# Convex Optimization M2

## Lecture 4

# Unconstrained minimization

# Unconstrained minimization

- terminology and assumptions

- gradient descent method

- steepest descent method

- Newton's method

- self-concordant functions

- implementation

# Unconstrained minimization

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

- $f$ convex, twice continuously differentiable (hence $\mathbf{dom}\, f$ open)

- we assume optimal value $p^\star = \inf_x f(x)$ is attained (and finite)

**unconstrained minimization methods**

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

$$f(x^{(k)}) \to p^\star$$

- can be interpreted as iterative methods for solving optimality condition

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

# Initial point and sublevel set

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

- $x^{(0)} \in \mathbf{dom}\, f$

- sublevel set $S = \{x \mid f(x) \le f(x^{(0)})\}$ 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 $\mathbf{dom}\, f = \mathbb{R}^n$

- true if $f(x) \to \infty$ as $x \to \mathbf{bd}\,\mathbf{dom}\, f$

examples of differentiable functions with closed sublevel sets:

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

# Strong convexity and implications

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

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

**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

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

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

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

# Descent methods

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

- other notations: $x^+ = x + t\Delta x$, $x := x + t\Delta x$

- $\Delta x$ is the *step*, or *search direction*; $t$ is the *step size*, or *step length*

- from convexity, $f(x^+) < f(x)$ implies $\nabla f(x)^T \Delta x < 0$
  (*i.e.*, $\Delta x$ is a *descent direction*)

*General descent method.*

**given** a starting point $x \in \mathbf{dom}\, f$.
**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
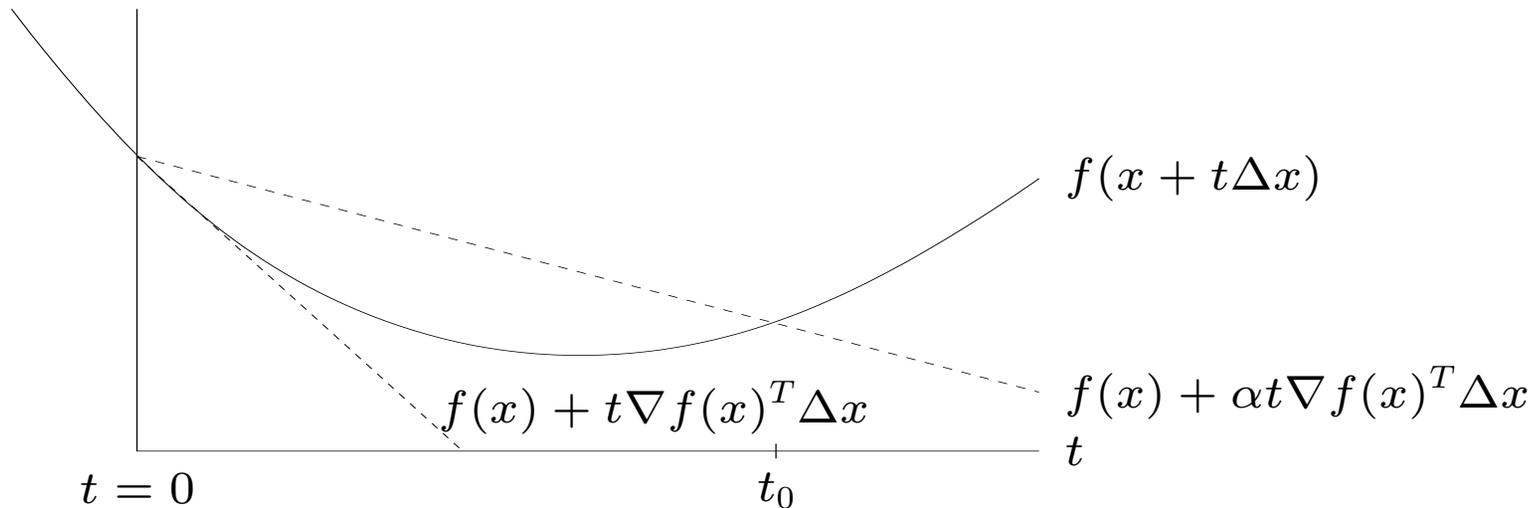
**exact line search:** $t = \text{argmin}_{t>0} f(x + t\Delta x)$

**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

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

> **given** a starting point $x \in \mathbf{dom}\, f$.
> **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 $f$,

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

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

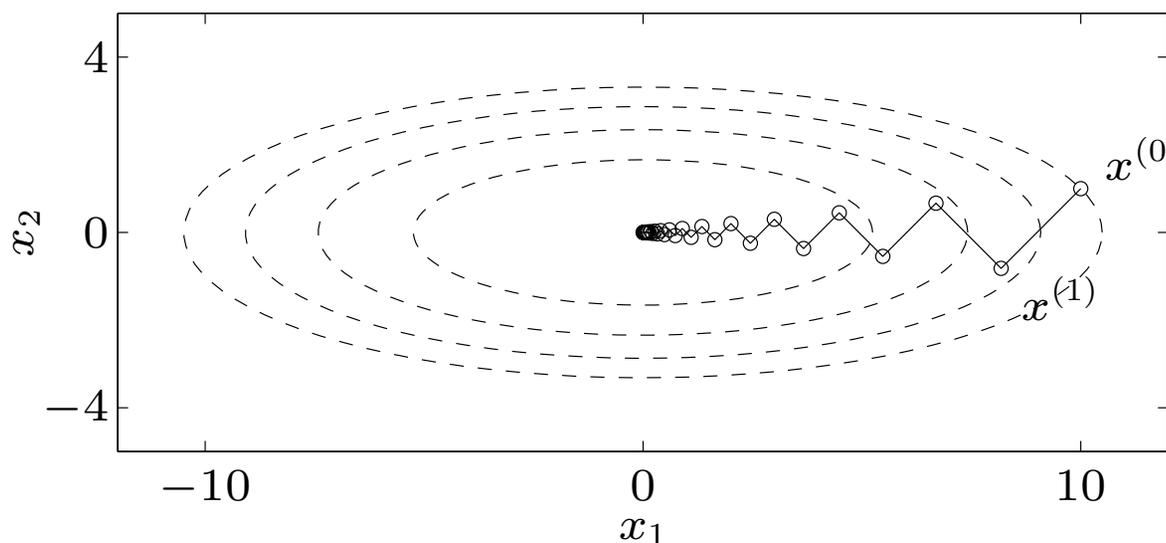- very simple, but often very slow; rarely used in practice

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

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

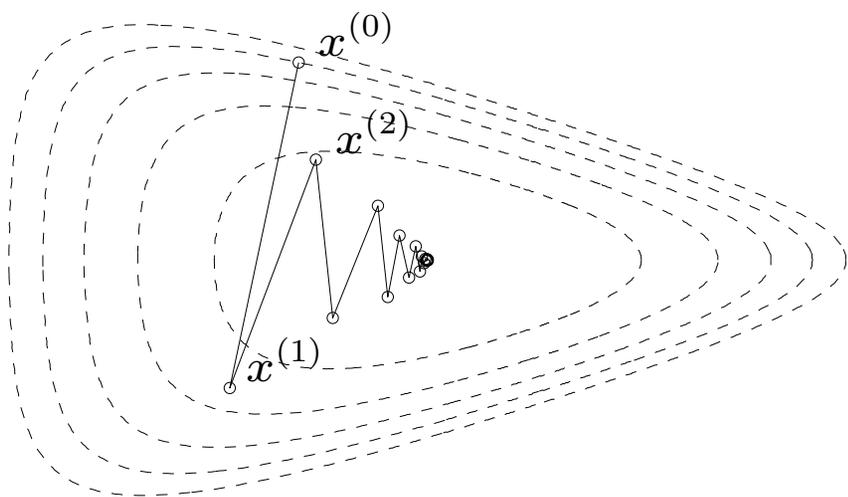with exact line search, starting at $x^{(0)} = (\gamma, 1)$:

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

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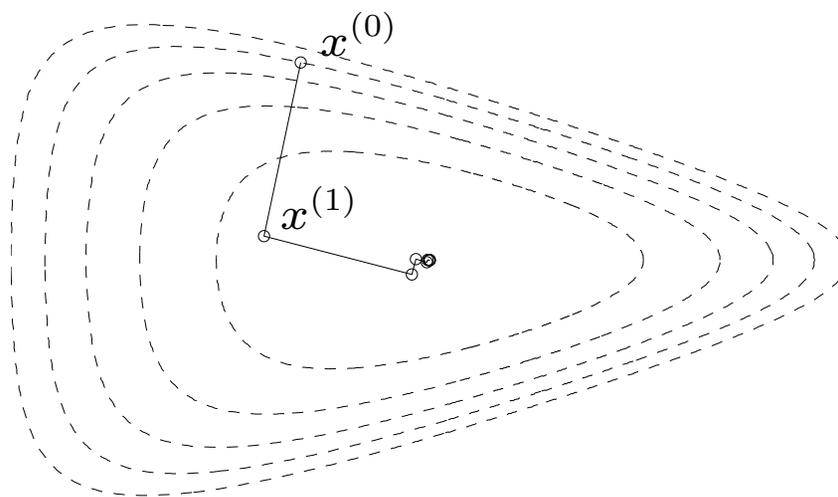- example for $\gamma = 10$:

# nonquadratic example

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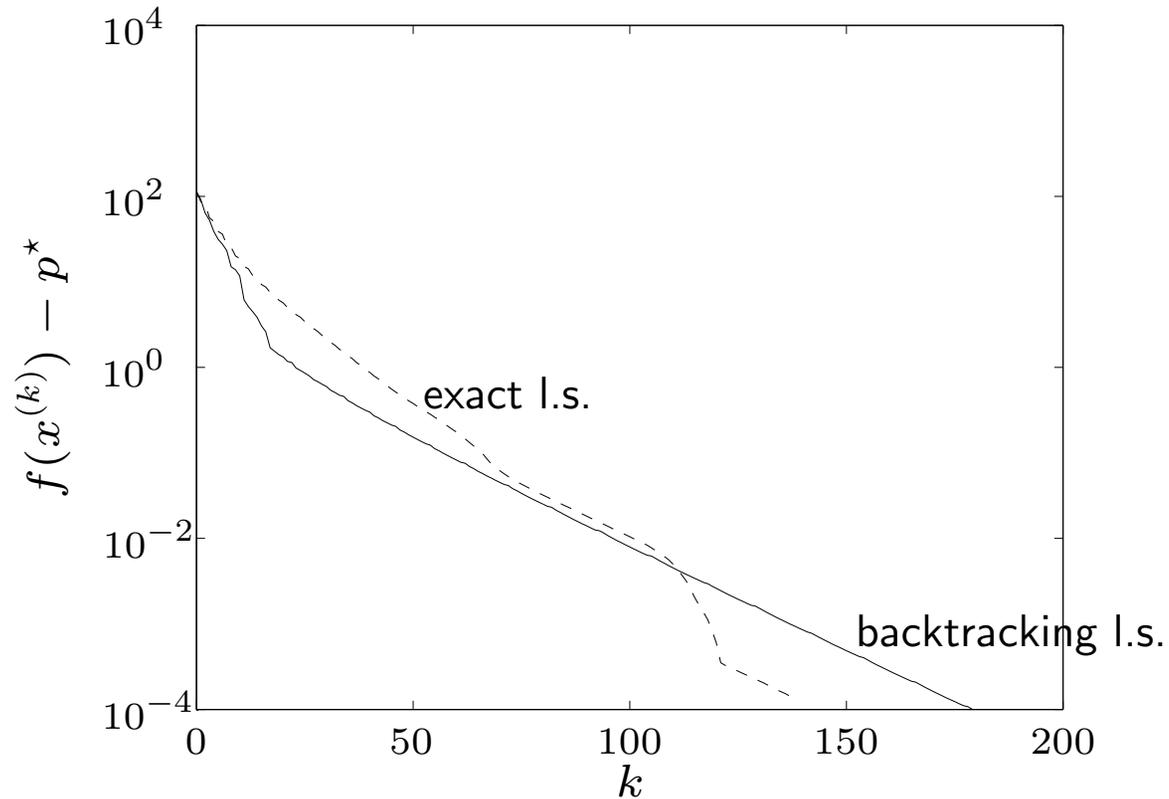


backtracking line search                    exact line search

# a problem in $\mathbb{R}^{100}$

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



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

# Steepest descent method

**normalized steepest descent direction** (at $x$, for norm $\|\cdot\|$):

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

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

**(unnormalized) steepest descent direction**

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

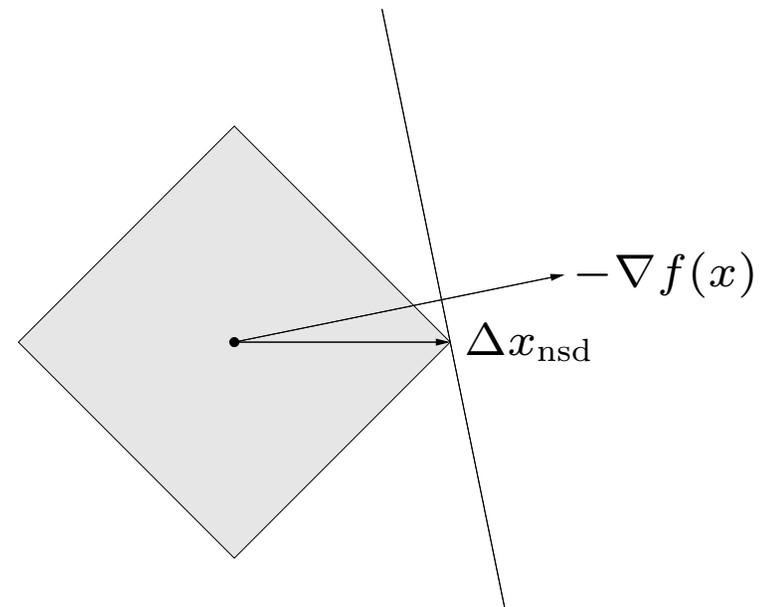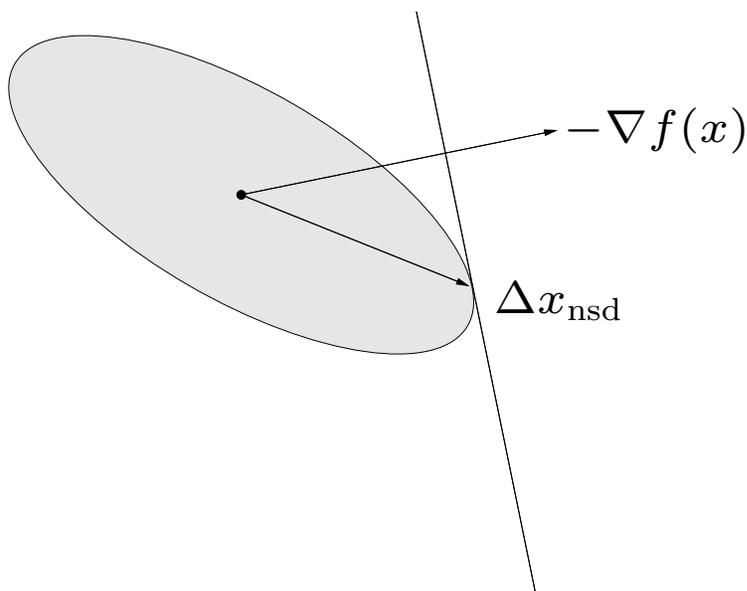satisfies $\nabla f(x)^T \Delta_{\mathrm{sd}} = -\|\nabla f(x)\|_*^2$

**steepest descent method**

- general descent method with $\Delta x = \Delta x_{\mathrm{sd}}$

- convergence properties similar to gradient descent

## examples

- Euclidean norm: $\Delta x_{\mathrm{sd}} = -\nabla f(x)$

- quadratic norm $\|x\|_P = (x^T P x)^{1/2}$ $(P \in \mathbf{S}_{++}^n)$: $\Delta x_{\mathrm{sd}} = -P^{-1}\nabla f(x)$

- $\ell_1$-norm: $\Delta x_{\mathrm{sd}} = -(\partial f(x)/\partial x_i)e_i$, 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 $\bar{x} = P^{1/2}x$

shows 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

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

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

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

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

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



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

$$\lambda(x) = \left(\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x)\right)^{1/2}$$
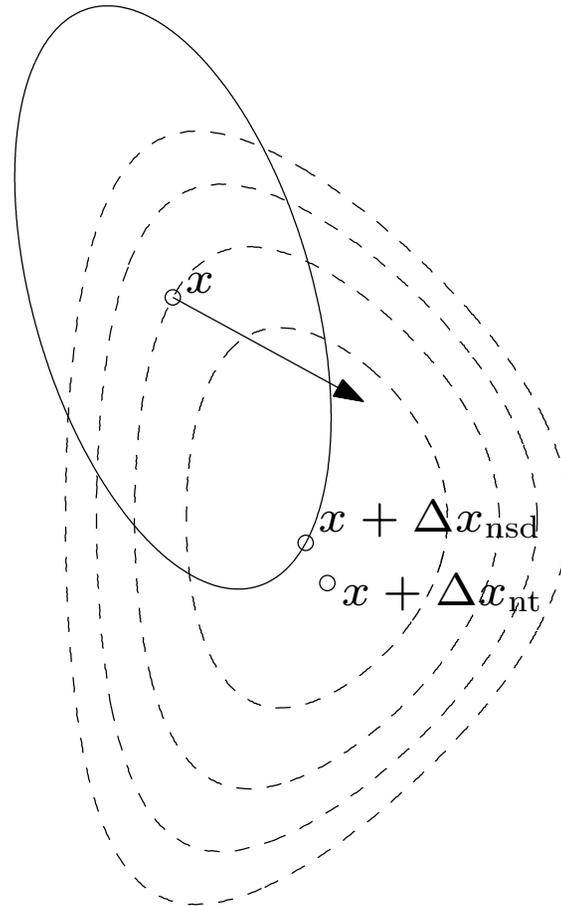
a measure of the proximity of $x$ to $x^\star$

**properties**

- gives an estimate of $f(x) - p^\star$, using quadratic approximation $\widehat{f}$:

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

- equal to the norm of the Newton step in the quadratic Hessian norm

$$\lambda(x) = \left(\Delta x_{\mathrm{nt}} \nabla^2 f(x) \Delta x_{\mathrm{nt}}\right)^{1/2}$$

- directional derivative in the Newton direction: $\nabla f(x)^T \Delta x_{\mathrm{nt}} = -\lambda(x)^2$

- affine invariant (unlike $\|\nabla f(x)\|_2$)

# Newton's method

**given** a starting point $x \in \mathbf{dom}\, f$, tolerance $\epsilon > 0$.
**repeat**

1. *Compute the Newton step and decrement.*
$$\Delta x_{\mathrm{nt}} := -\nabla^2 f(x)^{-1} \nabla f(x); \quad \lambda^2 := \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x).$$
2. *Stopping criterion.* **quit** if $\lambda^2/2 \le \epsilon$.
3. *Line search.* Choose step size $t$ by backtracking line search.
4. *Update.* $x := x + t\Delta x_{\mathrm{nt}}$.

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 $p^\star > -\infty$, this phase ends after at most $(f(x^{(0)}) - p^\star)/\gamma$ 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

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

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

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

- $\gamma$, $\epsilon_0$ 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** $\mathbb{R}^2$ (page 11)



- backtracking parameters $\alpha = 0.1$, $\beta = 0.7$

- converges in only 5 steps

- quadratic local convergence

# example in $\mathbb{R}^{100}$ (page 12)



- 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** $\mathbb{R}^{10000}$ (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

# 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

# Self-concordant functions

**definition**

- $f : \mathbb{R} \to \mathbb{R}$ is self-concordant if $|f'''(x)| \le 2f''(x)^{3/2}$ for all $x \in \mathbf{dom}\, f$

- $f : \mathbb{R}^n \to \mathbb{R}$ is self-concordant if $g(t) = f(x + tv)$ is self-concordant for all $x \in \mathbf{dom}\, f$, $v \in \mathbb{R}^n$

**examples on** $\mathbb{R}$

- linear and quadratic functions

- negative logarithm $f(x) = -\log x$

- negative entropy plus negative logarithm: $f(x) = x \log x - \log x$

**affine invariance:** if $f : \mathbb{R} \to \mathbb{R}$ is s.c., then $\tilde{f}(y) = f(ay + b)$ is s.c.:

$$\tilde{f}'''(y) = a^3 f'''(ay + b), \qquad \tilde{f}''(y) = a^2 f''(ay + b)$$

# Self-concordant calculus

**properties**

- preserved under positive scaling $\alpha \geq 1$, and sum

- preserved under composition with affine function

- if $g$ is convex with $\mathbf{dom}\, g = \mathbb{R}_{++}$ and $|g'''(x)| \leq 3g''(x)/x$ then

$$f(x) = \log(-g(x)) - \log x$$

  is self-concordant

**examples**: properties can be used to show that the following are s.c.

- $f(x) = -\sum_{i=1}^{m} \log(b_i - a_i^T x)$ on $\{x \mid a_i^T x < b_i,\ i = 1, \ldots, m\}$
- $f(X) = -\log \det X$ on $\mathbf{S}_{++}^n$
- $f(x) = -\log(y^2 - x^T x)$ on $\{(x, y) \mid \|x\|_2 < y\}$

# Convergence analysis for self-concordant functions

**summary**: there exist constants $\eta \in (0, 1/4]$, $\gamma > 0$ such that

- if $\lambda(x) > \eta$, then

$$f(x^{(k+1)}) - f(x^{(k)}) \leq -\gamma$$

- if $\lambda(x) \leq \eta$, then

$$2\lambda(x^{(k+1)}) \leq \left(2\lambda(x^{(k)})\right)^2$$

($\eta$ and $\gamma$ only depend on backtracking parameters $\alpha$, $\beta$)

**complexity bound:** number of Newton iterations bounded by

$$\frac{f(x^{(0)}) - p^\star}{\gamma} + \log_2 \log_2(1/\epsilon)$$

for $\alpha = 0.1$, $\beta = 0.8$, $\epsilon = 10^{-10}$, bound evaluates to $375(f(x^{(0)}) - p^\star) + 6$

**numerical example:** 150 randomly generated instances of

$$\text{minimize} \quad f(x) = -\sum_{i=1}^{m} \log(b_i - a_i^T x)$$

$\bigcirc$: $m = 100$, $n = 50$
$\square$: $m = 1000$, $n = 500$
$\diamondsuit$: $m = 1000$, $n = 50$



- number of iterations much smaller than $375(f(x^{(0)}) - p^\star) + 6$

- bound of the form $c(f(x^{(0)}) - p^\star) + 6$ with smaller $c$ (empirically) valid

# 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

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

- 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

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

- assume $A \in \mathbb{R}^{p \times n}$, 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

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

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

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

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

# Equality Constraints

# Equality Constraints

- equality constrained minimization

- eliminating equality constraints

- 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

- $A \in \mathbb{R}^{p \times n}$ with $\mathbf{Rank}\, A = p$

- we assume $p^\star$ is finite and attained

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

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

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

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

optimality condition:

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

- coefficient matrix is called KKT matrix

- KKT matrix is nonsingular if and only if

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

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

# Eliminating equality constraints

represent solution of $\{x \mid Ax = b\}$ as

$$\{x \mid Ax = b\} = \{Fz + \hat{x} \mid z \in \mathbb{R}^{n-p}\}$$

- $\hat{x}$ is (any) particular solution

- range of $F \in \mathbb{R}^{n \times (n-p)}$ is nullspace of $A$ ($\mathbf{Rank}\, F = n - p$ and $AF = 0$)

**reduced or eliminated problem**

$$\text{minimize} \quad f(Fz + \hat{x})$$

- an unconstrained problem with variable $z \in \mathbb{R}^{n-p}$

- from solution $z^\star$, obtain $x^\star$ and $\nu^\star$ as

$$x^\star = Fz^\star + \hat{x}, \qquad \nu^\star = -(AA^T)^{-1} A \nabla f(x^\star)$$

**example:** optimal allocation with resource constraint

$$
\begin{array}{ll}
\text{minimize} & f_1(x_1) + f_2(x_2) + \cdots + f_n(x_n) \\
\text{subject to} & x_1 + x_2 + \cdots + x_n = b
\end{array}
$$

eliminate $x_n = b - x_1 - \cdots - x_{n-1}$, $i.e.$, choose

$$
\hat{x} = be_n, \qquad F = \begin{bmatrix} I \\ -\mathbf{1}^T \end{bmatrix} \in \mathbb{R}^{n \times (n-1)}
$$

reduced problem:

$$
\text{minimize} \quad f_1(x_1) + \cdots + f_{n-1}(x_{n-1}) + f_n(b - x_1 - \cdots - x_{n-1})
$$

(variables $x_1, \ldots, x_{n-1}$)

# Newton step

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

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

**interpretations**

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

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

- equations follow from linearizing optimality conditions

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

# Newton decrement

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

**properties**

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

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

- directional derivative in Newton direction:

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

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

# Newton's method with equality constraints

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

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

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

- affine invariant

# Newton's method and elimination

**Newton's method for reduced problem**

$$\text{minimize} \quad \tilde{f}(z) = f(Fz + \hat{x})$$

- variables $z \in \mathbb{R}^{n-p}$

- $\hat{x}$ satisfies $A\hat{x} = b$; $\textbf{Rank}\, F = n - p$ and $AF = 0$

- Newton's method for $\tilde{f}$, started at $z^{(0)}$, generates iterates $z^{(k)}$

**Newton's method with equality constraints**

when started at $x^{(0)} = Fz^{(0)} + \hat{x}$, iterates are

$$x^{(k+1)} = Fz^{(k)} + \hat{x}$$

hence, don't need separate convergence analysis

# Newton step at infeasible points

2nd interpretation of page 39 extends to infeasible $x$ (*i.e.*, $Ax \neq b$)

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

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

**primal-dual interpretation**

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

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

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

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

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

# Infeasible start Newton method

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

**repeat**

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

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

        $t := 1$.

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

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

**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 $\Delta y = (\Delta x_{\mathrm{nt}}, \Delta \nu_{\mathrm{nt}})$ is

$$\frac{d}{dt} \|r(y + \Delta y)\|_2 \bigg|_{t=0} = -\|r(y)\|_2$$

# 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**

- $\mathrm{LDL}^\mathsf{T}$ factorization

- elimination (if $H$ nonsingular)

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

- elimination with singular $H$: write as

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

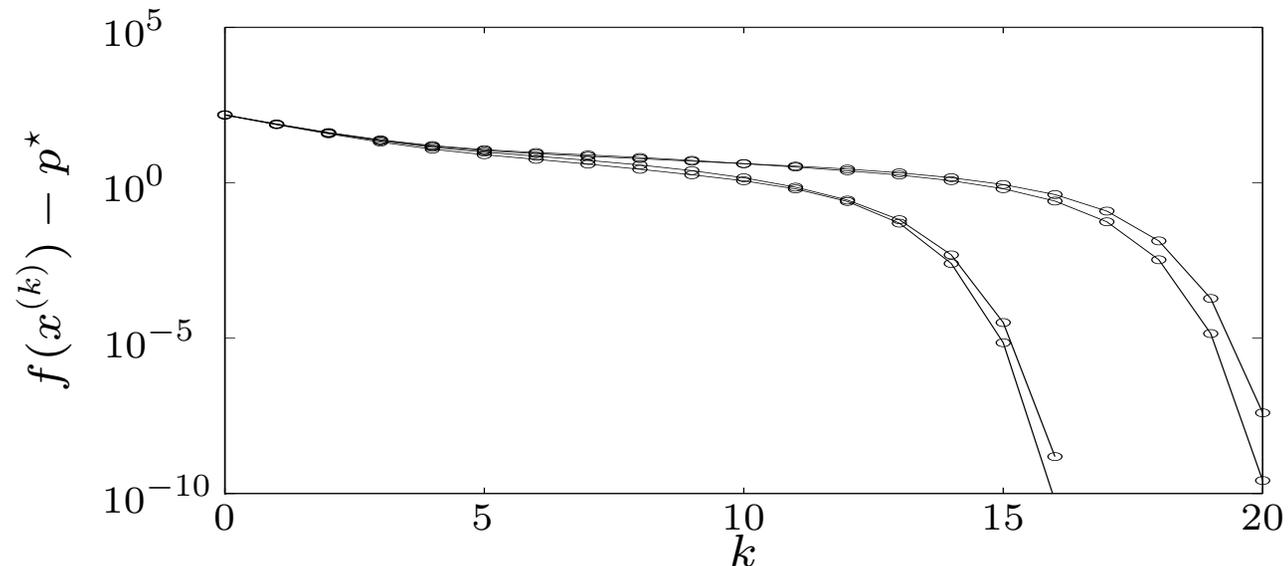with $Q \succeq 0$ for which $H + A^T Q A \succ 0$, and apply elimination

# Equality constrained analytic centering

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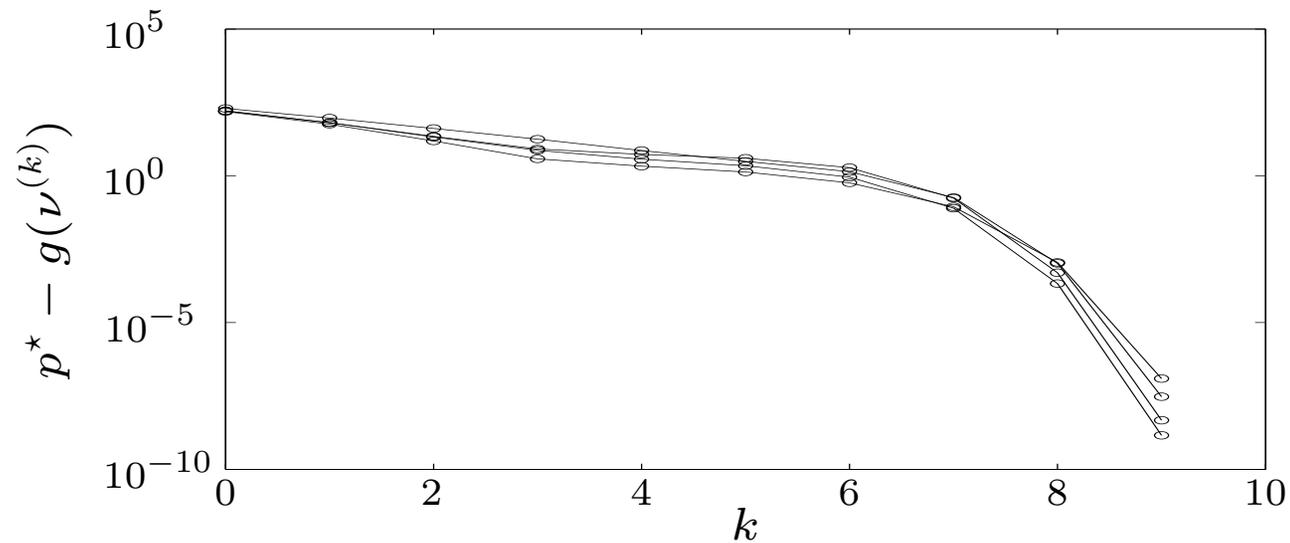
**dual problem:** maximize $-b^T \nu + \sum_{i=1}^{n} \log(A^T \nu)_i + n$

three methods for an example with $A \in \mathbb{R}^{100 \times 500}$, different starting points
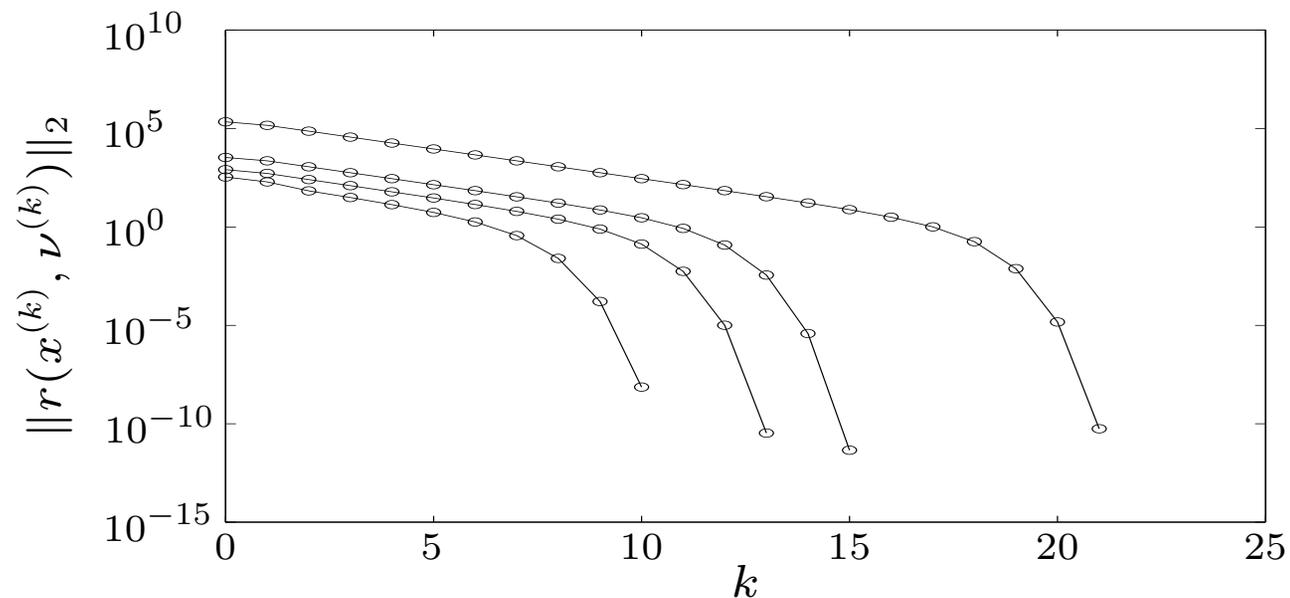
1. Newton method with equality constraints (requires $x^{(0)} \succ 0$, $Ax^{(0)} = b$)

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

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

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

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

3. use block elimination to solve KKT system

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

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 $\tilde{A} \in \mathbb{R}^{(p+1) \times n}$ defined as

$$
\tilde{A}_{ij} = \left\{ \begin{array}{rl}
1 & \text{arc } j \text{ leaves node } i \\
-1 & \text{arc } j \text{ enters node } i \\
0 & \text{otherwise}
\end{array} \right.
$$

■ reduced node-incidence matrix $A \in \mathbb{R}^{p \times n}$ is $\tilde{A}$ with last row removed

■ $b \in \mathbb{R}^p$ is (reduced) source vector

■ **Rank** $A = p$ if graph is connected

# KKT system

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

- $H = \mathbf{diag}(\phi_1''(x_1), \ldots, \phi_n''(x_n))$, positive diagonal

- solve via elimination:

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

sparsity pattern of coefficient matrix is given by graph connectivity

$$(AH^{-1}A^T)_{ij} \neq 0 \iff (AA^T)_{ij} \neq 0$$

$$\iff \text{nodes } i \text{ and } j \text{ are connected by an arc}$$

# Analytic center of linear matrix inequality

$$\begin{array}{ll} \text{minimize} & -\log\det X \\ \text{subject to} & \mathbf{Tr}(A_i X) = b_i, \quad i = 1, \ldots, p \end{array}$$

variable $X \in \mathbf{S}^n$

**optimality conditions**

$$X^\star \succ 0, \qquad -(X^\star)^{-1} + \sum_{j=1}^{p} \nu_j^\star A_i = 0, \qquad \mathbf{Tr}(A_i X^\star) = b_i, \quad i = 1, \ldots, p$$

**Newton equation at feasible $X$:**

$$X^{-1}\Delta X X^{-1} + \sum_{j=1}^{p} w_j A_i = X^{-1}, \qquad \mathbf{Tr}(A_i \Delta X) = 0, \quad i = 1, \ldots, p$$

- follows from linear approximation $(X + \Delta X)^{-1} \approx X^{-1} - X^{-1}\Delta X X^{-1}$
- $n(n+1)/2 + p$ variables $\Delta X$, $w$

## solution by block elimination

- eliminate $\Delta X$ from first equation: $\Delta X = X - \sum_{j=1}^{p} w_j X A_j X$

- substitute $\Delta X$ in second equation

$$\sum_{j=1}^{p} \mathbf{Tr}(A_i X A_j X) w_j = b_i, \quad i = 1, \dots, p \tag{2}$$

a dense positive definite set of linear equations with variable $w \in \mathbb{R}^p$

flop count (dominant terms) using Cholesky factorization $X = LL^T$:

- form $p$ products $L^T A_j L$: $(3/2)pn^3$

- form $p(p+1)/2$ inner products $\mathbf{Tr}((L^T A_i L)(L^T A_j L))$: $(1/2)p^2 n^2$

- solve (2) via Cholesky factorization: $(1/3)p^3$