

15.093 Optimization Methods

Lecture 19: Line Searches
and Newton's Method

1 Last Lecture

SLIDE 1

- *Necessary Conditions for Optimality*
(identifies candidates)
 x^* local min $\Rightarrow \nabla f(x^*) = 0, \nabla^2 f(x^*)$ PSD
- *Sufficient Conditions for Optimality*
 $\nabla f(x^*) = 0 \nabla^2 f(x)$ psd $x \in B_\epsilon(x^*) \Rightarrow x^*$ loc. min.
- *Characterizations of Convexity*
 - a. f convex $\Leftrightarrow f(y) \geq f(x) + \nabla f(x)'(y - x) \forall x, y$
 - b. f convex $\Leftrightarrow \nabla^2 f(x)$ PSD $\forall x$
- f convex then global min $\Leftrightarrow \nabla f(x^*) = 0$

2 Steepest descent

2.1 The algorithm

SLIDE 2

- Step 0** Given x^0 , set $k := 0$.
- Step 1** $d^k := -\nabla f(x^k)$. If $\|d^k\| \leq \epsilon$, then stop.
- Step 2** Solve $\min_\lambda h(\lambda) := f(x^k + \lambda d^k)$ for the step-length λ^k , perhaps chosen by an exact or inexact line-search.
- Step 3** Set $x^{k+1} \leftarrow x^k + \lambda^k d^k$, $k \leftarrow k + 1$.
Go to **Step 1**.

3 Outline

SLIDE 3

1. Bisection Method - Armijo's Rule
2. Motivation for Newton's method
3. Newton's method
4. Quadratic rate of convergence
5. Modification for global convergence

4 Choices of step sizes

SLIDE 4

- $\text{Min}_\lambda f(x^k + \lambda d^k)$
- Limited Minimization: $\text{Min}_{\lambda \in [0, s]} f(x^k + \lambda d^k)$
- Constant stepsize $\lambda^k = s$ constant

- Diminishing stepsize: $\lambda^k \rightarrow 0$, $\sum_k \lambda^k = \infty$
- Armijo Rule

4.1 Bisection Line- Search Algorithm

4.1.1 Convex functions

SLIDE 5

$$\bar{\lambda} := \arg \min_{\lambda} h(\lambda) := \arg \min_{\lambda} f(\bar{\mathbf{x}} + \lambda \bar{\mathbf{d}})$$

If $f(\mathbf{x})$ is convex, $h(\lambda)$ is convex.

Find $\bar{\lambda}$ for which $h'(\lambda) = 0$

$$h'(\lambda) = \nabla f(\bar{\mathbf{x}} + \lambda \bar{\mathbf{d}})' \bar{\mathbf{d}}$$

4.1.2 Algorithm

SLIDE 6

Step 0. Set $k = 0$. Set $\lambda_L := 0$ and $\lambda_U := \hat{\lambda}$.

Step k. Set $\tilde{\lambda} = \frac{\lambda_U + \lambda_L}{2}$ and compute $h'(\tilde{\lambda})$.
 If $h'(\tilde{\lambda}) > 0$, re-set $\lambda_U := \tilde{\lambda}$. Set $k \leftarrow k + 1$.
 If $h'(\tilde{\lambda}) < 0$, re-set $\lambda_L := \tilde{\lambda}$. Set $k \leftarrow k + 1$.
 If $h'(\tilde{\lambda}) = 0$, stop.

4.1.3 Analysis

SLIDE 7

- At the k th iteration of the bisection algorithm, the current interval $[\lambda_L, \lambda_U]$ must contain a point $\bar{\lambda}$: $h'(\bar{\lambda}) = 0$
- At the k th iteration of the bisection algorithm, the length of the current interval $[\lambda_L, \lambda_U]$ is

$$\text{length} = \left(\frac{1}{2}\right)^k (\hat{\lambda}).$$

- A value of λ such that $|\lambda - \bar{\lambda}| \leq \epsilon$ can be found in at most

$$\left\lceil \log_2 \left(\frac{\hat{\lambda}}{\epsilon} \right) \right\rceil$$

steps of the bisection algorithm.

4.1.4 Example

SLIDE 8

$$h(\lambda) = (x_1 + \lambda d_1) - 0.6(x_2 + \lambda d_2) + 4(x_3 + \lambda d_3) + \\ + 0.25(x_4 + \lambda d_4) - \sum_{i=1}^4 \log(x_i + \lambda d_i) - \\ - \log\left(5 - \sum_{i=1}^4 x_i - \lambda \sum_{i=1}^4 d_i\right)$$

where $(x_1, x_2, x_3, x_4) = (1, 1, 1, 1)'$ and $(d_1, d_2, d_3, d_4) = (-1, 0.6, -4, -0.25)'$.

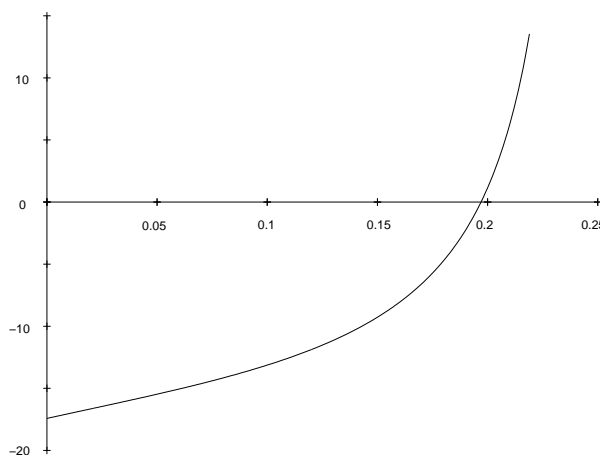
SLIDE 9

and

$$h(\lambda) = 4.65 - 17.4225\lambda - \log(1 - \lambda) - \log(1 + 0.6\lambda) - \\ - \log(1 - 4\lambda) - \log(1 - 0.25\lambda) - \log(1 + 4.65\lambda)$$

$$h'(\lambda) = -17.4225 + \frac{1}{1 - \lambda} - \frac{0.6}{1 + 0.6\lambda} + \frac{4}{1 - 4\lambda} + \\ + \frac{0.25}{1 - 0.25\lambda} - \frac{4.65}{1 + 4.65\lambda}$$

SLIDE 10



SLIDE 11

k	λ_l^k	λ_u^k	$h'(\tilde{\lambda})$
1	0.0000000	1.0000000	NaN
2	0.0000000	0.5000000	NaN
3	0.0000000	0.2500000	-11.520429338348
4	0.1250000	0.2500000	-2.952901763683
5	0.1875000	0.2500000	13.286386294218
6	0.1875000	0.2187500	2.502969022220
7	0.1875000	0.2031250	-0.605144021505
8	0.1953125	0.2031250	0.831883373151
9	0.1953125	0.1992188	0.087369215988
10	0.1953125	0.1972656	-0.265032213496

SLIDE 12

k	λ_l^k	λ_u^k	$h'(\lambda)$
20	0.1970253	0.1970272	-0.000184301091
30	0.1970268	0.1970268	-0.000000146531
40	0.1970268	0.1970268	0.000000000189
41	0.1970268	0.1970268	0.000000000023
42	0.1970268	0.1970268	-0.000000000059
43	0.1970268	0.1970268	-0.000000000018
44	0.1970268	0.1970268	0.000000000003
45	0.1970268	0.1970268	-0.000000000008
46	0.1970268	0.1970268	-0.000000000002
47	0.1970268	0.1970268	0.000000000000

4.2 Armijo Rule

SLIDE 13

Bisection is accurate but may be expensive in practice

Need cheap method guaranteeing sufficient accuracy

Inexact line search method. Requires two parameters: $\epsilon \in (0, 1)$, $\sigma > 1$.

$$\bar{h}(\lambda) = h(0) + \lambda \epsilon h'(0)$$

$\bar{\lambda}$ acceptable by Armijo's rule if:

- $h(\bar{\lambda}) \leq \bar{h}(\bar{\lambda})$
 - $h(\sigma \bar{\lambda}) \geq \bar{h}(\sigma \bar{\lambda})$ (prevents the step size be small)
- (i.e. $f(x^k + \bar{\lambda} d^k) \leq f(x^k) + \bar{\lambda} \epsilon \nabla f(x^k)' d^k$)

SLIDE 14

We get a range of acceptable stepsizes.

Step 0: Set $k = 0$, $\lambda^0 = \bar{\lambda} > 0$

Step k: If $h(\lambda^k) \leq \bar{h}(\lambda^k)$, choose λ^k stop. If $h(\lambda^k) > \bar{h}(\lambda^k)$ let $\lambda^{k+1} = \frac{1}{\sigma} \lambda^k$ (for example, $\sigma = 2$)

5 Newton's method

5.1 History

SLIDE 15

Steepest Descent is simple but slow

Newton's method complex but fast

Origins not clear

Raphson became member of the Royal Society in 1691 for his book "Analysis Aequationum Universalis" with Newton method.

Raphson published it 50 years before Newton.

6 Newton's method

6.1 Motivation

SLIDE 16

Consider

$$\min f(\mathbf{x})$$

- Taylor series expansion around $\bar{\mathbf{x}}$

$$f(\mathbf{x}) \approx g(\mathbf{x}) = f(\bar{\mathbf{x}}) + \nabla f(\bar{\mathbf{x}})'(\mathbf{x} - \bar{\mathbf{x}}) + \frac{1}{2}(\mathbf{x} - \bar{\mathbf{x}})'\nabla^2 f(\bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})$$

- Instead of $\min f(\mathbf{x})$, solve $\min g(\mathbf{x})$, i.e., $\nabla g(\mathbf{x}) = \mathbf{0}$

- $\nabla f(\bar{\mathbf{x}}) + \nabla^2 f(\bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}}) = \mathbf{0}$

$$\Rightarrow \mathbf{x} - \bar{\mathbf{x}} = -(\nabla^2 f(\bar{\mathbf{x}}))^{-1}\nabla f(\bar{\mathbf{x}})$$

- The direction $\mathbf{d} = -(\nabla^2 f(\bar{\mathbf{x}}))^{-1}\nabla f(\bar{\mathbf{x}})$ is the Newton direction

6.2 The algorithm

SLIDE 17

Step 0 Given \mathbf{x}^0 , set $k := 0$.

Step 1 $\mathbf{d}^k := -(\nabla^2 f(\mathbf{x}^k))^{-1}\nabla f(\mathbf{x}^k)$.
If $\|\mathbf{d}^k\| \leq \epsilon$, then stop.

Step 2 Set $\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k + \mathbf{d}^k$, $k \leftarrow k + 1$.
Go to **Step 1**.

6.3 Comments

SLIDE 18

- The method assumes that $\nabla^2 f(\mathbf{x}^k)$ is nonsingular at each iteration
- There is no guarantee that $f(\mathbf{x}^{k+1}) \leq f(\mathbf{x}^k)$
- We can augment the algorithm with a line-search:

$$\min_{\lambda} f(\mathbf{x}^k + \lambda\mathbf{d}^k)$$

6.4 Properties

SLIDE 19

Theorem If $\mathbf{H} = \nabla^2 f(\mathbf{x}^k)$ is PD, then \mathbf{d}^k is a descent direction: $\nabla f(\mathbf{x}^k)'\mathbf{d}^k < 0$

Proof:

$$\nabla f(\mathbf{x}^k)'\mathbf{d}^k = -\nabla f(\mathbf{x}^k)'\mathbf{H}^{-1}\nabla f(\mathbf{x}^k) < 0$$

if \mathbf{H}^{-1} is PD. But,

$$0 < \mathbf{v}'(\mathbf{H}^{-1})'\mathbf{H}\mathbf{H}^{-1}\mathbf{v} = \mathbf{v}'\mathbf{H}^{-1}\mathbf{v} \\ \Rightarrow \mathbf{H}^{-1} \text{ is PD}$$

6.5 Example 1

SLIDE 20

$$f(x) = 7x - \log x, \quad x^* = \frac{1}{7} = 0.14857143$$

$$f'(x) = 7 - \frac{1}{x}, \quad f''(x) = \frac{1}{x^2}$$

$$\Rightarrow d = -(f''(x))^{-1} f'(x) = -\left(\frac{1}{x^2}\right)^{-1} \left(7 - \frac{1}{x}\right) = x - 7x^2$$

$$x^{k+1} = x^k + (x^k - 7(x^k)^2) = 2x^k - 7(x^k)^2$$

SLIDE 21

k	x^k	x^k	x^k	x^k
0	1	0	0.01	0.1
1	-5	0	0.0193	0.13
2	-185	0	0.03599	0.1417
3	-239945	0	0.062917	0.14284777
4	-4E11	0	0.098124	0.142857142
5	-112E22	0	0.128849782	0.142857143
6		0	0.141483700	0.142857143
7		0	0.142843938	0.142857143
8		0	0.142857142	0.142857143

6.6 Example 2

SLIDE 22

$$f(x_1, x_2) = -\log(1 - x_1 - x_2) - \log x_1 - \log x_2$$

$$\nabla f(x_1, x_2) = \begin{bmatrix} \frac{1}{1 - x_1 - x_2} - \frac{1}{x_1} \\ \frac{1}{1 - x_1 - x_2} - \frac{1}{x_2} \end{bmatrix}$$

$$\nabla^2 f(x_1, x_2) = \begin{bmatrix} \left(\frac{1}{1 - x_1 - x_2}\right)^2 + \left(\frac{1}{x_1}\right)^2 & \left(\frac{1}{1 - x_1 - x_2}\right)^2 \\ \left(\frac{1}{1 - x_1 - x_2}\right)^2 & \left(\frac{1}{1 - x_1 - x_2}\right)^2 + \left(\frac{1}{x_2}\right)^2 \end{bmatrix}$$

SLIDE 23

$$(x_1^*, x_2^*) = \left(\frac{1}{3}, \frac{1}{3}\right), \quad f(x_1^*, x_2^*) = 3.295836867$$

k	x_1^k	x_2^k	$\ \mathbf{x} - \mathbf{x}^*\ $
0	0.85	0.05	0.58925565
1	0.7170068	0.09659864	0.45083106
2	0.5129752	0.17647971	0.23848325
3	0.35247858	0.27324878	0.06306103
4	0.33844902	0.32623807	0.00874717
5	0.33333772	0.33325933	7.4133E-05
6	0.33333334	0.33333333	1.1953E-08
7	0.33333333	0.33333333	1.5701E-16

7 Quadratic convergence

SLIDE 24

- Recall

$$z = \lim_{n \rightarrow \infty} z_n, \quad \limsup_{n \rightarrow \infty} \frac{|z_{n+1} - z|}{|z_n - z|^2} = \delta < \infty$$

- Example: $z_n = a^{2^n}$, $0 < a < 1$, $z = 0$

$$\frac{|z_{n+1} - z|}{|z_n - z|^2} = \frac{a^{2^{n+1}}}{a^{2^{n+1}}} = 1$$

7.1 Intuitive statement

SLIDE 25

Theorem Suppose $f(\mathbf{x}) \in C^3$ (thrice cont. diff/ble), \mathbf{x}^* : $\nabla f(\mathbf{x}^*) = 0$ and $\nabla^2 f(\mathbf{x}^*)$ is nonsingular. If Newton's method is started sufficiently close to \mathbf{x}^* , the sequence of iterates converges quadratically to \mathbf{x}^* .

SLIDE 26

Lemma:

Suppose $f(\mathbf{x}) \in C^3$ on \mathbb{R}^n and x a given pt.

Then $\forall \epsilon > 0, \exists \beta > 0: \|x - y\| \leq \epsilon \Rightarrow$

$$\|\nabla f(x) - \nabla f(y) - H(y)(x - y)\| \leq \beta \|x - y\|^2$$

7.2 Formal statement

SLIDE 27

Theorem

Suppose f twice cont. diff/ble. Let $\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x})$, $\mathbf{H}(\mathbf{x}) = \nabla^2 f(\mathbf{x})$, and $\mathbf{x}^* : \mathbf{g}(\mathbf{x}^*) = \mathbf{0}$. Suppose $f(\mathbf{x})$ is twice differentiable and its Hessian satisfies:

- $\|\mathbf{H}(\mathbf{x}^*)\| \geq h$
- $\|\mathbf{H}(\mathbf{x}) - \mathbf{H}(\mathbf{x}^*)\| \leq L \|\mathbf{x} - \mathbf{x}^*\|, \forall \mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\| \leq \beta$

(Reminder: If \mathbf{A} is a square matrix:

$$\|\mathbf{A}\| = \max_{\mathbf{x}: \|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\| = \max\{|\lambda_1|, \dots, |\lambda_n|\}$$

Suppose

SLIDE 28

$$\|\mathbf{x} - \mathbf{x}^*\| \leq \min\left(\beta, \frac{2h}{3L}\right) = \gamma.$$

Let $\bar{\mathbf{x}}$ be the next iterate in the Newton method. Then,

$$\|\bar{\mathbf{x}} - \mathbf{x}^*\| \leq \|\mathbf{x} - \mathbf{x}^*\|^2 \left(\frac{2h}{3L}\right) \quad \text{and}$$

$$\|\bar{\mathbf{x}} - \mathbf{x}^*\| < \|\mathbf{x} - \mathbf{x}^*\| < \gamma.$$

7.3 Proof

SLIDE 29

$\bar{\mathbf{x}} - \mathbf{x}^*$

$$\begin{aligned} &= \mathbf{x} - \mathbf{H}(\mathbf{x})^{-1}\mathbf{g}(\mathbf{x}) - \mathbf{x}^* \\ &= \mathbf{x} - \mathbf{x}^* + \mathbf{H}(\mathbf{x})^{-1}\left(\mathbf{g}(\mathbf{x}^*) - \mathbf{g}(\mathbf{x})\right) \\ &= \mathbf{x} - \mathbf{x}^* + \mathbf{H}(\mathbf{x})^{-1} \\ &\quad \int_0^1 \mathbf{H}(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x}))(\mathbf{x}^* - \mathbf{x}) dt \quad (\text{Lemma 1}) \\ &= \mathbf{H}(\mathbf{x})^{-1} \int_0^1 (\mathbf{H}(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x})) - \mathbf{H}(\mathbf{x})) (\mathbf{x}^* - \mathbf{x}) dt \end{aligned}$$

SLIDE 30

$\|\bar{\mathbf{x}} - \mathbf{x}^*\|$

$$\begin{aligned} &\leq \|\mathbf{H}(\mathbf{x})^{-1}\| \\ &\quad \int_0^1 \left\| \mathbf{H}(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x})) - \mathbf{H}(\mathbf{x}) \right\| \cdot \|\mathbf{x} - \mathbf{x}^*\| dt \\ &\leq \|\mathbf{H}(\mathbf{x})^{-1}\| \cdot \|\mathbf{x} - \mathbf{x}^*\| \int_0^1 L \|\mathbf{x} - \mathbf{x}^*\| t dt \\ &= \frac{1}{2} L \|\mathbf{H}(\mathbf{x})^{-1}\| \cdot \|\mathbf{x} - \mathbf{x}^*\|^2 \end{aligned}$$

SLIDE 31

Now

$$\begin{aligned} h &\leq \|\mathbf{H}(\mathbf{x}^*)\| \\ &= \|\mathbf{H}(\mathbf{x}) + \mathbf{H}(\mathbf{x}^*) - \mathbf{H}(\mathbf{x})\| \\ &\leq \|\mathbf{H}(\mathbf{x})\| + \|\mathbf{H}(\mathbf{x}^*) - \mathbf{H}(\mathbf{x})\| \\ &\leq \|\mathbf{H}(\mathbf{x})\| + L\|\mathbf{x}^* - \mathbf{x}\| \\ &\Rightarrow \|\mathbf{H}(\mathbf{x})\| \geq h - L\|\mathbf{x} - \mathbf{x}^*\| \end{aligned}$$

SLIDE 32

$$\Rightarrow \|\mathbf{H}(\mathbf{x})^{-1}\| \leq \frac{1}{\|\mathbf{H}(\mathbf{x})\|} \leq \frac{1}{h - L\|\mathbf{x} - \mathbf{x}^*\|}$$

$$\begin{aligned}
\Rightarrow \|\bar{\mathbf{x}} - \mathbf{x}^*\| &\leq \|\mathbf{x} - \mathbf{x}^*\|^2 \frac{L}{2(h - L\|\mathbf{x} - \mathbf{x}^*\|)} \\
&\leq \|\mathbf{x} - \mathbf{x}^*\|^2 \frac{L}{2(h - \frac{2}{3}h)} \\
&= \frac{3L}{2h} \|\mathbf{x} - \mathbf{x}^*\|^2
\end{aligned}$$

7.4 Lemma 1

SLIDE 33

- Fix \mathbf{w} . Let $\phi(t) = \mathbf{g}(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x}))' \mathbf{w}$
- $\phi(0) = \mathbf{g}(\mathbf{x})' \mathbf{w}$, $\phi(1) = \mathbf{g}(\mathbf{x}^*)' \mathbf{w}$
- $\phi'(t) = \mathbf{w}' \mathbf{H}(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x}))(\mathbf{x}^* - \mathbf{x})$
- $\phi(1) = \phi(0) + \int_0^1 \phi'(t) dt \Rightarrow$

SLIDE 34

$$\begin{aligned}
\forall \mathbf{w} : \quad \mathbf{w}'(\mathbf{g}(\mathbf{x}^*) - \mathbf{g}(\mathbf{x})) &= \\
\mathbf{w}' \int_0^1 \mathbf{H}(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x}))(\mathbf{x}^* - \mathbf{x}) dt & \\
\Rightarrow \mathbf{g}(\mathbf{x}^*) - \mathbf{g}(\mathbf{x}) &= \int_0^1 \mathbf{H}(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x}))(\mathbf{x}^* - \mathbf{x}) dt
\end{aligned}$$

7.5 Critical comments

SLIDE 35

- The iterates from Newton's method are equally attracted to local min and local max.
- We do not know β , h , L in general.
- Note, however, that they are only used in the proof, not in the algorithm.
- We do not assume convexity, only that $\mathbf{H}(\mathbf{x}^*)$ is nonsingular and not badly behaved near \mathbf{x}^* .

SLIDE 36

7.6 Properties of Convergence

Proposition:

Let $r^k = \|\mathbf{x}^k - \mathbf{x}^*\|$ and $C = \frac{3L}{2h}$. If $\|\mathbf{x}^0 - \mathbf{x}^*\| < \gamma$ then

$$r^k \leq \frac{1}{C} (Cr^0)^{2^k}$$

SLIDE 37

Proposition:

If $\|x^0 - x^*\| > \epsilon > 0 \Rightarrow \|x^k - x^*\| < \epsilon$

$$\forall k \geq \left\lceil \frac{\log\left(\frac{\log\left(\frac{1}{\epsilon}\right)}{\log\left(\frac{1}{Cr^0}\right)}\right)}{\log 2} \right\rceil$$

8 Solving systems of equations

SLIDE 38

$$\mathbf{g}(\mathbf{x}) = \mathbf{0}$$

$$\mathbf{g}(\mathbf{x}^{t+1}) \approx \mathbf{g}(\mathbf{x}^t) + \nabla \mathbf{g}(\mathbf{x}^t)(\mathbf{x}^{t+1} - \mathbf{x}^t) \Rightarrow$$

$$\mathbf{x}^{t+1} = \mathbf{x}^t - (\nabla \mathbf{g}(\mathbf{x}^t))^{-1} \mathbf{g}(\mathbf{x}^t)$$

Application in optimization: $\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x})$

9 Modifications for global convergence

SLIDE 39

- Perform line search
- When Hessian is singular or near singular, use:

$$\mathbf{d}^k = -\left(\nabla^2 f(\mathbf{x}^k) + \mathbf{D}^k\right)^{-1} \nabla f(\mathbf{x}^k)$$

- \mathbf{D}^k is a diagonal matrix:

$$\nabla^2 f(\mathbf{x}^k) + \mathbf{D}^k \text{ is PD}$$

10 Summary

SLIDE 40

1. Line search methods:
Bisection Method.
Armijo's Rule.
2. The Newton's method
3. Quadratic rate of convergence
4. Modification for global convergence

MIT OpenCourseWare
<http://ocw.mit.edu>

15.093J / 6.255J Optimization Methods
Fall 2009

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.