# Numerical Optimization with applications: Homework 01

104021601 林俊傑

104021602 吳彥儒

104021615 黃翊軒

September 28, 2016

**Exercise 1.** *Compute the gradient $\nabla f(x)$ and Hessian $\nabla^2 f(x)$ of the Rosenbrock function*

$$f(x) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

*Show that $x^* = (1,1)^T$ is the only local minimizer of this function, and that the Hessian matrix at that point is positive definite.*

*Proof.* Calculating the gradient of $f(x)$

$$\nabla f(x) = \begin{bmatrix} f_{x_1} \\ f_{x_2} \end{bmatrix} = \begin{bmatrix} 200(x_2 - x_1^2)(2x_1) + 2(1 - x_1)(-1) \\ 200(x_2 - x_1^2) \end{bmatrix} = \begin{bmatrix} -400x_1(x_2 - x_1^2) + 2(x_1 - 1) \\ 200(x_2 - x_1^2) \end{bmatrix}$$

Solving $\nabla f(x) = 0$, we obtain that $\nabla f(x) = 0$ if and only if $x$ equals to $x^* = (1,1)^T$.

On the other hand, calculating the Hessian of $f(x)$

$$\nabla^2 f(x) = \begin{bmatrix} f_{x_1 x_1} & f_{x_1 x_2} \\ f_{x_2 x_1} & f_{x_2 x_2} \end{bmatrix} = \begin{bmatrix} 400(x_2 - x_1^2) - 400x_1(2x_2) + 4 & 400x_1 \\ -400x_1 & 200 \end{bmatrix}$$

Observe that

$$\begin{aligned}
P^T \nabla^2 f(x^*) P &= \begin{bmatrix} p_1 & p_2 \end{bmatrix} \begin{bmatrix} 804 & 400 \\ -400 & 200 \end{bmatrix} \begin{bmatrix} p_1 \\ p_2 \end{bmatrix} \\
&= 804p_1^2 + 400p_1 p_2 - 400p_1 p_2 + 200p_2^2 \\
&= 804p_1^2 + 200p_2^2
\end{aligned}$$

Therefore, $P^T \nabla^2 f(x^*) P > 0$ for all nonzero vector $P$. *i.e. the Hessian matrix at $x^*$ is positive definite.* By Theorem2.4 (Second-Order Sufficient Conditions), $x^* = (1,1)^T$ is the only local minimizer of this function. $\square$

**Exercise 7.** *Suppose that $f(x) = x^T Q x$, where $Q$ is an $n \times n$ symmetric positive semidefinite matrix. Show using the definition (1.4) that $f(x)$ is convex on the domain $\mathbb{R}^n$. Hint: It may be convenient to prove the following equivalent inequality:*

$$f(y + \alpha(x - y)) - \alpha f(x) - (1 - \alpha)f(y) \leq 0$$

*for all $\alpha \in [0,1]$ and all $x, y \in \mathbb{R}$.*

*Proof.* By the definition of $f$ and (1.4), for any $x, y \in \mathbb{R}^n$

$$\begin{aligned}
&f(y + \alpha(x - y)) - \alpha f(x) - (1 - \alpha)f(y) \\
&= (y + \alpha(x - y))^T Q(y + \alpha(x - y)) - \alpha x^T Q x - (1 - \alpha)y^T Q y \\
&= y^T Q y + \alpha y^T Q(x - y) + \alpha(x - y)^T Q y + \alpha^2(x - y)^T Q(x - y) - \alpha x^T Q x - (1 - \alpha)y^T Q y \\
&= \alpha[y^T Q y + y^T Q(x - y) - x^T Q x + (x - y)^T Q y] + \alpha^2(x - y)^T Q(x - y) \\
&= \alpha[y^T Q x - x^T Q x + (x - y)^T Q y] + \alpha^2(x - y)^T Q(x - y) \\
&= \alpha[-(x - y)^T Q x + (x - y)^T Q y] + \alpha^2(x - y)^T Q(x - y) \\
&= -\alpha(x - y)^T Q(x - y) + \alpha^2(x - y)^T Q(x - y) \\
&= (\alpha - \alpha^2)(x - y)^T Q(x - y) \leq 0
\end{aligned}$$

since $\alpha \in [0,1]$ and Q is positive semidefinite. This completes the proof. $\square$

**Exercise 8.** *Suppose that $f$ is a convex function. Show that the set of global minimizer of $f$ is a convex set.*

*Proof.* Let $E$ denotes the set of global minimizer of $f$.
For any $x^*, y^* \in E$, $f(x^*) \le f(y^*)$ and $f(y^*) \le f(x^*)$. i.e. $f(x^*) = f(y^*)$
Then for any $\alpha \in [0, 1]$

$$f(\alpha x^* + (1 - \alpha)y^*) \le \alpha f(x^*) + (1 - \alpha)f(y^*) = f(x^*) \le f(x) \qquad \forall x \in \mathbb{R}^n$$

since $x^*$ is a global minimizer.
By above, $f(\alpha x^* + (1 - \alpha)y^*) \le f(x), \forall x \in \mathbb{R}, \alpha \in [0, 1]$.
$E$ is a convex set. $\qquad\qquad\square$

**Exercise 16.** *Consider the sequence $x_k$ defined by*

$$x_k = \begin{cases} \left(\frac{1}{4}\right)^{2^k}, & k \text{ even}, \\ (x_{k-1})/k, & k \text{ odd}. \end{cases}$$

*Is this sequence Q-superlinearly convergent? Q-quadratically convergent? R-quadratically convergent?*

*Proof.* Clearly, $x_k$ converges to $x^* = 0$

(i) Q-superlinearly:

If k is even: $\displaystyle\lim_{k\to\infty} \frac{|x_{k+1} - x^*|}{|x_k - x*|} = \lim_{k\to\infty} \frac{\frac{x_k}{k+1}}{x_k} = \lim_{k\to\infty} \frac{1}{k+1} = 0$

If k is odd: $\displaystyle\lim_{k\to\infty} \frac{|x_{k+1} - x^*|}{|x_k - x*|} = \lim_{k\to\infty} \frac{x_{k+1}}{\frac{x_{k-1}}{k}} = \lim_{k\to\infty} \frac{k(\frac{1}{4})^{2^k}}{(\frac{1}{4})^{2^{k-1}}} = \lim_{k\to\infty} k(\frac{1}{4})^{2^{k-1}} = 0$

This implies $x_k$ is Q-superlinearly convergence.

(ii) Q-quadratically:

If k is even: $\displaystyle\lim_{k\to\infty} \frac{|x_{k+1} - x^*|}{|x_k - x*|^2} = \lim_{k\to\infty} \frac{\frac{x_k}{k+1}}{x_k^2} = \lim_{k\to\infty} \frac{1}{kx_k} = \lim_{k\to\infty} \frac{1}{k(\frac{1}{4})^{2^k}} = +\infty$ This implies $x_k$ is not Q-quadratically convergent.

(iii) R-quadratically:

Let $\epsilon_k = \begin{cases} \left(\frac{1}{4}\right)^{2^k}, & k \text{ even}, \\ \left(\frac{1}{4}\right)^{2^{k-1}}, & k \text{ odd}. \end{cases}$

When k is even, $|x_k - x^*| = |x_k| = (\frac{1}{4})^{2^k} \le \epsilon_k$
When k is odd, $|x_k - x^*| = |x_k| = x_{k-1}/k = (\frac{1}{4})^{2^{k-1}}(\frac{1}{k}) \le (\frac{1}{4})^{2^{k-1}} = \epsilon_k$
By the above, $|x_k - x^*| \le \epsilon_k \quad \forall k$

If k is even: $\displaystyle\lim_{k\to\infty} \frac{|\epsilon_{k+1} - 0|}{|\epsilon_k - 0|^2} = \lim_{k\to\infty} \frac{(\frac{1}{4})^{2^k}}{(\frac{1}{4})^{2^k}} = 1$

If k is odd: $\displaystyle\lim_{k\to\infty} \frac{|\epsilon_{k+1} - 0|}{|\epsilon_k - 0|^2} = \lim_{k\to\infty} \frac{(\frac{1}{4})^{2^{k+1}}}{(\frac{1}{4})^{2^{k-1}}} = \lim_{k\to\infty} [(\frac{1}{4})^{2^{k-1}}]^3 = 0$

This implies $\epsilon_k$ is Q-quadratically convergent.
And therefore, $x_k$ is R-quadratically convergent.

$\qquad\qquad\square$

# Numerical Optimization with applications: Homework 02

104021601 林俊傑

104021602 吳彥儒

104021615 黃翊軒

**Exercise 5.** *Prove that $\|Bx\| \geq \frac{\|x\|}{\|B^{-1}\|}$ for any nonsingular matrix $B$. Use this fact to establish (3.19).*

*Proof.* For simplicity, we drop the iteration index $k$ in the proof.

Note that from (3.2) we use the fact $P = -B^{-1}\nabla f$. Thus by multiplying both sides by $B$ and taking transport, we have $BP = -\nabla f$ and $P^T B^T = -\nabla f^T$. We are now prepared to estimate $\cos \theta$ :

$$
\begin{aligned}
\cos \theta &= \frac{-\nabla f^T P}{\|\nabla f\| \|P\|} \\
&= \frac{\left(P^T B^T\right) P}{\|\nabla f\| \|B^{-1} BP\|} \\
&= \frac{P^T BP}{\|BP\| \|B^{-1} BP\|} \\
&\geq \frac{P^T BP}{\|B\| \|P\| \|B^{-1}\| \|BP\|} \\
&= \left(\frac{P^T BP}{\|P\| \|BP\|}\right) \frac{1}{\|B\| \|B^{-1}\|} \\
&\geq \frac{1}{\|B\| \|B^{-1}\|} \\
&\geq \frac{1}{M},
\end{aligned}
$$

where the last two inequality hold by the assumption that $B$ is positive definite and has a uniformly bounded condition number. Therefore, (3.19) follows. $\square$

**Exercise 7.** *Prove the result (3.28) by working through the following steps. First, use (3.26) to show that*

$$
\|x_k - x^*\|_Q^2 - \|x_{k+1} - x^*\|_Q^2 = 2\alpha_k \nabla f_k^T Q(x_k - x^*) - \alpha_k^2 \nabla f_k^T Q \nabla f_k
$$

*where $\|\cdot\|_Q$ is defined by (3.27). Second, use the fact that $\nabla f_k = Q(x_k - x^*)$ to obtain*

$$
\|x_k - x^*\|_Q^2 - \|x_{k+1} - x^*\|_Q^2 = \frac{2(\nabla f_k^T \nabla f_k)^2}{(\nabla f_k^T Q \nabla f_k)} - \frac{(\nabla f_k^T \nabla f_k)^2}{(\nabla f_k^T Q \nabla f_k)}
$$

*and*

$$
\|x_k - x^*\|_Q^2 = \nabla f_k^T Q^{-1} \nabla f_k.
$$

*Proof.* (1)

$$
\begin{aligned}
&\|x_k - x^*\|_Q^2 - \|x_{k+1} - x^*\|_Q^2 \\
&= 2f(x_k) - 2f(x_{k+1}) \\
&= x_k^T Q x_k - 2b^T x_k - (x_k - \alpha_k \nabla f_k)^T Q(x_k - \alpha_k \nabla f_k) + 2b^T(x_k - \alpha_k \nabla f_k) \\
&= x_k^T Q(\alpha_k \nabla f_k) + (\alpha_k \nabla f_k)^T Q x_k - \alpha_k^2 \nabla f_k^T Q \nabla f_k - 2\alpha_k b^T \nabla f_k \\
&= 2\alpha_k \nabla f_k^T Q x_k - \alpha_k^2 \nabla f_k^T Q \nabla f_k - 2\alpha_k (Qx^*)^T \nabla f_k \\
&= 2\alpha_k \nabla f_k^T Q(x_k - x^*) - \alpha_k^2 \nabla f_k^T Q \nabla f_k
\end{aligned}
$$

(2) Combining the result of (1) and the fact that $\nabla f_k = Q(x_k - x^*)$, we have

$$
\begin{aligned}
&\|x_k - x^*\|_Q^2 - \|x_{k+1} - x^*\|_Q^2 \\
&= 2\alpha_k \nabla f_k^T \nabla f_k - \alpha_k^2 \nabla f_k^T Q \nabla f_k \\
&= 2\left(\frac{\nabla f_k^T \nabla f_k}{\nabla f_k^T Q \nabla f_k}\right)\nabla f_k^T \nabla f_k - \left(\frac{\nabla f_k^T \nabla f_k}{\nabla f_k^T Q \nabla f_k}\right)^2 \nabla f_k^T Q \nabla f_k \\
&= \frac{2(\nabla f_k^T \nabla f_k)^2}{(\nabla f_k^T Q \nabla f_k)} - \frac{(\nabla f_k^T \nabla f_k)^2}{(\nabla f_k^T Q \nabla f_k)}
\end{aligned}
$$

(3) The definition of the weight norm: $\|x\|_Q^2 = x^T Q x$, therefore

$$
\begin{aligned}
\|x_k - x^*\|_Q^2 &= (x_k - x^*)^T Q(x_k - x^*) \\
&= (x_k - x^*)T(QQ^{-1})Q(x_k - x^*) \\
&= (Q(x_k - x^*))^T Q^{-1} Q(x_k - x^*) \\
&= \nabla f_k^T Q^{-1} \nabla f_k
\end{aligned}
$$

Since $Q$ is symmetric and nonsingular.

(4) Now we turn to prove (3.28) by using the result of (2) and (3).

$$
\begin{aligned}
\|x_{k+1} - x^*\|_Q^2 &= \|x_k - x^*\|_Q^2 - \frac{(\nabla f_k^T \nabla f_k)^2}{(\nabla f_k^T Q \nabla f_k)} \\
&= \|x_k - x^*\|_Q^2 - \frac{(\nabla f_k^T \nabla f_k)^2}{(\nabla f_k^T Q \nabla f_k)(\nabla f_k^T Q^{-1} \nabla f_k)}(\nabla f_k^T Q^{-1} \nabla f_k) \\
&= \left(1 - \frac{(\nabla f_k^T \nabla f_k)^2}{(\nabla f_k^T Q \nabla f_k)(\nabla f_k^T Q^{-1} \nabla f_k)}\right)\|x_k - x^*\|_Q^2
\end{aligned}
$$

Therefore the proof is completed.

$\square$

**Exercise 8.** *Let $Q$ be a positive definite symmetric matrix. Prove that for any vector $x$, we have*

$$
\frac{(x^T x)^2}{(x^T Q x)(x^T Q^{-1} x)} \geq \frac{4\lambda_n \lambda_1}{(\lambda_n + \lambda_1)^2}
$$

*where $\lambda_n$ and $\lambda_1$ are, respectively the largest and smallest eigenvalues of $Q$. (This relation, which is known as the Kantorovich inequality, can be used to deduce (3.29) from (3.28).)*

*Proof.* Since $Q$ is positive definite and symmetric, we have eigenvalue decompsition $Q = U\Lambda U^T$. Let $x = Uy$. Then

$$
\frac{(x^T x)^2}{(x^T Q x)(x^T Q^{-1} x)} = \frac{(y^T y)^2}{(y^T \Lambda y)(y^T \Lambda^{-1} y)} = \frac{(\sum_{i=1}^n y_i^2)^2}{(\sum_{i=1}^n \lambda_i y_i^2)(\sum_{i=1}^n y_i^2/\lambda_i)}
$$

Let $\eta_i = \dfrac{y_i^2}{\sum_{j=1}^n y_j^2}$ and $f(\lambda) = \frac{1}{\lambda}$. Then

$$
\frac{(x^T x)^2}{(x^T Q x)(x^T Q^{-1} x)} = \frac{1}{(\sum_{i=1}^n \lambda_i \eta_i)(\sum_{i=1}^n f(\lambda_i)\eta_i)}
$$

Let $\lambda = \sum_{i=1}^n \lambda_i \eta_i$ , $\lambda_f = \sum_{i=1}^n f(\lambda_i)\eta_i$

Since $\eta_i \geq 0 \quad \forall i$ and $\sum_{i=1}^n \eta_i = 1$, $\lambda_1 \leq \lambda \leq \lambda_n$

Write $\lambda_i = \dfrac{\lambda_n - \lambda_i}{\lambda_n - \lambda_1}\lambda_1 + \dfrac{\lambda_i - \lambda_1}{\lambda_n - \lambda_1}\lambda_n$. This shows $\lambda_i$ is a convex combination of $\lambda_1$ and $\lambda_n$ $\quad \forall i$

$\because f$ is convex $\quad \therefore f(\lambda_i) \le \dfrac{\lambda_n - \lambda_i}{\lambda_n - \lambda_1}f(\lambda_1) + \dfrac{\lambda_i - \lambda_1}{\lambda_n - \lambda_1}f(\lambda_n)$

Therefore,

$$\lambda_f \le \sum_{i=1}^{n}\left[\frac{\lambda_n - \lambda_i}{\lambda_n - \lambda_1}f(\lambda_1) + \frac{\lambda_i - \lambda_1}{\lambda_n - \lambda_1}f(\lambda_n)\right]\eta_i = \sum_{i=1}^{n}\left[\frac{\lambda_n - \lambda_i}{\lambda_n - \lambda_1}\frac{1}{\lambda_1} + \frac{\lambda_i - \lambda_1}{\lambda_n - \lambda_1}\frac{1}{\lambda_n}\right]\eta_i$$

$$= \sum_{i=1}^{n}\frac{\eta_i}{\lambda_n - \lambda_1}\left[\frac{\lambda_n - \lambda_i}{\lambda_1} + \frac{\lambda_i - \lambda_1}{\lambda_n}\right] = \sum_{i=1}^{n}\frac{\eta_i}{\lambda_n - \lambda_1}\left[\frac{\lambda_n^2 - \lambda_i\lambda_n + \lambda_i\lambda_1 - \lambda_1^2}{\lambda_n\lambda_1}\right]$$

$$= \sum_{i=1}^{n}\frac{\eta_i}{\lambda_n - \lambda_1}\left[\frac{(\lambda_n + \lambda_1)(\lambda_n - \lambda_1) - \lambda_i(\lambda_n - \lambda_1)}{\lambda_n\lambda_1}\right] = \sum_{i=1}^{n}\frac{\lambda_n + \lambda_1 - \lambda_i}{\lambda_n\lambda_1}\eta_i$$

$$= \frac{\lambda_n \sum_{i=1}^{n}\eta_i + \lambda_1\sum_{i=1}^{n}\eta_i - \sum_{i=1}^{n}\lambda_i\eta_i}{\lambda_n\lambda_1} = \frac{\lambda_n + \lambda_1 - \lambda}{\lambda_n\lambda_1}$$

We conclude that

$$\frac{(x^Tx)^2}{(x^TQx)(x^TQ^{-1}x)} = \frac{1}{\lambda\lambda_f} \ge \frac{\lambda_n\lambda_1}{\lambda(\lambda_n + \lambda_1 - \lambda)} \ge \frac{\lambda_n\lambda_1}{\max_{\lambda\in[\lambda_1,\lambda_n]}\lambda(\lambda_n + \lambda_1 - \lambda)}$$

Let $g(\lambda) = \lambda(\lambda_n + \lambda_1 - \lambda) = -\lambda^2 + (\lambda_n + \lambda_1)\lambda$. Then $g(\lambda)$ has maxmun at $\bar\lambda = \dfrac{\lambda_n + \lambda_1}{2} \in [\lambda_1, \lambda_n]$.

$$g(\bar\lambda) = -\frac{(\lambda_n + \lambda_1)^2}{4} + \frac{(\lambda_n + \lambda_1)^2}{2} = \frac{(\lambda_n + \lambda_1)^2}{4}$$

This implies that

$$\frac{(x^Tx)^2}{(x^TQx)(x^TQ^{-1}x)} \ge \frac{\lambda_n\lambda_1}{g(\bar\lambda)} = \frac{4\lambda_n\lambda_1}{(\lambda_n + \lambda_1)^2}$$

$\square$

**Exercise 13.** *Show that the quadratic function that interpolates $\phi(0)$, $\phi'(0)$, and $\phi(\alpha_0)$ is given by (3.57). Then, make use of the fact that the sufficient decrease condition (3.6a) is not satisfied at $\alpha(0)$ to show that this quadratic has positive curvature and that the minimizer satisfies*

$$\alpha_1 < \frac{\alpha_0}{2(1 - c_1)}.$$

*Since $c_1$ is chosen to be quite small in practice, this inequality indicates that $\alpha_1$ cannot be much greater than $\frac{1}{2}$ (and may be smaller), which gives us an idea of the new step length.*

*Proof.* By assuming a quadratic function

$$\phi_q(\alpha) = a\alpha^2 + b\alpha + c$$

and solving coefficients through standard calculations, we have

$$\phi_q(0) = c = \phi(0).$$

Also, since

$$\phi_q'(\alpha) = 2a\alpha + b,$$

we find

$$\phi_q'(0) = b = \phi'(0).$$

On the other hand,

$$\phi_q(\alpha_0) = \phi(\alpha_0) = a\alpha_0^2 + \phi'(0)\alpha_0 + \phi(0),$$

3

we obtain
$$a = \frac{\phi(\alpha_0) - \phi(0) - \phi'(0)\alpha_0}{\alpha_0^2}.$$

By assumptions, we now discuss a senario that $(3.6a)$ is not satisfied at $\alpha_0$. Thus we have

$$\phi(\alpha_0) > \phi(0) + c_1\alpha_0\phi'(0).$$

Hence, a minor manipulation of this inequality gives

$$a = \frac{\phi(\alpha_0) - \phi(0) - \phi'(0)\alpha_0}{\alpha_0} > 0,$$

which guarantees that the quadratic has positive curvature. Moreover, since the minimizer of a quadratic is $-b/2a$, we obtain

$$\begin{aligned}
\alpha_1 = \frac{-b}{2a} &= \frac{-\alpha_0^2\phi'(0)}{2[\phi(\alpha_0) - \phi(0) - \phi'(0)\alpha_0]} \\
&= \frac{\alpha_0}{2[1 - \frac{\phi(\alpha_0)-\phi(0)}{\alpha_0\phi'(0)}]} \\
&\leq \frac{\alpha_0}{2(1 - c_1)},
\end{aligned}$$

where the second equality holds by dividing $-\alpha_0\phi'(0)$ both upper and lower sides, and the last inequality follows from the given senario.

$\square$

# Numerical Optimization with applications: Homework 03

104021601 林俊傑
104021602 吳彥儒
104021615 黃翊軒

November 2, 2016

**Exercise 6.** *The Cauchy-Schwarz inequality states that for any vectors u and v, we have*

$$|u^T v|^2 \le (u^T u)(v^T v),$$

*with equality only when u and v are parallel. When B is positive definite, use this inequality to show that*

$$\gamma := \frac{\|g\|^4}{(g^T B g)(g^T B^{-1} g)} \le 1,$$

*with equality only if g and Bg (and $B^{-1}g$) are parallel.*

*Proof.* B is a positive definite matrix, so there exists an orthonormal matrix $Q$ and a diagonal matrix

$$\Lambda = \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{pmatrix} \text{ s.t. } B = Q\Lambda Q^T.$$

Define the matrix $\sqrt{B} = Q\sqrt{\Lambda}Q^T$ where $\sqrt{\Lambda} = \begin{pmatrix} \sqrt{\lambda_1} & & & \\ & \sqrt{\lambda_2} & & \\ & & \ddots & \\ & & & \sqrt{\lambda_n} \end{pmatrix}$

Obviously, $\sqrt{B}$ is also symmetric.
Claim: $(\sqrt{B})^{-1} = \sqrt{B^{-1}}$
proof of claim:
$(\sqrt{B})^{-1} = Q(\sqrt{\Lambda})^{-1}Q^T = Q\sqrt{\Lambda^{-1}}Q^T = \sqrt{B^{-1}}$
We proved the claim.
Now we use the claim above, the symmetricity of $\sqrt{B}$ and Cauchy-Schwarz inequality. We have the following statement:

$$\begin{aligned}
\|g\|^4 = (g^T g)^2 = (g^T \sqrt{B}(\sqrt{B})^{-1}g)^2 &= ((\sqrt{B}g)^T(\sqrt{B^{-1}}g))^2 \\
&\le (\sqrt{B}g)^T(\sqrt{B}g)(\sqrt{B^{-1}}g)^T(\sqrt{B^{-1}}g) \\
&= (g^T\sqrt{B}\sqrt{B}g)(g^T\sqrt{B^{-1}}\sqrt{B^{-1}}g) \\
&= (g^T B g)(g^T B^{-1} g)
\end{aligned}$$

When the equality holds only if $\sqrt{B}g$ and $\sqrt{B^{-1}}g$ are parallel.
i.e. $\sqrt{B}g = k\sqrt{B^{-1}}g$ for some constant $k$.
1. Multiplying both side by $\sqrt{B}$.
   $Bg = kg \implies Bg$ and $g$ are parallel.

2. Multiplying both side by $\sqrt{B^{-1}}$.
   $g = kB^{-1}g \implies B^{-1}g$ and $g$ are parallel.
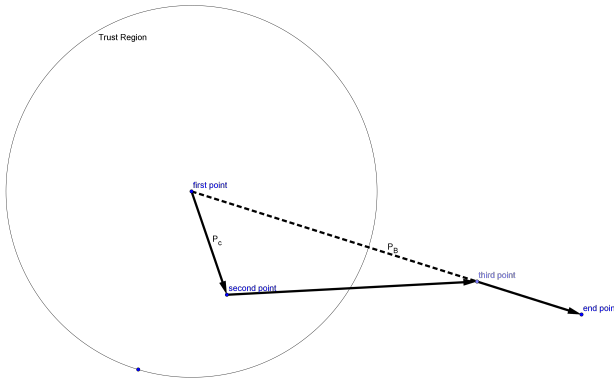
$\square$

**Exercise 7.** *When $B$ is positive definite, the double-dog leg method constructs a path with three line segments from the origin to the full step. The four points that define the path are*

- *the origin;*

- *the unconstrained Cauchy step $p^c = -(g^T g)/(g^T B g)g$;*

- *a fraction of the full step $\bar{\gamma}p^B = -\bar{\gamma}B^{-1}g$, for some $\bar{\gamma} \in (\gamma, 1]$, where $\gamma$ is defined in the previous question; and*

- *the full step $p^B = -B^{-1}g$*

*Show that $\|p\|$ increases monotonically along this path.*
*(Note: The double-dogleg method, as discussed in Dennis and Schnabel [92, Section 6.4.2], was for some time thought to be superior to the standard dogleg method, but later testing has not shown much difference in performance.)*



*Proof.* It is obviously that $\|p\|$ increases monotonically along the first segment and the last segment because $\alpha\|v\|$ increases as $\alpha$ increases, where $\alpha \in (0,1)$. Now we consider the second segment. Let $P^A = -\bar{\gamma}B^{-1}g$, and $P^U = -(g^T g)/(g^T B g)g$, then the parametrization of the second segment is

$$P(\alpha) = \alpha(P^A - P^U) + P^U.$$

Define

$$
\begin{aligned}
h(\alpha) &= (1/2)\|P(\alpha)\|^2 \\
&= (1/2)\|\alpha(P^A - P^U) + P^U\|^2 \\
&= (1/2)\|P^U\|^2 + \alpha(P^U)^T(P^A - P^U) + (1/2)\alpha^2\|P^A - P^U\|^2
\end{aligned}
$$

Then we have

$$
\begin{aligned}
h'(\alpha) &= -(P^U)^T(P^U - P^B) + \alpha\|P^U - P^B\|^2 \\
&\geq -(P^U)^T(P^U - P^A) \\
&= \frac{g^T g}{g^T B g}g^T\left(-\frac{g^T g}{g^T B g}g + \bar{\gamma}B^{-1}g\right) \\
&= g^T g\frac{gB^{-1}g}{gBg}\left(\bar{\gamma} - \frac{(g^T g)^2}{(g^T B g)(g^T B^{-1}g)}\right) \\
&> 0
\end{aligned}
$$

Since $h'(\alpha) > 0$ for all $\alpha \in (0,1)$, $h(\alpha)$ is increasing monotonically on $(0,1)$, that is, $\|p(\alpha)\|$ is increasing monotonically on $(0,1)$. Therefore, the $\|p\|$ increases monotonically along this segament. $\qquad\square$

**Exercise 8.** *Show that*

$$\lambda^{(l+1)} = \lambda^{(l)} - \frac{\phi_2(\lambda^{(l)})}{\phi_2'(\lambda^{(l)})}, \quad and \quad \lambda^{(l+1)} = \lambda^{(l)} + \left(\frac{\|p_l\|}{\|q_l\|}\right)^2 \left(\frac{\|p_l\| - \Delta}{\Delta}\right)$$

*are equivalents.*

*Proof.* First, we caculate

$$\phi_2'(\lambda) = \frac{d}{d\lambda}\left(\frac{1}{\|p(\lambda)\|}\right) = \frac{d}{d\lambda}\left(\|p(\lambda)\|^2\right)^{-1/2} = -\frac{1}{2}\left(\|p(\lambda)\|^2\right)^{-3/2}\frac{d}{d\lambda}\|p(\lambda)\|^2$$

Since B is symmetric, there is an orthonormal matrix $U$ and a diagonal matrix $\Lambda$ such that $B = U\Lambda U^T$, where

$$\Lambda = diag(\lambda_1, \lambda_2, \cdots, \lambda_n)$$

Then, $B + \lambda I = U(\Lambda + \lambda I)U^T$. We have,

$$p(\lambda) = -U(\Lambda + \lambda I)U^T g = -\sum_{j=1}^{n} -\frac{u_j^T g}{\lambda_j + \lambda}u_j$$

where $u_j$ denotes the jth column of $U$. Therefore, by orthonormality of $u_1, u_2, \cdots, u_n$, we have

$$\|p(\lambda)\|^2 = \sum_{j=1}^{n} \frac{(u_j^T g)^2}{(\lambda_j + \lambda)^2}$$

Hence, we can caculate

$$\frac{d}{d\lambda}\|p(\lambda)\|^2 = -2\sum_{j=1}^{n} \frac{(u_j^T g)^2}{(\lambda_j + \lambda)^3}$$

On the other hand, we have
$$\|q_l\|^2 = \left\|R^{-T}p_l\right\|^2 = p_l^T R^{-1}R^{-T}p_l = [-(R^T R)^{-1}g]^T (R^T R)^{-1}[-(R^T R)^{-1}g] = g^T[(R^T R)^{-1}]^3 g$$
$$= g^T[(B + \lambda^{(l)}I)^{-1}]^3 g = g^T U(\Lambda + \lambda^{(l)}I)^{-3}U^T g = \sum_{j=1}^{n} \frac{(u_j^T g)^2}{(\lambda_j + \lambda^{(l)})^3}$$

We conclude that

$$\phi_2'(\lambda^{(l)}) = -\frac{1}{2}\left\|p(\lambda^{(l)})\right\|^{-3}\left(-2\sum_{j=1}^{n} \frac{(u_j^T g)^2}{(\lambda_j + \lambda^{(l)})^3}\right) = \|p_l\|^{-3}\|q_l\|^2$$

Finally, we get

$$-\frac{\phi_2(\lambda^{(l)})}{\phi_2'(\lambda^{(l)})} = \left(\frac{1}{\Delta} - \frac{1}{\|p(\lambda^{(l)})\|}\right)\left(\frac{\|p_l\|^3}{\|q_l\|^2}\right) = \left(\frac{\|p_l\| - \Delta}{\Delta\|p_l\|}\right)\left(\frac{\|p_l\|^3}{\|q_l\|^2}\right) = \left(\frac{\|p_l\|}{\|q_l\|}\right)^2\left(\frac{\|p_l\| - \Delta}{\Delta}\right)$$

Therefore, the two equations above are equivalent. □

# Numerical Optimization with applications: Homework 04

104021601 林俊傑
104021602 吳彥儒
104021615 黃翊軒

November 9, 2016

**Exercise 1.** *Implement Algorithm 5.2 and use to it solve linear systems in which $A$ is the Hilbert matrix, whose elements are $A_{i,j} = 1/(i+j-1)$. Set the right-hand-side to $b = (1, 1, ..., 1)^T$ and the initial point to $x_0 = 0$. Try dimensions $n = 5, 8, 12, 20$ and report the number of iterations required to reduce the residual below $10^{-6}$.*

**Solution.** The numbers of iterations as the table below.

| n | 5 | 8 | 12 | 20 |
|---|---|---|---|---|
| number of iteration | 6 | 19 | 38 | 73 |
| condition number | 4.766E+05 | 1.526E+10 | 1.633E+16 | 2.596E+18 |

Observe that the condition number in the case $n = 20$ is greater than the others. By(5.36), the rate of convergence should be less than the others. ◄

**Exercise 2.** *Show that if the nonzero vectors $p_0, p_1, ..., p_l$ satisfy (5.5), where $A$ is symmetric and positive definite, then these vectors are linearly independent. (This result implies that $A$ has at most $n$ conjugate direction.)*

*Proof.* Suppose $a_0 p_0 + a_1 p_1 + ... + a_l p_l = 0$. For any $p_j$, we have the following argument.

$$
\begin{aligned}
0 &= p_j^T A(a_0 p_0 + a_1 p_1 + ... + a_l p_l) \\
&= a_0(p_j^T A p_0) + a_1(p_j^T A p_1) + ... + a_j(p_j^T A p_j) + ... + a_l(p_j^T A p_l) \\
&= a_0 \cdot 0 + a_1 \cdot 0 + ... + a_j \cdot (p_j^T A p_j) + ... + a_l \cdot 0 \\
&= a_j \cdot (p_j^T A p_j)
\end{aligned}
$$

Since $A$ is positive definite, $p_j^T A p_j > 0$, this implies $a_j = 0. \quad \forall j$
Consequently, $p_0, p_1, ..., p_l$ are linearly independent.

$\square$

**Exercise 4.** *Show that if $f(x)$ is a strictly convex quadratic, then the function $h(\sigma) \stackrel{def}{=} f(x_0 + \sigma_0 p_0 + \cdots + \sigma_{k-1} p_{k-1})$ also is a strictly convex quadratic in the variable $\sigma = (\sigma_0, \sigma_1, \cdots, \sigma_{k-1})^T$.*

*Proof.* By the definition of strictly convex quadratic function, we can assume

$$
f(x) = \frac{1}{2} x^T A x - b^T x,
$$

where $A$ is a positive definite symmetric matrix and $b$ is a constant vector. We want prove that $h(\sigma)$ is also a strictly convex quadratic function by showing

$$
h(\sigma) = \frac{1}{2} \sigma^T B \sigma - c^T \sigma + d,
$$

where $B$ is a positive definite symmetric matrix and $c, d$ are constant vectors. Since $p_i^T A p_j = 0$ for all

$i \neq j$, we obtain that

$$
\begin{aligned}
h(\sigma) &= f(x_0 + \sigma_0 p_0 + \cdots + \sigma_{k-1} p_{k-1}) \\
&= \frac{1}{2}(x_0 + \sigma_0 p_0 + \cdots + \sigma_{k-1} p_{k-1})^T A(x_0 + \sigma_0 p_0 + \cdots + \sigma_{k-1} p_{k-1}) - b^T(x_0 + \sigma_0 p_0 + \cdots + \sigma_{k-1} p_{k-1}) \\
&= \frac{1}{2}x_0^T A(x_0 + \sigma_0 p_0 + \cdots + \sigma_{k-1} p_{k-1}) + \frac{1}{2}(\sigma_0 p_0)^T A(x_0 + \sigma_0 p_0 + \cdots + \sigma_{k-1} p_{k-1}) + \cdots \\
&\quad + \frac{1}{2}(\sigma_{k-1} p_{k-1})^T A(x_0 + \sigma_0 p_0 + \cdots + \sigma_{k-1} p_{k-1}) - b^T(x_0 + \sigma_0 p_0 + \cdots + \sigma_{k-1} p_{k-1}) \\
&= \frac{1}{2}(\sigma_0 p_0)^T A(\sigma_0 p_0) + \frac{1}{2}(\sigma_1 p_1)^T A(\sigma_1 p_1) + \cdots + \frac{1}{2}(\sigma_{k-1} p_{k-1})^T A(\sigma_{k-1} p_{k-1}) \\
&\quad + \frac{1}{2}x_0^T A(x_0 + \sigma_0 p_0 + \cdots + \sigma_{k-1} p_{k-1}) - b^T(x_0 + \sigma_0 p_0 + \cdots + \sigma_{k-1} p_{k-1}) \\
&= \frac{1}{2}\sigma^T B \sigma + \frac{1}{2}x_0^T A P \sigma - b^T P \sigma + \frac{1}{2}x_0^T A x_0 - b^T x_0 \\
&= \frac{1}{2}\sigma^T B \sigma + (\frac{1}{2}x_0^T A P - b^T P)\sigma + \frac{1}{2}x_0^T A x_0 - b^T x_0
\end{aligned}
$$

where $B = \begin{bmatrix} p_0^T A p_0 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & p_{k-1}^T A p_{k-1} \end{bmatrix}$ is positive definite sysmmetric matrix and $P = (p_0, p_1, \cdots, p_{k-1})$
and $\sigma = (\sigma_0, \sigma_1, \cdots, \sigma_{k-1})^T$. Hence, $h(\sigma)$ is also a strictly convex quadratic function. $\square$

**Exercise 7.** *Let $\{\lambda_i, v_i\}$ $i = 1, 2, \cdots, n$ be the eigenpairs of the symmetric matrix $A$. Show that the eigenvalues and eigenvectors of $[I + P_k(A)A]^T A[I + P_k(A)A]$ are $\lambda_i[1 + \lambda_i P_k(\lambda_i)]^2$ and $v_i$, respectively.*

*Proof.* We first show that
$$
P_k(A)v_i = P_k(\lambda_i)v_i
$$
for any polynomials $P_k(x)$ of degree $k$.

Let $P_k(x) = \sum_{j=0}^{k} a_j x^j$. Then
$$
P_k(A)v_i = \sum_{j=0}^{k} a_j A^j v_i = \sum_{j=0}^{k} a_j A^{j-1}(\lambda_i v_i) = \sum_{j=0}^{k} a_j A^{j-2}(\lambda_i^2 v_i) = \cdots = \sum_{j=0}^{k} a_j \lambda_i^j v_i = P_k(\lambda_i)v_i
$$
Since $[I + P_k(x)x]$ is a ploynomial, we have

$$
[I + P_k(A)A]v_i = [1 + \lambda_i P_k(\lambda_i)]v_i
$$

$A$ is symmetric, therefore, $[I + P_k(A)A]^T = [I + P_k(A)A]$
Now, we are ready to compute

$$
\begin{aligned}
[I + P_k(A)A]^T A[I + P_k(A)A]v_i &= [I + P_k(A)A]A[I + P_k(A)A]v_i \\
&= [I + P_k(A)A]A[1 + \lambda_i P_k(\lambda_i)]v_i \\
&= [I + P_k(A)A](Av_i)[1 + \lambda_i P_k(\lambda_i)] \\
&= [I + P_k(A)A]\lambda_i v_i[1 + \lambda_i P_k(\lambda_i)] \\
&= [I + P_k(A)A]v_i \lambda_i[1 + \lambda_i P_k(\lambda_i)] \\
&= [1 + \lambda_i P_k(\lambda_i)]v_i \lambda_i[1 + \lambda_i P_k(\lambda_i)] \\
&= \lambda_i[1 + \lambda_i P_k(\lambda_i)]^2 v_i
\end{aligned}
$$

We conclude that $\{\lambda_i[1 + \lambda_i P_k(\lambda_i)]^2, v_i\}$ $i = 1, 2, \cdots, n$ are the eigenpairs of
$[I + P_k(A)A]^T A[I + P_k(A)A]$. $\square$

# Numerical Optimization with applications: Homework 05

104021601 林俊傑

104021602 吳彥儒

104021615 黃翊軒

November 23, 2016

**Exercise 6.** *The square root of a matrix $A$ is a matrix $A^{1/2}$ such that $A^{1/2}A^{1/2} = A$. Show that any symmetric positive definite matrix $A$ has a square root, and that this square root is itself symmetric and positive definite.(Hint: factorization $A = UDU^T$ (A.16), where $U$ is orthogonal and $D$ is diagonal with positive diagonal elements.)*

*Proof.* First, we show that a real symmetric matrix $A$ is diagonalizable. Prove it by contradiction, which means there is a generalized eigenvector $v$ of order 2, that is $(A - \lambda I)v \neq 0$ and $(A - \lambda I)^2 = 0$, and we have the following statement.

$$0 = v^T(A - \lambda I)^2 v = v^T(A - \lambda I)^T(A - \lambda I)v$$
$$= \|(A - \lambda I)v\|^2 \neq 0 \rightarrow\leftarrow$$

Thus, every eigenvector of $A$ is of order 1 and $A$ is diagonalizable. We may Assume $A = UDU^T$, where $U$ is orthogonal and by $A > 0$, $D = diag(\lambda_1, \lambda_2, ..., \lambda_n)$ is diagonal with positive diagonal elements. Define the square root of $A$

$$A^{1/2} := U\sqrt{D}U^T = Udiag(\sqrt{\lambda_1}, \sqrt{\lambda_2}, ...\sqrt{\lambda_n})U^T$$

Obviously, $A^{1/2}$ is symmetric, and positive number $\sqrt{\lambda_i}$ is the eigenvalue correspond to the $i$th column vector of $U$. Hence $A^{1/2}$ is also positive definite. □

**Exercise 10.** *(a) Show that $\det(I + xy^T) = 1 + y^Tx$, where $x$ and $y$ are n-vectors.*

*(b) Using similar technique to prove that*

$$\det(I + xy^T + uv^T) = (1 + y^Tx)(1 + v^Tu) - (x^Tv)(y^Tu).$$

*(c) Use this relation to establish*

$$\det(B_{k+1}) = \det(B_k)\frac{y_k^Ts_k}{s_k^TB_ks_k}.$$

*Proof.* (a) Assuming $x \neq 0$, we can find vectors $q_1, q_2, \cdots, q_{n-1}$ such that the matrix Q defined by

$$Q = [x, q_1, q_2, \cdots, q_{n-1}]$$

is nonsingular and $x = Qe_1$. If we define

$$y^TQ = (w_1, w_2, \cdots, w_n)$$

then

$$w_1 = y^TQe_1 = y^Tx$$

and
$$\det(I + xy^T) = \det(Q^{-1}(I + xy^T)Q) = \det(I + Q^{-1}xy^TQ) = \det(I + e_1y^TQ)$$

$$= \det\left(I + \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}(w_1, w_2, \cdots, w_n)\right) = \det\left(\begin{bmatrix} 1 + w_1 & w_2 & \cdots & w_{n-1} & w_n \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & \cdots & 0 & 1 \end{bmatrix}\right)$$

$$= \det\left(\begin{bmatrix} 1+w_1 & w_2 & \cdots & w_{n-1} \\ 0 & 1 & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}\right) = \cdots = \det\left(\begin{bmatrix} 1+w_1 & w_2 \\ 0 & 1 \end{bmatrix}\right)$$

$$= 1 + w_1 = 1 + y^T x$$

(b) Assuming $x, u \neq 0$, we can find vectors $q_1, q_2, \cdots, q_{n-2}$ such that the matrix Q defined by

$$Q = [x, u, q_1, q_2, \cdots, q_{n-2}]$$

is nonsingular and $x = Qe_1$, $u = Qe_2$. If we define

$$y^T Q = (w_1, w_2, \cdots, w_n)$$
$$v^T Q = (z_1, z_2, \cdots, z_n)$$

then

$$w_1 = y^T Q e_1 = y^T x \qquad w_2 = y^T Q e_2 = y^T u$$
$$z_1 = v^T Q e_1 = v^T x \qquad z_2 = v^T Q e_2 = v^T u$$

and

$$\det(I + xy^T + uv^T) = \det\left(I + [x \quad u]\begin{bmatrix} y^T \\ v^T \end{bmatrix}\right) = \det\left(Q^{-1}(I + [x \quad u]\begin{bmatrix} y^T \\ v^T \end{bmatrix})Q\right)$$

$$= \det\left(I + [Q^{-1}x \quad Q^{-1}u]\begin{bmatrix} y^T Q \\ v^T Q \end{bmatrix}\right) = \det\left(I + [e_1 \quad e_2]\begin{bmatrix} w_1 & w_2 & \cdots & w_n \\ z_1 & z_2 & \cdots & z_n \end{bmatrix}\right)$$

$$= \det\left(I + \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ \vdots & \vdots \\ 0 & 0 \end{bmatrix}\begin{bmatrix} w_1 & w_2 & \cdots & w_n \\ z_1 & z_2 & \cdots & z_n \end{bmatrix}\right) = \det\left(\begin{bmatrix} 1+w_1 & w_2 & \cdots & w_{n-1} & w_n \\ z_1 & 1+z_2 & \cdots & z_{n-1} & z_n \\ \vdots & & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & \cdots & 0 & 1 \end{bmatrix}\right)$$

$$= \det\left(\begin{bmatrix} 1+w_1 & w_2 & \cdots & w_{n-1} \\ z_1 & 1+z_2 & \cdots & z_{n-1} \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}\right) = \cdots = \det\left(\begin{bmatrix} 1+w_1 & w_2 \\ z_1 & 1+z_2 \end{bmatrix}\right)$$

$$= (1+w_1)(1+z_2) - z_1 w_2 = (1 + y^T x)(1 + v^T u) - (x^T v)(y^T u)$$

(c) We have $B_{k+1} = B_k - \dfrac{B_k s_k s_k^T B_k}{s_k^T B_k s_k} + \dfrac{y_k y_k^T}{y_k^T s_k}$. So,

$$\det(B_{k+1}) = \det(B_k)\det\left(I + \left(\frac{-s_k}{s_k^T B_k s_k}\right)(s_k^T B_k)\left(\frac{B_k^{-1} y_k}{y_k^T s_k}\right)(y_k^T)\right)$$

Let

$$x = \left(\frac{-s_k}{s_k^T B_k s_k}\right) \quad y^T = (s_k^T B_k) \quad u = \left(\frac{B_k^{-1} y_k}{y_k^T s_k}\right) \quad v^T = (y_k^T)$$

Using (b), we can caculate

$$\det\left(I + \left(\frac{-s_k}{s_k^T B_k s_k}\right)(s_k^T B_k) + \left(\frac{B_k^{-1} y_k}{y_k^T s_k}\right)(y_k^T)\right)$$

$$= \left[1 + (s_k^T B_k)\left(\frac{-s_k}{s_k^T B_k s_k}\right)\right]\left[1 + (y_k^T)\left(\frac{B_k^{-1} y_k}{y_k^T s_k}\right)\right] - \left[(y_k^T)\left(\frac{-s_k}{s_k^T B_k s_k}\right)\right]\left[(s_k^T B_k)\left(\frac{B_k^{-1} y_k}{y_k^T s_k}\right)\right]$$

$$= 0 \times \left[1 + (y_k^T)\left(\frac{B_k^{-1} y_k}{y_k^T s_k}\right)\right] - \left[\frac{-y_k^T s_k}{s_k^T B_k s_k}\right] \times 1 = \frac{y_k^T s_k}{s_k^T B_k s_k}$$

We conclude that

$$\det(B_{k+1}) = \det(B_k)\frac{y_k^T s_k}{s_k^T B_k s_k}$$

□

**Exercise 12.** *Show that if $f$ satisfies Assumption 6.1 and if the sequence of gradients satisfies* $\liminf \|\nabla f_k\| = 0$, , *then the whole sequence of iterates $x$ converges to the solution $x^*$.*

*Proof.* Since $f(x_k)$ deceases at each step and by Assumption 6.1(ii) the convexity of the set $\mathcal{L} = \{x | f(x) \le f(x_0)\}$ , the fact $\liminf \|\nabla f_k\| = 0$ implies there exists a subsequence $\{x_{n_j}\}$ converges to the unique minimizer $x^*$. We are now proving the whole sequence $\{x_k\}$ converges to $x^*$. By Taylor's thm, for all $x \in \mathbb{R}^n$ we have

$$f(x) = f(x^* + (x - x^*)) = f(x^*) + \nabla f(x^*)^T(x - x^*) + \frac{1}{2}(x - x^*)^T \nabla^2 f(\xi)(x - x^*)$$

If $x$ belongs to $\mathcal{L} = \{x | f(x) \le f(x_0)\}$ and satisfies

$$f(x) \le f(x^*) + \varepsilon$$

for some given $\varepsilon > 0$, we obtain following by using the fact $\nabla f(x^*) = 0$

$$\frac{1}{2}(x - x^*)^T \nabla^2 f(\xi)(x - x^*) \le \varepsilon.$$

By Assumption 6.1(ii) again, we conclude that

$$m\|x - x^*\|_2^2 \le (x - x^*)^T \nabla^2 f(\xi)(x - x^*) \le 2\varepsilon.$$

So,

$$\|x - x^*\|_2^2 \le (2\varepsilon/m)$$

On the other hand, the whole sequence $\{f(x_k)\}$ is nonincreasing by any descent direction Algorithm, and we already know that there exists a subsequence $\{f(x_{n_j})\}$ converges to the $f(x^*)$. So given $\varepsilon > 0$, we can find a $N \in \mathbb{N}$ such that

$$f(x_k) \le f(x_{n_j}) \le f(x*) + \varepsilon$$

for all $k \ge n_j \ge N$. Hence, combining the two inequality gives

$$\|x_k - x^*\|_2^2 \le (2\varepsilon/m)$$

for for all $k \ge N$. So the whole sequence $\{x_k\}$ converges to $x^*$. □

# Numerical Optimization with applications: Homework 06

104021601 林俊傑
104021602 吳彥儒
104021615 黃翊軒

**Exercise 1.** *The following example from [268] with a single variable $x \in \mathbb{R}$ and a single equality constraint shows that strict local solutions are not necessarily isolated. Consider*

$$\min_x x^2 \quad \text{subject to } c(x) = 0, \text{ where } c(x) = \begin{cases} x^6 \sin(1/x) & \text{if } x \neq 0 \\ 0 & \text{if } x = 0 \end{cases} \quad (12.96)$$

*(a) Show that the constraint function is twice continuously differentiable at all $x$ (including at $x = 0$) and that the feasible points are $x = 0$ and $x = 1/(k\pi)$ for all nonzero integers $k$.*

*(b) Verify that each feasible point except $x = 0$ is an isolated local solution by showing that there is a neighborhood $\mathcal{N}$ around each such point within which it is the only feasible point.*

*Proof.* (a) We first show directly that constraint function is twice continuously differentiable at all $x$.

If $x \neq 0$, then

$$c(x) = x^6 \sin(1/x)$$
$$c'(x) = 6x^5 \sin(1/x) - x^4 \cos(1/x)$$
$$c''(x) = (30x^4 - x^2)\sin(1/x) - 10x^3 \cos(1/x)$$

If $x = 0$, then by definition we obtain

$$c'(0) = \lim_{h \to \infty} \frac{h^6 \sin(1/h) - 0}{h} = 0$$
$$c''(0) = \lim_{h \to \infty} \frac{[6h^5 \sin(1/h) - h^4 \cos(1/h)] - 0}{h} = 0$$

Hence, the constraint function is twice continuously differentiable at all $x$.

Second, we show that the feasible points are $x = 0$ and $x = 1/(k\pi)$ for all nonzero integers $k$.

If $x = 0$, then $c(x) = 0$ by definition. If $x = 1/(k\pi)$, then $\sin(1/x) = 0$ and thus we have $c(x) = 0$.

(b) If $x = 1/(k\pi)$ for some fixed nonzero integers $k$ and we choose $r = \left| \frac{1}{k\pi} - \frac{1}{(k+1\pi)} \right|$, then the open interval $\mathcal{N} = (x - r, x + r)$ contains only a feasible point, which is $x$ itself.

On the other hand, if $x = 0$, then for all $r > 0$ there exists nonzero positive integers $k$ such that $x_k = 1/(k\pi) < r$. However, $x_k$ are feasible points in the neighborhood $(-r, r)$.

Therefore, combining disscusion above, we have that each feasible point except $x = 0$ is an isolated local solution by showing that there is a neighborhood $\mathcal{N}$ around each such point within which it is the only feasible point. $\qquad \square$

**Exercise 15.** *Consider the following modification of (12.36), where $t$ is a parameter to be fixed prior to solving the problem:*

$$\min_x (x_1 - \frac{3}{2})^2 + (x_2 - t)^4 \quad s.t. \quad \begin{bmatrix} 1 - x_1 - x_2 \\ 1 - x_1 + x_2 \\ 1 + x_1 - x_2 \\ 1 + x_1 + x_2 \end{bmatrix} \geq 0$$

*(a)For what value of $t$ does the point $x^* = (1,0)^T$ satisfy the KKT conditions?*
*(b)Show that when $t = 1$, only the first constraint is active at the solution, and find the solution.*
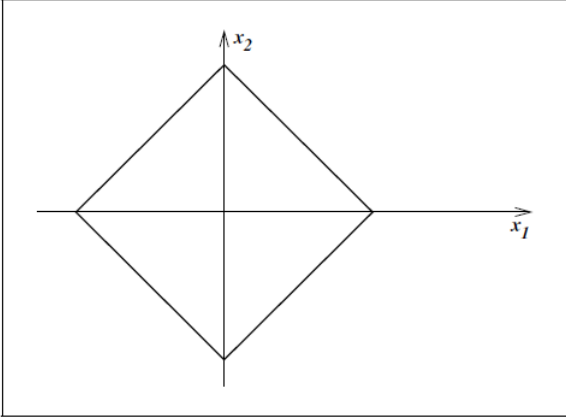
*Proof.* (a) First, we check the complementary condition of KKT, i.e. $\lambda_i c_i(x^*) = 0$ for $i = 1, 2, 3, 4$.

$$\begin{cases} \lambda_1(1 - 1 - 0) = 0 \\ \lambda_2(1 - 1 + 0) = 0 \\ \lambda_3(1 + 1 - 0) = 0 \\ \lambda_4(1 + 1 + 0) = 0 \end{cases} \Rightarrow \begin{cases} \lambda_3 = 0 \\ \lambda_4 = 0 \end{cases} \tag{1}$$

Obviously, $c(x^*) \geq 0$ holds. Consider

$$\nabla_x L(x^*, \lambda) = \nabla_x L((0, 1)^T, \lambda)$$
$$= \begin{bmatrix} -1 + \lambda_1 + \lambda_2 \\ -4t^3 + \lambda_1 - \lambda_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Since $\lambda_1, \lambda_2 \geq 0$ and $\lambda_1 + \lambda_2 = 1$, we have $\lambda_1 - \lambda_2 \in [-1, 1]$. Hence, $4t^3 \in [-1, 1]$, then $t \in [-\sqrt[3]{4}, \sqrt[3]{4}]$.



(b) We know the feasible set $E = \{x \in \mathbb{R}^2 \mid \|x\|_1 = 1\}$. Compute the gradient of $f(x)$

$$\nabla f(x) = \begin{bmatrix} 2(x_1 - \frac{3}{2}) \\ 4(x_2 - 1)^3 \end{bmatrix} \leq 0 \quad \forall x \in E$$

From above, $\forall x, y \in E$, if $x_1 \leq y_1$ and $x_2 \leq y_2$ then $f(y) \leq f(x)$. Therefore, we only have to consider the case that the first constraint is active: $1 - x_1 - x_2 = 0$. Substituting $x_2 = 1 - x_1$ into $f(x)$ and find the minimum of $f$:

$$f(x) = (x_1 - \frac{3}{2})^2 + (-x_1)^4$$
$$f'(x) = 2(x_1 - \frac{3}{2}) + 4x_1^3$$
$$= 4x_1^3 + 2x_1 - 3$$

$f'(x^*) = 0$ if $x_1^* = \frac{\sqrt[3]{27+\sqrt{753}}}{2 \times 3^{2/3}} - \frac{1}{\sqrt[3]{3(27+\sqrt{753})}}$. Consequently, the minimizer of $f$ is $(x_1^*, 1 - x_1^*)$.

$\square$

2

**Exercise 19.** *Consider the problem*

$$\min_{x\in R^2} = -2x_1 + x_2 \quad subject\ to \quad \begin{cases} (1-x_1)^3 - x_2, & \geq 0 \\ x_2 + 0.25x_1^2 - 1, & \geq 0 \end{cases}$$

*The optimal solution is $x^* = (0,1)^T$, where both constraints are active.*

*(a) Do the LICQ hold at this point?*

*(b) Are the KKT conditions satisfied?*

*(c) Write down the set $\mathcal{F}(x^*)$ and $\mathcal{C}(x^*, \lambda^*)$.*

*(d) Are the second-order necessary conditions satisfied? Are the second-order sufficient conditions satisfied?*

*Proof.* (a)

$$A(x^*) = [\nabla C_i(x^*)]_{i\in\mathcal{A}(x^*)} = \begin{bmatrix} -3(1-x_1)^2 & 0.5x_1 \\ -1 & 1 \end{bmatrix}_{x=x*} = \begin{bmatrix} -3 & 0 \\ -1 & 1 \end{bmatrix}$$

$A(x^*)$ is nonsingular. Therefore, the LICQ holds.

(b)

$$\nabla f(x^*) = \begin{bmatrix} -2 \\ 1 \end{bmatrix}, \quad \nabla c_1(x^*) = \begin{bmatrix} -3 \\ -1 \end{bmatrix}, \quad \nabla c_2(x^*) = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

Therefore, the KKT conditions (12.34a)-(12.34e) are satisfied when we set

$$\lambda^* = \left(\frac{2}{3}, \frac{5}{3}\right)^T$$

(c)

$$\mathcal{F}(x^*) = \{d|\nabla c_i(x^*)^T d \geq 0\} = \{(d_1, d_2)^T| \begin{bmatrix} -3 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \end{bmatrix} \geq 0\} = \{(d_1, d_2)^T|d_2 \geq 0,\ 3d_1 + d_2 \leq 0\}$$

$$\begin{aligned} \mathcal{C}(x^*, \lambda^*) &= \{w \in \mathcal{F}(x^*)| \nabla c_i(x^*)^T w = 0\ \forall i \in \mathcal{A}(x^*) \cap I \text{ with } \lambda_i^* > 0\} \\ &= \{w \in \mathcal{F}(x^*)| \nabla c_i(x^*)^T w = 0 \text{ for } i = 1, 2\} \\ &= \{(w_1, w_2)^T \in \mathcal{F}(x^*)| \begin{bmatrix} -3 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = 0\} \\ &= \{(0,0)^T\} \end{aligned}$$

(d)

$$\forall w \in \mathcal{C}(x^*, \lambda^*) = \{(0,0)^T\} \qquad w^T \nabla_{xx} \mathcal{L}(x^*, \lambda^*) w = 0 \tag{2}$$

Since at $x^*$ LICQ holds and $(x^*, \lambda^*)$ satisfies KKT. By (2), the second-order necessary conditions is satisfied.

Since $x^*$ is a feasible solution and $(x^*, \lambda^*)$ satisfies KKT. By (2), the second-order sufficient conditions is satisfied.

<div style="text-align: right;">□</div>

# Numerical Optimization with applications: Homework 07

104021601 林俊傑

104021602 吳彥儒

104021615 黃翊軒

December 14, 2016

**Exercise 2** (Chapter 7). *Show that the matrix $\widehat{H}_{k+1} = (I - \frac{s_k y_k^T}{y_k^T s_k})$ is singular.*

*Proof.* Consider $\widehat{H}_{k+1} s_k$, then we have

$$
\begin{aligned}
\widehat{H}_{k+1} s_k &= (I - \frac{s_k y_k^T}{y_k^T s_k}) s_k \\
&= s_k - \frac{s_k (y_k^T s_k)}{y_k^T s_k} \\
&= s_k - s_k \\
&= 0
\end{aligned}
$$

Since $s_k = x_{k+1} - x_k \neq 0$, thus $\widehat{H}_{k+1}$ is singular. $\qquad \square$

**Exercise 5** (Chapter 10). *Suppose that each residual function $r_j$ and its gradient are Lipschitz continuous with Lipschitz constant $L$, that is ,*

$$
\|r_j(x) - r_j(\widehat{x})\| \leq L\|x - \widehat{x}\|, \quad \|\bigtriangledown r_j(x) - \bigtriangledown r_j(\widehat{x})\| \leq L\|x - \widehat{x}\|
$$

*for all $j = 1, 2, ..., m$ and all $x, \widehat{x} \in \mathcal{D}$, where $\mathcal{D}$ is a compact subset of $\mathbb{R}^n$. Assume also that the $r_j$ are bounded on $\mathcal{D}$, that is there exist $M > 0$ such that $|r_j(x)| \leq M$ for all $j = 1, 2, ..., m$ and all $x \in \mathcal{D}$. Find Lipschitz constant for the Jacobian $J$ and the gradient $\bigtriangledown f$ over $\mathcal{D}$.*

$$
J(x) = \left[\frac{\partial r_j}{\partial x_i}\right]_{\substack{j=1,2,...,m \\ i=1,2,...,n}} = \begin{bmatrix} \bigtriangledown r_1(x)^T \\ \bigtriangledown r_2(x)^T \\ \vdots \\ \bigtriangledown r_m(x)^T \end{bmatrix}
$$

$$
\bigtriangledown f(x) = \Sigma_{j=1}^m r_j(x) \bigtriangledown r_j(x) = J(x)^T r(x)
$$

*Proof.* Since all norms in $\mathbb{R}^n$ are equivalent.

$$
\exists \alpha > 0 \quad \text{such that} \quad \|x\| \leq \alpha \|x\|_\infty \quad \forall x \in \mathbb{R}^n
$$

We have,

$$
\begin{aligned}
\|J(x_1) - J(x_2)\| &= \max_{\|y\|=1} \|(J(x_1) - J(x_2))y\| \\
&= \max_{\|y\|=1} \left\| \begin{bmatrix} (\bigtriangledown r_1(x_1) - \bigtriangledown r_1(x_2))^T y \\ \vdots \\ (\bigtriangledown r_m(x_1) - \bigtriangledown r_m(x_2))^T y \end{bmatrix} \right\| \\
&\leq \max_{\|y\|=1} \alpha \left\| \begin{bmatrix} (\bigtriangledown r_1(x_1) - \bigtriangledown r_1(x_2))^T y \\ \vdots \\ (\bigtriangledown r_m(x_1) - \bigtriangledown r_m(x_2))^T y \end{bmatrix} \right\|_\infty \\
&= \alpha \max_{\|y\|=1} \max_{1 \leq j \leq m} |(\bigtriangledown r_j(x_1) - \bigtriangledown r_j(x_2))^T y| \\
&\leq \alpha \max_{\|y\|=1} \max_{1 \leq j \leq m} |(\bigtriangledown r_j(x_1) - \bigtriangledown r_j(x_2))| \; |y| \\
&\leq \alpha \max_{\|y\|=1} \max_{1 \leq j \leq m} L\|x_1 - x_2\| \; |y| \\
&= \alpha L \|x_1 - x_2\|
\end{aligned}
$$

We conclude that $J$ is Lipschitz continuous with constant $\tilde{L} = \alpha L$.

On the other hand, Given $x, \tilde{x}$ in $\mathcal{D}$, we estimate

$$
\begin{aligned}
\|\nabla f(x) - \nabla f(\tilde{x})\| &= \|J(x)^T r(x) - J(\tilde{x})^T r(\tilde{x})\| \\
&= \| \left[ J(x)^T r(x) - J(\tilde{x})^T r(x) \right] + \left[ J(\tilde{x})^T r(x) - J(\tilde{x})^T r(\tilde{x}) \right] \| \\
&= \| \left( J(x)^T - J(\tilde{x})^T \right) r(x) + J(\tilde{x})^T \left( r(x) - r(\tilde{x}) \right) \| \\
&\leq \|J(x)^T - J(\tilde{x})^T\| \|r(x)\| + \|J(\tilde{x})^T\| \|r(x) - r(\tilde{x})\| \\
&\leq M\alpha L \|x - \tilde{x}\| + M' L \|x - \tilde{x}\| \\
&= \mathcal{L} \|x - \tilde{x}\|
\end{aligned}
$$

where $\mathcal{L} = M\alpha L + M'L$ and $\|J(\tilde{x})^T\|$ is bounded since it is Lipschitz continuous on a compact set $\mathcal{D}$. $\qquad\square$

# Numerical Optimization with applications: Homework 07

104021601 林俊傑

104021602 吳彥儒

104021615 黃翊軒

December 14, 2016

*Proof.* Since all norms in $\mathbb{R}^n$ are equivalent.

$$\exists \alpha > 0 \quad \text{such that} \quad \|x\| \leq \alpha \|x\|_\infty \quad \forall x \in \mathbb{R}^n$$

We have,

$$
\begin{aligned}
\|J(x_1) - J(x_2)\| &= \max_{\|y\|=1} \|(J(x_1) - J(x_2))y\| \\
&= \max_{\|y\|=1} \left\| \begin{bmatrix} (\bigtriangledown r_1(x_1) - \bigtriangledown r_1(x_2))^T y \\ \vdots \\ (\bigtriangledown r_m(x_1) - \bigtriangledown r_m(x_2))^T y \end{bmatrix} \right\| \\
&\leq \max_{\|y\|=1} \alpha \left\| \begin{bmatrix} (\bigtriangledown r_1(x_1) - \bigtriangledown r_1(x_2))^T y \\ \vdots \\ (\bigtriangledown r_m(x_1) - \bigtriangledown r_m(x_2))^T y \end{bmatrix} \right\|_\infty \\
&= \alpha \max_{\|y\|=1} \max_{1 \leq j \leq m} |(\bigtriangledown r_j(x_1) - \bigtriangledown r_j(x_2))^T y| \\
&\leq \alpha \max_{\|y\|=1} \max_{1 \leq j \leq m} |(\bigtriangledown r_j(x_1) - \bigtriangledown r_j(x_2))| \, |y| \\
&\leq \alpha \max_{\|y\|=1} \max_{1 \leq j \leq m} L \|x_1 - x_2\| \, |y| \\
&= \alpha L \|x_1 - x_2\|
\end{aligned}
$$

We conclude that $J$ is Lipschitz continuous with constant $\tilde{L} = \alpha L$. $\square$

*Proof.* Given $x, \tilde{x}$ in $\mathcal{D}$, we estimate

$$
\begin{aligned}
\|\nabla f(x) - \nabla f(\tilde{x})\| &= \|J(x)^T r(x) - J(\tilde{x})^T r(\tilde{x})\| \\
&= \| [J(x)^T r(x) - J(\tilde{x})^T r(x)] + [J(\tilde{x})^T r(x) - J(\tilde{x})^T r(\tilde{x})] \| \\
&= \| (J(x)^T - J(\tilde{x})^T) r(x) + J(\tilde{x})^T (r(x) - r(\tilde{x})) \| \\
&\leq \|J(x)^T - J(\tilde{x})^T\| |r(x)| + \|J(\tilde{x})^T\| \|r(x) - r(\tilde{x})\| \\
&\leq M \alpha L \|x - \tilde{x}\| + M' L \|x - \tilde{x}\| \\
&= \mathcal{L} \|x - \tilde{x}\|
\end{aligned}
$$

where $\mathcal{L} = M \alpha L + M' L$ and $\|J(\tilde{x})^T\|$ is bounded since it is Lipschitz continuous on a compact set $\mathcal{D}$. $\square$

# Numerical Optimization with applications
## CHAPTER 17: Penalty and Augmented Lagrangian Methods

104021601 林俊傑

104021602 吳彥儒

104021615 黃翊軒

January 15, 2017

## 17.1 THE QUADRATIC PENALTY METHOD

Given the original constrained optimization problem

$$\min_{x} f(x) \quad \text{subject to } c_i(x) = 0, \quad i \in \mathcal{E}, \quad c_i(x) \geq 0, \quad i \in \mathcal{I}, \tag{1}$$

we could define the corresponding quadratic penalty function as

$$Q(x;\mu) := f(x) + \frac{\mu}{2} \sum_{i \in \mathcal{E}} c_i^2(x) + \frac{\mu}{2} \sum_{i \in \mathcal{I}} \left( [c_i(x)]^- \right)^2, \tag{2}$$

where $\mu > 0$ and $[y]^-$ is the abbreviated symbol of $\max(-y, 0)$.

If we take a sequence $\mu_k \nearrow \infty$ into the quadratic penalty function, we could find that $Q(x;\mu_k)$ deverges if $x$ is infeasible. The larger $\mu_k$ is, the severer constraint violations we panalize. As a result, the minimizer of the quadratic penalty function $Q(x;\mu_k)$ is closer to the feasible region as $k$ increases.

**Framework 17.1** (Quadratic Penalty Method).

Given $\mu_0 > 0$, a nonnegative sequence $\{\tau_k\}$ with $\tau_k \to 0$, and a starting point $x_0^s$;

**for** $k = 0, 1, 2, \ldots$

    Find an approximate minimizer $x_k$ of $Q(\cdot; \mu_k)$, starting at $x_k^s$,

        and terminating when $\|\nabla_x Q(x; \mu_k)\| \leq \tau_k$;

    **if** final convergence test satisfied

        **stop** with approximate solution $x_k$;

    **end (if)**

    Choose new penalty parameter $\mu_{k+1} > \mu_k$;

    Choose new starting point $x_{k+1}^s$;

**end (for)**

We have two theorems to support the convergence of Framework 17.1.

**Theorm 17.1** states that the global minimizer $x_k$ of quadratic penalty function $Q(x;\mu_k)$ converges to the constrained optimization problem $x$, i.e. $x_k \to x$.

**Theorm 17.2** states that if $\tau_k \to 0$ and $x_k$ only satisfies

$$\|\nabla_x Q(x;\mu_k)\| \leq \tau_k,$$

then

$$x_k \to x^*,$$

where $x^*$ is a stationary point of $\|c(x)\|^2$. Besides, if $\nabla c_i(x^*)$ is linearly independent, then

$$\lim_{k \to \infty} -\mu_k c_i(x_k) = \lambda_i^* \quad \forall i \in \mathcal{E}$$

and $(X^*, \lambda^*)$ satisfy the KKT conditions.

### Practical problems

Even if $\nabla^2 f(x^*)$ is well-conditioned, the Hessian $\nabla_{xx}^2 Q(x;\mu_k)$ might become ill-conditioned as $\mu_k \to \infty$.

By defining

$$A(x)^T = (\nabla c_i(x))_{i \in \mathcal{E}}$$

and considering equality constraints only, we have

$$\nabla^2_{xx} Q(x; \mu_k) = \nabla^2 f(x) + \sum_{i \in \mathcal{E}} \mu_k c_i(x) \nabla^2 c_i(x) + \mu_k A(X)^T A(X).$$

From **Theorem 17.2**, we have

$$\mu_k c_i(x) \approx -\lambda_i^*$$

for $x$ near a minimizer. Hence, we obtain

$$\nabla^2_{xx} Q(x; \mu_k) \approx \nabla^2_{xx} \mathcal{L}(x, \lambda^*) + \mu_k A(X)^T A(X).$$

We find that $\nabla^2_{xx} Q(x; \mu_k)$ have problems with ill-conditioning since the second term diverges as $\mu_k \to \infty$.

For Newton's method step

$$\nabla^2_{xx} Q(x; \mu_k) p = \nabla_x Q(x; \mu),$$

we can apply a reformulation

$$\begin{pmatrix} \nabla^2 f(x) + \sum_{i \in \mathcal{E}} \mu_k c_i(x) \nabla^2 c_i(x) & A(x)^T \\ A(x) & -(1/\mu_k) I \end{pmatrix} \begin{pmatrix} p \\ \mu A(x) p \end{pmatrix} = \begin{pmatrix} -\nabla_x Q(x; \mu_k) \\ 0 \end{pmatrix}$$

to avoid the ill-conditioning since $p$ solves both systems. Note that this system has dimension $n + |\mathcal{E}|$ rather than n.

## 17.2 NONSMOOTH PENALTY FUNCTIONS

A penalty function ia called *exact* if, for certains coice of penalty parameters, the minimizer $x^\star$ is the exact solution of the original constrained optimization problem. Nevertheless, the quadratical penalty function is not exact. In this section, we introduce the *nonsmooth* penalty functions.
A popular nonsmooth penalty function is the $l_1$ *penalty function* defined by

$$\phi_1(x; \mu) = f(x) + \mu \sum_{i \in \mathcal{E}} |c_i(x)| + \mu \sum_{i \in \mathcal{I}} [c_i(x)]^-. \tag{3}$$

The next two theorems establish the *exactness* of (3).
**Theorm 17.3** states that if $x^\star$ is a strictly local minimizer of (1), with Lagrange miltipliers $\lambda^\star$. Then $x^\star$ is a local minimizer of (3) $\quad \forall \mu > \mu^\star$, where

$$\mu^\star = \|\lambda^\star\|_\infty \tag{4}$$

**Theorm 17.4** states that if $\hat{x}$ is a stationary points of $\phi_1(x; \mu)$ for all $\mu$ large enouth. Then, $\hat{x}$ is either satisfying KKT conditions for (1) or it is an infeasible stationary points.
Define the measure of infeasibility

$$h(x) = sum_{i \in \mathcal{E}} |c_i(x)| + sum_{i \in \mathcal{I}} [c_i(x)]^- \tag{5}$$

Then, we can develope an algorithm framwork via the $l_1$ penalty funtion.
**Framework 17.2** (Classical $\ell_1$ Penalty Method).
   Given $\mu_0 > 0$, tolerance $\tau > 0$, starting point $x_0^s$;
   **for** $k = 0, 1, 2, \ldots$
        Find an approximate minimizer $x_k$ of $\phi_1(x; \mu_k)$, starting at $x_k^s$;
        **if** $h(x_k) \leq \tau$
            **stop** with approximate solution $x_k$;
        **end (if)**
        Choose new penalty parameter $\mu_{k+1} > \mu_k$;
        Choose new starting point $x_{k+1}^s$;
   **end (for)**
Since $\phi_1(x; \mu)$ is nonsmooth, the minimization will be difficule. However, we can transform $\phi_1(x; \mu)$ into a smooth model.

## A PRATICAL $l_1$ PENALTY METHOD

As we did for the unconstrained optimization problem, we can transform (3) into a smooth model by replacing $f$ by its Taylor expension and $c_i$ by its linearization,as follows:

$$q(p;\mu) = f(x) + \bigtriangledown f(x)^T p + \frac{1}{2} p^T W p + \mu \sum_{i\in\mathcal{E}} |c_i(x) + \bigtriangledown c_i(x)^T p| + \mu \sum_{i\in\mathcal{I}} [c_i(x) + \bigtriangledown c_i(x)^T p]^-$$

where $W$ is an approximation of Hessian about $f$ and $c_i$. The function $q(p;\mu)$ is still not smooth, but we can reformulate it into a smooth quadratic optimization problem by introducing some new variables, as follows:

$$
\begin{aligned}
\min_{p,r,s,t} \quad & f(x) + \frac{1}{2} p^T W p + \bigtriangledown f(x)^T p + \mu \sum_{i\in\mathcal{E}} |r_i + s_i| + \mu \sum_{i\in\mathcal{I}} t_i \\
subject\ to \quad & \bigtriangledown c_i(x)^T p + c_i(x) = r_i - s_i, \quad i \in \mathcal{E} \\
& \bigtriangledown c_i(x)^T p + c_i(x) \geq -t_i, \quad i \in \mathcal{I} \\
& r,s,t \geq 0
\end{aligned}
\tag{6}
$$

Even after adding a trust region constraint $\|p\|_\infty \leq \triangle$, (6) is still a quadratic problem. It can be solved by a quadratic programming solver.

## A GENERAL CLASS OF NONSMOOTH PENALTY METHODS

Exact nonsmooth penalty funtions can use other norms.

$$\phi(x;\mu) = f(x) + \mu \|c_\mathcal{E}(x)\| + \mu \|[c_\mathcal{I}(x)]^-\| \tag{7}$$

Framework 17.2 can work on these penalty functions by simply redefinind the measure of infeasibility (5) as $h(x) = \|c_\mathcal{E}(x)\| + \|[c_\mathcal{I}(x)]^-\|$.

The properties garguaranteed by Theorem 17.3 and Theorem 17.4 can be extended to the general class (7).In Theorem 17.3, we replace $\mu^\star$ in (4) by

$$\mu^\star = \|\lambda^\star\|_D,$$

where $\|\bullet\|_D$ is the dual norm of $\|\bullet\|$. Theorem 17.4 applies without modification.

# 17.3 AUGMENTED LAGRANGIAN METHOD: EQUALITY CONSTRAINTS

In section 17.1, we know that even $\mu_k$ is large, the approximate minimizer $x_k$ of the quadratic penalty function $Q(x;\mu_k)$ may be infeasible, the violation of $c_i(x) \approx -\lambda_i^*/\mu_k$. To make the approximate solution $x_k$ closer to the feasible region, we introduce the Augmented Lagrangian function:

$$\mathcal{L}_A(x,\lambda;\mu) := f(x) - \sum_{i\in\mathcal{E}} \lambda_i c_i(x) + \frac{\mu}{2} \sum_{i\in\mathcal{E}} c_i^2(x) \tag{8}$$

Use the fact of Theorem 2.2 and (17.17), and rearranging the expression, we have $c_i(x_k) \approx -\frac{1}{\mu_k}(\lambda_i^* - \lambda_i^k)$, the violent of $x_k$ is much smaller than $\frac{1}{\mu_k}$. We can set the Lagrangian multiplier vector of the next step $\lambda_i^{k+1} = \lambda_i^k - \mu_k c_i(x_k)$, for all $i \in \mathcal{E}$.

**Framework 17.3** (Augmented Lagrangian Method-Equality Constraints).

Given $\mu_0 > 0$, tolerance $\tau_0 > 0$, starting points $x_0^s$ and $\lambda^0$;

**for** $k = 0, 1, 2, \ldots$

        Find an approximate minimizer $x_k$ of $\mathcal{L}_A(\cdot, \lambda^k; \mu_k)$, starting at $x_k^s$,

            and terminating when $\|\nabla_x \mathcal{L}_A(x_k, \lambda^k; \mu_k)\| \le \tau_k$;

        **if** a convergence test for (17.1) is satisfied

            **stop** with approximate solution $x_k$;

        **end (if)**

        Update Lagrange multipliers using (17.39) to obtain $\lambda^{k+1}$;

        Choose new penalty parameter $\mu_{k+1} \ge \mu_k$;

        Set starting point for the next iteration to $x_{k+1}^s = x_k$;

        Select tolerance $\tau_{k+1}$;

**end (for)**

**Theorm 17.5** states that if we know the exact Lagrangian multiplier $\lambda^*$, then the solution of (1) is a strict minimizer of $\mathcal{L}_A(x, \lambda; \mu)$ for $\mu$ large enough. Even though we only have a "good" estimate of $\lambda^*$, we can still get a good estimate of $x^*$ by minimizing $\mathcal{L}_A(x, \lambda; \mu)$ with large $\mu$.

**Theorem 17.6** states the advantage of the augmented Lagrangian method. Different from the quadratic penalty method, we can get a good approximation of $x^*$ if $\lambda_k$ is close to $\lambda^*$ or if the penalty parameter $\mu_k$ is large. On the other hand, by (17.46), we can improve the accuracy of $\lambda^*$ by choosing a large $\mu_k$.

# 17.4 PRACTICAL AUGMENTED LAGRANGIAN METHOD

In section 17.3, we only discuss the problem with equality constrains. Now for the general case, there are three useful formulations.

**Bound-Constrained Formulation**

Use the slack variable $s_i$ to turn inequalities into equalities. That is

$$c_i(x) - s_i = 0, \quad s_i \ge 0, \quad \forall i \in \mathcal{I}$$

We can reformulate the problem into

$$\min_{x \in \mathbb{R}^n} f(x) \quad s.t. \quad c_i(x) = 0, \quad i = 1, 2, \ldots, m, \quad l \le x \le u$$

The Bounded-constrained Lagrangian will be:

$$\min_x \mathcal{L}_A(x, \lambda; \mu) = f(x) - \sum_{i=1}^m \lambda_i c_i(x) + \frac{\mu}{2} \sum_{i=1}^m c_i^2(x) \quad s.t. \quad l \le x \le u$$

Solve this problem and update $\lambda$ and $\mu$ repeatedly.

**Linearly Constrained Formulation**

LCL method is to solve the subproblem of minimizing the augmented Lagrangian function subject to linearization of the constrains.

$$\min_x \quad F_k(x)$$

$$s.t. \quad c(x_k) + A_k(x - x_k) = 0, \quad l \le x \le u.$$

where

$$c_i^{-k}(x) = c_i(x) - c_i(x_k) - \nabla c_i(x_k)^T(x - x_k).$$

Current Augmented Lagrangian funciton

$$F_k(x) = f(x) - \sum_{i=1}^m \lambda_i^k c_i^{-k}(x) + \frac{\mu}{2} \sum_{i=1}^m [c_i^{-k}(x)]^2$$

4

**Unconstrained Formulation**

Suppose the problem has no equality constrain, i.e. $\mathcal{E} = \emptyset$, then we can rewrite the problem as

$$\min_{x \text{ feasible}} f(x) = \min_{x \in \mathbb{R}^n} F(x)$$

where

$$F(x) = \max_{\lambda \geq 0}\{f(x) - \sum_{i \in \mathcal{I}} \lambda_i c_i(x)\}$$

Note that if x is feasible, $F(x) = f(x)$ and $\lambda_i$ should be zero. Otherwise $F(x)$ turns to infinity, and $\lambda_i$ can be chosen arbitrary large. Consequently, $F$ is not smooth, so it is not practical to minimize directly. We replace F by a smooth approximated function

$$\widehat{F}(x; \lambda^k, \mu_k) = \max_{\lambda \geq 0}\{f(x) - \sum_{i \in \mathcal{I}} \lambda_i c_i(x) - \frac{1}{2\mu_k} \sum_{i \in \mathcal{I}} (\lambda_i - \lambda_i^k)^2\}$$

where the last term can enforce the mew maximizer $\lambda$ close to the previous estimate $\lambda^k$.
By above, we can obtain the explicit maximization of $\lambda$. Then we have

$$\widehat{F}(x; \lambda^k, \mu_k) = f(x) + \sum_{i \in \mathcal{I}} \psi(c_i(x), \lambda_i^k; \mu_k)$$

where the function$\psi$ is defined as

$$\psi(t, \sigma; \mu) := \begin{cases} -\sigma t + \frac{\mu}{2} t^2 & \text{if} t - \sigma/\mu \leq 0, \\ -\frac{1}{2\mu} \sigma^2 & \text{otherwise,} \end{cases}$$

Hence, we can obtain $x_k$ by minimizing $\widehat{F}$, and update Lagrange multiplier estimates repeatedly.