

# Network Optimization

## Chapter 2: Mathematical Tools

Yiyin Wang  
yiyinwang@sjtu.edu.cn  
Department of Automation  
Shanghai Jiao Tong University

Nov. 4, 2014

# Acknowledgements

Slides borrow heavily from lectures by Stephen Boyd (Stanford), Lieven Vandenberghe (UCLA), and Mung Chiang (Princeton)

## Part 3: Algorithms

- ▶ 2.4 Descent algorithm
- ▶ 2.5 Interior-point algorithm

# Unconstrained minimization

- ▶ Given  $f : \mathbf{R}^n \rightarrow \mathbf{R}$  convex and twice differentiable:

$$\text{minimize} \quad f(x)$$

Optimizer  $x^*$ . Optimized value  $p^* = f(x^*)$

- ▶ Necessary and sufficient condition of optimality:

$$\nabla f(x^*) = 0$$

# Unconstrained minimization

- ▶ Iterative algorithm: computes a sequence of points  $x^{(k)} \in \text{dom } f$ ,  $k = 0, 1, \dots$ , such that

$$\lim_{k \rightarrow \infty} f(x^{(k)}) = p^*$$

- ▶ Terminate algorithm when  $f(x^{(k)}) - p^* \leq \epsilon$  for a specified  $\epsilon > 0$
- ▶ Can be interpreted as iterative methods for solving optimality condition

$$\nabla f(x^*) = 0$$

- ▶ The iterative algorithm require a starting point  $x^{(0)}$  such that
  - ▶  $x^{(0)} \in \text{dom } f$
  - ▶ Sublevel set  $S = \{x \mid f(x) \leq f(x^{(0)})\}$  is closed

# Examples

- ▶ Least-squares: minimize

$$\|Ax - b\|_2^2 = x^T (A^T A)x - 2(A^T b)^T x + b^T b$$

Optimality condition ( $\nabla f(x^*) = 0$ ) is system of linear equations:

$$A^T A x^* = A^T b$$

called normal equations for least squares

- ▶ Unconstrained geometric programming: minimize

$$f(x) = \log \left( \sum_{i=1}^m \exp(a_i^T x + b_i) \right)$$

Optimality condition ( $\nabla f(x^*) = 0$ ) has no analytic solution:

$$\nabla f(x^*) = \frac{1}{\sum_{j=1}^m \exp(a_j^T x^* + b_j)} \sum_{i=1}^m \exp(a_i^T x^* + b_i) a_i = 0$$

## Strong convexity

- ▶  $f$  assumed to be strongly convex: there exists  $m > 0$  such that

$$\nabla^2 f(x) \succeq mI$$

which also implies that there exists  $M \geq m$  such that

$$\nabla^2 f(x) \preceq MI$$

- ▶ Bound optimal value:

$$f(x) - \frac{1}{2m} \|\nabla f(x)\|_2^2 \leq p^* \leq f(x) - \frac{1}{2M} \|\nabla f(x)\|_2^2$$

- ▶ Suboptimality condition:

$$\|\nabla f(x)\|_2^2 \leq (2m\epsilon)^{1/2} \rightarrow f(x) - p^* \leq \epsilon$$

- ▶ Distance between  $x$  and optimal  $x^*$

$$\|x - x^*\|_2 \leq \frac{2}{m} \|\nabla f(x)\|_2$$

# Descent Methods

- ▶ Minimizing sequence  $x^{(k)}$ ,  $k = 1, \dots$ , (where  $t^{(k)} > 0$ )

$$x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x^{(k)}$$

- ▶  $\Delta x^{(k)}$ : the search direction or step
  - ▶  $t^{(k)}$ : the step length or step size
- ▶ Descent methods:

$$f(x^{(k+1)}) < f(x^{(k)})$$

By convexity of  $f$ , search direction must make an acute angle with negative gradient:

$$\nabla f(x^{(k)})^T \Delta x^{(k)} < 0$$

Because otherwise,  $f(x^{(k+1)}) \geq f(x^{(k)})$  since  
 $f(x^{(k+1)}) \geq f(x^{(k)}) + \nabla f(x^{(k)})^T (x^{(k+1)} - x^{(k)})$



# General descent method

**Given** a starting point  $x^{(0)} \in \text{dom } f$

**Repeat**

1. Determine a descent direction  $\Delta x^{(k)}$
2. Line search: choose a step size  $t > 0$
3. Update:  $x^{(k+1)} = x^{(k)} + t\Delta x^{(k)}$

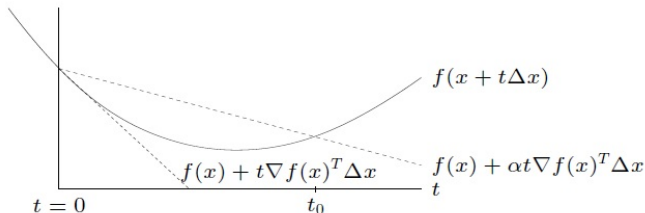
**Until** stopping criterion satisfied

# Line search types

- ▶ Exact line search:  $t = \operatorname{argmin}_{s \geq 0} f(x + s\Delta x)$
- ▶ Backtracking line search (with parameters  $\alpha \in (0, 1/2), \beta \in (0, 1)$ )
  - ▶ Starting at  $t = 1$ , repeat  $t := \beta t$  until

$$f(x + t\Delta x) < f(x) + \alpha t \nabla f(x)^T \Delta x$$

- ▶ graphical interpretation: backtrack until  $t \leq t_0$



# Gradient descent method

**Given** a starting point  $x \in \text{dom } f$

**Repeat**

1.  $\Delta x := -\nabla f(x)$
2. Line search: choose a step size  $t > 0$
3. Update:  $x := x + t\Delta x$

**Until** stopping criterion satisfied

- ▶ Stopping criterion usually of the form  $\|\nabla f(x)\|_2 \leq \epsilon$
- ▶ Convergence result: for strong convex  $f$ ,

$$f(x^{(k)}) - p^* \leq c^k(f(x^{(0)}) - p^*)$$

$c \in (0, 1)$  depends on  $m, x^{(0)}$ , line search type

- ▶ Very simple, but often very slow; rarely used in practice

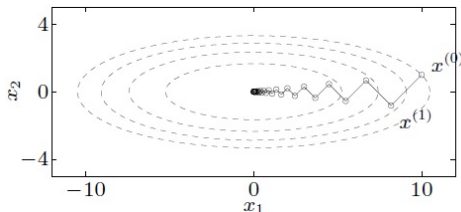
## Examples: quadratic problem in $\mathbf{R}^2$

$$\text{minimize } f(x) = \frac{1}{2}(x_1^2 + \gamma x_2^2) \quad (\gamma > 0)$$

with exact line search, starting at  $x^{(0)} = (\gamma, 1)$ :

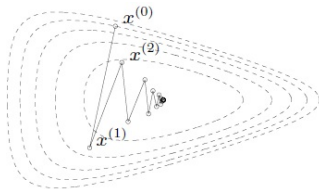
$$x_1^{(k)} = \gamma \left( \frac{\gamma - 1}{\gamma + 1} \right)^k, \quad x_2^{(k)} = \left( -\frac{\gamma - 1}{\gamma + 1} \right)^k$$

- ▶ Very slow if  $\gamma \gg 1$  or  $\gamma \ll 1$
- ▶ Example for  $\gamma = 10$

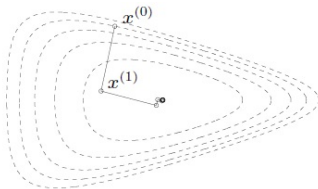


## Examples: nonquadratic example

$$f(x_1, x_2) = e^{x_1+3x_2-0.1} + e^{x_1-3x_2-0.1} + e^{-x_1-0.1}$$



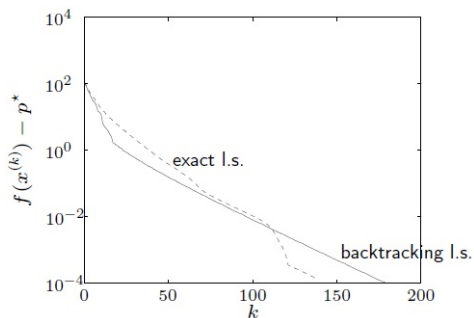
Backtracking line search



Exact line search

## Examples: a problem in $\mathbf{R}^{100}$

$$f(x) = c^T x - \sum_{i=1}^{500} \log(b_i - a_i^T x)$$



'Linear' convergence, i.e. a straight line on a semilog plot

# Steepest descent method

- ▶ Normalized steepest descent direction (at  $x$ , for norm  $\|\cdot\|$ ):

$$\Delta x_{\text{nsd}} = \operatorname{argmin}\{\nabla f(x)^T \nu \mid \|\nu\| = 1\}$$

Interpretation: for small  $\nu$ ,  $f(x + \nu) \approx f(x) + \nabla f(x)^T \nu$ ;  
direction  $\Delta x_{\text{nsd}}$  is unit-norm step with most negative  
directional derivative

- ▶ (unnormalized) steepest descent direction

$$\Delta x_{\text{sd}} = \|\nabla f(x)\|_* \Delta x_{\text{nsd}}$$

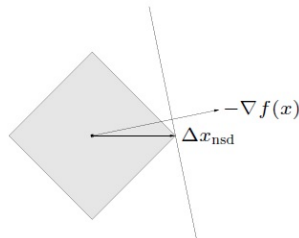
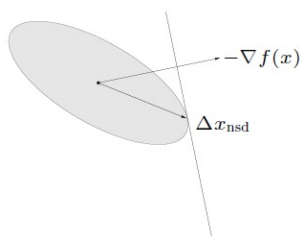
satisfies  $\nabla f(x)^T \Delta x_{\text{sd}} = -\|\nabla f(x)\|_*^2$

- ▶ Steepest decent method
  - ▶ General descent method with  $\Delta x = \Delta x_{\text{sd}}$
  - ▶ Convergence properties similar to gradient descent

# Examples

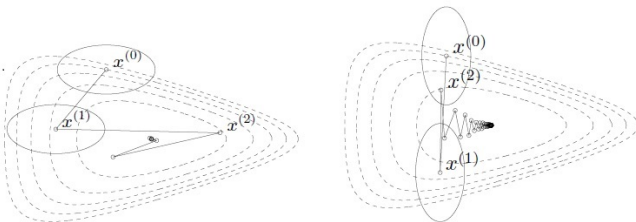
- ▶ Euclidean norm:  $\Delta x_{\text{sd}} = -\nabla f(x)$
- ▶ Quadratic norm  
 $\|x\|_P = (x^T P x)^{1/2} (P \in \mathbf{S}_{++}^n) : \Delta x_{\text{sd}} = -P^{-1} \nabla f(x)$
- ▶  $\ell_1$ -norm:  $\Delta x_{\text{sd}} = -(\partial f(x)/\partial x_i)e_i$ , where  
 $|\partial f(x)/\partial x_i| = \|\nabla f(x)\|_\infty$

Unit balls and normalized steepest descent directions for a quadratic norm and the  $\ell_1$ -norm





# Choice of norm for steepest descent



- ▶ Steepest descent with backtracking line search for two quadratic norms
- ▶ Ellipses show  $\{x \mid \|x - x^{(k)}\|_P = 1\}$
- ▶ Equivalent interpretation of steepest descent with quadratic norm  $\|\cdot\|_P$ : gradient descent after change of variables  $\bar{x} = P^{1/2}x$   
Show choice of  $P$  has strong effect on speed of convergence

# Newton step

$$\Delta x_{\text{nt}} = -\nabla^2 f(x)^{-1} \nabla f(x)$$

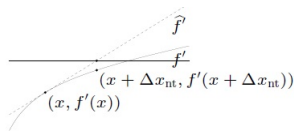
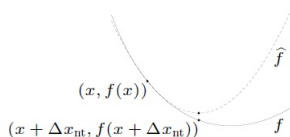
## Interpretations

- ▶  $x + \Delta x_{\text{nt}}$  minimizes second order approximation

$$\hat{f}(x + \nu) = f(x) + \nabla f(x)^T \nu + \frac{1}{2} \nu^T \nabla^2 f(x) \nu$$

- ▶  $x + \Delta x_{\text{nt}}$  solves linearized optimality condition

$$\nabla f(x + \nu) \approx \nabla \hat{f}(x + \nu) = \nabla f(x) + \nabla^2 f(x) \nu = 0$$



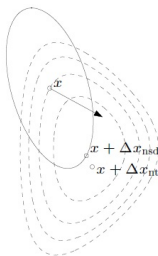
## Newton step

- ▶  $\Delta x_{\text{nt}}$  is independent of linear (or affine) changes of coordinates. It is affine invariant

$$x = Ty, \quad x + \Delta x_{\text{nt}} = T(y + \Delta y_{\text{nt}})$$

- ▶  $\Delta x_{\text{nt}}$  is steepest descent direction at  $x$  in local Hessian norm

$$\|u\|_{\nabla^2 f(x)} = (u^T \nabla^2 f(x) u)^{1/2}$$



Dash lines are contour lines of  $f$ ; ellipse is  $\{x + \nu \mid \nu^T \nabla^2 f(x) \nu = 1\}$  arrow shows  $-\nabla f(x)$

# Newton decrement

$$\lambda(x) = (\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x))^{1/2}$$

a measure of the proximity of  $x$  to  $x^*$

Properties

- ▶ Gives an estimate of  $f(x) - p^*$ , using quadratic approximation  $\hat{f}$ :

$$f(x) - \inf_y \hat{f}(y) = f(x) - \hat{f}(x + \Delta x_{\text{nt}}) = \frac{1}{2} \lambda(x)^2$$

- ▶ Equal to the norm of the Newton step in the quadratic Hessian norm

$$\lambda(x) = (\Delta x_{\text{nt}}^T \nabla^2 f(x) \Delta x_{\text{nt}})^{1/2}$$

- ▶ Directional derivative in the Newton direction,  $\nabla f(x)^T \Delta x_{\text{nt}} = -\lambda(x)^2$ , can be used in the backtracking line search
- ▶ Affine invariant (unlike  $\|\nabla f(x)\|_2$ )

# Newton's method

- ▶ Given a starting point  $x \in \text{dom } f$ , tolerance  $\epsilon > 0$
- ▶ Repeat

1. Compute the Newton step and decrement

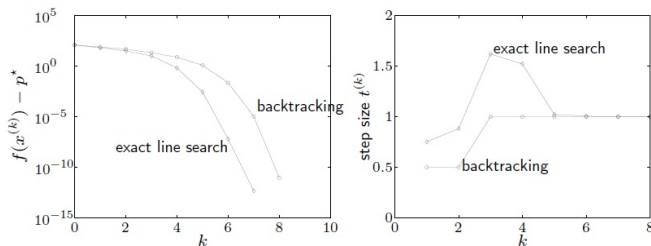
$$\Delta x_{\text{nt}} := -\nabla^2 f(x)^{-1} \nabla f(x); \quad \lambda^2 := \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x)$$

2. Stopping criterion. Quit if  $\lambda^2/2 \leq \epsilon$
3. Line search. Choose step size  $t$  by backtracking line search.
4. Update.  $x := x + t\Delta x_{\text{nt}}$

Advantages of Newton method: fast, robust, scalable

## Example in $\mathbf{R}^{100}$

$$f(x) = c^T x - \sum_{i=1}^{500} \log(b_i - a_i^T x)$$



- ▶ Backtracking parameters  $\alpha = 0.01, \beta = 0.5$
- ▶ Backtracking line search almost as fast as exact l.s. (and much simpler)
- ▶ Clearly shows two phases in algorithm

# Summary about Newton's method

- ▶ Convergence of Newton's method is rapid in general, and quadratic near  $x^*$
- ▶ Newton's method is affine invariant
- ▶ Newton's method scales well with problem size
- ▶ The good performance of Newton's method is not dependent on the choice of algorithm parameters
- ▶ The main disadvantage of Newton's method is the cost of forming and storing the Hessian ( $\nabla^2 f(x)$ ), and the cost of computing the Newton step