

Nonlinear Optimization and Differential Equations

Jorge Nocedal

with Frank Curtis

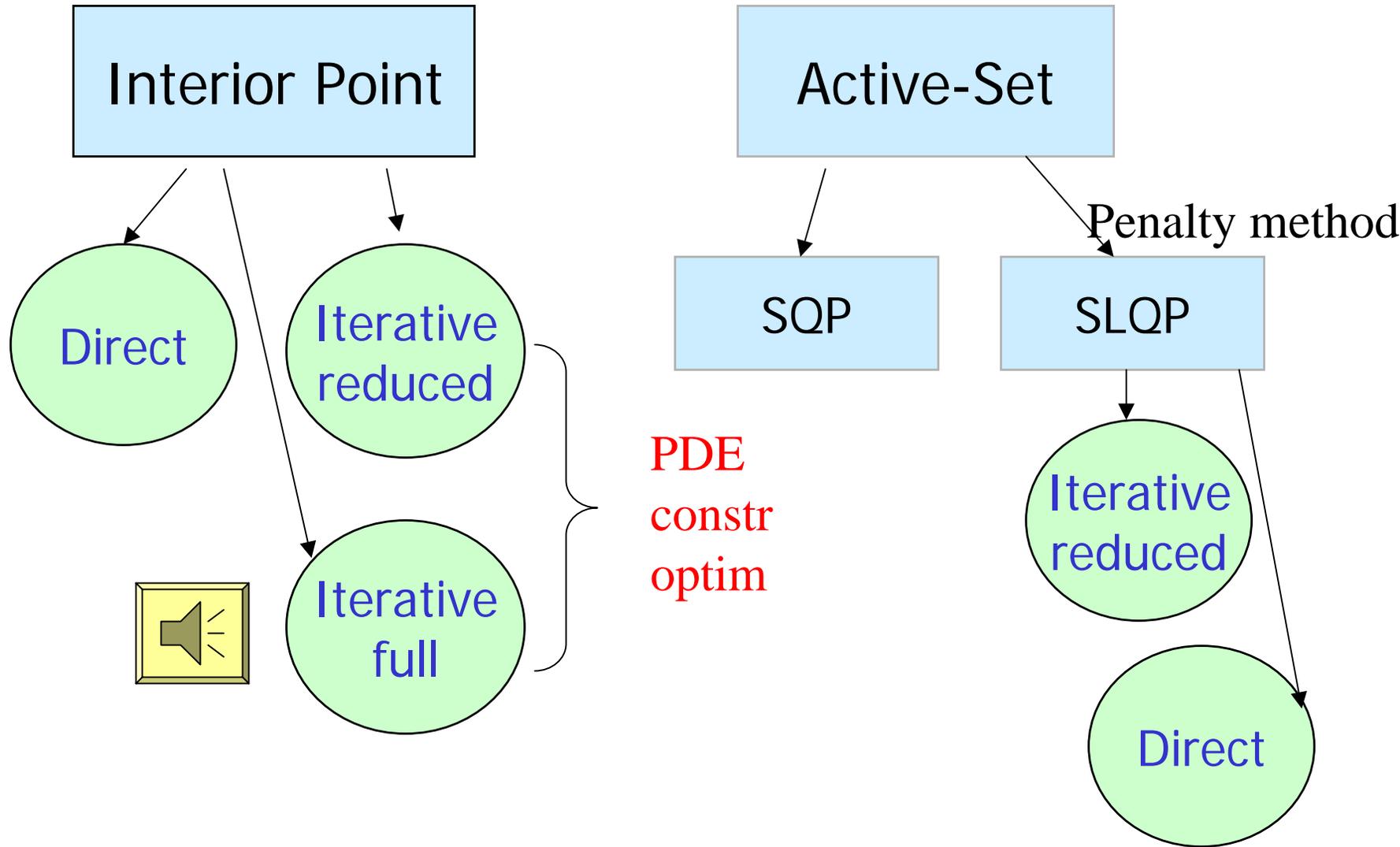
Northwestern University

Livermore, May 2007

Overview

- Review of *some* research contributions by the nonlinear optimization community
- New results: questions in the design of algorithms for PDE-optimization

Nonlinear Optimization, 2007



Looking back at progress made: An Example

- OPF: whose objective is reactive power injection in all network buses as a diagnosis of lack of reactive support in system
- Network represents the Brazilian high voltage generation/transmission system with about 3500 buses and 5000 circuits
- Nonlinear, nonconvex, written in AMPL

Number of variables:	14873
Number of nonlinear equality constraints:	6892
Number of nonzeros in Jacobian:	57971
Number of nonzeros in Hessian:	31501

SQP (SNOPT 7.2):

35 mins

Interior point (KNITRO 5.0):

30 seconds

Why?

- Interior-points much point better in these problems
- Exact second derivatives vs quasi-Newton
- Large reduced space
- Active set method with exact Hessians? See next



Main point: could not get this performance 10 years ago!

In operational models ($n=100,000$) Hessian **NOT** available
Quasi-Newton needed?

$$\begin{bmatrix} W & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} d \\ \delta \end{bmatrix} = - \begin{bmatrix} \nabla f + A^T \lambda \\ c \end{bmatrix}$$

Alternative to Quasi-Newton:

Compute Wd by finite difference of gradients

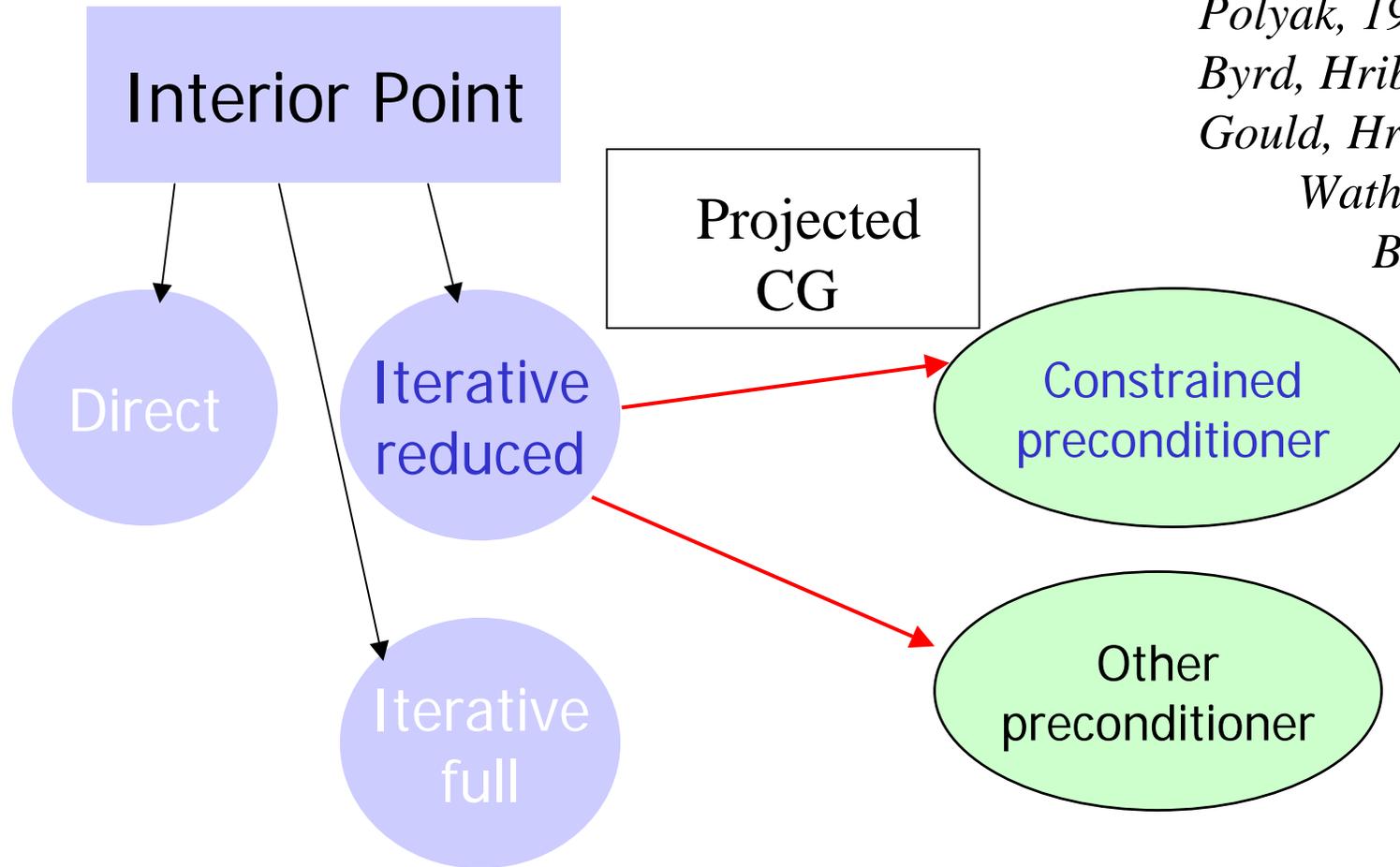
Iterative Method: reduced space interior point method

Projected CG with **constraint preconditioner** (to remove barrier ill conditioning)

Software implementation: 2001

$$\begin{bmatrix} D & A^T \\ A & 0 \end{bmatrix}$$

Reduced space iterative solve



Polyak, 1969

Byrd, Hribar, N 1999

Gould, Hribar, N. 2001

Wathen, Gould,

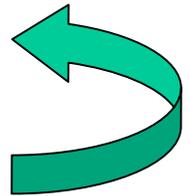
Benzi, Golub

Gill et al

Performance of various algorithms in KNITRO

Algorithm	Iterations	CPU (secs)
Interior/Default	53	32
Interior/Direct, barrule=4	29	17
Interior/Direct, LBFGS	185	350
Active Set/Hessian	***	> 20 mins (CPLEX)
Interior/CG, FinDiffHess	34	48

Projected CG with constraint preconditioner; **NO Hessian**



One more comment about recent advances
in nonlinear programming

There has been important research done in the last
five years on **exact** penalty functions

$$\phi(x) = f(x) + \pi \|c(x)\|_1$$

- Expands frontiers of NLP, addresses robustness
- Complementarity constraints (theory, algorithms)
Scheel-Scholtes, Anitescu, Ralph, Leyffer, Pang
- General NLP: new active set methods, *NU, Wright*
- General NLP: achieving robustness, *Chen-Goldfarb*
- Dynamic (non-heuristic) rules for updating penalty parameter, *Byrd, Waltz, Nocedal*
- Design of inexact Newton methods for PDE optimization (today!)

Anitescu et al., 00, 04, 07

Benson, Vanderbei, Shanno 04

Byrd, Nocedal, Waltz, 06

Chen and Goldfarb 05, 06

Fletcher and Chin 03

Gould, Orban, Toint 03

Hu and Ralph, 02

Leyffer, Lopez, Nocedal, 04

Scholtes et al 02, 05

Partial list of
references

PDE-Optimization

- Inexact Newton Methods for constrained optimization
- Negative Curvature
- Models of Penalty functions

For Simplicity: Equality Constrained Optimization

$$\begin{array}{ll} \min_x & f(x) \\ \text{s.t.} & c(x) = 0 \end{array}$$

No. of variables in the millions
Jacobian A not formed but

$$Av \quad A^T v$$

available

Nonlinear Elimination

Reduced Space

Step-Decomposition

Full Space (Primal-Dual) SQP ←

Geophysics, Meteorology

Biros-Ghattas, Haber,...

Heinkenschloss, Ridzal

SQP: Approximate Solution of

$$\begin{bmatrix} W & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} d \\ \delta \end{bmatrix} = - \begin{bmatrix} \nabla f + A^T \lambda \\ c \end{bmatrix}$$

$$x_{k+1} = x_k + \alpha_k d \quad \lambda_{k+1} = \lambda_k + \alpha_k \delta$$

Iterative Methods

- Symmetric QMR
- GMRES
- other ...

We impose no structure on the linear solver

When do we stop iterative solver?

Use non-smooth penalty function as a guide!

$$\Phi(x; \pi) = f(x) + \pi \|c(x)\|$$

Negative curvature?

Borrow from trust region methods

retro-1980s

Algorithm: Newton's method

$$\begin{bmatrix} W & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} d \\ \delta \end{bmatrix} = - \begin{bmatrix} \nabla f + A^T \lambda \\ c \end{bmatrix}$$

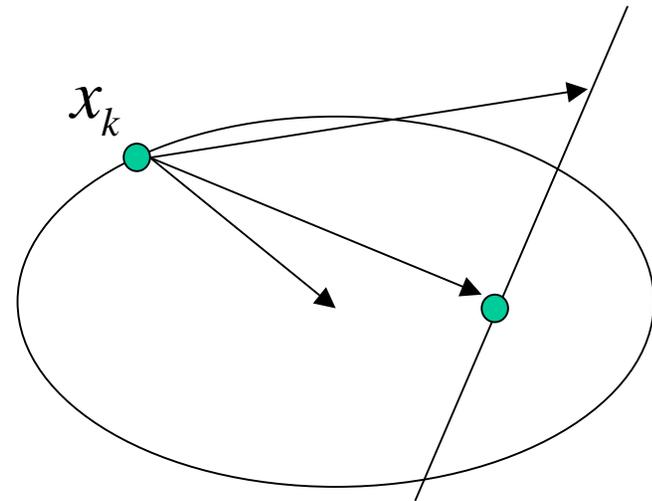
Algorithm: SQP

$$\begin{array}{ll} \min_d & \nabla f^T d + \frac{1}{2} d^T W d \\ \text{s.t.} & c + A d = 0 \end{array}$$

Question: can we ensure convergence with a

- step to constraints?
- step to reduce objective?

Preferably both, but if we can't do both?



(Heinkenschloss and Vicente, 2001)

Our approach

$$W > 0$$



$$\begin{bmatrix} W & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} d \\ \delta \end{bmatrix} = - \begin{bmatrix} \nabla f + A^T \lambda \\ c \end{bmatrix} + \begin{bmatrix} \rho \\ r \end{bmatrix}$$

control residual
components **separately**

Use a model of the merit function

$$m(d) = f + \nabla f^T d + \frac{1}{2} d^T W d + \pi (\|c + Ad\|)$$

to determine conditions for ρ and r

Require :

$$\Delta m(d) \geq 0.1\pi \max\{\|c\|, \|r\| - \|c\|\}$$

Model condition implies descent

Instead of imposing descent directly (KNITRO experience)

Algorithm Outline ($W > 0$)

- Iteratively solve

$$\begin{bmatrix} W & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} d \\ \delta \end{bmatrix} = - \begin{bmatrix} g + A^T \lambda \\ c \end{bmatrix} + \begin{bmatrix} \rho \\ r \end{bmatrix}$$

- until

$$\begin{aligned} \|r\| &\leq \varepsilon \|c\|, & 0 < \varepsilon < 1 \\ \|\rho\| &\leq \beta \|c\|, & 0 < \beta \end{aligned} \quad \text{or}$$

$$\begin{aligned} \|\rho\| &\leq \varepsilon \|g + A^T \lambda\|, & 0 < \varepsilon < 1 \\ \Delta m(d) &\geq 0.1\pi \max\{\|c\|, \|r\| - \|c\|\} \end{aligned}$$

- **Update** penalty parameter
- Perform backtracking line search
- Update iterate

If W is positive definite only on the null space of constraints

Only one change: modify the model

$$m(d) = f + \nabla f^T d + \frac{\omega}{2} d^T W d + \pi(\|c + Ad\|)$$
$$\omega = \begin{cases} 1 & \text{if } d^T W d \geq \theta \|d\|^2 \\ 0 & \text{otherwise} \end{cases}$$

- linear model for negative curvature direction

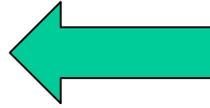
Global Convergence

Byrd, Curtis, N., 2006

$$\lim_{k \rightarrow \infty} \|c_k\| = 0$$
$$\lim_{k \rightarrow \infty} \|g_k + A_k^T \lambda_k\| = 0$$

W is positive definite only on the null space of constraints

Negative curvature



*W has negative eigenvalues
on null space of A*

$$\begin{bmatrix} W & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} d \\ \delta \end{bmatrix} = - \begin{bmatrix} \nabla f + A^T \lambda \\ c \end{bmatrix} + \begin{bmatrix} \rho \\ r \end{bmatrix}$$

Our approach

1. First identify conditions under which the step d is acceptable **even** if negative curvature is present
2. Introduce a modification of W if conditions cannot be fulfilled

$$(W + \gamma I)$$

Algorithm Sequential Model Reduction

Repeat until convergence

Set $\gamma = 0$

repeat

Apply 1 or more steps of linear solver

If Test 1 or Test 2 hold **break**

else increase γ

end repeat

update penalty

perform backtracking line search

End repeat

$$\begin{bmatrix} W + \gamma I & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} d \\ \delta \end{bmatrix} = - \begin{bmatrix} g + A^T \lambda \\ c \end{bmatrix} + \begin{bmatrix} \rho \\ r \end{bmatrix}$$

Model Reduction Condition

$$\Delta m(d) \geq (1 - \omega)\theta \|u\|^2 + 0.1\pi \max\{\|c\|, \|r\| - \|c\|\}$$

$\omega = 0,1$ $\theta =$ small constant

$u =$ approximation of tangential component of step

How?

Test II:

Model reduction condition

plus

$$\|\rho\| \leq \varepsilon \|g + A^T \lambda\|,$$

as before...

Test I:

$$\|r\| \leq \varepsilon \|c\|,$$

$$0 < \varepsilon < 1$$

$$\|\rho\| \leq \beta \|c\|,$$

$$0 < \beta$$

$$\|u\|^2 \leq \beta \|v\|^2 \quad \text{if } \omega = 0$$

How?

$d = u + v$ tangential, normal components

$$\|v\| \geq \|Ad\| / \|d\|$$

$$\|u\|^2 \leq \|d\|^2 - \|Ad\|^2 / \|d\|^2$$

Motivation, Justification

Is this all:

- Incremental?
- Recycled?
- Just plain weird?
- .. Or simply wrong headed!

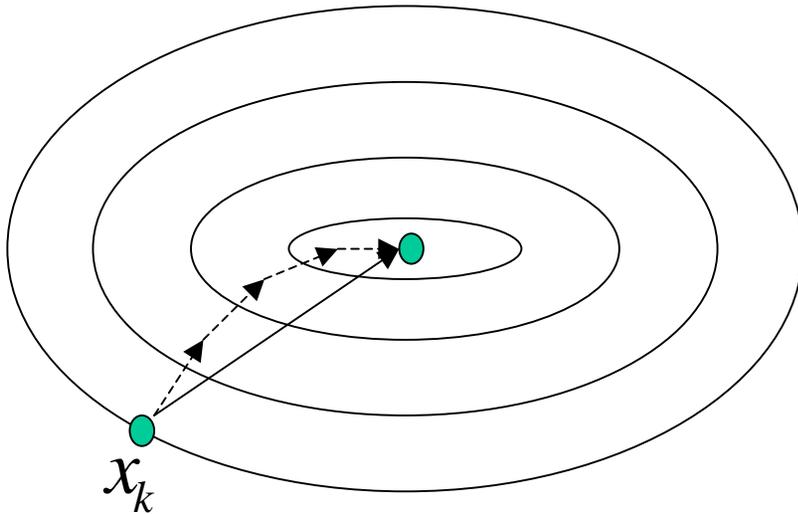
Unconstrained optimization - Inexactness

$$\min_x f(x)$$

Algorithm: Inexact Newton method (CG)

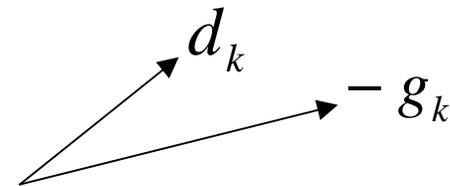
$$\nabla^2 f(x_k) d_k = -\nabla f(x_k)$$

Choosing *any* intermediate step ensures global convergence; sufficient (Cauchy) **reduction in model**



Suppose we **cannot use CG...**

Negative curvature:
angle condition:



If $\nabla^2 f_k > 0$ (positive def) inexactness condition:

$$\nabla^2 f_k d = -\nabla f_k + \rho \quad \|\rho_k\| \leq \varepsilon \|\nabla f_k\| \quad \varepsilon < 1$$

But this condition does not imply descent if f . Define model

$$m(d) = f + \nabla f^T d + d^T W d$$

and require

$$\Delta m(d) = m(0) - m(d) \geq 0$$

$$-\nabla f^T d - \frac{1}{2} d^T W d \geq 0 \quad \text{Easy}$$

Unconstrained optimization: Negative curvature

Exact Newton Method

$$\nabla^2 f(x_k) d_k = -\nabla f(x_k)$$

$$(\nabla^2 f(x_k) + \gamma I) d_k = -\nabla f(x_k)$$

Why modify so quickly? Step could point downhill (toward saddle point)

Crucial question for inexact Newton case

If Hessian not positive def:

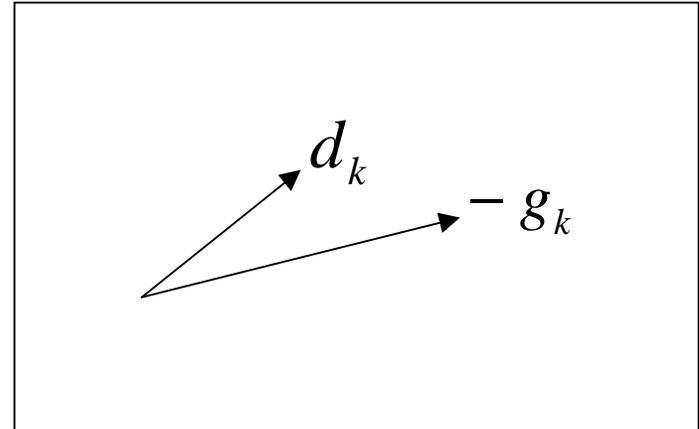
modified Cholesky

shift modification

That is, convex models

or

Trust region approach



Unconstrained optimization: Negative Curvature

$$\nabla^2 f_k d = -\nabla f \quad \text{Newton step}$$

Simple idea: step is acceptable if

$$d^T \nabla^2 f d \geq \theta \|d\|^2 \quad (*)$$

For scale invariance choose $\theta = 10^{-8} \|W\|$

- We prefer this to an angle test, which is not practical in the constrained setting
- Express (*) using a model

$$m(d) = f + \nabla f^T d + \frac{\omega}{2} d^T W d \quad \omega = 0,1$$

Unconstrained optimization: Model Reduction

Focus on model reduction, not spectrum

Form of model depends
on the step d

The new model for Newton's method

$$m(d) = f + \nabla f^T d + \frac{\omega}{2} d^T W d$$

$$\omega = \begin{cases} 1 & \text{if } d^T W d \geq \theta \|d\|^2 \\ 0 & \text{otherwise} \end{cases}$$

Model Reduction Condition: Step d is acceptable if

$$\Delta m(d) = m(0) - m(d) \geq (1 - \omega)\theta \|d\|^2$$

If not acceptable, modify W , or add trust region, or...

Model reduction condition ensures descent

A Numerical Test:

Exact Newton step, negative curvature

Algorithm I

$$(\nabla^2 f(x_k) + \gamma I)d_k = -\nabla f(x_k) \quad (\text{loop})$$

with $\gamma > 0$ whenever $\text{neig}(\nabla^2 f) > 0$ (margin)

Perform line search

Algorithm II

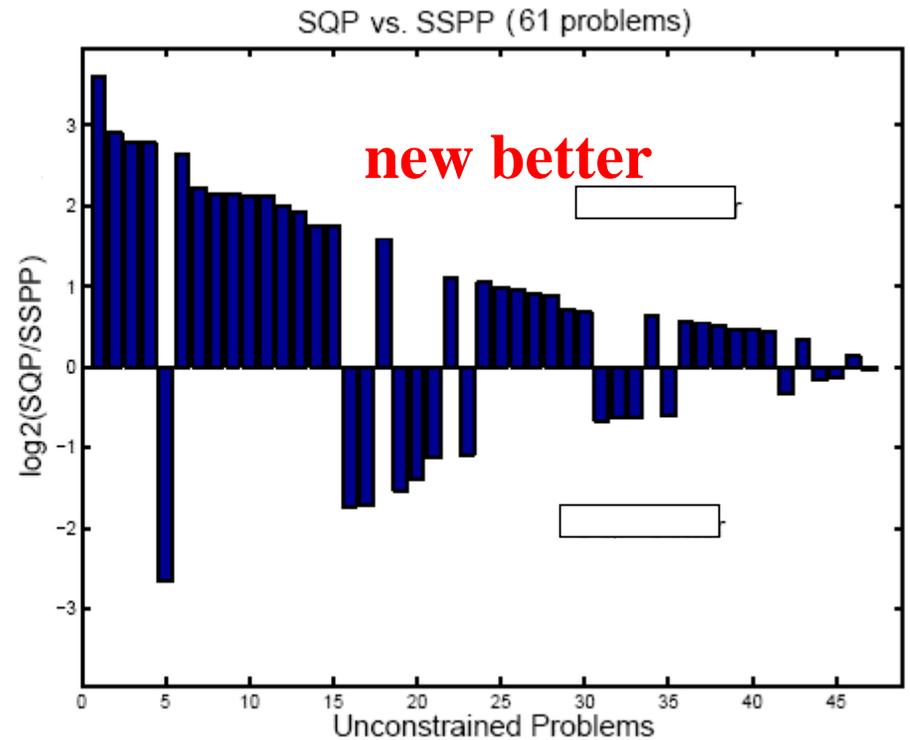
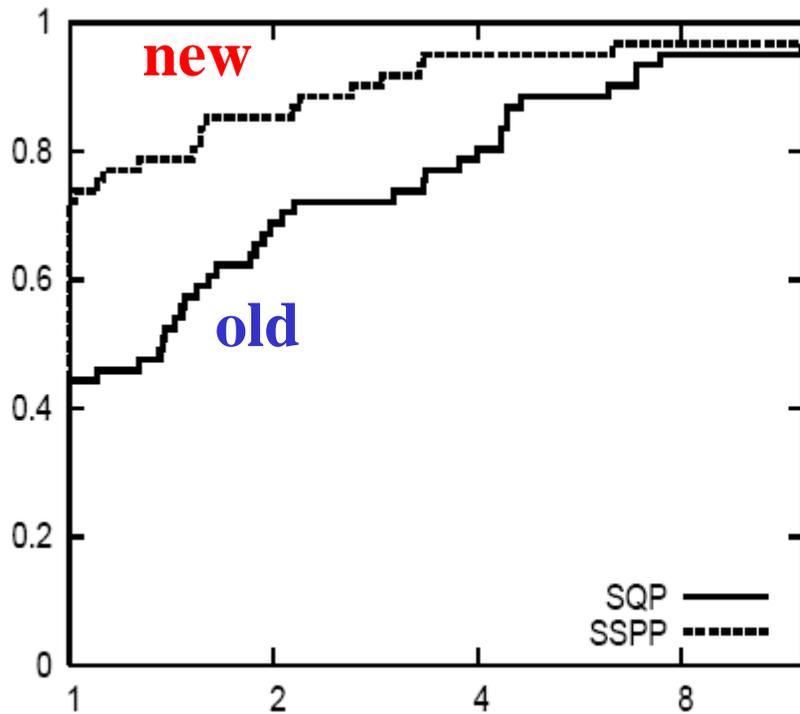
Shift only if model decrease

condition does not hold

$$\Delta m(d) \geq (1 - \omega)10^{-6} \|d\|^2$$

(loop)

Number of factorizations (number of iterations similar for both approaches)
71 problems CUTeR, COPS



Performance profiles for matrix factorizations

Repeat experiments with **iterative solution** of Newton equations
Inertia information not available

$$d^T \nabla^2 f d \geq \theta \|d\|^2$$

Use QMR and SYMQR

Clear-cut advantage of model reduction approach

Thanks to:

Richard Byrd

Eldad Haber

Richard Waltz

Todd Plantenga

Constrained Optimization

$W > 0$ A full rank

Algorithm: Newton's method

$$\begin{bmatrix} W & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} d \\ \delta \end{bmatrix} = - \begin{bmatrix} \nabla f + A^T \lambda \\ c \end{bmatrix}$$

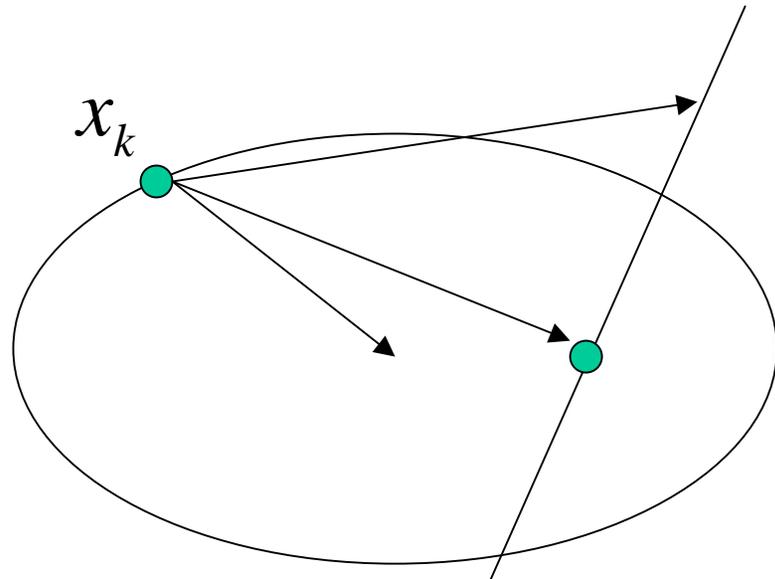
Algorithm: SQP

$$\begin{aligned} \min_d \quad & \nabla f^T d + \frac{1}{2} d^T W d \\ \text{s.t.} \quad & c + A d = 0 \end{aligned}$$

Question: can we ensure convergence with a

- step to constraints?
- step to reduce objective?

Preferably both, but if we can't do both?



(Heinkenschloss and Vicente, 2001)

Model Decrease Condition

$$\Phi(x; \pi) = f(x) + \pi \|c(x)\|$$

$$\begin{array}{ll} \min_d & \nabla f^T d + \frac{1}{2} d^T W d \\ \text{s.t.} & c + Ad = 0 \end{array}$$

$$m(d) = f + \nabla f^T d + \frac{1}{2} d^T W d + \pi (\|c + Ad\|)$$

Quantify reduction obtained from step

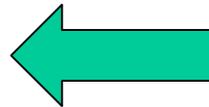
$$\Delta(d) = m(0) - m(d)$$

$$= -\nabla f^T d - \frac{1}{2} d^T W d + \pi (\|c\| - \|c + Ad\|)$$

\curvearrowright
 $cred(d)$

Require :

$$\Delta m(d) \geq 0.1\pi \|cred\|$$



Algorithm Outline

- Iteratively solve

$$\begin{bmatrix} W & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} d \\ \delta \end{bmatrix} = - \begin{bmatrix} g + A^T \lambda \\ c \end{bmatrix} + \begin{bmatrix} \rho \\ r \end{bmatrix}$$

- until

$$\begin{aligned} \|\rho\| &\leq \varepsilon \|g + A^T \lambda\|, & 0 < \varepsilon < 1 \\ \Delta m(d) &\geq \sigma \pi \|c\| \end{aligned}$$

or

$$\begin{aligned} \|r\| &\leq \varepsilon \|c\|, & 0 < \varepsilon < 1 \\ \|\rho\| &\leq \beta \|c\|, & 0 < \beta \end{aligned}$$

- Update penalty parameter
- Perform backtracking line search
- Update iterate