_p = v_d$ and both optimal values are attained when finite.

If $\exists \hat{x}$ or \hat{y} obey constraints $\Rightarrow v_p = v_d$ and finite.

If $\nexists \hat{x}$ obey constraints, and $\exists \hat{y}$ obey constraints

$$\Rightarrow v_p = \max_x c^T x = -\infty$$

$$\text{and } v_d = \min_y b^T y = -\infty$$

If $\nexists \hat{x}, \hat{y} \Rightarrow v_p = v_d = \infty$.

	Maximize $c^T x$ subject to $Ax \trianglelefteq b$ $x \triangleleft$	Minimize $b^T y$ subject to $A^T y \trianglelefteq c$ $y \triangleleft$
♣	i th constraint \leq i th constraint \geq i th constraint $=$	i th variable ≥ 0 i th variable ≤ 0 i th variable unrestricted
◊	j th variable ≥ 0 j th variable ≤ 0 j th variable unrestricted	j th constraint \geq j th constraint \leq j th constraint $=$

日期:

Semi definite Programming

SDP

primal: $\underset{X \in S^n}{\text{Min}} \text{tr}(CX)$

s.t. $\text{tr}(A_i X) = b_i, i = 1, \dots, m$
 $X \succeq 0$.

dual: $\underset{y \in \mathbb{R}^m}{\text{Max}} b^T y$

s.t. $C - \sum_i y_i A_i \succeq 0$

Thm: if $\exists \bar{y}$, s.t. $C - \sum_i \bar{y}_i A_i \succ 0$

Strong duality $\left\{ (\text{tr}(CX), \text{tr}(A_i X) - \text{tr}(A_i X))^T \mid X \succ 0 \right\}$
 is closed, \Rightarrow primal optimal value
 is attainable.

Theorem 6.1 (Strong duality for SDPs)

Consider the following primal-dual SDP pairs:

Primal : $\begin{cases} \underset{X \in S^n}{\text{Minimize}} \text{tr}(CX) \\ \text{subject to } \text{tr}(A_i X) = b_i, i = 1, \dots, m, \\ X \succeq 0, \end{cases}$

Dual : $\begin{cases} \underset{y \in \mathbb{R}^m}{\text{Maximize}} b^T y \\ \text{subject to } C - \sum_{i=1}^m y_i A_i \succeq 0, \end{cases}$

where $C \in S^n$ and $A_i \in S^n$ for all i . Let v_p and v_d denote their optimal values respectively. Then the following statements hold.

1. If there exists $\bar{X} \succ 0$ such that $\text{tr}(A_i \bar{X}) = b_i$ for all i , then $v_p = v_d$ and v_d is attained when finite.
2. If there exists $\bar{y} \in \mathbb{R}^m$ such that $C - \sum_{i=1}^m \bar{y}_i A_i \succ 0$, then $v_p = v_d$ and v_p is attained when finite.

if X is unique
element in
feasible region

Def: (slater point): \bar{X} satisfies $(\succ 0)$

7/15

Thm: ① if $\exists \bar{X}$ obey primal constraints and $(\succ 0)$
 $\Rightarrow v_p = v_d$, and v_d is attainable when finite

② if $\exists \bar{y}$ obey $C - \sum_i \bar{y}_i A_i \succ 0$ not $\succeq 0$
 $\Rightarrow v_p = v_d$, and v_p is attained when finite

日期

To understand Theorem 6.1, we need the following.

Theorem 6.2

Let $A \in S_+^n$ and $C \in S_+^m$. Then $\text{tr}(AC) \geq 0$.



S_+^n : symmetric, positive semi-definite

Def: principal matrices.
(ordered)

Schur complement

The following result is crucial in reformulating problems into SDPs.



Theorem 7.1

Let $A \in S^m$, $C \in S^n$, $B \in \mathbb{R}^{m \times n}$, and $A \succ 0$. Then

$$\begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \succeq 0 \iff C - B^T A^{-1} B \succeq 0.$$

Note: We call $C - B^T A^{-1} B$ the Schur complement of A in $\begin{bmatrix} A & B \\ B^T & C \end{bmatrix}$.

Thm: $C - B^T A^{-1} B \succeq 0$, $A \succ 0$, $\iff A \succeq 0$, $\begin{pmatrix} A & B \\ B^T & C \end{pmatrix} \succeq 0$.

Use: primal: $\min_{\mathbf{x}} \mathbf{x}^\top Q \mathbf{x}$

$$\iff \min \text{tr}(\mathbf{x}^\top Q \mathbf{x}) = \text{tr}(Q \mathbf{x} \cdot \mathbf{x}^\top)$$

$$\iff \min \text{tr}(Q Y) \quad (\text{Thm: } \dots) \\ \text{s.t. } Y \succeq \mathbf{x} \cdot \mathbf{x}^\top$$

$$\iff \begin{array}{ll} \min & \text{tr}(Q Y) \\ \text{s.t.} & \begin{pmatrix} Y & \mathbf{x} \\ \mathbf{x}^\top & 1 \end{pmatrix} \succeq 0 \end{array}.$$

日期:

$$\text{Thm: } \quad \textcircled{1} \quad t \geq \|x\|_2, \quad (t \geq 0)$$

$$\Leftrightarrow \begin{pmatrix} t \cdot 1 & x \\ x^T & t \end{pmatrix} \succeq 0$$

$$\textcircled{2} \quad ab \geq \|x\|_2, \quad a \geq 0, b \geq 0$$

$$\Leftrightarrow \begin{pmatrix} a \cdot 1 & x \\ x^T & b \end{pmatrix} \succeq 0.$$

$$\textcircled{3} \quad y \geq x^6 \Leftrightarrow y \geq t^2 \geq x^6, \quad t \geq 0$$

$$\Leftrightarrow y \geq t^2, \quad t \geq |x|^3 \Leftrightarrow y \geq t^2, \quad t \geq z^3 \geq |x|^3, \quad z \geq 0$$

$$\Leftrightarrow y \geq t^2, \quad z \geq |x|, \quad z \geq 0, \quad t \cdot z \geq s^2 \geq z^4, \quad s \geq 0$$



Thm: convex QP are SDP

Pf. QP: $\min x^T Ax$

$$\Leftrightarrow \min \operatorname{tr}(x^T Ax) = \operatorname{tr}(Ax \cdot x^T)$$

since $(A \succeq 0), \quad \Leftrightarrow$

$$\begin{aligned} &\min \operatorname{tr}(AT) \\ \text{s.t. } &T \succeq x \cdot x^T \end{aligned}$$

$$\Leftrightarrow \min \operatorname{tr}(AT)$$

$$\text{s.t. } \begin{pmatrix} Y & x \\ x^T & 1 \end{pmatrix} \succeq 0$$

\therefore is SDP.

$$\text{Def: } \|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} \quad (\text{F norm})$$

KKT conditions for LP

Theorem 8.1: (Karush-Kuhn-Tucker conditions for the LP)

Consider the linear program

$$\begin{array}{ll} \text{primal} & \begin{array}{l} \text{Minimize}_{x \in \mathbb{R}^n} c^T x \\ \text{subject to } Bx = d, \\ Ax \leq b. \end{array} \end{array} \quad (2)$$

$$\text{dual: } (\mu, \lambda) = y$$

where $c \in \mathbb{R}^n$, $B \in \mathbb{R}^{p \times n}$ and $A \in \mathbb{R}^{q \times n}$. Then $x^* \in \mathbb{R}^n$ is an optimal solution of (2) if and only if there exist $\lambda^* \in \mathbb{R}^q$ and $\mu^* \in \mathbb{R}^p$ so that the following conditions hold:

- (Primal feasibility) $Bx^* = d$ and $Ax^* \leq b$; and
- (Dual feasibility) $B^T \mu^* + A^T \lambda^* + c = 0$ and $\lambda^* \geq 0$ and
- (Complementary slackness) $\lambda^{*T}(Ax^* - b) = 0$. \leftarrow LP means

Remark: Since $\lambda^* \geq 0$ and $Ax^* \leq b$, we note that $\lambda^{*T}(Ax^* - b) = 0$ means $\lambda_i^*(Ax^* - b)_i = 0$ for each $i = 1, \dots, q$.

$$\begin{array}{ll} \text{Minimize}_{x \in \mathbb{R}^n} & f(x) \\ \text{subject to} & h_j(x) = 0, \quad j \in J, \\ & g_i(x) \leq 0, \quad i \in I. \end{array}$$



$$\begin{cases} I = \text{all } i \\ I(x) = \{i \mid g_i(x) = 0\} \end{cases}$$

MFCQ is for explanation to

complementary slackness

When is the approximation good?

Definition: (Mangasarian-Fromovitz constraint qualification, MFCQ)

Consider the feasible set of (1) and let x^* be feasible. We say that the **Mangasarian-Fromovitz constraint qualification (MFCQ)** holds at x^* if the following conditions hold:

- if

$$\sum_{i \in J} \mu_j \nabla h_j(x^*) + \sum_{i \in I(x^*)} \lambda_i \nabla g_i(x^*) = 0 \text{ and } \lambda_i \geq 0 \quad \forall i \in I(x^*),$$

then $\lambda_i = 0$ for all $i \in I(x^*)$ and $\mu_j = 0$ for all $j \in J$.



$$I(x^*) = \{i \mid g_i(x^*) = 0\}$$

Remark:

- If $g_i(x^*) < 0$ for all $i \in I$ so that $I(x^*) = \emptyset$, then MFCQ means $\{\nabla h_j(x^*) : j \in J\}$ is linearly independent.

- If $J = \emptyset$, then MFCQ means $\{\nabla g_i(x^*) : i \in I(x^*)\}$ is positively independent.

$\nabla h_j(x^*)$ means tangent

since $\nabla g_i(x^*) \leq 0$ or ≥ 0



: Need to consider if $g_i(x^*) = 0$, $g_i(x^*) \neq 0$.

KKT conditions for NLP

$L \neq I(x^*)$

Theorem 8.2: (KKT conditions for NLP) ①

Consider (1) and let x^* be a local minimizer. Suppose that MFCQ holds at x^* . Then there exist $\lambda^* \in \mathbb{R}^m$ and $\mu^* \in \mathbb{R}^p$ so that

$$\nabla f(x^*) + \sum_{i \in I} \lambda_i^* \nabla g_i(x^*) + \sum_{j \in J} \mu_j^* \nabla h_j(x^*) = 0; \text{ and } \quad \text{KKT}$$

$$\bullet \lambda_i^* \geq 0 \text{ and } \lambda_i^* g_i(x^*) = 0 \text{ for all } i \in I. \quad \text{written as}$$

$$c + A^T \mu + B^T \lambda = 0.$$

Remarks:

- Under MFCQ, the approximating LP is "good" around x^* . We hence look at the KKT of this LP.
- The first bullet point follows from the dual feasibility condition of the approximating LP, and by defining $\lambda_i^* = 0$ for $i \notin I(x^*)$.
- The second bullet point follows from the definition of $I(x^*)$ and by defining $\lambda_i^* = 0$ for $i \notin I(x^*)$.

Remarks cont.:

- The λ^* and μ^* are called Lagrange multipliers at x^* .
- We consider the so-called stationary points:



Definition: (Stationary points)

Consider (1). An \bar{x} is called a stationary point of (1) if it is feasible and there exist $\bar{\lambda} \in \mathbb{R}^m$ and $\bar{\mu} \in \mathbb{R}^p$ so that

$$\star \nabla f(\bar{x}) + \sum_{i \in I} \bar{\lambda}_i \nabla g_i(\bar{x}) + \sum_{j \in J} \bar{\mu}_j \nabla h_j(\bar{x}) = 0; \text{ and } \quad \text{KKT}$$

$$\star \bar{\lambda} \geq 0 \text{ and } \bar{\lambda}_i g_i(\bar{x}) = 0 \text{ for all } i \in I.$$

- According to Theorem 8.2, if x^* is a local minimizer and if the MFCQ holds at x^* , then x^* is a stationary point.

Stationary point $\bar{x} \Leftrightarrow \bar{x}$ satisfies KKT

日期:



$$\begin{aligned} & \text{Min } f(x) \\ \text{s.t. } & h_j(x) = 0, j \in J \\ & g_i(x) \leq 0, i \in I \end{aligned}$$



KKT conditions for (1):

- $g_i(x) \leq 0$ for all $i \in I$ and $h_j(x) = 0$ for all $j \in J$; and
- $\nabla f(x) + \sum_{i \in I} \lambda_i \nabla g_i(x) + \sum_{j \in J} \mu_j \nabla h_j(x) = 0$; and
- $\lambda \geq 0$ and $\lambda_i g_i(x) = 0$ for all $i \in I$.

x satisfy KKT $\Leftrightarrow x$ is stationary point



Thm: x^* is local minimizer, MFCQ \Rightarrow stationary point

18/32

- ① check MFCQ holds for $\forall x$, $\Rightarrow x^*$ is stationary.
- ② f can be not convex Slater's condition (\Rightarrow KKT satisfies)



Theorem 8.3: (MFCQ from Slater)

Consider the set defined by

$$I = (\text{all } i), I(x^*) \neq I$$



$$S := \{x \in \mathbb{R}^n : g_i(x) \leq 0 \ \forall i \in I\},$$

where g_i are convex. Suppose that there exists \bar{x} satisfying

$$g_i(\bar{x}) < 0 \ \forall i \in I.$$

- ① convex
② slater point \bar{x} .

Then MFCQ holds at every point in S .

Show in exam

Remark: $\lambda_i \cdot \nabla g_i(x^*) = 0, \lambda_i \geq 0 \Rightarrow \lambda_i = 0$ (MFCQ)

- The set S in the above theorem is closed and convex.
- The condition that "there exists \bar{x} satisfying $g_i(\bar{x}) < 0$ for all $i \in I$ " is called the **Slater's condition**. The \bar{x} is called a **Slater point**.
- One can indeed show that for the above S , the MFCQ holds at every point in S if and only if Slater's condition holds.



日期: /

Generalized Slater's condition

Theorem 8.4: (MFCQ from generalized Slater)

Consider the set defined by

\bar{x} is called slater point

$$S := \{x \in \mathbb{R}^n : g_i(x) \leq 0 \ \forall i \in I, Ax = b\},$$

where g_i are convex C^1 and $A \in \mathbb{R}^{p \times n}$. Suppose that there exists \bar{x} satisfying

$$A\bar{x} = b, \quad g_i(\bar{x}) < 0 \quad \forall i \in I,$$

and A has full row rank. Then MFCQ holds at every point in S .

Remark: $\sum \mu_j \triangleright h_j(x^*) = \sum \mu \cdot A = 0$, A full row rank $\Rightarrow \mu = 0$.

- The set S in the above theorem is closed and convex.

Role of convexity

Consider the following special instance of (1)

$$\begin{array}{ll} \text{Minimize} & f(x) \\ \text{subject to} & \begin{array}{l} Ax = b \\ g_i(x) \leq 0, \quad i \in I \end{array} \end{array} \quad (3)$$

here f and g_i are all convex C^1 functions, $A \in \mathbb{R}^{p \times n}$.

Theorem 8.5: (Sufficiency under convexity)

Consider (3). Suppose that there exist $x^* \in \mathbb{R}^n$, $\lambda^* \in \mathbb{R}^m$ and $\mu^* \in \mathbb{R}^p$ so that

- $Ax^* = b$ and $g_i(x^*) \leq 0$ for all $i \in I$; and
- $\nabla f(x^*) + \sum_{i \in I} \lambda_i^* \nabla g_i(x^*) + A^T \mu^* = 0$; and
- $\lambda_i^* \geq 0$ and $\lambda_i^* g_i(x^*) = 0$ for all $i \in I$.

KKT

Then x^* is a global minimizer of (3).

Thm: f, g convex, stationary point (KKT in I) \Rightarrow global minimizer

Only one minimizer

日期:



Thm: local minimizer + MFCQ

\Rightarrow stationary point?

Thm: $k \in \mathbb{I} \Leftrightarrow$ stationary point

- * Penalty method (exterior type);
- * Barrier method (interior type).

Penalty method (exterior type)

- Initial idea: Define

$$P(x) := \begin{cases} 0 & \text{if } g_i(x) \leq 0 \quad \forall i \in I, \\ \infty & \text{otherwise.} \end{cases}$$

Then we can consider the **unconstrained** problem:

$$\underset{x \in \mathbb{R}^n}{\text{Minimize}} \quad f(x) + P(x).$$

Examples:

- $P(x) = \sum_{i=1}^m \max\{g_i(x), 0\}.$
- $P(x) = \frac{1}{2} \sum_{i=1}^m \max\{g_i(x), 0\}^2.$ Courant-Beltrami penalty function: C^1 function

{ numerator
denominator

Penalty method: basic version

Penalty method for (1): basic version

Let $x^0 \in \mathbb{R}^n$, $c > 0$ and $\eta > 1$. Set $c_1 = c$. For $k = 1, \dots,$

- Find a minimizer (x^k) of

$$q_{c_k}(x) := f(x) + \frac{c_k}{2} \sum_{i=1}^m (\max\{g_i(x), 0\})^2,$$

using x^{k-1} as the initial point for the iterative method.

- Update $c_{k+1} = \eta c_k$.

$$\begin{pmatrix} x^0 \\ q_{c_0} \\ x^1 \\ \vdots \\ x^k \end{pmatrix} \rightarrow \begin{pmatrix} x^0 \\ q_{c_0} \\ x^1 \\ \vdots \\ x^k \end{pmatrix}$$

Remark:

- As c increases, q_c becomes more ill-conditioned. The choice of x^{k-1} as a starting point for the iterative method helps alleviate the ill-conditioning.
- The above algorithm is only conceptual because finding global minimizers can be challenging if q_{c_k} is not convex. Global minimizers also may

if $c_k \rightarrow +\infty$, compute x^*

Penalty method: basic version cont.

Theorem 9.1: (Convergence of penalty method: basic version)

Consider (1) and suppose that $\inf f > -\infty$. Let $\{x^k\}$ be generated by the penalty method on Slide 7. Then any accumulation point x^* of $\{x^k\}$ is a globally optimal solution of (1).

日期:



Barrier method

(interior)

Recall that

$$\begin{array}{ll} \text{Minimize}_{x \in \mathbb{R}^n} & f(x) \\ \text{subject to} & g_i(x) \leq 0, \quad i \in I. \end{array} \quad (2)$$

- One standard way is to make use of the log-barrier function:

$$B(x) := - \sum_{i=1}^m \ln[-g_i(x)]. \quad -g_i(x) \geq 0$$

Then one minimizes $f_\mu(x) := f(x) + \mu B(x)$ for some $\mu > 0$.

$$S^0 \triangleq \{x \mid g_i(x) \geq 0\}$$

slater point $g_i(x^0) < 0$

Barrier method: basic version

Barrier method for (2): basic version

Let $x^0 \in S^0$, $\mu > 0$ and $\eta > 1$. Set $\mu_1 = \mu$. For $k = 1, \dots$,

- Find a minimizer x^k of

$$f_{\mu_k}(x) := f(x) - \mu_k \sum_{i=1}^m \ln[-g_i(x)],$$

using x^{k-1} as the **initial point** for the iterative method.

- Update $\mu_{k+1} = \mu_k / \eta$.

$$\{\mu_k\} \downarrow \text{to } 0$$

日期:

Gram-Schmidt process

Theorem 10.1: (Gram-Schmidt process)

Given a set of linearly independent vectors $\{v^0, \dots, v^k\} \subset \mathbb{R}^n$.
Set $w^0 = v^0$ and for each $j = 1, \dots, k$

$$w^k = v^k - \sum_{j=0}^{k-1} \frac{v^{kT} w^j}{\|w^j\|_2^2} w^j$$



Then $w^i \neq 0$ for all i . Moreover, $w^{iT} w^j = 0$ whenever $i \neq j$, and for each $i = 0, 1, \dots, k$, it holds that

$$\text{Span}\{v^0, \dots, v^i\} = \text{Span}\{w^0, \dots, w^i\}.$$

seen as inner product $\langle v, w \rangle = v^T A \cdot w$

(Generalized) Gram-Schmidt process

Theorem 10.2: ((Generalized) Gram-Schmidt process)

Given $A \in \mathbb{R}^{n \times n}$ with $A \succ 0$ and a set of linearly independent vectors $\{v^0, \dots, v^k\} \subset \mathbb{R}^n$. Set $w^0 = v^0$ and for each $j = 1, \dots, k$

$$w^k = v^k - \sum_{j=0}^{k-1} \frac{v^{kT} Aw^j}{w^{jT} Aw^j} w^j \Rightarrow \text{that means, } w^i \perp w^j$$

in $\langle \cdot, \cdot \rangle$ inner product def

Then $w^i \neq 0$ for all i . Moreover, $w^{iT} Aw^j = 0$ whenever $i \neq j$, and for each $i = 0, 1, \dots, k$, it holds that

$$\text{Span}\{v^0, \dots, v^i\} = \text{Span}\{w^0, \dots, w^i\}.$$

Conjugate gradient method: Conceptual version

Conjugate gradient method: Conceptual version

Start at $x^0 \in \mathbb{R}^n$ and $d^0 = -\nabla f(x^0) = b - Ax^0$.

(see note)

For each $k = 0, 1, 2, \dots$,

- If $d^k = 0$, terminate.
- Pick α_k so that

$$\underline{\alpha_k \in \text{Arg min}\{f(x^k + \alpha d^k) : \alpha \geq 0\}}.$$

- Set $x^{k+1} = x^k + \alpha_k d^k$ and

$$\underline{d^{k+1} = -\nabla f(x^{k+1}) - \sum_{j=0}^k \frac{[-\nabla f(x^{k+1})]^T Ad^j}{d^{jT} Ad^j} d^j}$$

To do: upper to n orthogonal $\{d^k\}_{k=1}^n$.

- Prove the correctness: i.e. when $d^k = 0$ what happens?



Thm: $\nabla f(x^0)^T \nabla f(x^k) = 0$, $d^{iT} \nabla f(x^0) = 0$, $\forall i < j$

$$\Rightarrow \text{span}\{\nabla f(x^0), \dots, \nabla f(x^n)\} = \text{span}\{d^0, \dots, d^n\}$$

\Rightarrow at most, n iteration can solve the problem.

日期:

Theorem 10.3:

Let $A \succ 0$ and $x^0 \in \mathbb{R}^n$. Set $d^0 = -\nabla f(x^0)$. For $k = 0, 1, \dots$, suppose that $d^0, \dots, d^k \neq 0$, where for each $i = 0, \dots, k-1$,

$$d^{i+1} = -\nabla f(x^{i+1}) - \sum_{j=0}^i \frac{[-\nabla f(x^{i+1})]^T A d^j}{d^j T A d^j} d^i,$$

with $x^{i+1} = x^i + \alpha_i d^i$ and α_i coming from exact line search. Then $[\nabla f(x^j)]^T \nabla f(x^{k+1}) = 0$ and $d^j T \nabla f(x^{k+1}) = 0$ for $j < k+1$.

Proof: The proof is by induction. Let $k = 0$ and $d^0 \neq 0$. Then $d^0 = -\nabla f(x^0)$, and $d^0 T \nabla f(x^1) = 0$ holds because of exact line search.

$$\nabla f(x^0)^T \cdot \nabla f(x^1) = 0$$

Conjugate gradient method: Formal version

Start at $x^0 \in \mathbb{R}^n$ and $d^0 = -\nabla f(x^0) = b - Ax^0$.

For each $k = 0, 1, 2, \dots$,

- If $d^k = 0$, terminate.
- Pick α_k so that

$$\alpha_k \in \text{Arg min}\{f(x^k + \alpha d^k) : \alpha \geq 0\}.$$

- Set $x^{k+1} = x^k + \alpha_k d^k$ and

$$d^{k+1} = -\nabla f(x^{k+1}) + \frac{\|\nabla f(x^{k+1})\|_2^2}{\|\nabla f(x^k)\|_2^2} d^k$$

Problem : $\text{Min } f(x) = \frac{1}{2} x^T A x - b^T x$



Conjugate gradient method: Actual version

Conjugate gradient method: Actual version

Start at $x^0 \in \mathbb{R}^n$ and $r^0 = d^0 = b - Ax^0$.

For each $k = 0, 1, 2, \dots$,

- If $\|r^k\|$ (or, less commonly, $\|d^k\|$) is below a tolerance, terminate.
- (Exact line search) Compute

$$\alpha_k = \frac{r^k T r^k}{d^k T A d^k}, \quad x^{k+1} = x^k + \alpha_k d^k, \quad r^{k+1} = r^k - \alpha_k A d^k.$$

- (Update d^{k+1}) Compute

$$\beta_k = \frac{r^{k+1} T r^{k+1}}{r^k T r^k}, \quad d^{k+1} = r^{k+1} + \beta_k d^k.$$

Remark: β_k not saved

- One matrix-vector multiplication per iteration if Ad^k is saved.
- Keeping track of four vectors, x^k , r^k , d^k and the Ad^k saved.

Convergence rate

日期:



Conjugate gradient method must terminate in at most n iterations:
because $\{d^0, \dots, d^{n-1}\}$ must be a basis of \mathbb{R}^n , if $d^k \neq 0$ for each k .
Then necessarily, $d^n = 0$.

In practice, conjugate gradient method can converge much more quickly. We state the following result without proof. See Theorem 5.5 of Ref 2.

Theorem 10.4: (Luenberger)

Consider the conjugate gradient method for minimizing $f(x) = \frac{1}{2}x^T Ax - b^T x$ for some $b \in \mathbb{R}^n$ and $A \succ 0$. Let $\{x^k\}$ be the sequence generated and let x^* be the minimizer of f . If A has eigenvalues $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$, then

$$(x^{k+1} - x^*)^T A(x^{k+1} - x^*) \leq \left(\frac{\lambda_{n-k} - \lambda_1}{\lambda_{n-k} + \lambda_1} \right)^2 (x^0 - x^*)^T A(x^0 - x^*)$$

$$\hookrightarrow \lambda_1 \cdot \|x^{k+1} - x^*\|_2^2$$

so if $\lambda_{n-k} - \lambda_1 = 0$, we have $x^{k+1} = x^*$.

15/20

Thm:



that means $\{x^k\} = x^*$
from x^{k+1} .