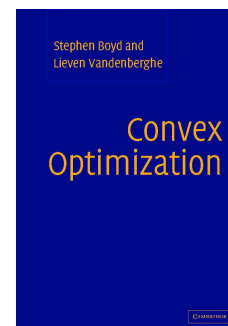


# Convex Optimization & Machine Learning

## Algorithms

[mcuturi@i.kyoto-u.ac.jp](mailto:mcuturi@i.kyoto-u.ac.jp)



Most slides in this lecture are taken from

---

# Unconstrained Convex Optimization Algorithms

- terminology and assumptions
- gradient descent method
- steepest descent method
- Newton's method
- self-concordant functions
- implementation

# Unconstrained minimization

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

- $f$  convex, twice continuously differentiable (hence  $\text{dom } f$  open)
- we assume optimal value  $p^* = \inf_x f(x)$  is attained (and finite)

## unconstrained minimization methods

- produce sequence of points  $x^{(k)} \in \text{dom } f$ ,  $k = 0, 1, \dots$  with

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

- can be interpreted as iterative methods for solving optimality condition

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

# Initial point and sublevel set

algorithms in this lecture require a starting point  $x^{(0)}$  such that

- $x^{(0)} \in \mathbf{dom} f$
- sublevel set  $S = \{x \mid f(x) \leq f(x^{(0)})\}$  is closed

2nd condition is hard to verify, except when *all* sublevel sets are closed:

- equivalent to condition that  $\mathbf{epi} f$  is closed
- true if  $\mathbf{dom} f = \mathbf{R}^n$
- true if  $f(x) \rightarrow \infty$  as  $x \rightarrow \mathbf{d} \mathbf{dom} f$

examples of differentiable functions with closed sublevel sets:

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

# Strong convexity and implications

$f$  is strongly convex on  $S$  if there exists an  $m > 0$  such that

$$\nabla^2 f(x) \succeq mI \quad \text{for all } x \in S$$

## implications

- for  $x, y \in S$ ,

$$f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{m}{2} \|x - y\|_2^2$$

hence,  $S$  is bounded

- $p^* > -\infty$ , and for  $x \in S$ ,

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

useful as stopping criterion (if you know  $m$ )

## Strong convexity: Proof of $f(x) - p^* \leq \frac{1}{2m} \|\nabla f(x)\|_2^2$

- By convexity  $f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{m}{2} \|x - y\|_2^2$
- the RHS is a quadratic form in  $\mathbf{u} \stackrel{\text{def}}{=} y - x$ :  $\frac{m}{2} \mathbf{u}^T \mathbf{u} + \nabla f(x)^T \mathbf{u} + f(x)$
- Recall that a quadratic form  $Q(\mathbf{u}) = \mathbf{u}^T P \mathbf{u} + \mathbf{q}^T \mathbf{u} + c$ 
  - has gradient  $\nabla Q(\mathbf{u}) = 2P\mathbf{u} + \mathbf{q}$
  - is thus minimal for  $2P\mathbf{u} = -\mathbf{q}$ , that is  $\mathbf{u} = -\frac{1}{2}P^{-1}\mathbf{q}$
- In this case, the r.h.s is thus minimal for  $y - x = -\frac{1}{2} \times \frac{2}{m} \nabla f(x) = -\frac{1}{m} \nabla f(x)$
- We obtain the bound

$$f(y) \geq f(x) + \nabla f(x)^T \left( -\frac{1}{m} \nabla f(x) \right) + \frac{m}{2} \times \left\| -\frac{1}{m} \nabla f(x) \right\|_2^2$$

$$f(y) \geq f(x) - \frac{1}{2m} \|\nabla f(x)\|_2^2$$

- This bound is valid for any couple of points  $(x, y)$ , in particular if  $f(y) = p^*$ :

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

---

# Descent methods

# Descent methods

$$x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x^{(k)} \quad \text{with } f(x^{(k+1)}) < f(x^{(k)})$$

- lighter notations:  $x^+ = x + t\Delta x$  ;  $x := x + t\Delta x$  (insist on iterative process)
- $\Delta x$  is the *step*, or *search direction*. **Can be any vector**
- $t$  is the *step size*, or *step length* which scales the step.



# Descent methods

A **descent** method means that  $f(x^{(k+1)}) < f(x^{(k)})$

- from convexity, we have that  $f(y) \geq f(x) + \nabla f(x)^T (y - x)$ .
- Using this inequality with  $y = x + t\Delta x$ , we get

$$f(y = x + t\Delta x) \geq f(x) + t\nabla f(x)^T (\Delta x)$$

- If we need  $f(x + t\Delta x) < f(x)$  then **necessarily**

$$\nabla f(x)^T \Delta x < 0.$$

(i.e.,  $\Delta x$  is a descent direction)

# General Descent methods

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

**repeat**

1. Determine a descent direction  $\Delta x$ .
2. **Line search**: Choose a step size  $t > 0$ .
3. **Update**:  $x := x + t\Delta x$ .

**until** stopping criterion is satisfied.

# Line search types

- **exact line search**: set  $t$  with the rule

$$t = \operatorname{argmin}_{u>0} f(x + u\Delta x)$$

- each gradient step involves another optimization problem!
- only one variable, but **usually** too costly to solve exactly.
- in most cases, better to look for just a “*good enough*” step.

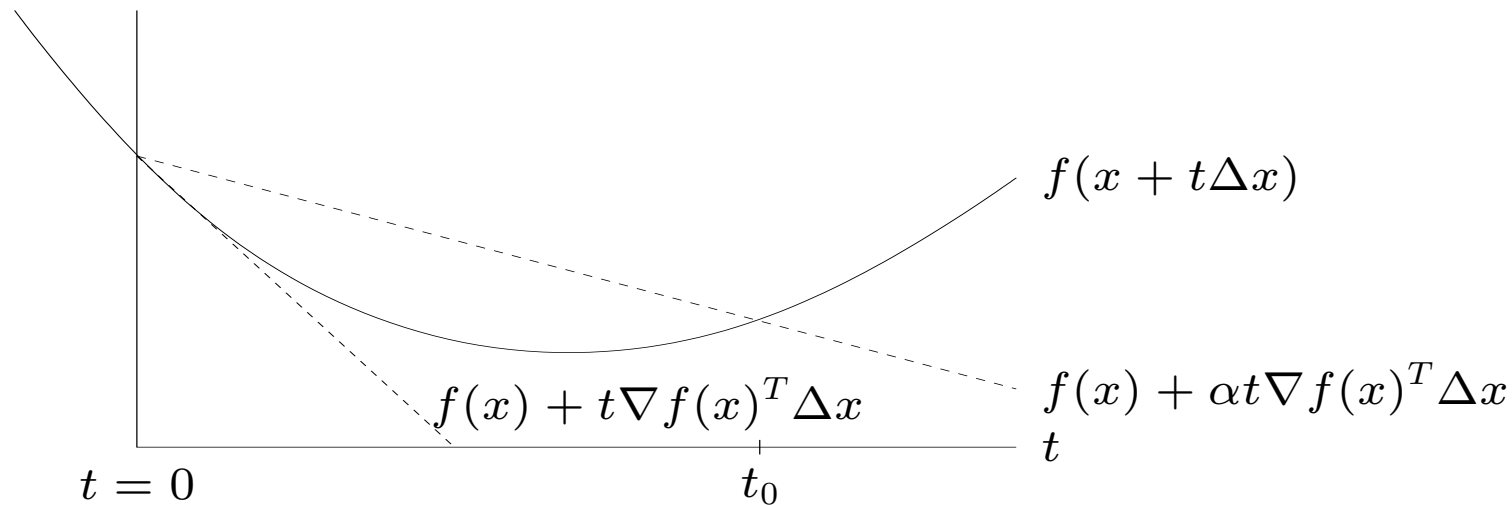
# Line search types

- **backtracking line search**

- **two** 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

general descent method with  $\Delta x = -\nabla f(x)$

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

**repeat**

1.  $\Delta x := -\nabla f(x)$ .
2. *Line search*. Choose step size  $t$  via exact or backtracking line search.
3. *Update*.  $x := x + t\Delta x$ .

**until** stopping criterion is satisfied.

- stopping criterion usually of the form  $\|\nabla f(x)\|_2 \leq \epsilon$
- convergence result: for strongly 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

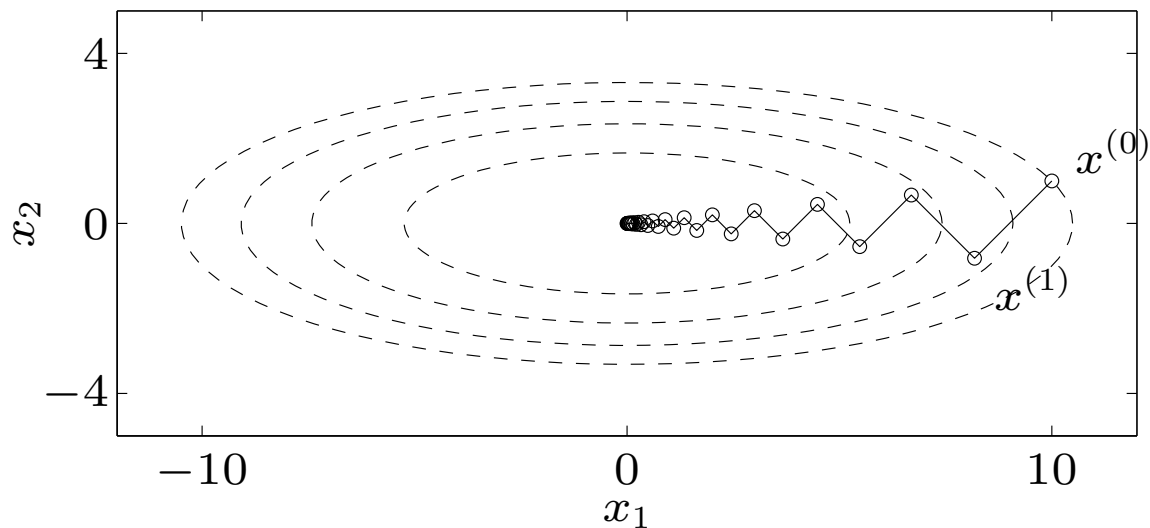
## quadratic problem in $\mathbf{R}^2$

$$f(x) = (1/2)(x_1^2 + \gamma x_2^2) \quad (\gamma > 0)$$

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

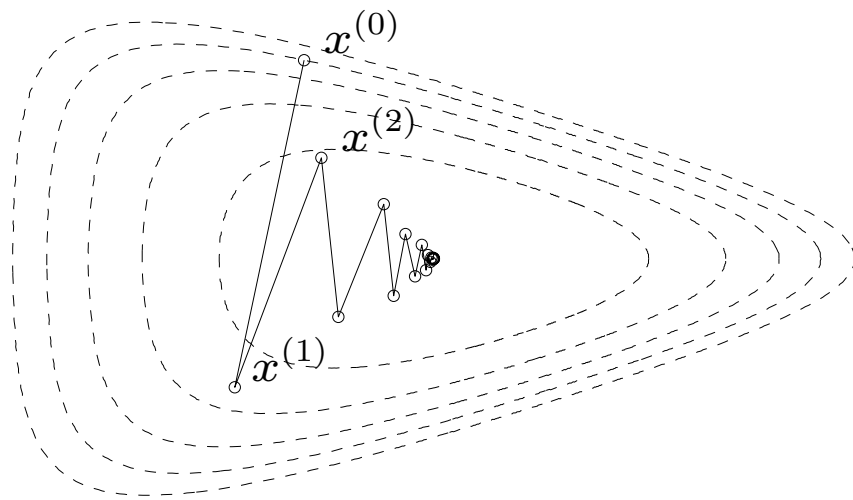
$$x_1^{(k)} = \gamma \left( \frac{\gamma - 1}{\gamma + 1} \right)^k, \quad x_2^{(k)} = \left( -\frac{\gamma - 1}{\gamma + 1} \right)^k$$

- very slow if  $\gamma \gg 1$  or  $\gamma \ll 1$
- example for  $\gamma = 10$ :

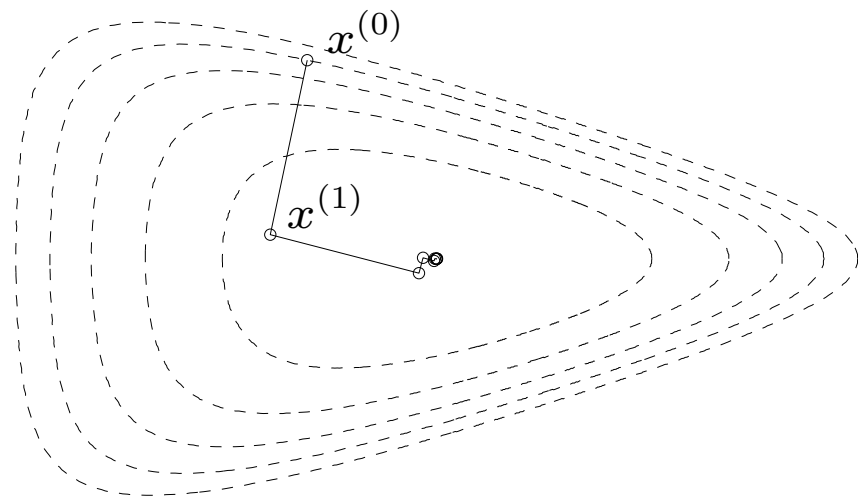


## nonquadratic example

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



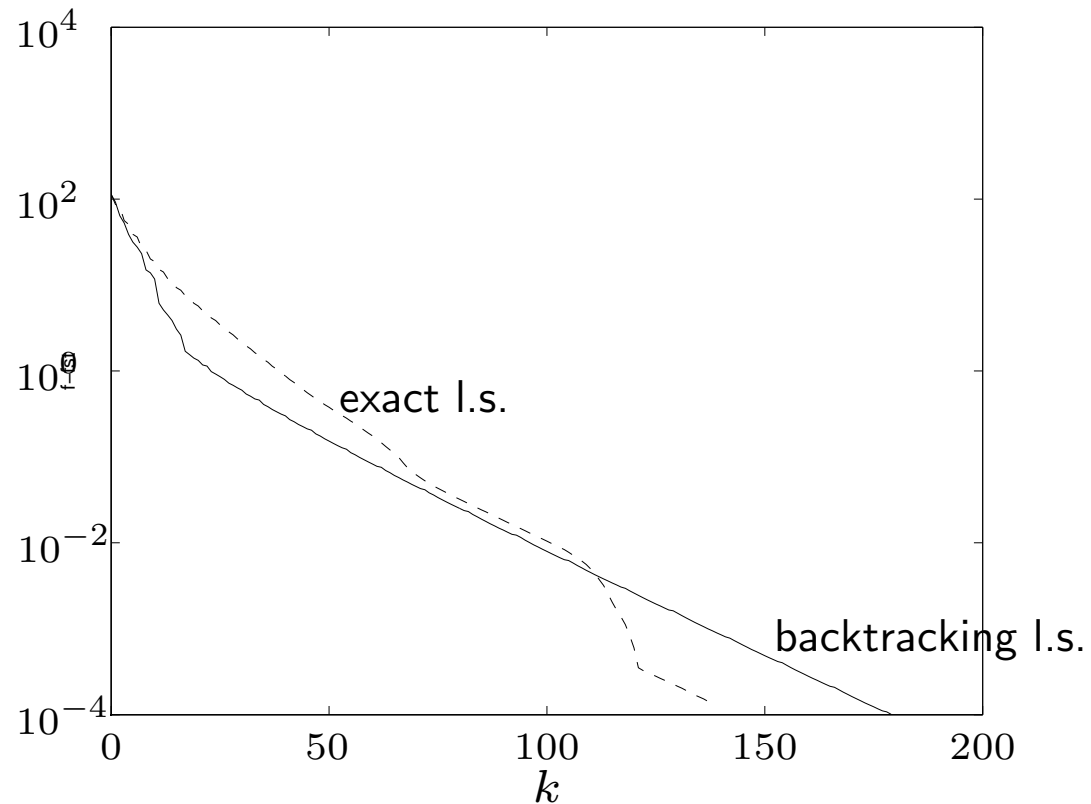
backtracking line search



exact line search

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  $\|\cdot\|$ ):

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

interpretation: for small  $v$ ,

$$f(x + v) \approx f(x) + \nabla f(x)^T v$$

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

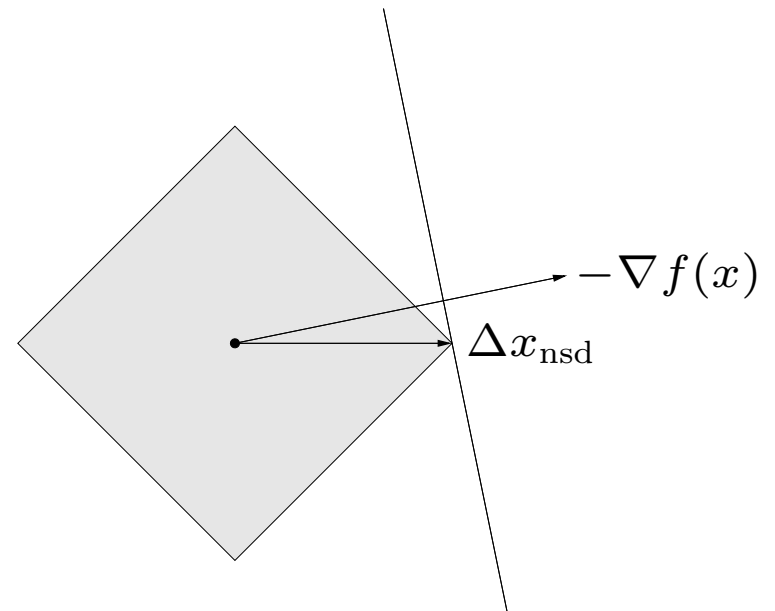
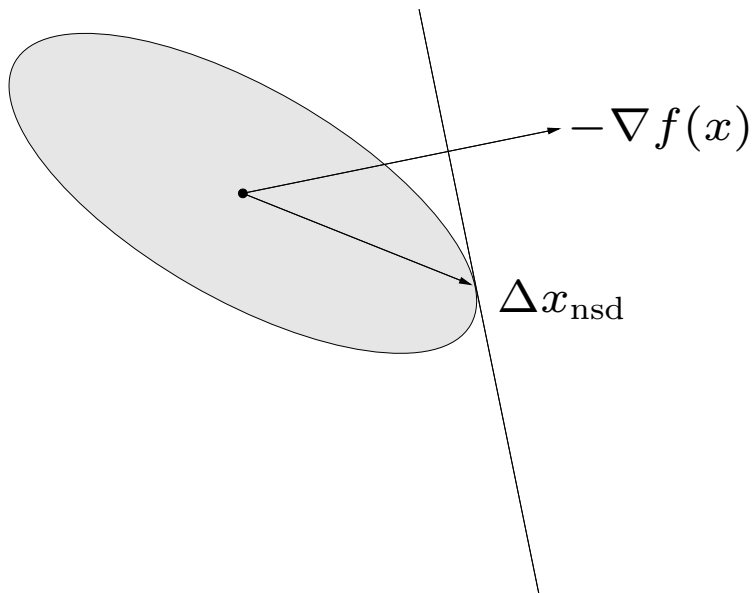
**steepest descent method**

- general descent method with  $\Delta x = \Delta x_{\text{sd}}$
- convergence properties similar to gradient descent

## examples

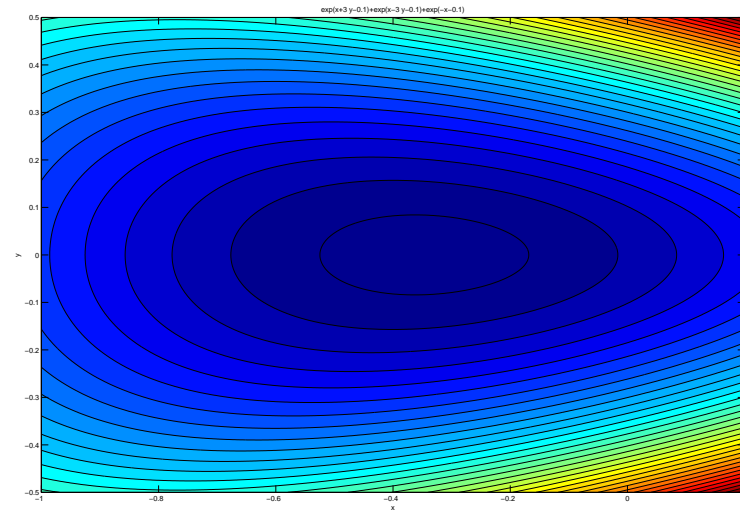
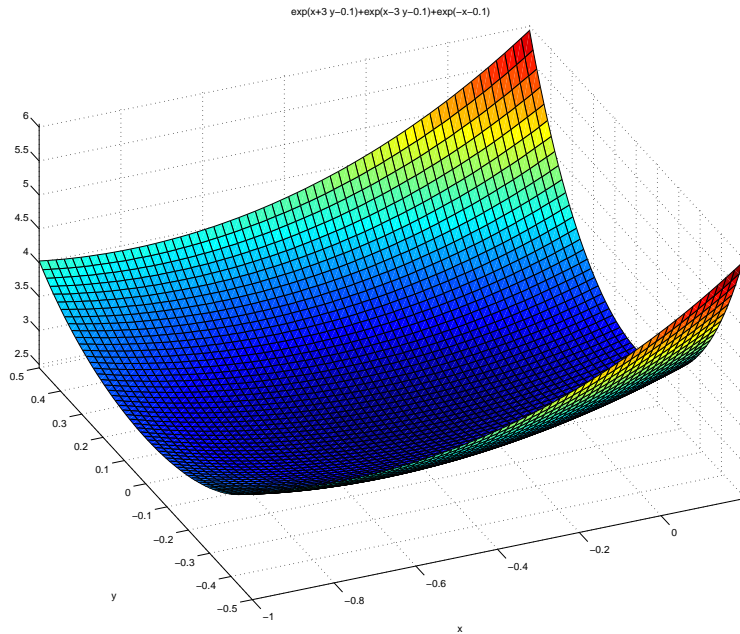
- Euclidean norm:  $\Delta x_{\text{sd}} = -\nabla f(x)$
- quadratic norm  $\|x\|_P = (x^T P x)^{1/2}$  ( $P \in \mathbf{S}_{++}^n$ ):  $\Delta x_{\text{sd}} = -P^{-1} \nabla f(x)$
- $\ell_1$ -norm:  $\Delta x_{\text{sd}} = -(\partial f(x)/\partial x_i)e_i$ , where  $|\partial f(x)/\partial x_i| = \|\nabla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

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



# choice of norm for steepest descent

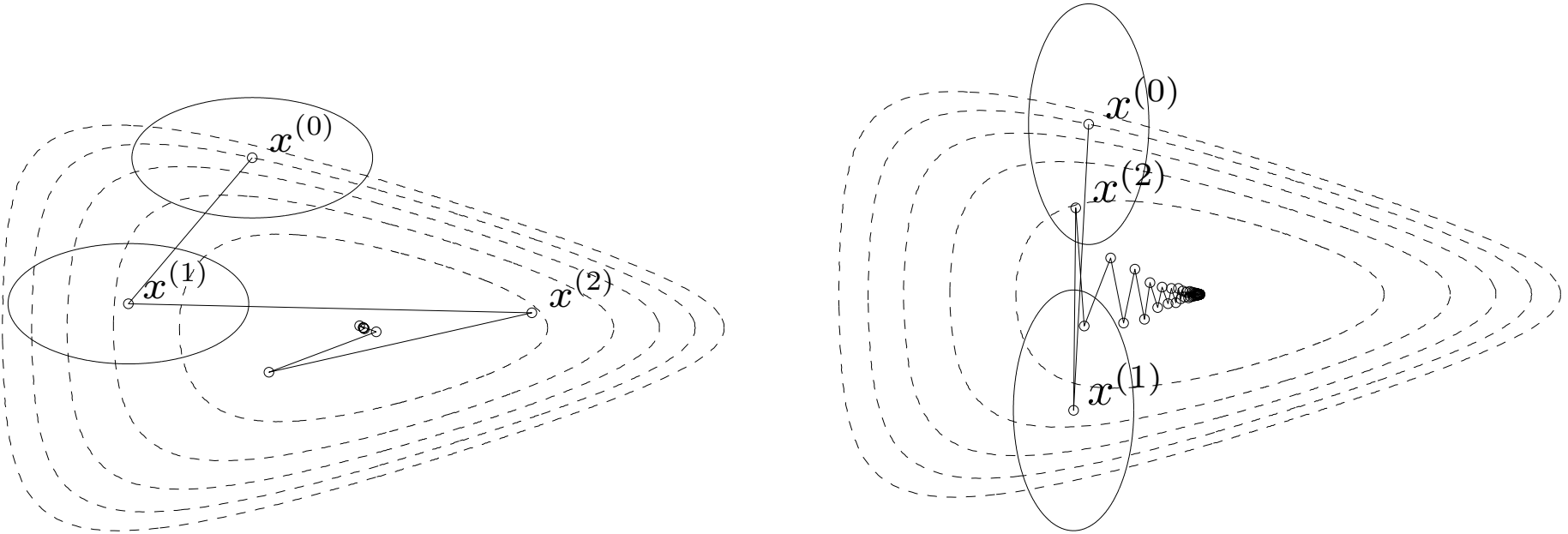
Consider again the function  $f(x_1, x_2) = e^{x_1+3x_2-0.1} + e^{x_1-3x_2-0.1} + e^{-x_1-0.1}$



- We consider two different elliptic norms to define the nsd:

$$P_1 = \begin{bmatrix} 2 & 0 \\ 0 & 8 \end{bmatrix}, P_2 = \begin{bmatrix} 8 & 0 \\ 0 & 2 \end{bmatrix}$$

# choice of norm for steepest descent



- steepest descent with backtracking line search, left  $\|\cdot\|_{P_1}$ , right  $\|\cdot\|_{P_2}$ ,
- ellipses show  $\{x \mid \|x - x^{(k)}\|_P = 1\}$

steepest descent with quadratic norm  $\|\cdot\|_P$



simple gradient descent with  $\|\cdot\|_2$  after change of variables  $\bar{x} = P^{1/2}x$

- shows choice of  $P$  has strong effect on speed of convergence

# Newton step (2nd Order)

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

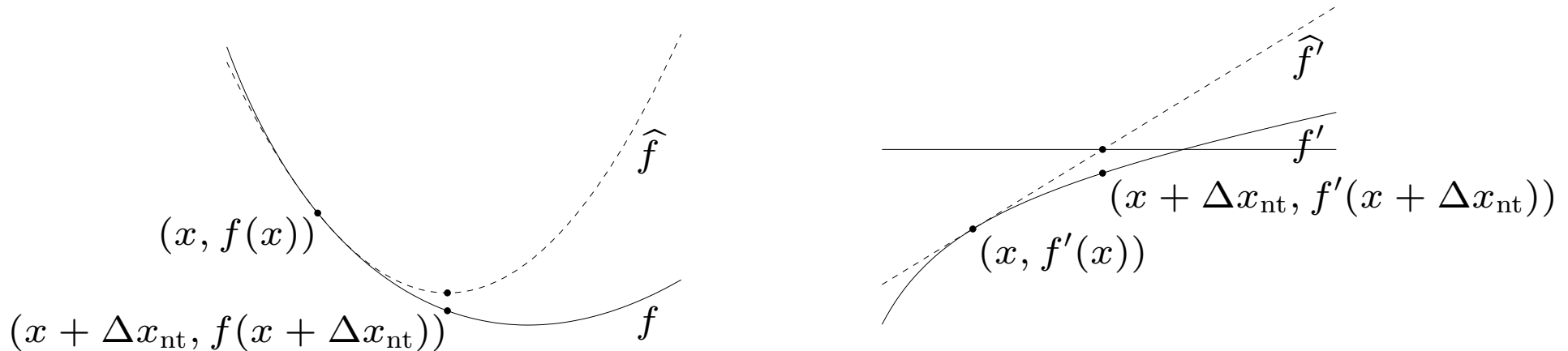
## interpretations

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

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

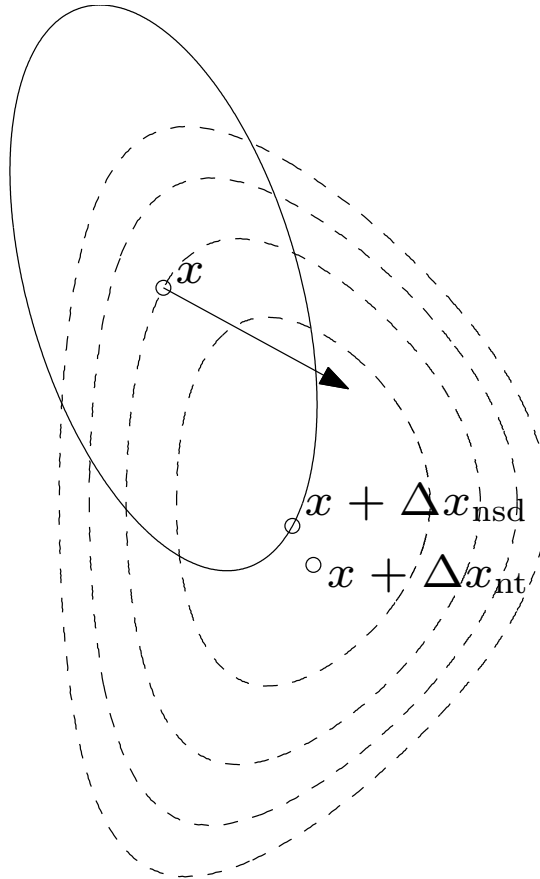
- $x + \Delta x_{\text{nt}}$  solves linearized optimality condition

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



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

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



dashed lines are contour lines of  $f$ ; ellipse is  $\{x + v \mid v^T \nabla^2 f(x) v = 1\}$

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  $\hat{f}$ :

$$f(x) - \inf_y \hat{f}(y) = \frac{1}{2} \lambda(x)^2$$

- 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:  $\nabla f(x)^T \Delta x_{\text{nt}} = -\lambda(x)^2$
- 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.*  $x := x + t\Delta x_{\text{nt}}$ .

affine invariant, *i.e.*, independent of linear changes of coordinates:

Newton iterates for  $\tilde{f}(y) = f(Ty)$  with starting point  $y^{(0)} = T^{-1}x^{(0)}$  are

$$y^{(k)} = T^{-1}x^{(k)}$$



# Classical convergence analysis

## assumptions

- $f$  strongly convex on  $S$  with constant  $m$
- $\nabla^2 f$  is Lipschitz continuous on  $S$ , with constant  $L > 0$ :

$$\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 \leq L\|x - y\|_2$$

( $L$  measures how well  $f$  can be approximated by a quadratic function)

**outline:** there exist constants  $\eta \in (0, m^2/L)$ ,  $\gamma > 0$  such that

- if  $\|\nabla f(x)\|_2 \geq \eta$ , then  $f(x^{(k+1)}) - f(x^{(k)}) \leq -\gamma$
- if  $\|\nabla f(x)\|_2 < \eta$ , then

$$\frac{L}{2m^2} \|\nabla f(x^{(k+1)})\|_2 \leq \left( \frac{L}{2m^2} \|\nabla f(x^{(k)})\|_2 \right)^2$$

## damped Newton phase ( $\|\nabla f(x)\|_2 \geq \eta$ )

- most iterations require backtracking steps
- function value decreases by at least  $\gamma$
- if  $p^* > -\infty$ , this phase ends after at most  $(f(x^{(0)}) - p^*)/\gamma$  iterations

## quadratically convergent phase ( $\|\nabla f(x)\|_2 < \eta$ )

- all iterations use step size  $t = 1$
- $\|\nabla f(x)\|_2$  converges to zero quadratically: if  $\|\nabla f(x^{(k)})\|_2 < \eta$ , then

$$\frac{L}{2m^2} \|\nabla f(x^l)\|_2 \leq \left( \frac{L}{2m^2} \|\nabla f(x^k)\|_2 \right)^{2^{l-k}} \leq \left( \frac{1}{2} \right)^{2^{l-k}}, \quad l \geq k$$

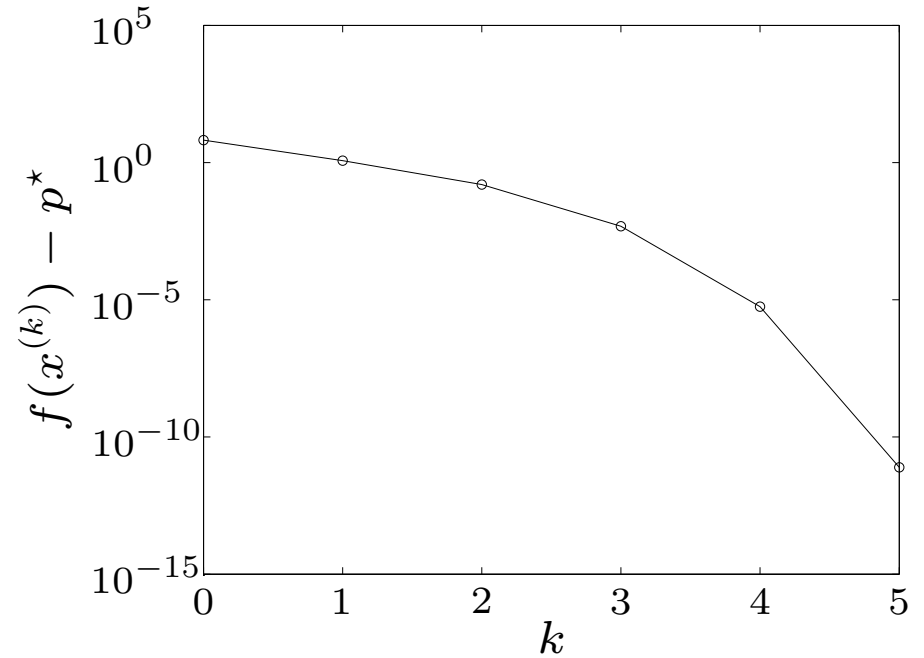
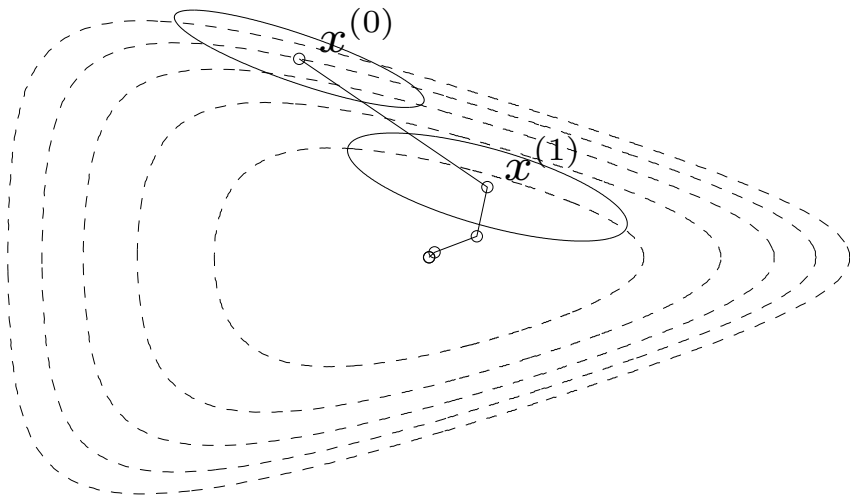
**conclusion:** number of iterations until  $f(x) - p^* \leq \epsilon$  is bounded above by

$$\frac{f(x^{(0)}) - p^*}{\gamma} + \log_2 \log_2(\epsilon_0/\epsilon)$$

- $\gamma, \epsilon_0$  are constants that depend on  $m, L, x^{(0)}$
- second term is small (of the order of 6) and almost constant for practical purposes
- in practice, constants  $m, L$  (hence  $\gamma, \epsilon_0$ ) are usually unknown
- provides qualitative insight in convergence properties (*i.e.*, explains two algorithm phases)

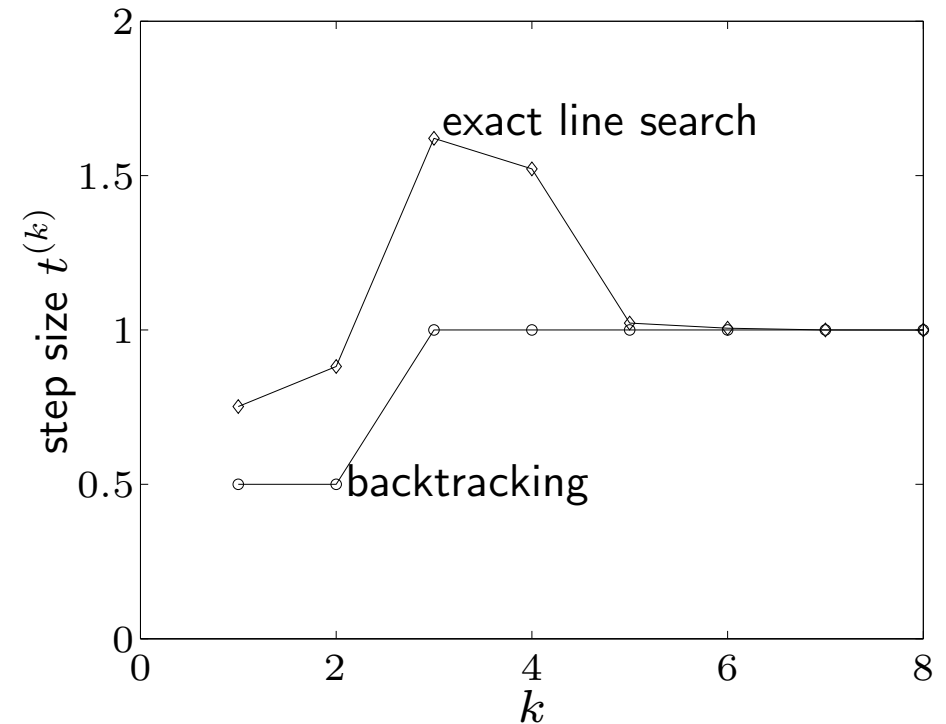
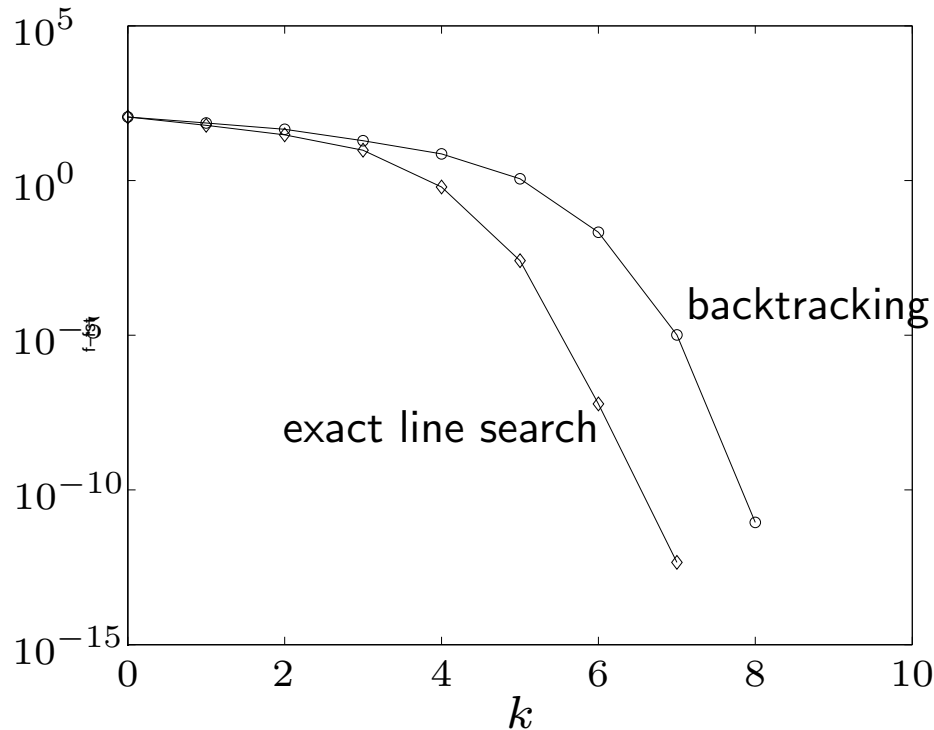
# Examples

## example in $\mathbb{R}^2$



- backtracking parameters  $\alpha = 0.1$ ,  $\beta = 0.7$
- converges in only 5 steps
- quadratic local convergence

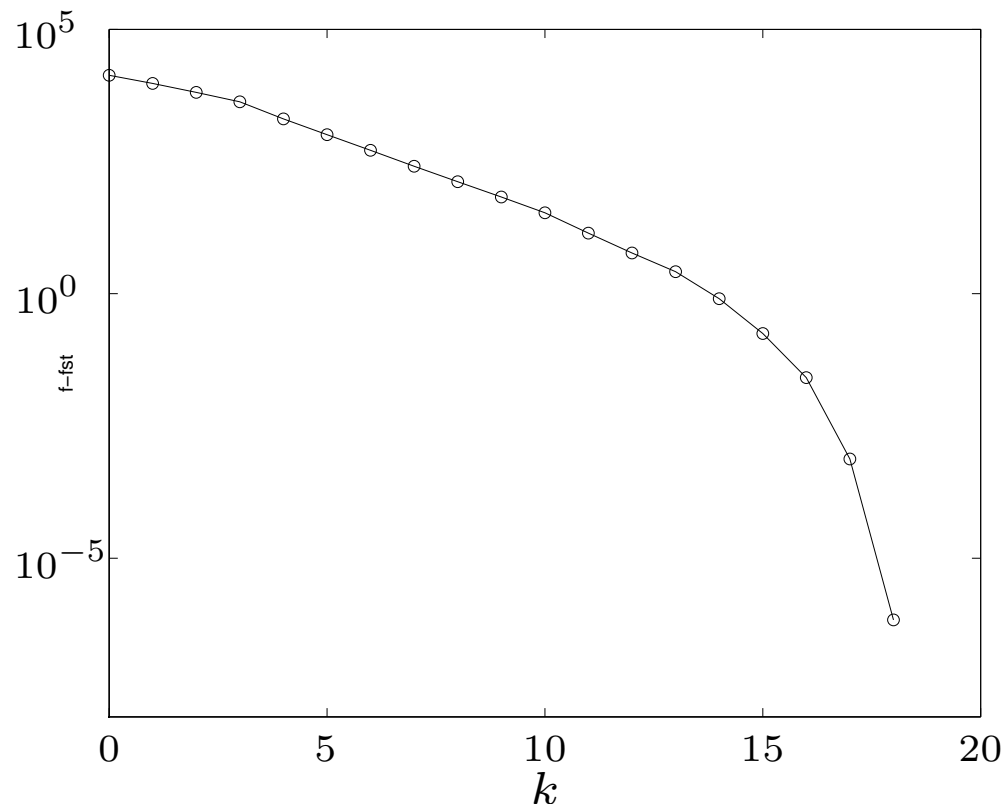
## example in $\mathbb{R}^{100}$ (page 16)



- 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

example in  $\mathbf{R}^{10000}$  (with sparse  $a_i$ )

$$f(x) = - \sum_{i=1}^{10000} \log(1 - x_i^2) - \sum_{i=1}^{100000} \log(b_i - a_i^T x)$$



- backtracking parameters  $\alpha = 0.01$ ,  $\beta = 0.5$ .
- performance similar as for small examples

# A few words on Self-concordance

## shortcomings of classical convergence analysis

- depends on unknown constants ( $m, L, \dots$ )
- bound is not affinely invariant, although Newton's method is

## convergence analysis via self-concordance (Nesterov and Nemirovski)

- does not depend on any unknown constants
- gives affine-invariant bound
- applies to special class of convex functions ('self-concordant' functions)
- developed to analyze polynomial-time interior-point methods for convex optimization
- Please check Boyd & Vandenberghe book for a review!

# Implementation

main effort in each iteration: evaluate derivatives and solve Newton system

$$H \Delta x = g$$

where  $H = \nabla^2 f(x)$ ,  $g = -\nabla f(x)$

**via Cholesky factorization**

$$H = LL^T, \quad \Delta x_{\text{nt}} = L^{-T} L^{-1} g, \quad \lambda(x) = \|L^{-1} g\|_2$$

- cost  $(1/3)n^3$  flops for unstructured system
- cost  $\ll (1/3)n^3$  if  $H$  sparse, banded



## example of dense Newton system with structure

$$f(x) = \sum_{i=1}^n \psi_i(x_i) + \psi_0(Ax + b), \quad H = D + A^T H_0 A$$

- assume  $A \in \mathbf{R}^{p \times n}$ , dense, with  $p \ll n$
- $D$  diagonal with diagonal elements  $\psi_i''(x_i)$ ;  $H_0 = \nabla^2 \psi_0(Ax + b)$

**method 1:** form  $H$ , solve via dense Cholesky factorization: (cost  $(1/3)n^3$ )

**method 2:** factor  $H_0 = L_0 L_0^T$ ; write Newton system as

$$D\Delta x + A^T L_0 w = -g, \quad L_0^T A \Delta x - w = 0$$

eliminate  $\Delta x$  from first equation; compute  $w$  and  $\Delta x$  from

$$(I + L_0^T A D^{-1} A^T L_0)w = -L_0^T A D^{-1} g, \quad D\Delta x = -g - A^T L_0 w$$

cost:  $2p^2 n$  (dominated by computation of  $L_0^T A D^{-1} A L_0$ )

---

# Convex Optimization Algorithms With Equality Constraints

- equality constrained minimization
- Newton's method with equality constraints
- infeasible start Newton method
- implementation

# Equality constrained minimization

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{subject to} & Ax = b \end{array}$$

- $f$  convex, twice continuously differentiable
- $A \in \mathbf{R}^{p \times n}$  with  $\mathbf{Rank} A = p$
- we assume  $p^*$  is finite and attained

**optimality conditions:**  $x^*$  is optimal iff there exists a  $\nu^*$  such that

$$\nabla f(x^*) + A^T \nu^* = 0, \quad Ax^* = b$$

## equality constrained quadratic minimization (with $P \in \mathbf{S}_+^n$ )

$$\begin{aligned} & \text{minimize} && (1/2)x^T P x + q^T x + r \\ & \text{subject to} && Ax = b \end{aligned}$$

optimality condition:

$$\begin{bmatrix} P & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} x^* \\ \nu^* \end{bmatrix} = \begin{bmatrix} -q \\ b \end{bmatrix}$$

- coefficient matrix is called KKT matrix
- KKT matrix is nonsingular if and only if

$$Ax = 0, \quad x \neq 0 \quad \implies \quad x^T P x > 0$$

- equivalent condition for nonsingularity:  $P + A^T A \succ 0$

# Newton step

Newton step of  $f$  at feasible  $x$  is given by (1st block) of solution of

$$\begin{bmatrix} \nabla^2 f(x) & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} \Delta x_{\text{nt}} \\ w \end{bmatrix} = \begin{bmatrix} -\nabla f(x) \\ 0 \end{bmatrix}$$

## interpretations

- $\Delta x_{\text{nt}}$  solves second order approximation (with variable  $v$ )

$$\begin{array}{ll} \text{minimize} & \hat{f}(x + v) = f(x) + \nabla f(x)^T v + (1/2)v^T \nabla^2 f(x)v \\ \text{subject to} & A(x + v) = b \end{array}$$

- equations follow from linearizing optimality conditions

$$\nabla f(x + \Delta x_{\text{nt}}) + A^T w = 0, \quad A(x + \Delta x_{\text{nt}}) = b$$

# Newton decrement

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

## properties

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

$$f(x) - \inf_{Ay=b} \hat{f}(y) = \frac{1}{2} \lambda(x)^2$$

- directional derivative in Newton direction:

$$\left. \frac{d}{dt} f(x + t \Delta x_{\text{nt}}) \right|_{t=0} = -\lambda(x)^2$$

- in general,  $\lambda(x) \neq (\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x))^{1/2}$

# Newton's method with equality constraints

**given** starting point  $x \in \text{dom } f$  with  $Ax = b$ , tolerance  $\epsilon > 0$ .

**repeat**

1. Compute the Newton step and decrement  $\Delta x_{\text{nt}}, \lambda(x)$ .
2. *Stopping criterion.* **quit** if  $\lambda^2/2 \leq \epsilon$ .
3. *Line search.* Choose step size  $t$  by backtracking line search.
4. *Update.*  $x := x + t\Delta x_{\text{nt}}$ .

- a feasible descent method:  $x^{(k)}$  feasible and  $f(x^{(k+1)}) < f(x^{(k)})$
- affine invariant

# Newton step at infeasible points

extends to infeasible  $x$  (*i.e.*,  $Ax \neq b$ )

linearizing optimality conditions at infeasible  $x$  (with  $x \in \text{dom } f$ ) gives

$$\begin{bmatrix} \nabla^2 f(x) & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} \Delta x_{\text{nt}} \\ w \end{bmatrix} = - \begin{bmatrix} \nabla f(x) \\ Ax - b \end{bmatrix} \quad (1)$$

## primal-dual interpretation

- write optimality condition as  $r(y) = 0$ , where

$$y = (x, \nu), \quad r(y) = (\nabla f(x) + A^T \nu, Ax - b)$$

- linearizing  $r(y) = 0$  gives  $r(y + \Delta y) \approx r(y) + Dr(y)\Delta y = 0$ :

$$\begin{bmatrix} \nabla^2 f(x) & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} \Delta x_{\text{nt}} \\ \Delta \nu_{\text{nt}} \end{bmatrix} = - \begin{bmatrix} \nabla f(x) + A^T \nu \\ Ax - b \end{bmatrix}$$

same as (1) with  $w = \nu + \Delta \nu_{\text{nt}}$



## Infeasible start Newton method

**given** starting point  $x \in \text{dom } f$ ,  $\nu$ , tolerance  $\epsilon > 0$ ,  $\alpha \in (0, 1/2)$ ,  $\beta \in (0, 1)$ .

**repeat**

1. Compute primal and dual Newton steps  $\Delta x_{\text{nt}}$ ,  $\Delta \nu_{\text{nt}}$ .

2. *Backtracking line search* on  $\|r\|_2$ .

$t := 1$ .

**while**  $\|r(x + t\Delta x_{\text{nt}}, \nu + t\Delta \nu_{\text{nt}})\|_2 > (1 - \alpha t)\|r(x, \nu)\|_2$ ,  $t := \beta t$ .

3. *Update*.  $x := x + t\Delta x_{\text{nt}}$ ,  $\nu := \nu + t\Delta \nu_{\text{nt}}$ .

**until**  $Ax = b$  and  $\|r(x, \nu)\|_2 \leq \epsilon$ .

- not a descent method:  $f(x^{(k+1)}) > f(x^{(k)})$  is possible
- directional derivative of  $\|r(y)\|_2^2$  in direction  $\Delta y = (\Delta x_{\text{nt}}, \Delta \nu_{\text{nt}})$  is

$$\left. \frac{d}{dt} \|r(y + \Delta y)\|_2^2 \right|_{t=0} = -\|r(y)\|_2^2$$

# Solving KKT systems

$$\begin{bmatrix} H & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} v \\ w \end{bmatrix} = - \begin{bmatrix} g \\ h \end{bmatrix}$$

## solution methods

- LDL<sup>T</sup> factorization
- elimination (if  $H$  nonsingular)

$$AH^{-1}A^T w = h - AH^{-1}g, \quad Hv = -(g + A^T w)$$

- elimination with singular  $H$ : write as

$$\begin{bmatrix} H + A^T Q A & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} v \\ w \end{bmatrix} = - \begin{bmatrix} g + A^T Q h \\ h \end{bmatrix}$$

with  $Q \succeq 0$  for which  $H + A^T Q A \succ 0$ , and apply elimination

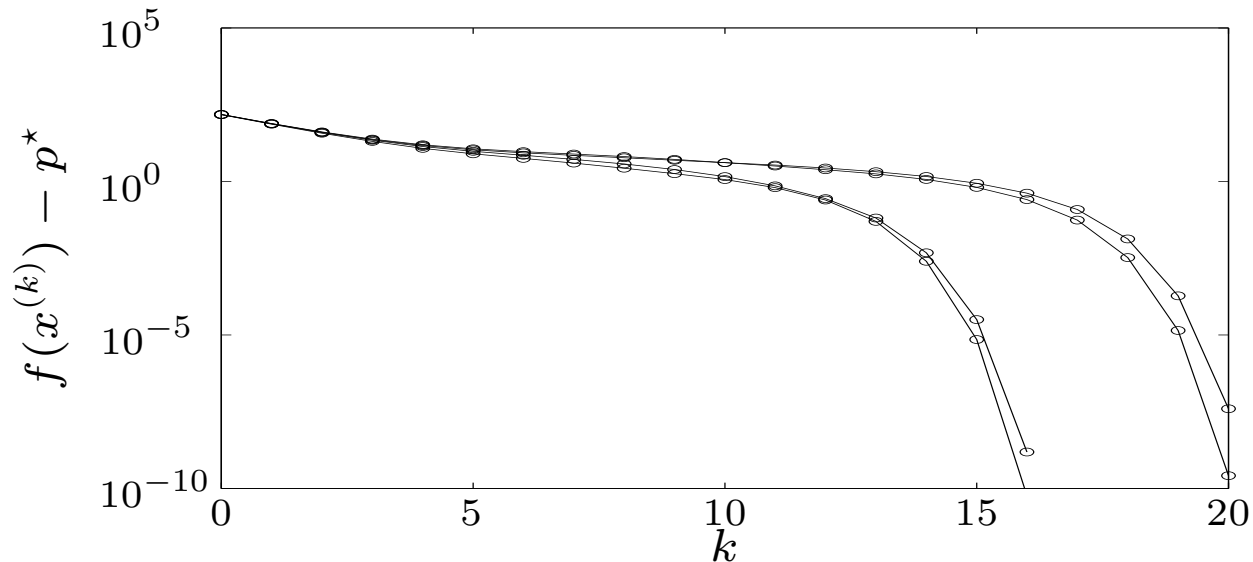
# Equality constrained analytic centering

**primal problem:** minimize  $-\sum_{i=1}^n \log x_i$  subject to  $Ax = b$

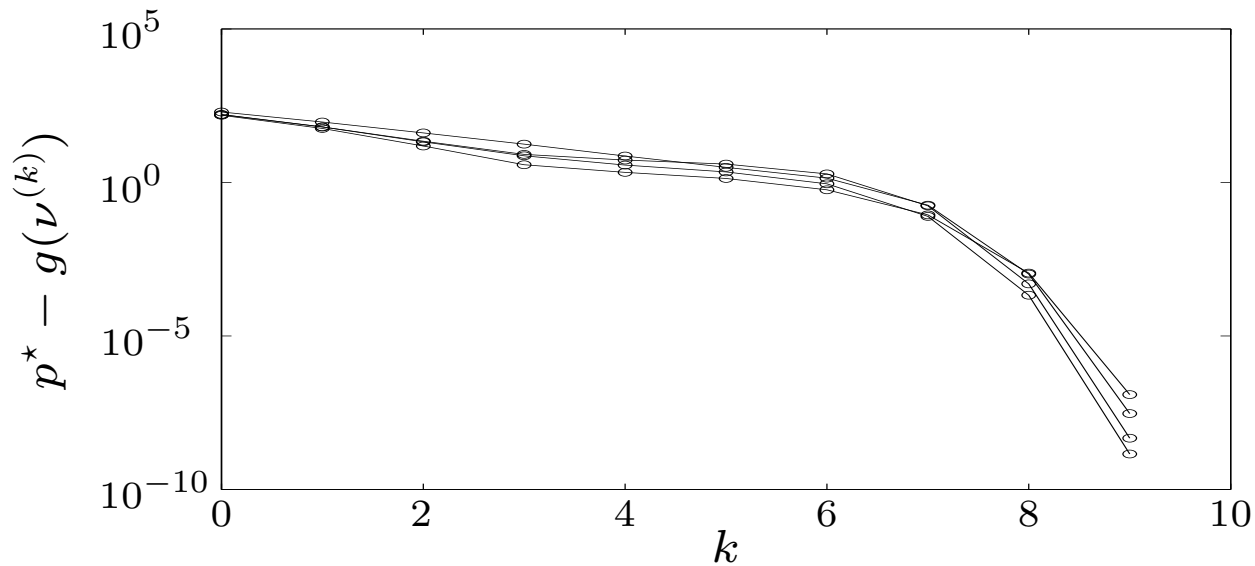
**dual problem:** maximize  $-b^T \nu + \sum_{i=1}^n \log(A^T \nu)_i + n$

three methods for an example with  $A \in \mathbf{R}^{100 \times 500}$ , different starting points

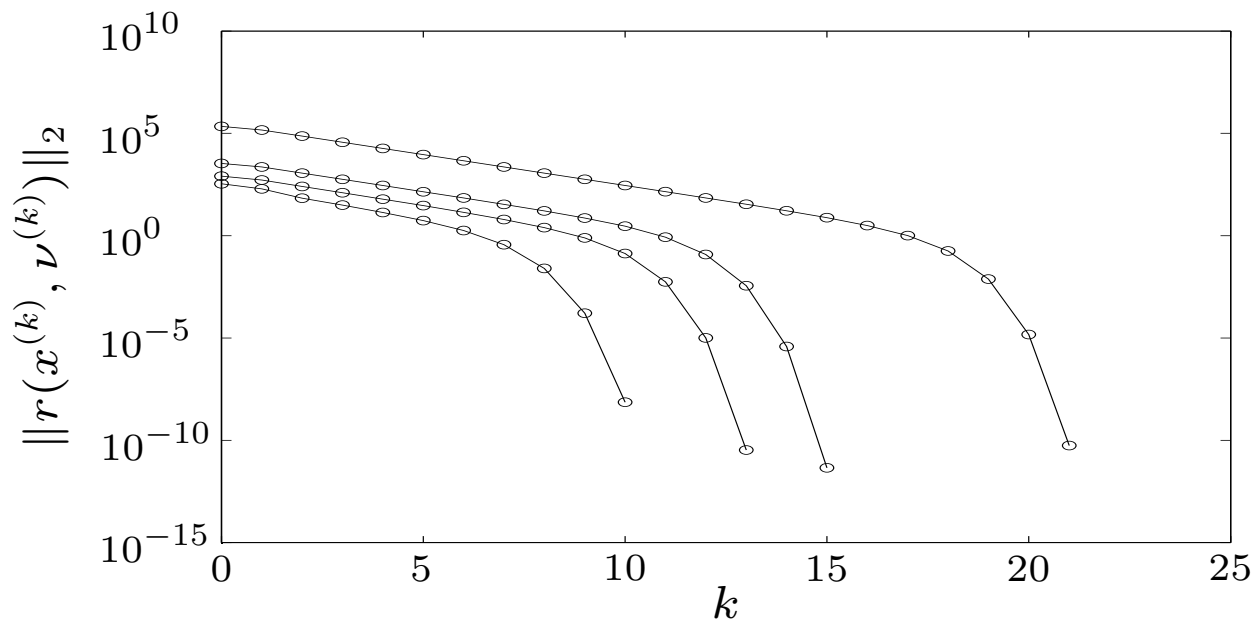
1. Newton method with equality constraints (requires  $x^{(0)} \succ 0$ ,  $Ax^{(0)} = b$ )



## 2. Newton method applied to dual problem (requires $A^T \nu^{(0)} \succ 0$ )



## 3. infeasible start Newton method (requires $x^{(0)} \succ 0$ )



## complexity per iteration of three methods is identical

1. use block elimination to solve KKT system

$$\begin{bmatrix} \mathbf{diag}(x)^{-2} & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} \Delta x \\ w \end{bmatrix} = \begin{bmatrix} \mathbf{diag}(x)^{-1} \mathbf{1}_{d,d} \\ 0 \end{bmatrix}$$

reduces to solving  $A \mathbf{diag}(x)^2 A^T w = b$

2. solve Newton system  $A \mathbf{diag}(A^T \nu)^{-2} A^T \Delta \nu = -b + A \mathbf{diag}(A^T \nu)^{-1} \mathbf{1}_{d,d}$

3. use block elimination to solve KKT system

$$\begin{bmatrix} \mathbf{diag}(x)^{-2} & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta \nu \end{bmatrix} = \begin{bmatrix} \mathbf{diag}(x)^{-1} \mathbf{1}_{d,d} \\ Ax - b \end{bmatrix}$$

reduces to solving  $A \mathbf{diag}(x)^2 A^T w = 2Ax - b$

conclusion: in each case, solve  $ADA^T w = h$  with  $D$  positive diagonal

# Network flow optimization

$$\begin{array}{ll} \text{minimize} & \sum_{i=1}^n \phi_i(x_i) \\ \text{subject to} & Ax = b \end{array}$$

- directed graph with  $n$  arcs,  $p + 1$  nodes
- $x_i$ : flow through arc  $i$ ;  $\phi_i$ : cost flow function for arc  $i$  (with  $\phi_i''(x) > 0$ )
- node-incidence matrix  $\tilde{A} \in \mathbf{R}^{(p+1) \times n}$  defined as

$$\tilde{A}_{ij} = \begin{cases} 1 & \text{arc } j \text{ leaves node } i \\ -1 & \text{arc } j \text{ enters node } i \\ 0 & \text{otherwise} \end{cases}$$

- reduced node-incidence matrix  $A \in \mathbf{R}^{p \times n}$  is  $\tilde{A}$  with last row removed
- $b \in \mathbf{R}^p$  is (reduced) source vector
- **Rank**  $A = p$  if graph is connected

## KKT system

$$\begin{bmatrix} H & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} v \\ w \end{bmatrix} = - \begin{bmatrix} g \\ h \end{bmatrix}$$

- $H = \mathbf{diag}(\phi_1''(x_1), \dots, \phi_n''(x_n))$ , positive diagonal
- solve via elimination:

$$AH^{-1}A^T w = h - AH^{-1}g, \quad Hv = -(g + A^T w)$$

sparsity pattern of coefficient matrix is given by graph connectivity

$$\begin{aligned} (AH^{-1}A^T)_{ij} \neq 0 &\iff (AA^T)_{ij} \neq 0 \\ &\iff \text{nodes } i \text{ and } j \text{ are connected by an arc} \end{aligned}$$

---

# The real deal: General Convex Problems

- inequality constrained minimization
- logarithmic barrier function and central path
- barrier method
- feasibility and phase I methods
- complexity analysis via self-concordance
- generalized inequalities



# Inequality constrained minimization

$$\begin{aligned} & \text{minimize} && f_0(x) \\ & \text{subject to} && f_i(x) \leq 0, \quad i = 1, \dots, m \\ & && Ax = b \end{aligned} \tag{1}$$

- $f_i$  convex, twice continuously differentiable
- $A \in \mathbf{R}^{p \times n}$  with  $\mathbf{Rank} A = p$
- we assume  $p^*$  is finite and attained
- we assume problem is strictly feasible: there exists  $\tilde{x}$  with

$$\tilde{x} \in \mathbf{dom} f_0, \quad f_i(\tilde{x}) < 0, \quad i = 1, \dots, m, \quad A\tilde{x} = b$$

hence, strong duality holds and dual optimum is attained

# Examples

- LP, QP, QCQP, GP
- entropy maximization with linear inequality constraints

$$\begin{array}{ll} \text{minimize} & \sum_{i=1}^n x_i \log x_i \\ \text{subject to} & Fx \preceq g \\ & Ax = b \end{array}$$

with  $\text{dom } f_0 = \mathbf{R}_{++}^n$

- differentiability may require reformulating the problem, *e.g.*, piecewise-linear minimization or  $\ell_\infty$ -norm approximation via LP

# Logarithmic barrier

reformulation of (1) via indicator function:

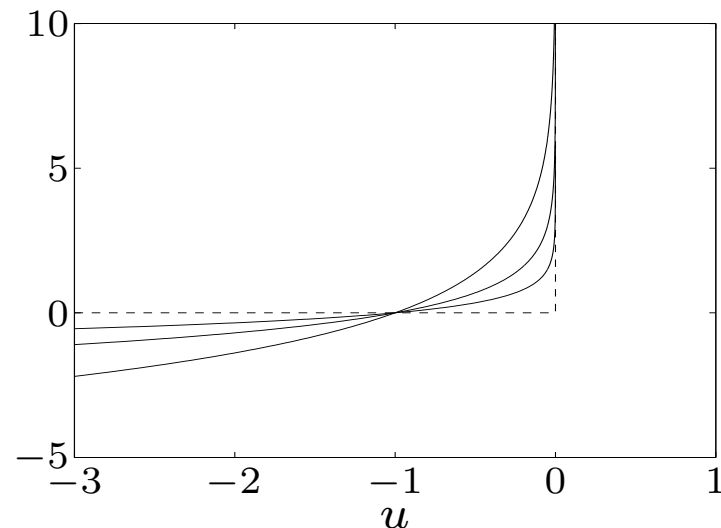
$$\begin{aligned} & \text{minimize} && f_0(x) + \sum_{i=1}^m I_-(f_i(x)) \\ & \text{subject to} && Ax = b \end{aligned}$$

where  $I_-(u) = 0$  if  $u \leq 0$ ,  $I_-(u) = \infty$  otherwise (indicator function of  $\mathbf{R}_-$ )

approximation via logarithmic barrier

$$\begin{aligned} & \text{minimize} && f_0(x) - (1/t) \sum_{i=1}^m \log(-f_i(x)) \\ & \text{subject to} && Ax = b \end{aligned}$$

- an equality constrained problem
- for  $t > 0$ ,  $-(1/t) \log(-u)$  is a smooth approximation of  $I_-$
- approximation improves as  $t \rightarrow \infty$



## logarithmic barrier function

$$\phi(x) = -\sum_{i=1}^m \log(-f_i(x)), \quad \mathbf{dom} \phi = \{x \mid f_1(x) < 0, \dots, f_m(x) < 0\}$$

- convex (follows from composition rules)
- twice continuously differentiable, with derivatives

$$\nabla \phi(x) = \sum_{i=1}^m \frac{1}{-f_i(x)} \nabla f_i(x)$$

$$\nabla^2 \phi(x) = \sum_{i=1}^m \frac{1}{f_i(x)^2} \nabla f_i(x) \nabla f_i(x)^T + \sum_{i=1}^m \frac{1}{-f_i(x)} \nabla^2 f_i(x)$$

# Central path

- for  $t > 0$ , define  $x^*(t)$  as the solution of

$$\begin{aligned} & \text{minimize} && t f_0(x) + \phi(x) \\ & \text{subject to} && Ax = b \end{aligned}$$

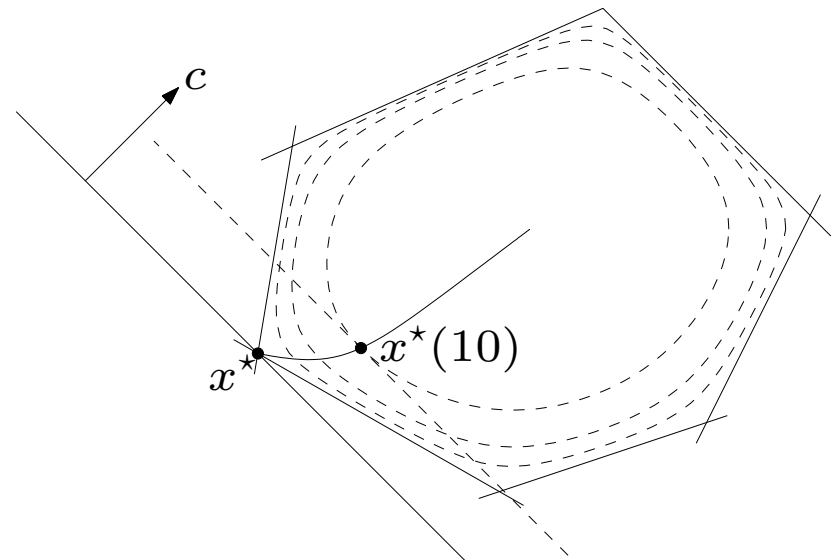
(for now, assume  $x^*(t)$  exists and is unique for each  $t > 0$ )

- central path is  $\{x^*(t) \mid t > 0\}$

**example:** central path for an LP

$$\begin{aligned} & \text{minimize} && c^T x \\ & \text{subject to} && a_i^T x \leq b_i, \quad i = 1, \dots, 6 \end{aligned}$$

hyperplane  $c^T x = c^T x^*(t)$  is tangent to level curve of  $\phi$  through  $x^*(t)$



## Dual points on central path

$x = x^*(t)$  if there exists a  $w$  such that

$$t\nabla f_0(x) + \sum_{i=1}^m \frac{1}{-f_i(x)} \nabla f_i(x) + A^T w = 0, \quad Ax = b$$

- therefore,  $x^*(t)$  minimizes the Lagrangian

$$L(x, \lambda^*(t), \nu^*(t)) = f_0(x) + \sum_{i=1}^m \lambda_i^*(t) f_i(x) + \nu^*(t)^T (Ax - b)$$

where we define  $\lambda_i^*(t) = 1/(-t f_i(x^*(t)))$  and  $\nu^*(t) = w/t$

- this confirms the intuitive idea that  $f_0(x^*(t)) \rightarrow p^*$  if  $t \rightarrow \infty$ :

$$\begin{aligned} p^* &\geq g(\lambda^*(t), \nu^*(t)) \\ &= L(x^*(t), \lambda^*(t), \nu^*(t)) \\ &= f_0(x^*(t)) - m/t \end{aligned}$$

# Interpretation via KKT conditions

$x = x^*(t)$ ,  $\lambda = \lambda^*(t)$ ,  $\nu = \nu^*(t)$  satisfy

1. primal constraints:  $f_i(x) \leq 0$ ,  $i = 1, \dots, m$ ,  $Ax = b$
2. dual constraints:  $\lambda \succeq 0$
3. approximate complementary slackness:  $-\lambda_i f_i(x) = 1/t$ ,  $i = 1, \dots, m$
4. gradient of Lagrangian with respect to  $x$  vanishes:

$$\nabla f_0(x) + \sum_{i=1}^m \lambda_i \nabla f_i(x) + A^T \nu = 0$$

difference with KKT is that condition 3 replaces  $\lambda_i f_i(x) = 0$

# Force field interpretation

**centering problem** (for problem with no equality constraints)

$$\text{minimize } tf_0(x) - \sum_{i=1}^m \log(-f_i(x))$$

## force field interpretation

- $tf_0(x)$  is potential of force field  $F_0(x) = -t\nabla f_0(x)$
- $-\log(-f_i(x))$  is potential of force field  $F_i(x) = (1/f_i(x))\nabla f_i(x)$

the forces balance at  $x^*(t)$ :

$$F_0(x^*(t)) + \sum_{i=1}^m F_i(x^*(t)) = 0$$



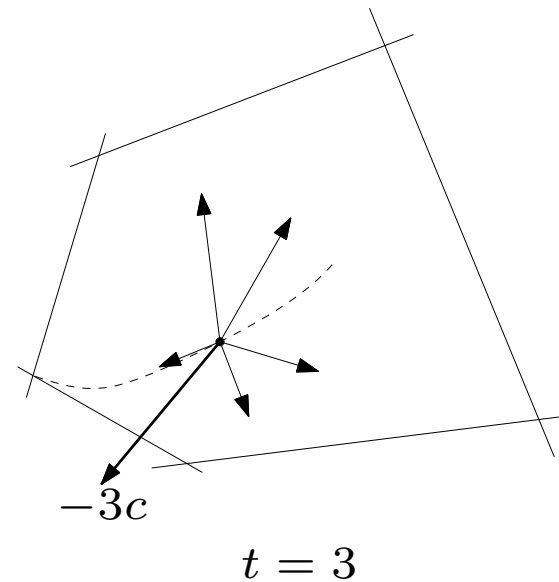
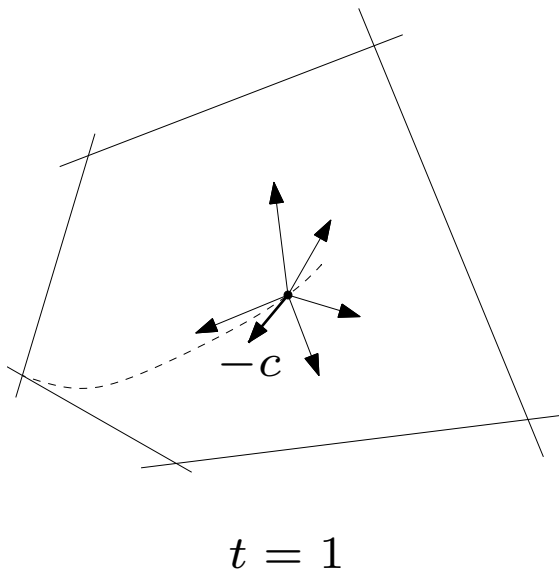
## example

$$\begin{aligned} & \text{minimize} && c^T x \\ & \text{subject to} && a_i^T x \leq b_i, \quad i = 1, \dots, m \end{aligned}$$

- objective force field is constant:  $F_0(x) = -tc$
- constraint force field decays as inverse distance to constraint hyperplane:

$$F_i(x) = \frac{-a_i}{b_i - a_i^T x}, \quad \|F_i(x)\|_2 = \frac{1}{\mathbf{dist}(x, \mathcal{H}_i)}$$

where  $\mathcal{H}_i = \{x \mid a_i^T x = b_i\}$



# Barrier method

**given** strictly feasible  $x$ ,  $t := t^{(0)} > 0$ ,  $\mu > 1$ , tolerance  $\epsilon > 0$ .

**repeat**

1. *Centering step.* Compute  $x^*(t)$  by minimizing  $tf_0 + \phi$ , subject to  $Ax = b$ .
2. *Update.*  $x := x^*(t)$ .
3. *Stopping criterion.* **quit** if  $m/t < \epsilon$ .
4. *Increase  $t$ .*  $t := \mu t$ .

- terminates with  $f_0(x) - p^* \leq \epsilon$  (stopping criterion follows from  $f_0(x^*(t)) - p^* \leq m/t$ )
- centering usually done using Newton's method, starting at current  $x$
- choice of  $\mu$  involves a trade-off: large  $\mu$  means fewer outer iterations, more inner (Newton) iterations; typical values:  $\mu = 10\text{--}20$
- several heuristics for choice of  $t^{(0)}$

# Convergence analysis

**number of outer (centering) iterations:** exactly

$$\left\lceil \frac{\log(m/(\epsilon t^{(0)}))}{\log \mu} \right\rceil$$

plus the initial centering step (to compute  $x^*(t^{(0)})$ )

**centering problem**

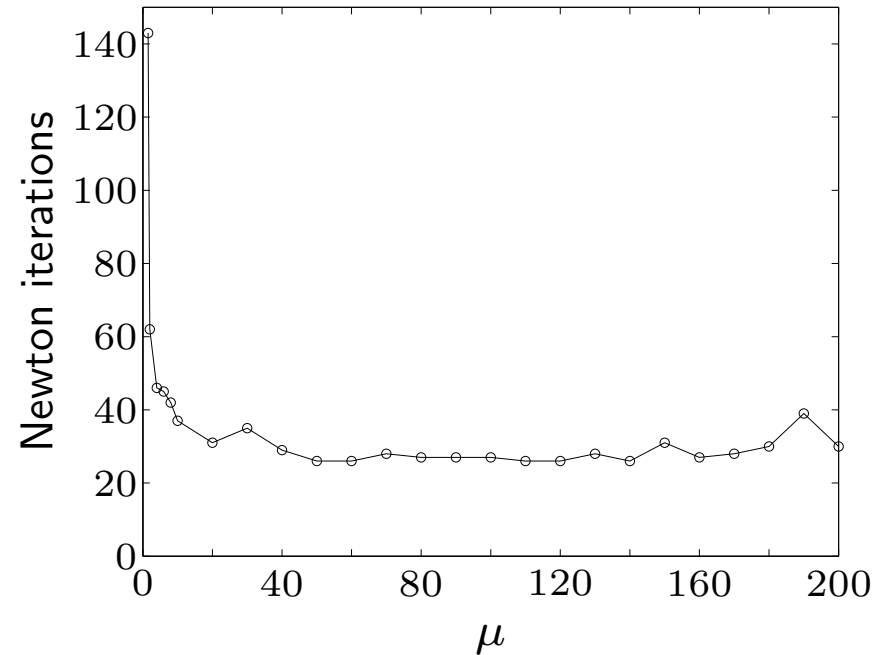
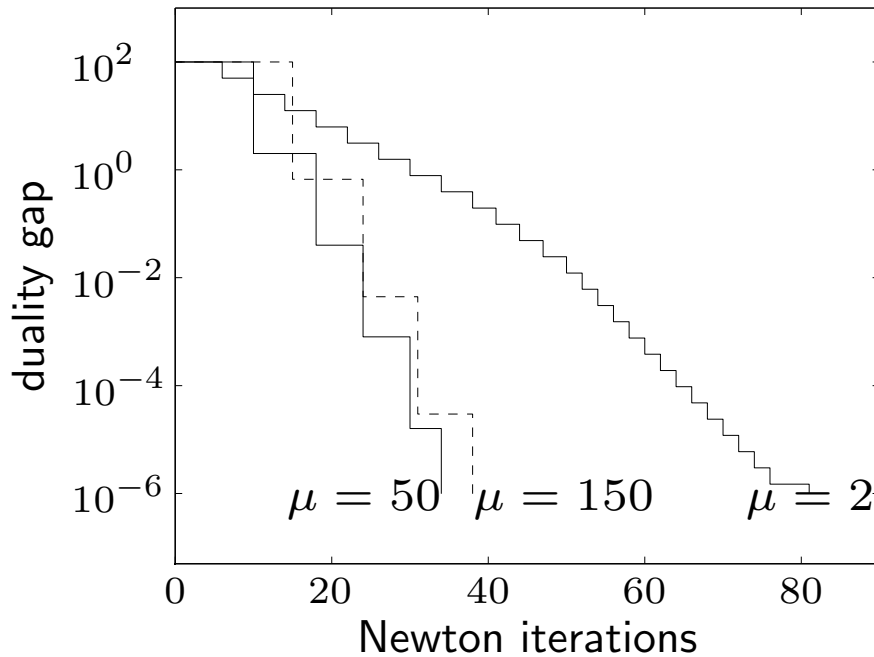
$$\text{minimize } tf_0(x) + \phi(x)$$

see convergence analysis of Newton's method

- $tf_0 + \phi$  must have closed sublevel sets for  $t \geq t^{(0)}$
- classical analysis requires strong convexity, Lipschitz condition
- analysis via self-concordance requires self-concordance of  $tf_0 + \phi$

# Examples

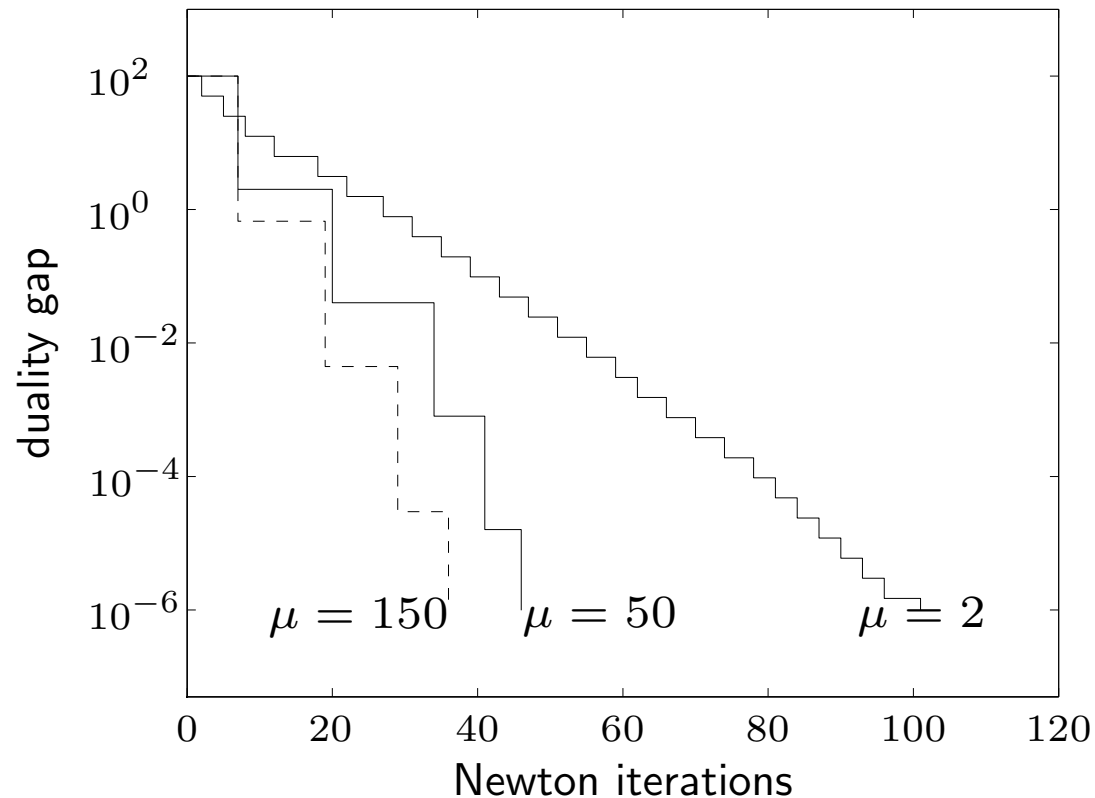
inequality form LP ( $m = 100$  inequalities,  $n = 50$  variables)



- starts with  $x$  on central path ( $t^{(0)} = 1$ , duality gap 100)
- terminates when  $t = 10^8$  (gap  $10^{-6}$ )
- centering uses Newton's method with backtracking
- total number of Newton iterations not very sensitive for  $\mu \geq 10$

**geometric program** ( $m = 100$  inequalities and  $n = 50$  variables)

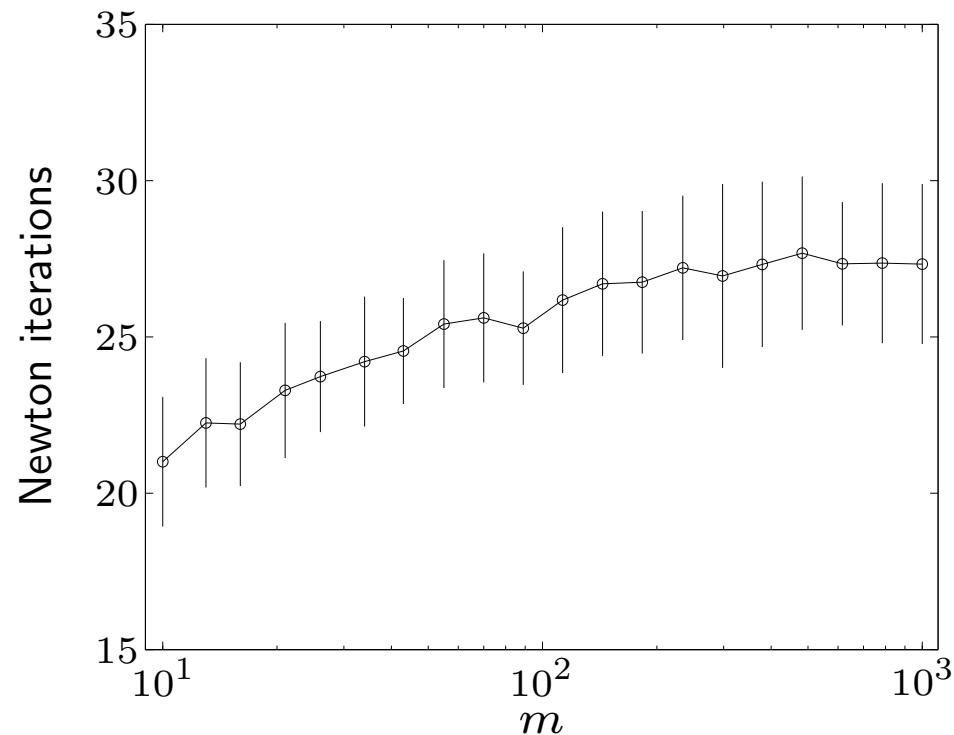
$$\begin{aligned} \text{minimize} \quad & \log \left( \sum_{k=1}^5 \exp(a_{0k}^T x + b_{0k}) \right) \\ \text{subject to} \quad & \log \left( \sum_{k=1}^5 \exp(a_{ik}^T x + b_{ik}) \right) \leq 0, \quad i = 1, \dots, m \end{aligned}$$



family of standard LPs ( $A \in \mathbf{R}^{m \times 2m}$ )

$$\begin{aligned} & \text{minimize} && c^T x \\ & \text{subject to} && Ax = b, \quad x \succeq 0 \end{aligned}$$

$m = 10, \dots, 1000$ ; for each  $m$ , solve 100 randomly generated instances



number of iterations grows very slowly as  $m$  ranges over a 100 : 1 ratio

# Feasibility and phase I methods

**feasibility problem:** find  $x$  such that

$$f_i(x) \leq 0, \quad i = 1, \dots, m, \quad Ax = b \quad (2)$$

**phase I:** computes strictly feasible starting point for barrier method

**basic phase I method**

$$\begin{array}{ll} \text{minimize (over } x, s) & s \\ \text{subject to} & f_i(x) \leq s, \quad i = 1, \dots, m \\ & Ax = b \end{array} \quad (3)$$

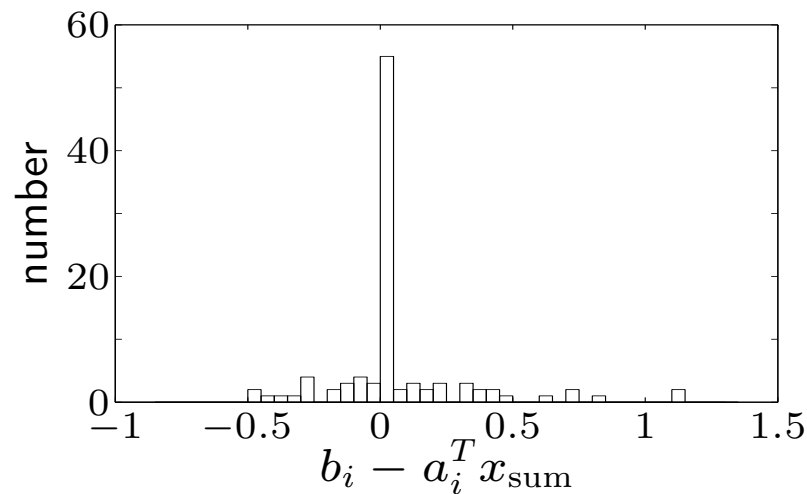
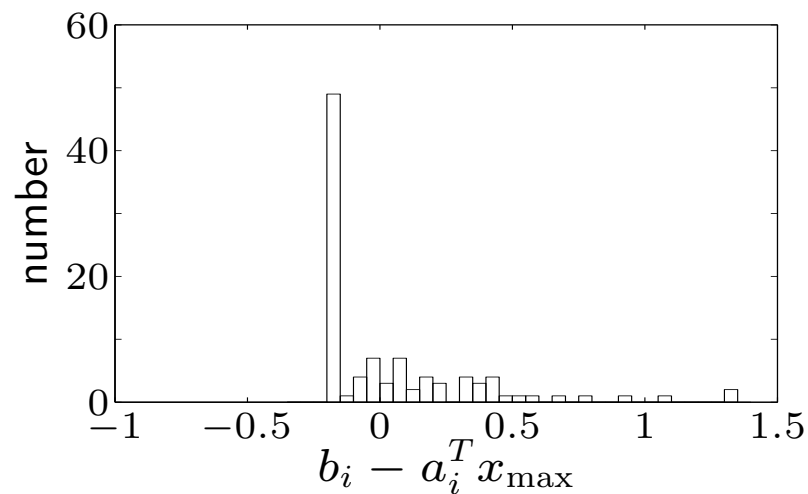
- if  $x, s$  feasible, with  $s < 0$ , then  $x$  is strictly feasible for (2)
- if optimal value  $\bar{p}^*$  of (3) is positive, then problem (2) is infeasible
- if  $\bar{p}^* = 0$  and attained, then problem (2) is feasible (but not strictly);  
if  $\bar{p}^* = 0$  and not attained, then problem (2) is infeasible

## sum of infeasibilities phase I method

$$\begin{aligned} & \text{minimize} && \mathbb{1}_{d,d}^T s \\ & \text{subject to} && s \succeq 0, \quad f_i(x) \leq s_i, \quad i = 1, \dots, m \\ & && Ax = b \end{aligned}$$

for infeasible problems, produces a solution that satisfies many more inequalities than basic phase I method

**example** (infeasible set of 100 linear inequalities in 50 variables)



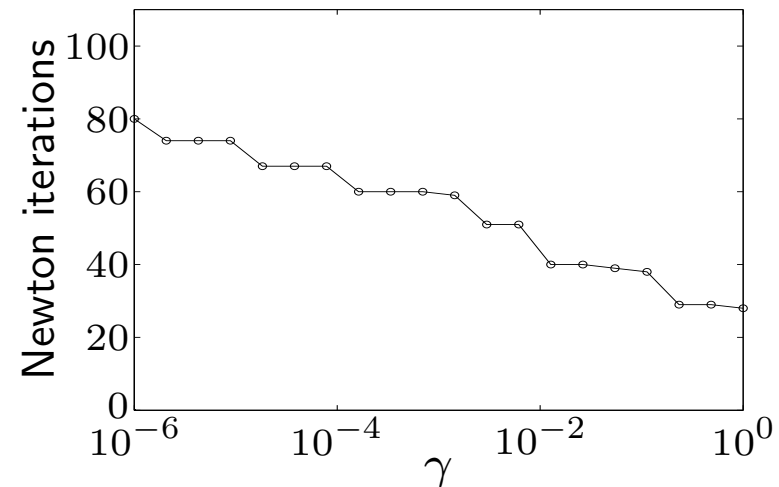
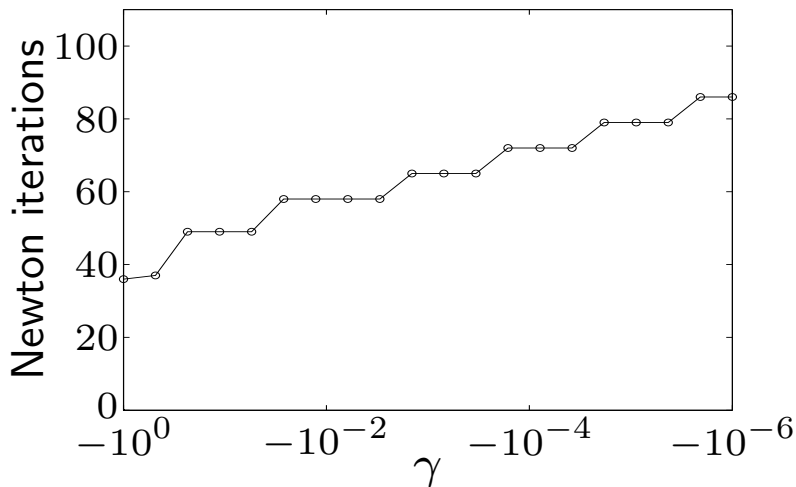
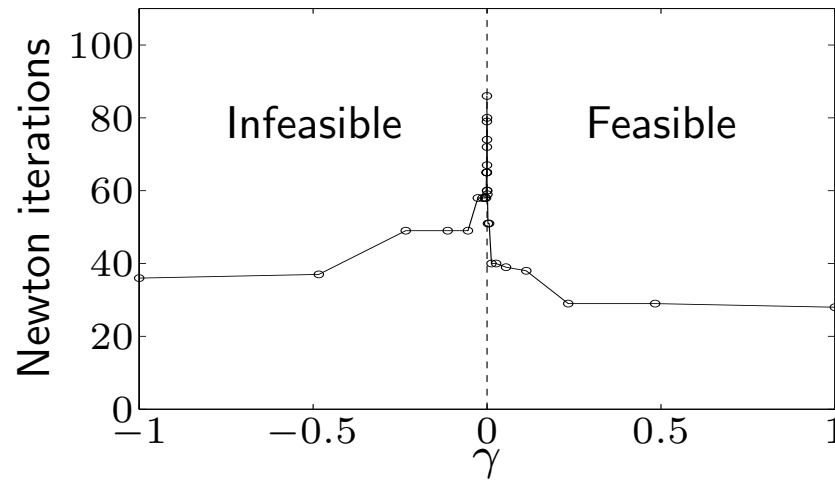
left: basic phase I solution; satisfies 39 inequalities

right: sum of infeasibilities phase I solution; satisfies 79 solutions



**example:** family of linear inequalities  $Ax \preceq b + \gamma \Delta b$

- data chosen to be strictly feasible for  $\gamma > 0$ , infeasible for  $\gamma \leq 0$
- use basic phase I, terminate when  $s < 0$  or dual objective is positive



number of iterations roughly proportional to  $\log(1/|\gamma|)$

# Complexity analysis via self-concordance

same assumptions as on page 49, plus:

- sublevel sets (of  $f_0$ , on the feasible set) are bounded
- $tf_0 + \phi$  is self-concordant with closed sublevel sets

second condition

- holds for LP, QP, QCQP
- may require reformulating the problem, *e.g.*,

$$\begin{array}{ll} \text{minimize} & \sum_{i=1}^n x_i \log x_i \\ \text{subject to} & Fx \preceq g \end{array} \quad \longrightarrow \quad \begin{array}{ll} \text{minimize} & \sum_{i=1}^n x_i \log x_i \\ \text{subject to} & Fx \preceq g, \quad x \succeq 0 \end{array}$$

- needed for complexity analysis; barrier method works even when self-concordance assumption does not apply

## Newton iterations per centering step: from self-concordance theory

$$\# \text{Newton iterations} \leq \frac{\mu t f_0(x) + \phi(x) - \mu t f_0(x^+) - \phi(x^+)}{\gamma} + c$$

- bound on effort of computing  $x^+ = x^*(\mu t)$  starting at  $x = x^*(t)$
- $\gamma, c$  are constants (depend only on Newton algorithm parameters)
- from duality (with  $\lambda = \lambda^*(t), \nu = \nu^*(t)$ ):

$$\begin{aligned} & \mu t f_0(x) + \phi(x) - \mu t f_0(x^+) - \phi(x^+) \\ &= \mu t f_0(x) - \mu t f_0(x^+) + \sum_{i=1}^m \log(-\mu t \lambda_i f_i(x^+)) - m \log \mu \\ &\leq \mu t f_0(x) - \mu t f_0(x^+) - \mu t \sum_{i=1}^m \lambda_i f_i(x^+) - m - m \log \mu \\ &\leq \mu t f_0(x) - \mu t g(\lambda, \nu) - m - m \log \mu \\ &= m(\mu - 1 - \log \mu) \end{aligned}$$

## total number of Newton iterations (excluding first centering step)

$$\# \text{Newton iterations} \leq N = \left\lceil \frac{\log(m/(t^{(0)}\epsilon))}{\log \mu} \right\rceil \left( \frac{m(\mu - 1 - \log \mu)}{\gamma} + c \right)$$

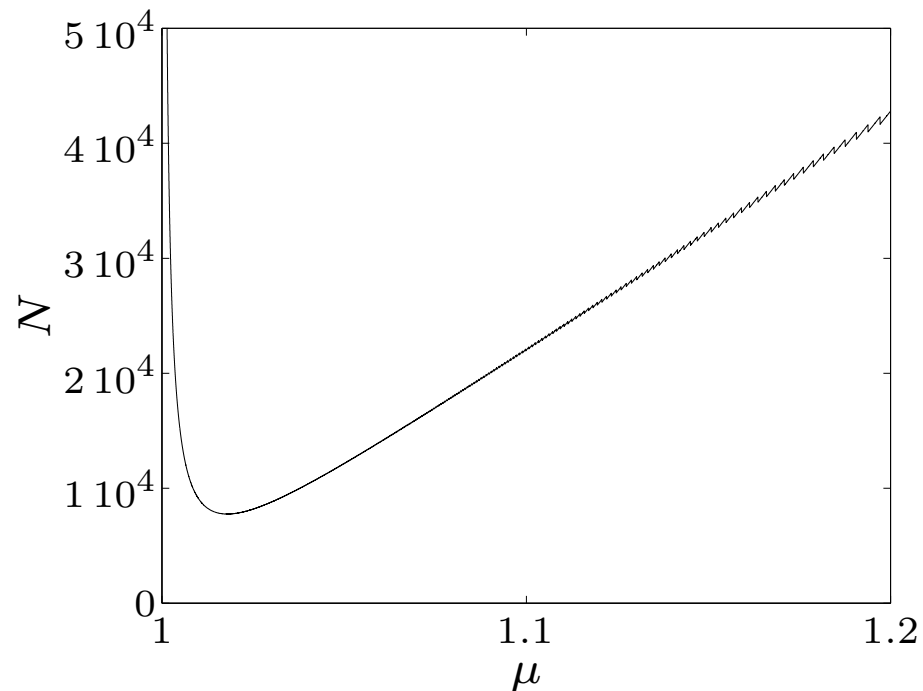


figure shows  $N$  for typical values of  $\gamma, c$ ,

$$m = 100, \quad \frac{m}{t^{(0)}\epsilon} = 10^5$$

- confirms trade-off in choice of  $\mu$
- in practice, #iterations is in the tens; not very sensitive for  $\mu \geq 10$

## polynomial-time complexity of barrier method

- for  $\mu = 1 + 1/\sqrt{m}$ :

$$N = O\left(\sqrt{m} \log\left(\frac{m/t^{(0)}}{\epsilon}\right)\right)$$

- number of Newton iterations for fixed gap reduction is  $O(\sqrt{m})$
- multiply with cost of one Newton iteration (a polynomial function of problem dimensions), to get bound on number of flops

this choice of  $\mu$  optimizes worst-case complexity; in practice we choose  $\mu$  fixed ( $\mu = 10, \dots, 20$ )

# Barrier method

**given** strictly feasible  $x$ ,  $t := t^{(0)} > 0$ ,  $\mu > 1$ , tolerance  $\epsilon > 0$ .

**repeat**

1. *Centering step.* Compute  $x^*(t)$  by minimizing  $tf_0 + \phi$ , subject to  $Ax = b$ .
2. *Update.*  $x := x^*(t)$ .
3. *Stopping criterion.* **quit** if  $(\sum_i \theta_i)/t < \epsilon$ .
4. *Increase  $t$ .*  $t := \mu t$ .

- only difference is duality gap  $m/t$  on central path is replaced by  $\sum_i \theta_i/t$
- number of outer iterations:

$$\left\lceil \frac{\log((\sum_i \theta_i)/(\epsilon t^{(0)}))}{\log \mu} \right\rceil$$