

AMS 553.761: Nonlinear Optimization I
Midterm, Fall 2020

- There are 4 questions on this test.
- You have to upload your answers on Blackboard by 4:30 pm on Tuesday, October 13, 2020. If something goes wrong with the upload you may email it to me before 4:30pm on Tuesday, October 13, 2020. Please use the email option only if the upload does not work. No answers will be accepted after this deadline.

Please hand in ONE submission - multiple submissions will not be tolerated.

- You are not allowed to discuss any problem with any other human being, except the instructor.
- You can use a computer only as a word processor; in particular, you cannot consult the internet in regards to this midterm. You CAN use the slides from class and books from the library.
- You CAN cite any result we have mentioned in class or from the HWs without proof. If you cite a result (e.g., from a book) that was NOT mentioned in class, you should include a complete proof of this fact.
- The level of rigor expected is the same as the HW solutions. Make sure you justify all your answers.

1. **(25 pts)** Give complete proofs of Lemmas 2.2, 2.3 and 2.6 from the lecture notes on “Conjugate Gradient”.
2. **(25 pts)** Prove that the bisection method for Wolfe linesearch, i.e., Algorithm 13 from slide 76 in the “Line Search Methods” slides on the course webpage, terminates with a steplength satisfying the weak Wolfe conditions. Assume that the function f is bounded from below and has a gradient that is Lipschitz continuous with a global Lipschitz constant γ .
3. **(20 pts)** Prove Theorem 4.8 from slide 103 in the “Line Search Methods” slides on the course webpage, that establishes a global convergence rate of $O\left(\left(\frac{1}{\epsilon}\right)^2\right)$ for a modified or quasi Newton method with *Wolfe linesearch*. You may use Theorem 3.4 on slide 71 (due to Zoutendijk) without proof. [Recall that the condition number of B_k is given by $\frac{\lambda_{\max}(B_k)}{\lambda_{\min}(B_k)}$ which is the ratio of the largest eigenvalue of B_k to the smallest eigenvalue of B_k .]
4. **(10 pts)** Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function, not necessarily differentiable. Show that any local minimum is also a global minimum.

1. (25 pts) Give complete proofs of Lemmas 2.2, 2.3 and 2.6 from the lecture notes on "Conjugate Gradient".

Lemmas 2.2 9/10

Consider L is span of a set of orthogonal basis $\{b_1, b_2, \dots, b_k\}$ and $f(\bar{x})$ is the minimizer of $f(x)$ restricted to $x_0 + L$, while $\nabla f(\bar{x}) = g(\bar{x})$ is not orthogonal to L w.r.t standard inner product. Thus

$$\langle g(\bar{x}), \sum_{i=1}^k \lambda_i b_i \rangle \neq 0$$

Then we can compute the projection of $g(\bar{x})$ into L :

$$\text{projection of } g(\bar{x}) = s' = \sum_{i=1}^k \frac{\langle g(\bar{x}), b_i \rangle}{\langle b_i, b_i \rangle} b_i$$

Thus, we can compute direction derivate of $f(\bar{x})$ w.r.t direction s' . $f(p) = \frac{1}{2} p^T A p - b^T p$ is quadratic function, which is differential continuously. Thus

$$f'(\bar{x}; s') = \nabla f(\bar{x}) \cdot s' = g(\bar{x})^T \cdot s' \neq 0$$

if $g(\bar{x})^T \cdot s' > 0$, then we set $s = -s'$ and consider $\phi(\alpha) = f(\bar{x} + \alpha s)$

$$\phi'(\alpha) = g(\bar{x})^T s = -g(\bar{x})^T \cdot s' < 0$$

since f is continuous differentiable, so is ϕ , thus there exists $0 < \delta < \varepsilon$ such that $\phi'(c) < 0$ for all $c \in (\bar{x}, \bar{x} + \delta)$

By the Mean Value theorem:

$$\phi(s) = \phi(0) + \phi'(t)(s-0)$$

for some $t \in (0, s)$ since $\phi'(t) < 0$ we have

$$f(\bar{x} + ss) = \phi(s) < \phi(0) = f(\bar{x})$$

If $g(\bar{x})^T s' < 0$, then we can set $s = s'$, by doing some proof, we can show that there must exist a point that make $f(x) < f(\bar{x})$, which is a contradiction to our hypothesis that $f(\bar{x})$ is the minimizer of $f(p)$ restricted to $x_0 + L$. Also. [if $\bar{x} \in x_0 + L$ and $s \in x_0 + L$, then $x_0 + ss \in x_0 + L$] (Due to the linearity of space). Thus, if $f(\bar{x})$ is the minimizer of $f(x)$ restricted to L , then $\nabla f(\bar{x})$ is orthogonal to L w.r.t the standard inner product.

Lemma 2.3

9/10

we assume lemma 2.6 is correct, which will be proved later, but we will uses it to prove lemma 2.3

The minimizer of $\phi(\alpha) = f(\bar{x} + \alpha s)$ is $\alpha = \frac{s^T(b - A\bar{x})}{\langle s, s \rangle_A}$

we consider \bar{x} is the minimizer of f_{cp} restricted to L^\perp

Then \bar{x} can be consider as a linear combination of bases of L which is s_1, \dots, s_k . Thus we can represent $\bar{x} = \sum_{i=1}^k \lambda_i s_i$ where s_i are mutually A -conjugate. Thus

$$\alpha = \frac{s^T(b - A\bar{x})}{\langle s, s \rangle_A} = \frac{s^T b - s^T A\bar{x}}{\langle s, s \rangle_A} = \frac{s^T b - \langle s, \bar{x} \rangle_A}{\langle s, s \rangle_A}$$

$$\langle s, \bar{x} \rangle_A = \langle s, \sum_{i=1}^k \lambda_i s_i \rangle_A = \sum_{i=1}^k \lambda_i \langle s, s_i \rangle_A$$

$\therefore s$ is any vector that is A -conjugate to L

$\therefore \langle s, s_i \rangle_A = 0 \quad \forall s_i \text{ for } i = 1, \dots, k$

$$\therefore \alpha_{\text{minimizer}} = \frac{s^T b}{\langle s, s \rangle_A}$$

And when we want to find the minimizer of f_{cp} restricted to $\text{span}(L \cup \{s\})$, and s is A -conjugate to L , then every point can be represented as linear combination of $\{s_1, \dots, s_k, s\}$

Thus we consider a global minimizer of f_{cp} restricted to $x^* + \text{span}(L \cup \{s\})$ is $x' = \sum_{i=1}^k \lambda_i s_i + \alpha' s + x^*$

According to the lecture note, For each coefficient of basis, [it can be computed independently due to the convexity of the quadratic function:]

$$\lambda_i = \frac{b^T s_i}{\langle s_i, s_i \rangle_A} \quad \text{and} \quad \lambda' = \frac{b^T s}{\langle s, s \rangle_A}$$

And it's see to see that the first k coordinate of λ' should be same as \bar{x} and

$$\lambda' = \frac{b^T s}{\langle s, s \rangle_A} = \frac{s^T b}{\langle s, s \rangle_A}$$

which is as same as λ minimizer and $\hat{x} = \bar{x} + \frac{s^T b}{\langle s, s \rangle_A} s$

Thus, it's clear to show that

$$x' \stackrel{\neq}{=} \hat{x} \quad x' \in \text{span}(L \cup \{s\}) \quad -1$$

By having such proof, we can say that \hat{x} is also the minimizer of fcp restricted to $\text{Span}\{L \cup \{s\}\}$

Lemma 2.6 415

$$\begin{aligned}
 \phi(\alpha) &= f(x+2s) = \frac{1}{2}(x+2s)^T A(x+2s) - b^T(x+2s) \\
 &= \frac{1}{2}(x+2s)^T(Ax+2As) - b^T(x+2s) \\
 &= \frac{1}{2}(x^T+2s^T)(Ax+2As) - b^Tx - 2b^Ts \\
 &= \frac{1}{2}x^TAx + \frac{1}{2}2x^TAs + \frac{1}{2}2s^TAx + \frac{1}{2}2^2s^TAs \\
 &\quad - b^Tx - 2b^Ts \\
 &= \frac{1}{2}x^TAx + 2x^TAs + \frac{1}{2}2^2s^TAs - b^Tx - 2b^Ts
 \end{aligned}$$

$$\frac{\partial \phi(\alpha)}{\partial \alpha} = x^TAs + 2s^TAs - b^Ts$$

$$\frac{\partial^2 \phi(\alpha)}{\partial \alpha^2} = s^TAs$$

To find minimum of $\phi(\alpha)$ w.r.t α , we set $\frac{\partial \phi(\alpha)}{\partial \alpha} = 0$

$$\therefore x^TAs + 2s^TAs - b^Ts = 0$$

$$\Rightarrow \alpha \cdot \langle s, s \rangle_A = b^Ts - x^TAs = s^Tb - s^TAx$$

$$\Rightarrow \alpha = \frac{s^T(b-Ax)}{\langle s, s \rangle_A}$$

also $\frac{\partial^2 \phi(\alpha)}{\partial \alpha^2} = s^TAs$, and A is a positive definite matrix

$$\therefore s^TAs \geq 0 \text{ for } \forall s \in \mathbb{R}^n \therefore \frac{\partial^2 \phi(\alpha)}{\partial \alpha^2} \geq 0$$

$\therefore \alpha = \frac{s^T(b-Ax)}{\langle s, s \rangle_A}$ is the minimizer of $\phi(\alpha)$

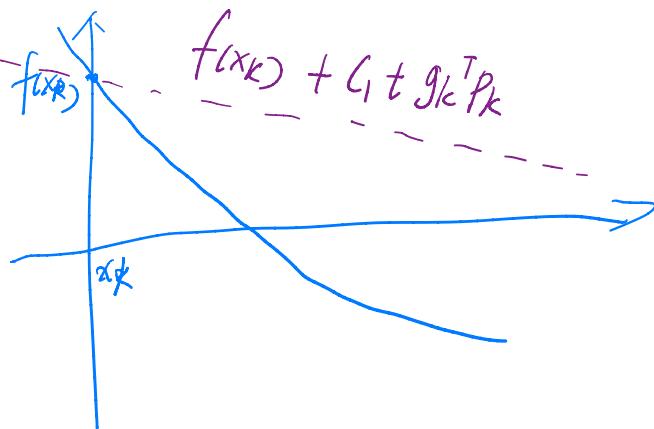
2. (25 pts) Prove that the bisection method for Wolfe linesearch, i.e., Algorithm 13 from slide 76 in the "Line Search Methods" slides on the course webpage, terminates with a steplength satisfying the weak Wolfe conditions. Assume that the function f is bounded from below and has a gradient that is Lipschitz continuous with a global Lipschitz constant γ .

Answer:

The only situation that the Algorithm won't terminate is that

$$f(x_k + t_p k) \leq f(x_k) + C_1 t_p g_k^T p_k \text{ but } g(x_k + t_p k)^T p_k < C_2 (g_k^T p_k)$$

In each iteration, α is doubled and β is always $= \infty$ which mean the function $f(x)$ is unbounded below and $g(x)^T p_k < 0$ for all x , which can be shown on following graph



But as said in the question, $f(x)$ is bounded below, this situation is impossible, $f(x_k) + C_1 t_p g_k^T p_k$ must intersect with $f(x)$ at least once.

If we assume $\alpha_k = 0$ for all k , which mean

$$f(x_k + t_k p_k) > f(x_k) + C_1 t_k g_k^T p_k \quad \forall k$$

$$\Rightarrow \frac{f(x_k + t_k p_k) - f(x_k)}{t_k} - C_1 g_k^T p_k > 0$$

When k become sufficiently large, due to each iteration

$$t_k = \frac{\alpha_k + \beta_k}{2} = \frac{\beta_k}{2} \quad \text{and} \quad \beta_k = \beta_{k-1}$$

Thus $\lim_{k \rightarrow \infty} t_k = 0$

$$\therefore \lim_{k \rightarrow \infty} \frac{f(x_k + t_k p_k) - f(x_k)}{t_k} = g_k^T p_k$$

Then $\lim_{k \rightarrow \infty} g_k^T p_k \geq C_1 g_k^T p_k$, which is contradiction. Thus we now can assume that for specific k_0 when $k \geq k_0$, $\alpha_k > 0$ and $0 < \alpha_k < t_k < \beta_k < \infty$, while α_k, t_k, β_k are approaching a specific value $\bar{\ell}$. Then we can compute following inequality:

$$f(x + \alpha_k p_k) \leq f(x) + C_1 \alpha_k g_k^T p_k \quad \text{--- (1)}$$

$$f(x + \beta_k p_k) \geq f(x) + C_1 \beta_k g_k^T p_k \quad \text{--- (2)}$$

$$\nabla f(x + \alpha_k p_k)^T p_k \leq C_2 \nabla f(x)^T p_k$$

$$\lim_{k \rightarrow \infty} \nabla f(x + \alpha_k p_k)^T p_k \leq C_2 \nabla f(x)^T p_k = \nabla f(x + \bar{\ell} p_k)^T p_k \leq C_2 \nabla f(x)^T p_k$$

By adding (1) and (2) together:

$$f(x + \alpha_k p_k) + f(x) + C_1 \beta_k g_k^T p_k \leq f(x) + C_1 \alpha_k g_k^T p_k + f(x + \beta_k p_k)$$

$$f(x + \beta_k p_k) - f(x + \alpha_k p_k) \geq C_1 (\beta_k - \alpha_k) g_k^T p_k \quad \text{--- (3)}$$

By using mean value theorem:

$$f(x + \beta_k p_k) - f(x + \alpha_k p_k) = (\beta_k - \alpha_k) \cdot \nabla f(x + \bar{\ell} p_k)^T p_k \quad \text{--- (4)}$$

When $\alpha_k \leq \bar{\ell} \leq \beta_k$

Combining ③ and ④, we can get:

$$(\beta_k - \alpha_k) \nabla f(x + \frac{1}{k} p_k)^T p_k \geq C_1 (\beta_k - \alpha_k) \nabla f(x)^T p_k$$

$$\Rightarrow \nabla f(x + \frac{1}{k} p_k)^T p_k \geq C_1 \nabla f(x)^T p_k$$

When $k \rightarrow \infty$, $\frac{1}{k} \approx 0$, which yield

$$\nabla f(x + \frac{1}{k} p_k)^T p_k \geq C_1 \nabla f(x)^T p_k$$

However, $\nabla f(x)^T p_k$ is supposed to be less than 0 for the reason that p_k is the descent direction of $f(x)$.

Thus with $0 < C_1 < C_2 < 1$:

$C_1 \nabla f(x)^T p_k$ is supposed to be larger than $C_2 \nabla f(x)^T p_k$ but in here:

$$\nabla f(x + \frac{1}{k} p_k)^T p_k \geq C_1 \nabla f(x)^T p_k$$

$$\nabla f(x + \frac{1}{k} p_k)^T p_k \leq C_2 \nabla f(x)^T p_k$$

which is a contradiction to our hypothesis. Thus we can say bisection method can find a step length to satisfy wolfe condition

3. (20 pts) Prove Theorem 4.8 from slide 103 in the “Line Search Methods” slides on the course webpage, that establishes a global convergence rate of $O\left(\left(\frac{1}{\epsilon}\right)^2\right)$ for a modified or quasi Newton method with Wolfe linesearch. You may use Theorem 3.4 on slide 71 (due to Zoutendijk) without proof. [Recall that the condition number of B_k is given by $\frac{\lambda_{\max}(B_k)}{\lambda_{\min}(B_k)}$ which is the ratio of the largest eigenvalue of B_k to the smallest eigenvalue of B_k .]

Answer:

Based on Wolfe linesearch, each iteration satisfy the following conditions

$$f(x_{k+1}) \leq f(x_k) - c_1 \alpha_k \|g_k^T p_k\| \quad \text{--- (1)}$$

$$g_{k+1}^T p_k \geq c_2 g_k^T p_k \quad \text{--- (2)}$$

while “ f is Lipschitz continuous on \mathbb{R}^n ” implies that

$$(g_{k+1} - g_k)^T p_k \leq \gamma^2 \|p_k\|_2^2 \quad \text{--- (3)}$$

From Wolfe condition (2), we can get

$$(g_{k+1} - g_k)^T p_k \geq (c_2 - 1) g_k^T p_k \quad \text{--- (4)}$$

Combining (3) and (4), we can get following relation

$$\alpha_k \geq \frac{c_2 - 1}{\gamma^2} \cdot \frac{g_k^T p_k}{\|p_k\|_2^2} \quad \text{--- (5)}$$

By substituting (5) back to Wolfe condition (1), we can obtain:

$$f(x_{k+1}) \leq f(x_k) - c_1 \frac{c_2 - 1}{\gamma^2} \frac{(g_k^T p_k)^2}{\|p_k\|_2^2} \quad \text{--- (6)}$$

By summing iteration from 0 to T and set $C = c_1 \frac{c_2 - 1}{\gamma^2}$, we can obtain following inequality:

$$\left. \begin{aligned} f(x_{T+1}) &\leq f(x_T) - C \frac{(g_T^T p_T)^2}{\|p_T\|_2^2} \\ &\vdots \\ f(x_1) &\leq f(x_0) - C \frac{(g_0^T p_0)^2}{\|p_0\|_2^2} \end{aligned} \right\} \xrightarrow{\text{sum}} f(x_{T+1}) \leq f(x_0) - C \sum_{k=0}^T \frac{(g_k^T p_k)^2}{\|p_k\|_2^2} \quad \text{--- (7)}$$

While $f(x_0)$ is bounded below, we can rewrite (7) as

$$C \sum_{k=0}^T \frac{(g_k^T p_k)^2}{\|p_k\|_2^2} \leq f(x_0) - f(x_{T+1}) \leq M' \quad \text{--- (8)}$$

Due to the modified Newton method ($B_k p_k = -g_k$), and B_k is a positive definite matrix, Thus

$$\lambda_{\min}(B_k) \leq \frac{s^T B_k s}{\|s\|^2} \leq \lambda_{\max}(B_k) \Leftrightarrow \lambda_{\max}^{-1}(B_k) \leq \frac{s^T B_k^{-1} s}{\|s\|^2} \leq \lambda_{\min}^{-1}(B_k)$$

Where $\lambda_{\min}(B_k)$ and $\lambda_{\max}(B_k)$ are the smallest and largest eigenvalue respectively of B_k , Thus:

$$|g_k^T p_k| = |g_k^T B_k^{-1} g_k| \geq \lambda_{\max}^{-1}(B_k) \|g_k\|_2^2 \quad \text{--- (9)}$$

$$\|p_k\|_2^2 = g_k^T B_k^{-2} g_k \leq \lambda_{\min}^{-2}(B_k) \|g_k\|_2^2$$

$$\therefore \|p_k\|_2 \leq \lambda_{\min}^{-1}(B_k) \|g_k\|_2 \quad \text{--- (10)}$$

Combining (9) and (10), we can obtain:

$$\frac{|g_k^T p_k|}{\|p_k\|_2} \geq \frac{\lambda_{\max}^{-1}(B_k) \|g_k\|_2^2}{\lambda_{\min}^{-1}(B_k) \|g_k\|_2} = \frac{\lambda_{\min}(B_k)}{\lambda_{\max}(B_k)} \|g_k\|_2 = \frac{1}{\text{cond}(B_k)} \|g_k\|_2$$

While $\text{cond}(B_k)$ is bounded by β , Thus:

$$\frac{|g_k^T p_k|}{\|p_k\|_2} \geq \frac{1}{\text{cond}(\beta_k)} \|g_k\|_2 \geq \frac{1}{\beta} \cdot \|g_k\|_2$$

Thus

$$\frac{(g_k^T p_k)^2}{\|p_k\|_2^2} \geq \frac{1}{\beta^2} \cdot \|g_k\|_2^2 \quad \text{--- (11)}$$

Combining (11) and (8), we can obtain following inequality:

$$C \sum_{k=1}^T \frac{1}{\beta^2} \|g_k\|_2^2 \leq C \sum_{k=0}^T \frac{(g_k^T p_k)^2}{\|p_k\|_2^2} \leq M'$$

Thus

$$\sum_{k=0}^T \|g_k\|_2^2 \leq \frac{\beta^2 M'}{C}, \quad \text{where } C = C_1 \frac{C_2 - 1}{r^2} \quad \text{--- (12)}$$

while it's intuitively to prove that

$$\sum_{k=1}^T \|g_k\|_2^2 \geq (T+1) \cdot \min_{k=0, \dots, T} \|g_k\|_2^2 \quad \text{--- (13)}$$

With (12) and (13), then

$$\min_{k=0, 1, \dots, T} \|g_k\| \leq \sqrt{\frac{\beta^2 M'}{C(T+1)}} = \frac{M}{\sqrt{T+1}}$$

$$\text{where } M = \sqrt{\frac{\beta^2 M'}{C}} = \beta r^2 \sqrt{\frac{M'}{C(C_2 - 1)}}$$

4. (10 pts) Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function, not necessarily differentiable. Show that any local minimum is also a global minimum.

Answer:

Consider $f(x^*)$ is a local minimum for $f: \mathbb{R}^n \rightarrow \mathbb{R}$ in x^* neighbor region $B(x, \varepsilon)$ where $\varepsilon \in \mathbb{R}^n$, Thus

$$f(x^*) < f(x), \forall x \in B(x, \varepsilon)$$

and we consider a point $x' = (1-t)x^* + ty$ where y can be any point within domain($f(x)$), $y \in \text{domain}(f)$, if we make t sufficiently small to make x' land in x^* neighbor, Thus

$$f(x^*) \leq f(x') = f((1-t)x^* + ty) \quad \text{--- (1)}$$

While $f(x)$ is a convex function, it implies that

$$f((1-t)x^* + ty) \leq (1-t)f(x^*) + t f(y) \quad \text{--- (2)}$$

Combining (1) and (2), we can get: 10/10

$$f(x^*) \leq (1-t)f(x^*) + t f(y)$$

$$\Rightarrow t f(x^*) \leq t f(y)$$

$$\Rightarrow f(x^*) \leq f(y) \quad (0 \leq t \leq 1)$$

With the fact that y can be any point within the domain(f) we can conclude that $f(x^*)$ is the global minimizer of f . Thus, it can be proved that local minimizer of convex function is also a global minimizer.