

15.093 Optimization Methods

Lecture 18: Optimality Conditions and
Gradient Methods
for Unconstrained Optimization

1 Outline

SLIDE 1

1. Necessary and sufficient optimality conditions
2. Gradient methods
3. The steepest descent algorithm
4. Rate of convergence
5. Line search algorithms

2 Optimality Conditions

SLIDE 2

Necessary Conds for Local Optima

“If \bar{x} is local optimum then \bar{x} must satisfy ...”

Identifies all candidates for local optima.

Sufficient Conds for Local Optima

“If \bar{x} satisfies ..., then \bar{x} must be a local optimum ”

3 Optimality Conditions

3.1 Necessary conditions

SLIDE 3

Consider

$$\min_{\mathbf{x} \in \mathfrak{R}^n} f(\mathbf{x})$$

Zero first order variation along all directions

Theorem

Let $f(\mathbf{x})$ be continuously differentiable.

If $\mathbf{x}^* \in \mathfrak{R}^n$ is a local minimum of $f(\mathbf{x})$, then

$$\nabla f(\mathbf{x}^*) = \mathbf{0} \text{ and } \nabla^2 f(\mathbf{x}^*) \text{ PSD}$$

3.2 Proof

SLIDE 4

Zero slope at local min \mathbf{x}^*

- $f(\mathbf{x}^*) \leq f(\mathbf{x}^* + \lambda \mathbf{d})$ for all $\mathbf{d} \in \mathfrak{R}^n, \lambda \in \mathfrak{R}$

- Pick $\lambda > 0$

$$0 \leq \frac{f(\mathbf{x}^* + \lambda \mathbf{d}) - f(\mathbf{x}^*)}{\lambda}$$

- Take limits as $\lambda \rightarrow 0$

$$0 \leq \nabla f(\mathbf{x}^*)' \mathbf{d}, \quad \forall \mathbf{d} \in \mathbb{R}^n$$

- Since \mathbf{d} arbitrary, replace with $-\mathbf{d} \Rightarrow \nabla f(\mathbf{x}^*) = \mathbf{0}$.

SLIDE 5

Nonnegative curvature at a local min \mathbf{x}^*

- $f(\mathbf{x}^* + \lambda \mathbf{d}) - f(\mathbf{x}^*) = \nabla f(\mathbf{x}^*)'(\lambda \mathbf{d}) + \frac{1}{2}(\lambda \mathbf{d})' \nabla^2 f(\mathbf{x}^*)(\lambda \mathbf{d}) + \|\lambda \mathbf{d}\|^2 R(\mathbf{x}^*; \lambda \mathbf{d})$
where $R(\mathbf{x}^*; \mathbf{y}) \rightarrow 0$ as $\mathbf{y} \rightarrow 0$. Since $\nabla f(\mathbf{x}^*) = 0$,

$$= \frac{1}{2} \lambda^2 \mathbf{d}' \nabla^2 f(\mathbf{x}^*) \mathbf{d} + \lambda^2 \|\mathbf{d}\|^2 R(\mathbf{x}^*; \lambda \mathbf{d}) \Rightarrow$$

$$\frac{f(\mathbf{x}^* + \lambda \mathbf{d}) - f(\mathbf{x}^*)}{\lambda^2} = \frac{1}{2} \mathbf{d}' \nabla^2 f(\mathbf{x}^*) \mathbf{d} + \|\mathbf{d}\|^2 R(\mathbf{x}^*; \lambda \mathbf{d})$$

If $\nabla^2 f(\mathbf{x}^*)$ is not PSD, $\exists \bar{\mathbf{d}}: \bar{\mathbf{d}}' \nabla^2 f(\mathbf{x}^*) \bar{\mathbf{d}} < 0 \Rightarrow f(\mathbf{x}^* + \lambda \bar{\mathbf{d}}) < f(\bar{\mathbf{x}}), \forall \lambda$
suff. small QED.

3.3 Example

SLIDE 6

$$f(x) = \frac{1}{2}x_1^2 + x_1 x_2 + 2x_2^2 - 4x_1 - 4x_2 - x_2^3$$

$$\nabla f(x) = (x_1 + x_2 - 4, x_1 + 4x_2 - 4 - 3x_2^2) \text{ Candidates } \mathbf{x}^* = (4, 0) \text{ and } \bar{\mathbf{x}} = (3, 1)$$

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} 1 & 1 \\ 1 & 4 - 6x_2 \end{bmatrix}$$

$$\nabla^2 f(\mathbf{x}^*) = \begin{bmatrix} 1 & 1 \\ 1 & 4 \end{bmatrix}$$

PSD

SLIDE 7

$$\bar{\mathbf{x}} = (3, 1)$$

$$\nabla^2 f(\bar{\mathbf{x}}) = \begin{bmatrix} 1 & 1 \\ 1 & -2 \end{bmatrix}$$

Indefinite matrix

\mathbf{x}^* is the only candidate for local min

3.4 Sufficient conditions

SLIDE 8

Theorem f twice continuously differentiable. If $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\nabla^2 f(x)$ PSD in $B(\mathbf{x}^*, \epsilon)$, then \mathbf{x}^* is a local minimum.

Proof: Taylor series expansion: For all $\mathbf{x} \in B(\mathbf{x}^*, \epsilon)$

$$f(\mathbf{x}) = f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)'(\mathbf{x} - \mathbf{x}^*) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)' \nabla^2 f(\mathbf{x}^* + \lambda(\mathbf{x} - \mathbf{x}^*))(\mathbf{x} - \mathbf{x}^*)$$

for some $\lambda \in [0, 1]$

$$\Rightarrow f(\mathbf{x}) \geq f(\mathbf{x}^*)$$

3.5 Example Continued...

SLIDE 9

At $x^* = (4, 0)$, $\nabla f(x^*) = 0$ and

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} 1 & 1 \\ 1 & 4 - 6x_2 \end{bmatrix}$$

is PSD for $\mathbf{x} \in B(\mathbf{x}^*, \epsilon)$

SLIDE 10

$f(x) = x_1^3 + x_2^2$ and $\nabla f(x) = (3x_1^2, 2x_2)$ $x^* = (0, 0)$

$$\nabla^2 f(x) = \begin{bmatrix} 6x_1 & 0 \\ 0 & 2 \end{bmatrix} \text{ is not PSD in } B(\mathbf{0}, \epsilon)$$

$$f(-\epsilon, 0) = -\epsilon^3 < 0 = f(\mathbf{x}^*)$$

3.6 Characterization of convex functions

SLIDE 11

Theorem Let $f(\mathbf{x})$ be continuously differentiable.

Then $f(\mathbf{x})$ is convex if and only if

$$\nabla f(\mathbf{x})'(\bar{\mathbf{x}} - \mathbf{x}) \leq f(\bar{\mathbf{x}}) - f(\mathbf{x})$$

3.7 Proof

SLIDE 12

By convexity

$$f(\lambda\bar{\mathbf{x}} + (1 - \lambda)\mathbf{x}) \leq \lambda f(\bar{\mathbf{x}}) + (1 - \lambda)f(\mathbf{x})$$

$$\frac{f(\mathbf{x} + \lambda(\bar{\mathbf{x}} - \mathbf{x})) - f(\mathbf{x})}{\lambda} \leq f(\bar{\mathbf{x}}) - f(\mathbf{x})$$

As $\lambda \rightarrow 0$,

$$\nabla f(\mathbf{x})'(\bar{\mathbf{x}} - \mathbf{x}) \leq f(\bar{\mathbf{x}}) - f(\mathbf{x})$$

3.8 Convex functions

SLIDE 13

Theorem Let $f(\mathbf{x})$ be a continuously differentiable convex function. Then \mathbf{x}^* is a minimum of f if and only if

$$\nabla f(\mathbf{x}^*) = \mathbf{0}$$

Proof: If f convex and $\nabla f(\mathbf{x}^*) = \mathbf{0}$

$$f(\mathbf{x}) - f(\mathbf{x}^*) \geq \nabla f(\mathbf{x}^*)'(\mathbf{x} - \mathbf{x}^*) = 0$$

3.9 Descent Directions

SLIDE 14

Interesting Observation

f diff/ble at \bar{x}

$\exists d: \nabla f(\bar{x})'d < 0 \Rightarrow \forall \lambda > 0$, suff. small, $f(\bar{x} + \lambda d) < f(\bar{x})$

(d : descent direction)

3.10 Proof

SLIDE 15

$$f(\bar{x} + \lambda d) = f(\bar{x}) + \lambda \nabla f(\bar{x})'d + \lambda \|d\| R(\bar{x}, \lambda d)$$

where $R(\bar{x}, \lambda d) \xrightarrow{\lambda \rightarrow 0} 0$

$$\frac{f(\bar{x} + \lambda d) - f(\bar{x})}{\lambda} = \nabla f(\bar{x})'d + \|d\| R(\bar{x}, \lambda d)$$

$\nabla f(\bar{x})'d < 0$, $R(\bar{x}, \lambda d) \xrightarrow{\lambda \rightarrow 0} 0 \Rightarrow$

$\forall \lambda > 0$ suff. small $f(\bar{x} + \lambda d) < f(\bar{x})$. QED

4 Algorithms for unconstrained optimization

4.1 Gradient Methods-Motivation

SLIDE 16

- Decrease $f(\mathbf{x})$ until $\nabla f(\mathbf{x}^*) = \mathbf{0}$

-

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

- If $\nabla f(\bar{\mathbf{x}})' \mathbf{d} < 0$, then for small $\lambda > 0$,

$$f(\bar{\mathbf{x}} + \lambda \mathbf{d}) < f(\bar{\mathbf{x}})$$

5 Gradient Methods

5.1 A generic algorithm

SLIDE 17

- $\mathbf{x}^{k+1} = \mathbf{x}^k + \lambda^k \mathbf{d}^k$

- If $\nabla f(\mathbf{x}^k) \neq \mathbf{0}$, direction \mathbf{d}^k satisfies:

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

- Step-length $\lambda^k > 0$

- Principal example:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \lambda^k \mathbf{D}^k \nabla f(\mathbf{x}^k)$$

\mathbf{D}^k positive definite symmetric matrix

5.2 Principal directions

SLIDE 18

- Steepest descent:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \lambda^k \nabla f(\mathbf{x}^k)$$

- Newton's method:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \lambda^k (\nabla^2 f(\mathbf{x}^k))^{-1} \nabla f(\mathbf{x}^k)$$

5.3 Other directions

SLIDE 19

- Diagonally scaled steepest descent

$$\mathbf{D}^k = \text{Diagonal approximation to } (\nabla^2 f(\mathbf{x}^k))^{-1}$$

- Modified Newton's method

$$\mathbf{D}^k = \text{Diagonal approximation to } (\nabla^2 f(\mathbf{x}^0))^{-1}$$

- Gauss-Newton method for least squares problems $f(\mathbf{x}) = \|g(\mathbf{x})\|^2$
 $(\nabla g(\mathbf{x}^k) \nabla g(\mathbf{x}^k)')^{-1}$ $\mathbf{D}^k =$

6 Steepest descent

6.1 The algorithm

SLIDE 20

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

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

Step 2 Solve $\min_{\lambda} h(\lambda) := f(\mathbf{x}^k + \lambda \mathbf{d}^k)$ for the step-length λ^k , perhaps chosen by an exact or inexact line-search.

Step 3 Set $\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k + \lambda^k \mathbf{d}^k$, $k \leftarrow k + 1$.
Go to **Step 1**.

6.2 An example

SLIDE 21

minimize $f(x_1, x_2) = 5x_1^2 + x_2^2 + 4x_1x_2 - 14x_1 - 6x_2 + 20$

$$\mathbf{x}^* = (x_1^*, x_2^*)' = (1, 1)'$$

$$f(\mathbf{x}^*) = 10$$

SLIDE 22

Given \mathbf{x}^k

$$\mathbf{d}^k = -\nabla f(x_1^k, x_2^k) = \begin{pmatrix} -10x_1^k - 4x_2^k + 14 \\ -2x_2^k - 4x_1^k + 6 \end{pmatrix} = \begin{pmatrix} d_1^k \\ d_2^k \end{pmatrix}$$

$$\begin{aligned} h(\lambda) &= f(x^k + \lambda \mathbf{d}^k) \\ &= 5(x_1^k + \lambda d_1^k)^2 + (x_2^k + \lambda d_2^k)^2 + 4(x_1^k + \lambda d_1^k)(x_2^k + \lambda d_2^k) - \\ &\quad - 14(x_1^k + \lambda d_1^k) - 6(x_2^k + \lambda d_2^k) + 20 \end{aligned}$$

$$\lambda^k = \frac{(d_1^k)^2 + (d_2^k)^2}{2(5(d_1^k)^2 + (d_2^k)^2 + 4d_1^k d_2^k)}$$

Start at $x = (0, 10)^T$

$\varepsilon = 10^{-6}$

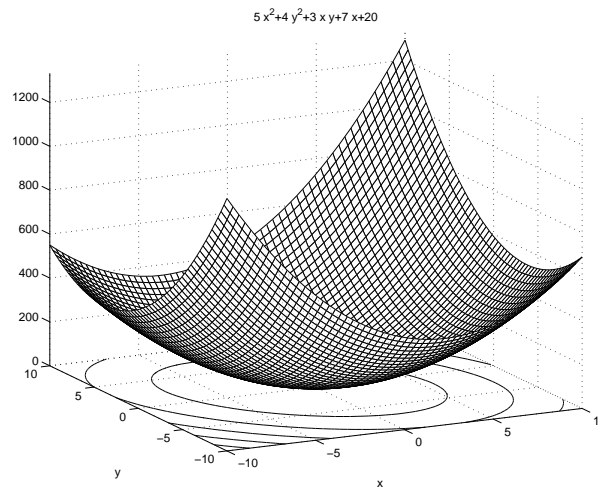
SLIDE 23

k	x_1^k	x_2^k	d_1^k	d_2^k	$\ d^k\ _2$	λ^k	$f(x^k)$
1	0.000000	10.000000	-26.000000	-14.000000	29.52964612	0.0866	60.000000
2	-2.252782	8.786963	1.379968	-2.562798	2.91071234	2.1800	22.222576
3	0.755548	3.200064	-6.355739	-3.422321	7.21856659	0.0866	12.987827
4	0.204852	2.903535	0.337335	-0.626480	0.71152803	2.1800	10.730379
5	0.940243	1.537809	-1.553670	-0.836592	1.76458951	0.0866	10.178542
6	0.805625	1.465322	0.082462	-0.153144	0.17393410	2.1800	10.043645
7	0.985392	1.131468	-0.379797	-0.204506	0.43135657	0.0866	10.010669
8	0.952485	1.113749	0.020158	-0.037436	0.04251845	2.1800	10.002608
9	0.996429	1.032138	-0.092842	-0.049992	0.10544577	0.0866	10.000638
10	0.988385	1.027806	0.004928	-0.009151	0.01039370	2.1800	10.000156

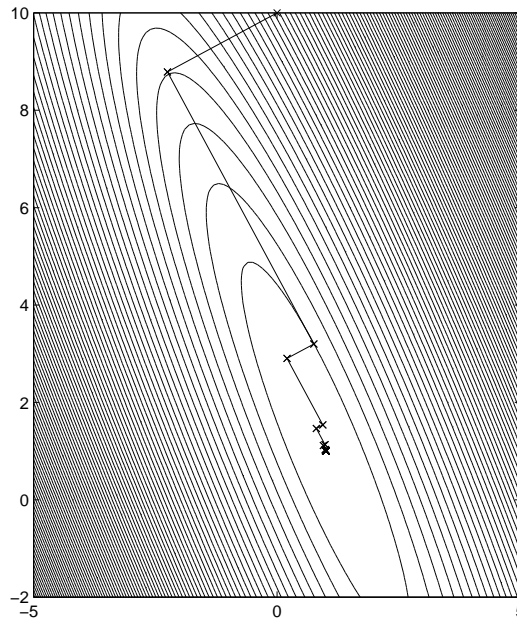
SLIDE 24

k	x_1^k	x_2^k	d_1^k	d_2^k	$\ d^k\ _2$	λ^k	$f(x^k)$
11	0.999127	1.007856	-0.022695	-0.012221	0.02577638	0.0866	10.000038
12	0.997161	1.006797	0.001205	-0.002237	0.00254076	2.1800	10.000009
13	0.999787	1.001920	-0.005548	-0.002987	0.00630107	0.0866	10.000002
14	0.999306	1.001662	0.000294	-0.000547	0.00062109	2.1800	10.000001
15	0.999948	1.000469	-0.001356	-0.000730	0.00154031	0.0866	10.000000
16	0.999830	1.000406	0.000072	-0.000134	0.00015183	2.1800	10.000000
17	0.999987	1.000115	-0.000332	-0.000179	0.00037653	0.0866	10.000000
18	0.999959	1.000099	0.000018	-0.000033	0.00003711	2.1800	10.000000
19	0.999997	1.000028	-0.000081	-0.000044	0.00009204	0.0866	10.000000
20	0.999990	1.000024	0.000004	-0.000008	0.00000907	2.1803	10.000000
21	0.999999	1.000007	-0.000020	-0.000011	0.00002250	0.0866	10.000000
22	0.999998	1.000006	0.000001	-0.000002	0.00000222	2.1817	10.000000
23	1.000000	1.000002	-0.000005	-0.000003	0.00000550	0.0866	10.000000
24	0.999999	1.000001	0.000000	-0.000000	0.00000054	0.0000	10.000000

SLIDE 25



SLIDE 26



6.3 Important Properties

SLIDE 27

- $f(\mathbf{x}^{k+1}) < f(\mathbf{x}^k) < \dots < f(\mathbf{x}^0)$ (because \mathbf{d}^k are descent directions)
- Under reasonable assumptions of $f(\mathbf{x})$, the sequence $\mathbf{x}^0, \mathbf{x}^1, \dots$, will have at least one cluster point $\bar{\mathbf{x}}$
- Every cluster point $\bar{\mathbf{x}}$ will satisfy $\nabla f(\bar{\mathbf{x}}) = \mathbf{0}$
- *Implication:* If $f(\mathbf{x})$ is a convex function, $\bar{\mathbf{x}}$ will be an optimal solution

7 Global Convergence Result

SLIDE 28

Theorem:

$f : R^n \rightarrow R$ is continuously diff/ble on $\mathcal{F} = \{x \in R^n : f(x) \leq f(x^0)\}$ closed, bounded set

Every cluster point \bar{x} of $\{x_k\}$ satisfies $\nabla f(\bar{x}) = 0$.

7.1 Work Per Iteration

SLIDE 29

Two computation tasks at each iteration of steepest descent:

- Compute $\nabla f(\mathbf{x}^k)$ (for quadratic objective functions, it takes $O(n^2)$ steps) to determine $\mathbf{d}^k = -\nabla f(\mathbf{x}^k)$

- Perform line-search of $h(\lambda) = f(\mathbf{x}^k + \lambda \mathbf{d}^k)$
to determine $\lambda^k = \arg \min_{\lambda} h(\lambda) = \arg \min_{\lambda} f(\mathbf{x}^k + \lambda \mathbf{d}^k)$

8 Rate of convergence of algorithms

Let $z_1, \dots, z_n, \dots \rightarrow z$ be a convergent sequence. We say that the order of convergence of this sequence is p^* if

SLIDE 30

$$p^* = \sup \left\{ p : \limsup_{k \rightarrow \infty} \frac{|z_{k+1} - z|}{|z_k - z|^p} < \infty \right\}$$

Let

$$\beta = \limsup_{k \rightarrow \infty} \frac{|z_{k+1} - z|}{|z_k - z|^{p^*}}$$

The larger p^* , the faster the convergence

8.1 Types of convergence

SLIDE 31

1. $p^* = 1$, $0 < \beta < 1$, then linear (or geometric) rate of convergence
2. $p^* = 1$, $\beta = 0$, super-linear convergence
3. $p^* = 1$, $\beta = 1$, sub-linear convergence
4. $p^* = 2$, quadratic convergence

8.2 Examples

SLIDE 32

- $z_k = a^k$, $0 < a < 1$ converges linearly to zero, $\beta = a$
- $z_k = a^{2^k}$, $0 < a < 1$ converges quadratically to zero
- $z_k = \frac{1}{k}$ converges sub-linearly to zero
- $z_k = \left(\frac{1}{k}\right)^k$ converges super-linearly to zero

8.3 Steepest descent

SLIDE 33

- $z_k = f(\mathbf{x}^k)$, $z = f(\mathbf{x}^*)$, where $\mathbf{x}^* = \arg \min f(\mathbf{x})$

- Then an algorithm exhibits *linear convergence* if there is a constant $\delta < 1$ such that

$$\frac{f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*)}{f(\mathbf{x}^k) - f(\mathbf{x}^*)} \leq \delta,$$

for all k sufficiently large, where \mathbf{x}^* is an optimal solution.

8.3.1 Discussion

SLIDE 34

$$\frac{f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*)}{f(\mathbf{x}^k) - f(\mathbf{x}^*)} \leq \delta < 1$$

- If $\delta = 0.1$, every iteration adds another digit of accuracy to the optimal objective value.
- If $\delta = 0.9$, every 22 iterations add another digit of accuracy to the optimal objective value, because $(0.9)^{22} \approx 0.1$.

9 Rate of convergence of steepest descent

9.1 Quadratic Case

9.1.1 Theorem

SLIDE 35

Suppose $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}' \mathbf{Q} \mathbf{x} - \mathbf{c}' \mathbf{x}$
 \mathbf{Q} is psd

λ_{\max} = largest eigenvalue of \mathbf{Q}

λ_{\min} = smallest eigenvalues of \mathbf{Q}

Linear Convergence Theorem: If $f(\mathbf{x})$ is a quadratic function and \mathbf{Q} is psd, then

$$\frac{f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*)}{f(\mathbf{x}^k) - f(\mathbf{x}^*)} \leq \left(\frac{\left(\frac{\lambda_{\max}}{\lambda_{\min}} \right) - 1}{\left(\frac{\lambda_{\max}}{\lambda_{\min}} \right) + 1} \right)^2$$

9.1.2 Discussion

SLIDE 36

$$\frac{f(x^{k+1}) - f(x^*)}{f(x^k) - f(x^*)} \leq \left(\frac{\left(\frac{\lambda_{\max}}{\lambda_{\min}} \right) - 1}{\left(\frac{\lambda_{\max}}{\lambda_{\min}} \right) + 1} \right)^2$$

- $\kappa(\mathbf{Q}) := \frac{\lambda_{\max}}{\lambda_{\min}}$ is the *condition number* of \mathbf{Q}

- $\kappa(\mathbf{Q}) \geq 1$
- $\kappa(\mathbf{Q})$ plays an extremely important role in analyzing computation involving \mathbf{Q}

SLIDE 37

$$\frac{f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*)}{f(\mathbf{x}^k) - f(\mathbf{x}^*)} \leq \left(\frac{\kappa(\mathbf{Q}) - 1}{\kappa(\mathbf{Q}) + 1} \right)^2$$

$\kappa(\mathbf{Q}) = \frac{\lambda_{\max}}{\lambda_{\min}}$	Upper Bound on Convergence Constant δ	Number of Iterations to Reduce the Optimality Gap by 0.10
1.1	0.0023	1
3.0	0.25	2
10.0	0.67	6
100.0	0.96	58
200.0	0.98	116
400.0	0.99	231

SLIDE 38

For $\kappa(\mathbf{Q}) \sim O(1)$ converges fast.

For large $\kappa(\mathbf{Q})$

$$\left(\frac{\kappa(\mathbf{Q}) - 1}{\kappa(\mathbf{Q}) + 1} \right)^2 \sim \left(1 - \frac{1}{\kappa(\mathbf{Q})} \right)^2 \sim 1 - \frac{2}{\kappa(\mathbf{Q})}$$

Therefore

$$(f(\mathbf{x}^k) - f(\mathbf{x}^*)) \leq \left(1 - \frac{2}{\kappa(\mathbf{Q})} \right)^k (f(\mathbf{x}^0) - f(\mathbf{x}^*))$$

In $k \sim \frac{1}{2}\kappa(\mathbf{Q})(-ln\epsilon)$ iterations, finds \mathbf{x}^k :

$$(f(\mathbf{x}^k) - f(\mathbf{x}^*)) \leq \epsilon(f(\mathbf{x}^0) - f(\mathbf{x}^*))$$

9.2 Example 2

SLIDE 39

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}'\mathbf{Q}\mathbf{x} - \mathbf{c}'\mathbf{x} + 10$$

$$\mathbf{Q} = \begin{bmatrix} 20 & 5 \\ 5 & 1 \end{bmatrix} \quad \mathbf{c} = \begin{pmatrix} 14 \\ 6 \end{pmatrix}$$

$$\kappa(\mathbf{Q}) = 30.234$$

$$\delta = \left(\frac{\kappa(\mathbf{Q}) - 1}{\kappa(\mathbf{Q}) + 1} \right)^2 = 0.8760$$

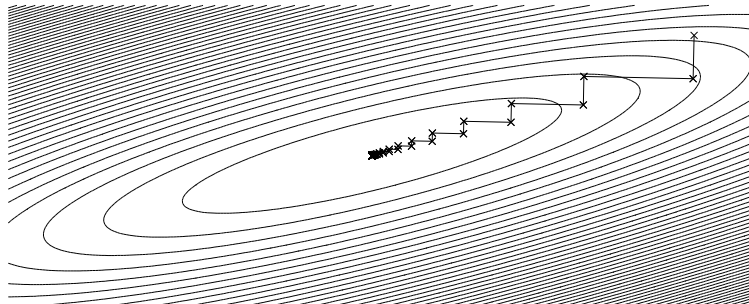
SLIDE 40

k	x_1^k	x_2^k	$\ \mathbf{d}^k\ _2$	λ^k	$f(\mathbf{x}^k)$	$\frac{f(\mathbf{x}^k) - f(\mathbf{x}^*)}{f(\mathbf{x}^{k-1}) - f(\mathbf{x}^*)}$
1	40.000000	-100.000000	286.06293014	0.0506	6050.000000	
2	25.542693	-99.696700	77.69702948	0.4509	3981.695128	0.658079
3	26.277558	-64.668130	188.25191488	0.0506	2620.587793	0.658079
4	16.763512	-64.468535	51.13075844	0.4509	1724.872077	0.658079
5	17.247111	-41.416980	123.88457127	0.0506	1135.420663	0.658079
6	10.986120	-41.285630	33.64806192	0.4509	747.515255	0.658079
7	11.304366	-26.115894	81.52579489	0.0506	492.242977	0.658079
8	7.184142	-26.029455	22.14307211	0.4509	324.253734	0.658079
9	7.393573	-16.046575	53.65038732	0.0506	213.703595	0.658079
10	4.682141	-15.989692	14.57188362	0.4509	140.952906	0.658079

SLIDE 41

k	x_1^k	x_2^k	$\ d^k\ _2$	λ^k	$f(x^k)$	$\frac{f(x^k) - f(x^*)}{f(x^{k-1}) - f(x^*)}$
20	0.460997	0.948466	1.79847660	0.4509	3.066216	0.658079
30	-0.059980	3.038991	0.22196980	0.4509	0.965823	0.658079
40	-0.124280	3.297005	0.02739574	0.4509	0.933828	0.658079
50	-0.132216	3.328850	0.00338121	0.4509	0.933341	0.658079
60	-0.133195	3.332780	0.00041731	0.4509	0.933333	0.658078
70	-0.133316	3.333265	0.00005151	0.4509	0.933333	0.658025
80	-0.133331	3.333325	0.00000636	0.4509	0.933333	0.654656
90	-0.133333	3.333332	0.00000078	0.0000	0.933333	0.000000

SLIDE 42



9.3 Example 3

SLIDE 43

$$f(x) = \frac{1}{2}x'Qx - c'x + 10$$

$$Q = \begin{bmatrix} 20 & 5 \\ 5 & 16 \end{bmatrix} \quad c = \begin{pmatrix} 14 \\ 6 \end{pmatrix}$$

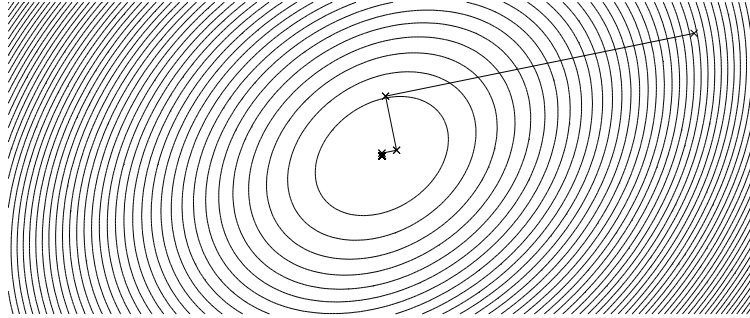
$$\kappa(Q) = 1.8541$$

$$\delta = \left(\frac{\kappa(Q) - 1}{\kappa(Q) + 1} \right)^2 = 0.0896$$

SLIDE 44

k	x_1^k	x_2^k	$\ d^k\ _2$	λ^k	$f(x^k)$	$\frac{f(x^k) - f(x^*)}{f(x^{k-1}) - f(x^*)}$
1	40.000000	-100.000000	1434.79336491	0.0704	76050.000000	
2	19.867118	-1.025060	385.96252652	0.0459	3591.615327	0.047166
3	2.513241	-4.555081	67.67315150	0.0704	174.058930	0.047166
4	1.563658	0.113150	18.20422450	0.0459	12.867208	0.047166
5	0.745149	-0.053347	3.19185713	0.0704	5.264475	0.047166
6	0.700361	0.166834	0.85861649	0.0459	4.905886	0.047166
7	0.661755	0.158981	0.15054644	0.0704	4.888973	0.047166
8	0.659643	0.169366	0.04049732	0.0459	4.888175	0.047166
9	0.657822	0.168996	0.00710064	0.0704	4.888137	0.047166
10	0.657722	0.169486	0.00191009	0.0459	4.888136	0.047166
11	0.657636	0.169468	0.00033491	0.0704	4.888136	0.047166
12	0.657632	0.169491	0.00009009	0.0459	4.888136	0.047161
13	0.657628	0.169490	0.00001580	0.0704	4.888136	0.047068
14	0.657627	0.169492	0.00000425	0.0459	4.888136	0.045002
15	0.657627	0.169491	0.00000075	0.0000	4.888136	0.000000

SLIDE 45



9.4 Empirical behavior

SLIDE 46

- The convergence constant bound is not just theoretical. It is typically experienced in practice.

- Analysis is due to Leonid Kantorovich, who won the Nobel Memorial Prize in Economic Science in 1975 for his contributions to optimization and economic planning.

SLIDE 47

- What about non-quadratic functions?
 - Suppose $\mathbf{x}^* = \arg \min_{\mathbf{x}} f(\mathbf{x})$

 - $\nabla^2 f(\mathbf{x}^*)$ is the Hessian of $f(\mathbf{x})$ at $\mathbf{x} = \mathbf{x}^*$

 - Rate of convergence will depend on $\kappa(\nabla^2 f(\mathbf{x}^*))$

10 Summary

SLIDE 48

1. Optimality Conditions
2. The steepest descent algorithm - Convergence
3. Rate of convergence of Steepest Descent

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>.