



# Derivative-free and Simulation-based Optimization

Stefan Wild

Argonne National Laboratory  
Mathematics and Computer Science Division

January 30, 2014

# The Plan

## Motivation

### Black-box Optimization

Random Search

Direct Search Methods

Model-Based Methods

Some Global Optimization

### Simulation-Based Optimization and Structure

Nonlinear Least Squares

Least Squares and Partials

Constraints

Nonsmoothness

### Computational Noise

Estimating Computational Noise

Noisy Finite Differences



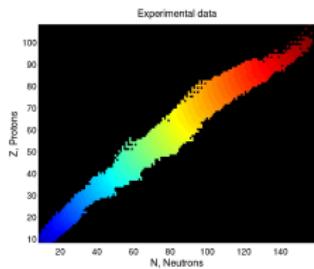
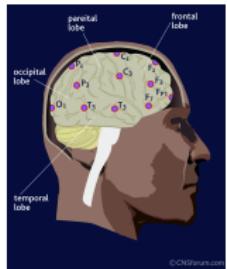
# Simulation-Based Optimization

$$\min_{x \in \mathbb{R}^n} \{f(x) = F[S(x)] : c(S(x)) \leq 0\}$$

- ◊  $S$  (numerical) simulation output, often “noisy” (even when deterministic)
- ◊ Derivatives  $\nabla_x S$  often unavailable or prohibitively expensive to obtain/approximate directly
- ◊  $S$  can contribute to objective and/or constraints
- ◊ Single evaluation of  $S$  could take seconds/minutes/hours/days

Evaluation is a bottleneck for optimization

Functions of complex numerical simulations arise everywhere



# Computing is Responsible for Pervasiveness of Simulations in Sci&Eng



Argonne's Blue Gene/P (2008: 163,840 cores)

Currently 67th fastest in the world



Argonne's Blue Gene/Q (2012: 786,432 cores)

Currently 5th fastest in the world

- ◊ Parallel/multi-core environments increasingly common
  - ◆ Small clusters/multi-core desktops/multi-core laptops pervasive
  - ◆ Leadership class machines increasingly parallel
- ◊ Simulations (the “forward problem”) become faster/more realistic/more complex

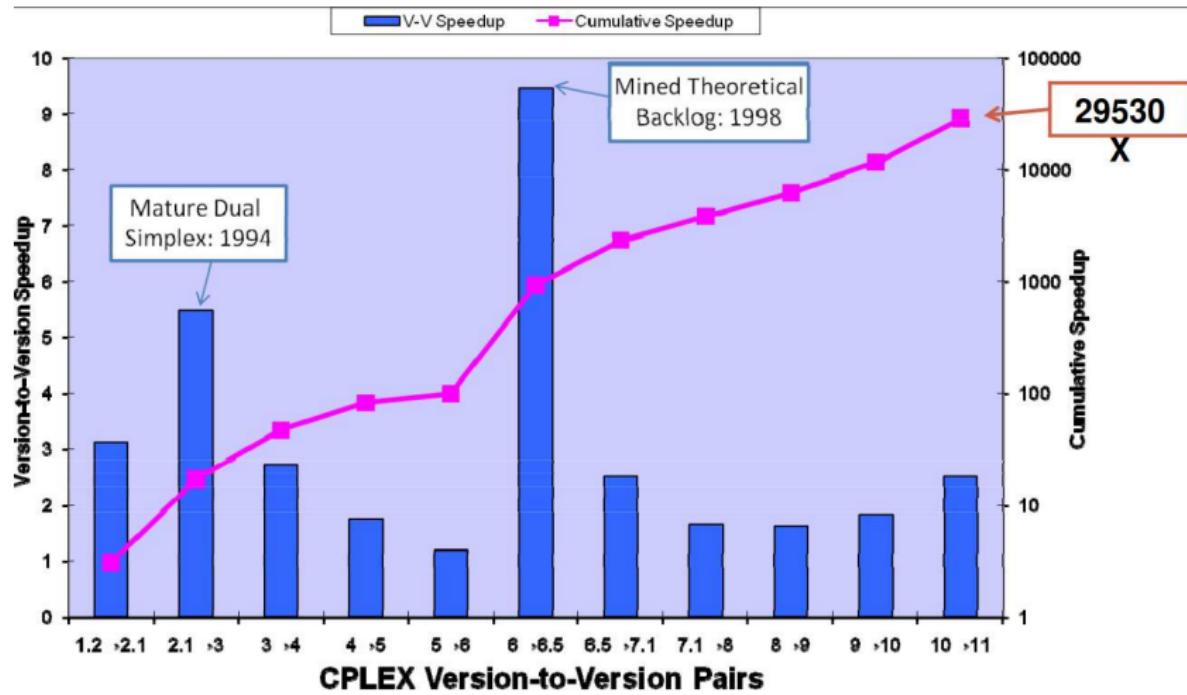
## Improvements from Algorithms Can Trump Those From Hardware

Martin Grötschel's production planning benchmark problem (a MIP):

1988	solve time using current computers and LP algorithms:	82 years
2003	solve time using current computers and LP algorithms:	1 minute

- ◊ Speed up of **43,000,000X**  
 $10^3$ X from processor improvements  
 $10^4$ X additional from **algorithmic improvements**

# Improvements from Algorithms Can Trump Those From Hardware



1991 (v1.2) to 2007 (v11.0): Moore's Law transistor speedup:  $\approx 256X$

[Slide from Bixby (CPLEX/GUROBI)]: Solves 1,852 MIPs

## Derivative-Free Optimization

"Some derivatives are unavailable for optimization purposes"

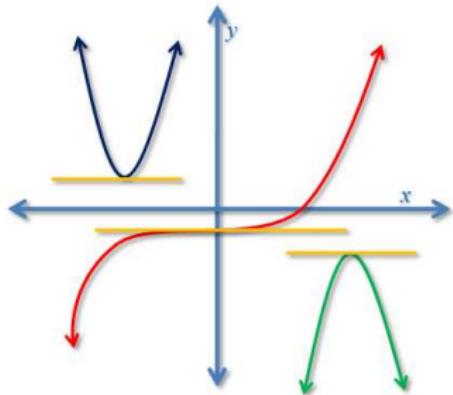
# Derivative-Free Optimization

"Some derivatives are unavailable for optimization purposes"

The Challenge: Optimization is tightly coupled with derivatives

Typical optimality (no noise, smooth functions)

$$\nabla_x f(x_*) + \lambda^T \nabla_x c_E(x_*) = 0, c_E(x_*) = 0$$



(sub)gradients  $\nabla_x f$ ,  $\nabla_x c$  enable:

- ◊ Faster feasibility
- ◊ Faster convergence
  - ◆ Guaranteed descent
  - ◆ Approximation of nonlinearities
- ◊ Better termination
  - ◆ Measure of criticality  
 $\|\nabla_x f\|$  or  $\|\mathcal{P}_\Omega(\nabla_x f)\|$
- ◊ Sensitivity analysis
  - ◆ Correlations, standard errors, UQ, ...

# Ways to Get Derivatives

(assuming they exist)

## Handcoding (HC)

"Army of students/programmers"

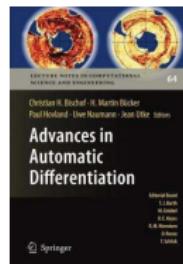
- ? Prone to errors/conditioning
- ? Intractable as number of ops increases



## Algorithmic/Automatic Differentiation (AD)

"Exact\* derivatives!"

- ? No black boxes allowed
- ? Not always automatic/cheap/well-conditioned



## Finite Differences (FD)

"Nonintrusive"

- ? Sensitive to stepsize choice/noise
- ? Expense grows with  $n$



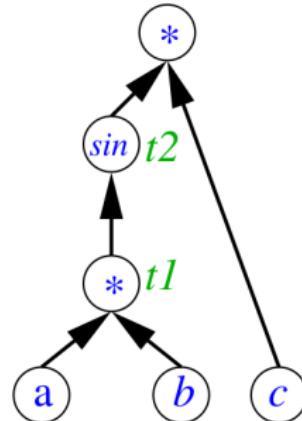
... then apply derivative-based method (that handles inexact derivatives)

# Algorithmic Differentiation

## Computational Graph

- ◊  $y = \sin(a * b) * c$
- ◊ Forward and reverse modes
- ◊ AD tool provides code for your derivatives

You should write codes and formulate problems with AD in mind!



Many tools (see [www.autodiff.org](http://www.autodiff.org)):

F/C [Tapenade](#), [Rapsodia](#)

C/C++ [ADOL-C](#), [ADIC](#)

Matlab [ADiMat](#), [INTLAB](#)

Also done in [AMPL](#) and [GAMS](#)!

## Caveats

You will pay a price for not having derivatives

We will focus primarily on:

- ◊ **Minimization** (sorry, economists, you'll need to stand on your head)
  - ◊ **Continuous** domains
  - ◊ **Unconstrained** problems
    - ◆ Simulation-based constraints bring additional challenges
  - ◊ Underlying **smooth** behavior (whatever that means)
  - ◊ **Deterministic** objective (or function of pseudorandom number generators)
  - ◊ Algorithms with convergence guarantees to **local** minima
    - ◆ GAs and other heuristics often require far too many evaluations
- Can address these in detail during lunch, dinner, TA/office hours

## Caveats

You will pay a price for not having derivatives

We will focus primarily on:

- ◊ **Minimization** (sorry, economists, you'll need to stand on your head)
  - ◊ **Continuous** domains
  - ◊ **Unconstrained** problems
    - ◆ Simulation-based constraints bring additional challenges
  - ◊ Underlying **smooth** behavior (whatever that means)
  - ◊ **Deterministic** objective (or function of pseudorandom number generators)
  - ◊ Algorithms with convergence guarantees to **local** minima
    - ◆ GAs and other heuristics often require far too many evaluations
- Can address these in detail during lunch, dinner, TA/office hours

Please: Interrupt, Ask Questions

## Global Optimization, $\min_{x \in \Omega} f(x)$

Careful:

- ◊ **Global convergence:** Convergence (to a local solution/stationary point) from anywhere in  $\Omega$
  - ◊ **Convergence to a global minimizer:** Obtain  $x_*$  with  $f(x_*) \leq f(x) \forall x \in \Omega$
-

## Global Optimization, $\min_{x \in \Omega} f(x)$

Careful:

- ◊ **Global convergence:** Convergence (to a local solution/stationary point) from anywhere in  $\Omega$
  - ◊ **Convergence to a global minimizer:** Obtain  $x_*$  with  $f(x_*) \leq f(x) \forall x \in \Omega$
- 

**Anyone selling you global solutions when derivatives are unavailable:**

either assumes more about your problem (e.g., convex  $f$ )

or expects you to wait forever

**Törn and Žilinskas:** An algorithm converges to the global minimum for any continuous  $f$  if and only if the sequence of points visited by the algorithm is dense in  $\Omega$ .

or cannot be trusted

---

## Global Optimization, $\min_{x \in \Omega} f(x)$

Careful:

- ◊ **Global convergence:** Convergence (to a local solution/stationary point) from anywhere in  $\Omega$
  - ◊ **Convergence to a global minimizer:** Obtain  $x_*$  with  $f(x_*) \leq f(x) \forall x \in \Omega$
- 

**Anyone selling you global solutions when derivatives are unavailable:**

either assumes more about your problem (e.g., convex  $f$ )

or expects you to wait forever

**Törn and Žilinskas:** An algorithm converges to the global minimum for any continuous  $f$  if and only if the sequence of points visited by the algorithm is dense in  $\Omega$ .

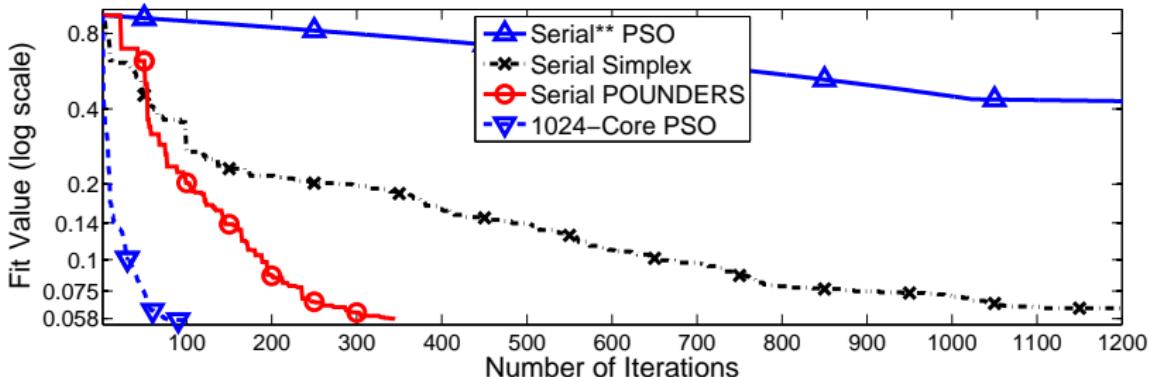
or cannot be trusted

---

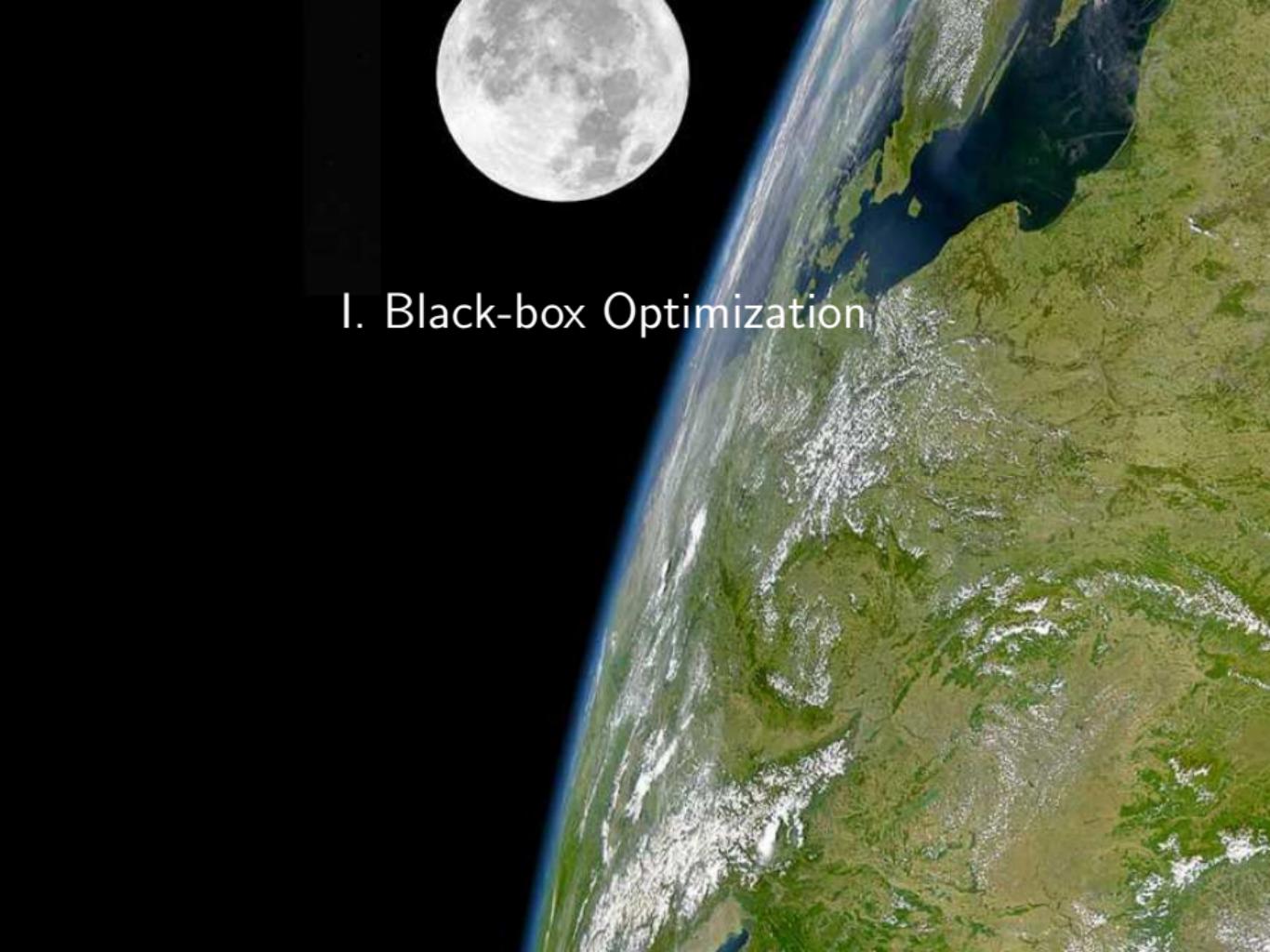
Instead:

- ◊ Rapidly find good local solutions and/or be robust to poor solutions
- ◊ Consider multistart approaches and/or structure of multimodality

## (One Reason) Why We Won't Be Talking About Heuristics



- ◊ Heuristics often “embarrassingly/naturally parallel”; PSO = particle swarm method
  - ◆ Typically through stochastic sampling/evolution
  - ◆ 1024 function evaluations per iteration
- ◊ Simplex is Nelder-Mead; POUNDERS is model-based trust-region algorithm
  - ◆ one function evaluation per iteration
- Is this an effective use of resources?
- How many cores would have sufficed?

A composite image featuring a large, detailed full moon in the upper left quadrant against a black background. The rest of the image shows a high-resolution satellite view of Earth's surface, focusing on the continent of Asia. The blue oceans and green continents are visible, with clouds appearing as white wisps.

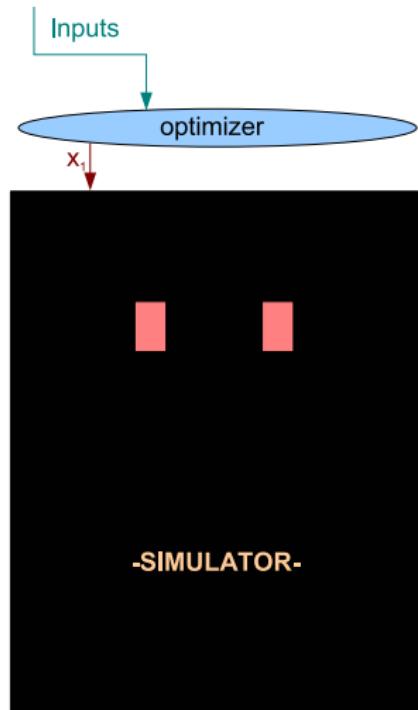
# I. Black-box Optimization

## Black-box Optimization Problems

NB- My “black box” is different from Ken’s “black box”

# Black-box Optimization Problems

NB- My “black box” is different from Ken’s “black box”



Only knowledge about  $f$  is obtained by sampling

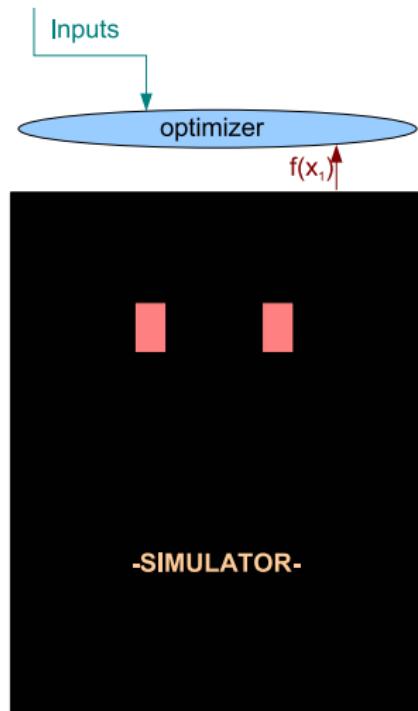
- ◊  $f = S$  a black box (running some executable-only code or performing an experiment in the lab)
- ◊ Only give a single output (no derivatives  $\nabla_x S(x)$ )

Good solutions guaranteed in the limit, but:

- ◊ Usually have computational budget (due to scheduling, finances, deadlines)
- ◊ Limited number of evaluations

# Black-box Optimization Problems

NB- My “black box” is different from Ken’s “black box”



Only knowledge about  $f$  is obtained by sampling

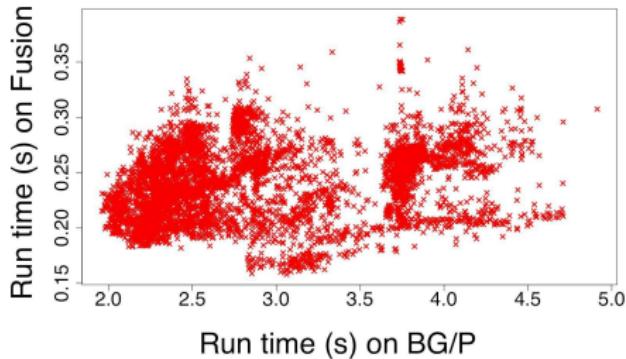
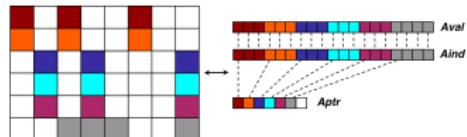
- ◊  $f = S$  a black box (running some executable-only code or performing an experiment in the lab)
- ◊ Only give a single output (no derivatives  $\nabla_x S(x)$ )

Good solutions guaranteed in the limit, but:

- ◊ Usually have computational budget (due to scheduling, finances, deadlines)
- ◊ Limited number of evaluations

# A Black Box: Automating Empirical Performance Tuning

Given semantically equivalent codes  
 $\mathcal{C}_1, \mathcal{C}_2, \dots$ , minimize run time subject to  
energy consumption



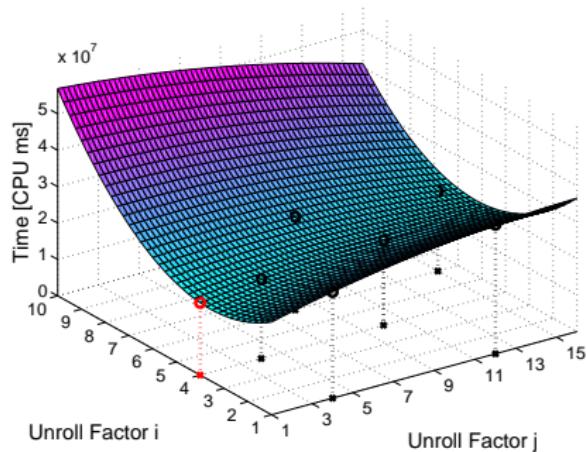
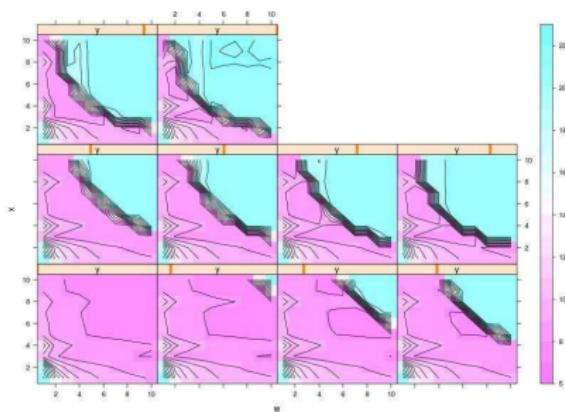
$$\min \{f(x) : (x_C, x_I, x_B) \in \Omega_C \times \Omega_I \times \Omega_B\}$$

- $x$  multidimensional parameterization (internal tolerances, unroll/tiling factors, compiler flags, compiler type, ...)
- $\Omega$  search domain (feasible transformation, no errors)
- $f$  quantifiable performance objective (requires a run)

# Optimization for Automatic Tuning of HPC Codes

Evaluation of  $f$  requires:

transforming source, compilation, (repeated?) execution, checking for correctness



## Challenges:

- Evaluating  $f(\Omega)$  prohibitively expensive (e.g.,  $10^{19}$  discrete decisions)
- $f$  noisy
- Discrete  $x$  unrelaxable
- $\nabla_x f$  unavailable/nonexistent
- Many distinct/local solutions

Solve general problems  $\min\{f(x) : x \in \mathbb{R}^n\}$ :

- ◊ Only require function values (no  $\nabla f(x)$ )
- ◊ Don't rely on finite-difference approximations to  $\nabla f(x)$
- ◊ Seek greedy and rapid decrease of function value
- ◊ Have asymptotic convergence guarantees
- ◊ Assume parallel resources are used within function evaluation

Main styles of DFO algorithms

- ◊ Random search methods
- ◊ Direct search methods (pattern search, Nelder-Mead, ...)
- ◊ Model-based methods (quadratics, radial basis functions, ...)

## General iteration

1. Generate direction  $d_k$
2. Evaluate gradient-free oracle  $g(x_k; h_k) = \frac{f(x_k + h_k d_k) - f(x_k)}{h_k} d_k$
3. Compute  $x_{k+1} = x_k - \delta_k g(x_k; h_k)$ , evaluate  $f(x_{k+1})$

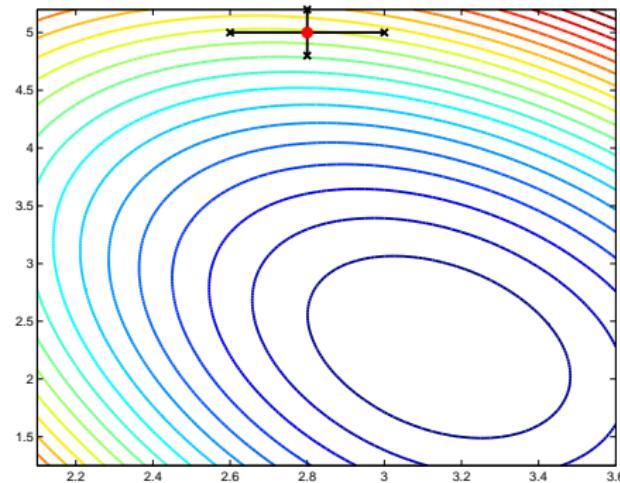
## Ex.- Nesterov's Random Gradient method

- ◊  $d_k$  Gaussian
- ◊  $h_k = h \leq \frac{5}{3(n+4)} \sqrt{\frac{\epsilon}{2L_1(f)}}$  to get accuracy  $\epsilon$  ( $= 10^{-5}$ )
- ◊  $\delta_k = \delta = \frac{1}{4(n+4)L_1(f)}$
- Requires Lipschitz constant  $L_1(f)$
- + Probabilistic guarantees in several cases (convex  $f$ , smooth  $f$ , ...)

# Pattern Search Methods

Choose a set of directions (pattern or mesh)  $\mathcal{D}_k$

Ex.-  $\pm$  coordinate directions ( $2n$  directions)



Basic iteration ( $k \geq 0$ ):

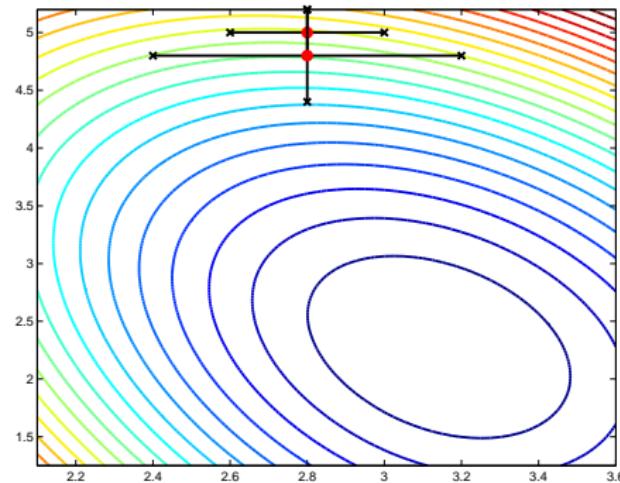
- ◊ Evaluate  $f(x_k + \Delta_k d_j)$ ,  $j = 1, \dots, |\mathcal{D}_k|$
- ◊ If  $[f(x_k + \Delta_k d_j) < f(x_k)]$ ,  
move to  $x_{k+1} = x_k + \Delta_k d_j$
- Otherwise shrink  $\Delta_k$
- ◊ Update  $\mathcal{D}_k$

This is an indicator function, does not say anything about the magnitude of  $f$  values, just the ordering

# Pattern Search Methods

Choose a set of directions (**pattern** or **mesh**)  $\mathcal{D}_k$

**Ex.-**  $\pm$  coordinate directions ( $2n$  directions)



Basic iteration ( $k \geq 0$ ):

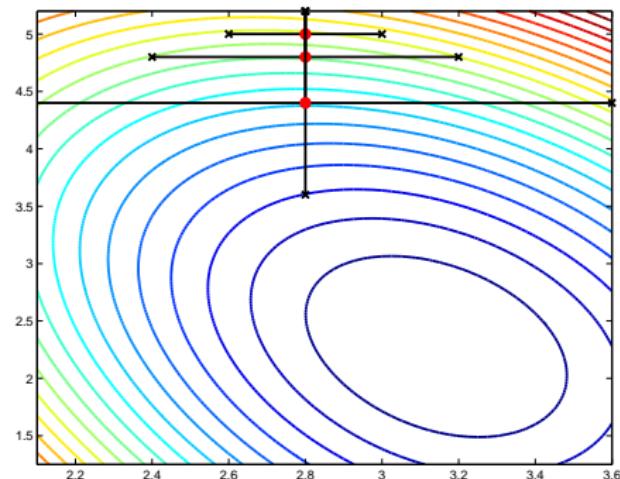
- ◊ Evaluate  $f(x_k + \Delta_k d_j)$ ,  $j = 1, \dots, |\mathcal{D}_k|$
- ◊ If  $[f(x_k + \Delta_k d_j) < f(x_k)]$ ,  
move to  $x_{k+1} = x_k + \Delta_k d_j$
- Otherwise shrink  $\Delta_k$
- ◊ Update  $\mathcal{D}_k$

This is an **indicator** function, does not say anything about the magnitude of  $f$  values, just the ordering

# Pattern Search Methods

Choose a set of directions (pattern or mesh)  $\mathcal{D}_k$

Ex.-  $\pm$  coordinate directions ( $2n$  directions)



Basic iteration ( $k \geq 0$ ):

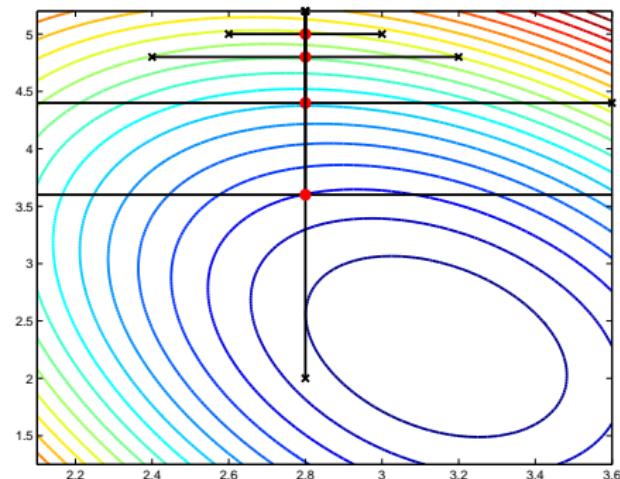
- ◊ Evaluate  $f(x_k + \Delta_k d_j)$ ,  $j = 1, \dots, |\mathcal{D}_k|$
- ◊ If  $[f(x_k + \Delta_k d_j) < f(x_k)]$ ,  
move to  $x_{k+1} = x_k + \Delta_k d_j$
- Otherwise shrink  $\Delta_k$
- ◊ Update  $\mathcal{D}_k$

This is an indicator function, does not say anything about the magnitude of  $f$  values, just the ordering

# Pattern Search Methods

Choose a set of directions (pattern or mesh)  $\mathcal{D}_k$

Ex.-  $\pm$  coordinate directions ( $2n$  directions)



Basic iteration ( $k \geq 0$ ):

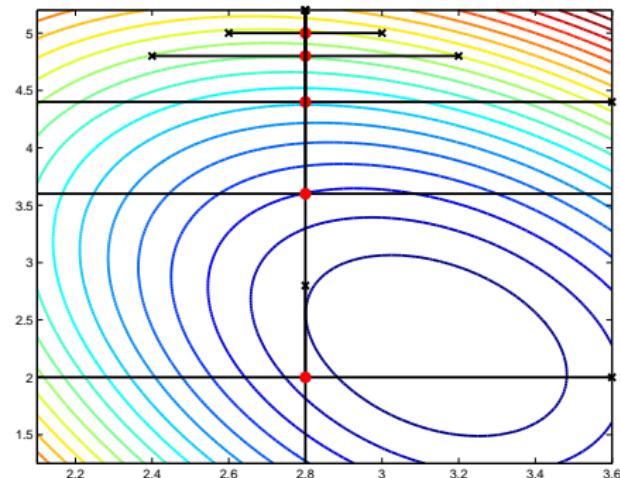
- ◊ Evaluate  $f(x_k + \Delta_k d_j)$ ,  $j = 1, \dots, |\mathcal{D}_k|$
- ◊ If  $[f(x_k + \Delta_k d_j) < f(x_k)]$ ,  
move to  $x_{k+1} = x_k + \Delta_k d_j$
- Otherwise shrink  $\Delta_k$
- ◊ Update  $\mathcal{D}_k$

This is an indicator function, does not say anything about the magnitude of  $f$  values, just the ordering

# Pattern Search Methods

Choose a set of directions (pattern or mesh)  $\mathcal{D}_k$

Ex.-  $\pm$  coordinate directions ( $2n$  directions)



Basic iteration ( $k \geq 0$ ):

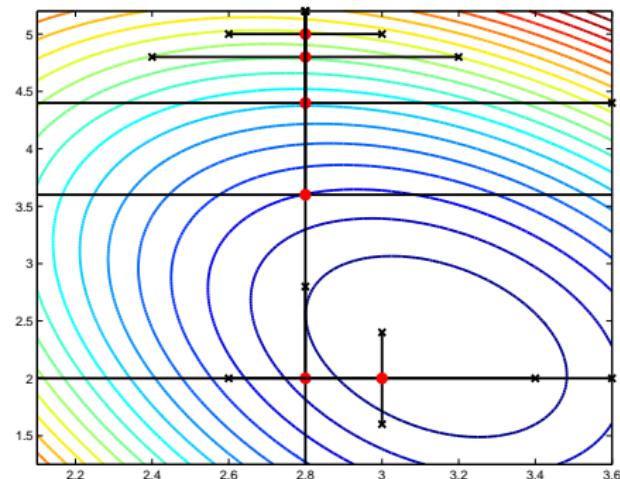
- ◊ Evaluate  $f(x_k + \Delta_k d_j)$ ,  $j = 1, \dots, |\mathcal{D}_k|$
- ◊ If  $[f(x_k + \Delta_k d_j) < f(x_k)]$ ,  
move to  $x_{k+1} = x_k + \Delta_k d_j$
- Otherwise shrink  $\Delta_k$
- ◊ Update  $\mathcal{D}_k$

This is an indicator function, does not say anything about the magnitude of  $f$  values, just the ordering

# Pattern Search Methods

Choose a set of directions (pattern or mesh)  $\mathcal{D}_k$

Ex.-  $\pm$  coordinate directions ( $2n$  directions)



Basic iteration ( $k \geq 0$ ):

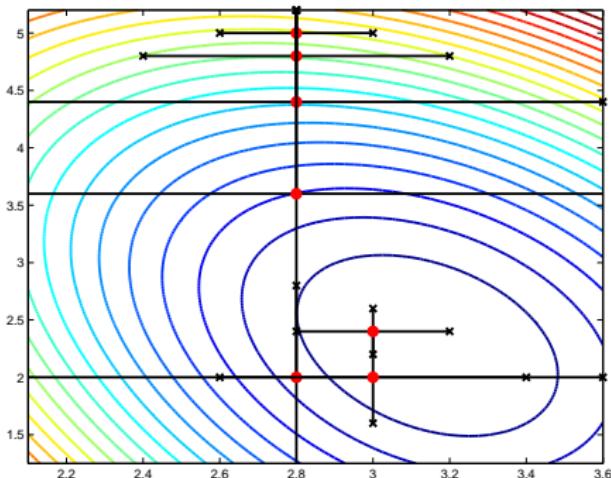
- ◊ Evaluate  $f(x_k + \Delta_k d_j)$ ,  $j = 1, \dots, |\mathcal{D}_k|$
- ◊ If  $[f(x_k + \Delta_k d_j) < f(x_k)]$ ,  
move to  $x_{k+1} = x_k + \Delta_k d_j$
- Otherwise shrink  $\Delta_k$
- ◊ Update  $\mathcal{D}_k$

This is an indicator function, does not say anything about the magnitude of  $f$  values, just the ordering

# Pattern Search Methods

Choose a set of directions (pattern or mesh)  $\mathcal{D}_k$

Ex.-  $\pm$  coordinate directions ( $2n$  directions)



Basic iteration ( $k \geq 0$ ):

- ◊ Evaluate  $f(x_k + \Delta_k d_j)$ ,  $j = 1, \dots, |\mathcal{D}_k|$
- ◊ If  $[f(x_k + \Delta_k d_j) < f(x_k)]$ ,  
move to  $x_{k+1} = x_k + \Delta_k d_j$
- Otherwise shrink  $\Delta_k$
- ◊ Update  $\mathcal{D}_k$

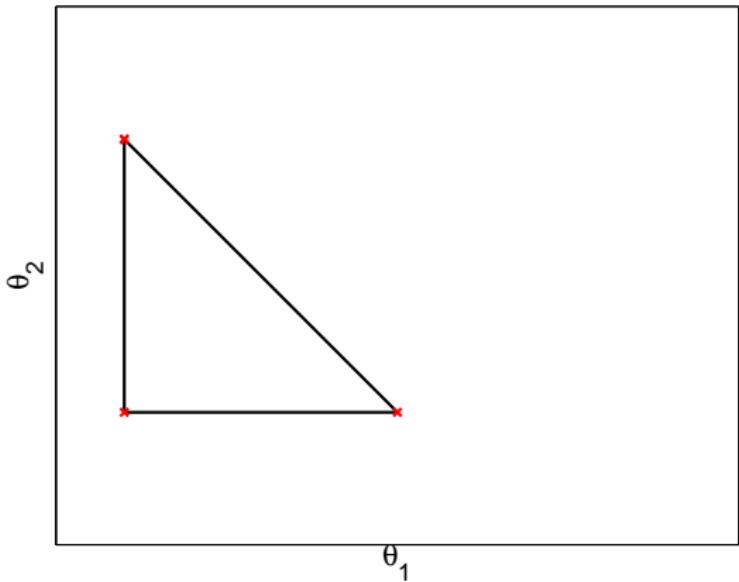
This is an indicator function, does not say anything about the magnitude of  $f$  values, just the ordering

- Lends itself well to doing concurrent function evaluations
- See also mesh-adaptive direct search methods (e.g., NOMAD)
- Can establish convergence for nonsmooth  $f$

# The Nelder-Mead Method [1965]

Basic iteration ( $k \geq 0$ ):

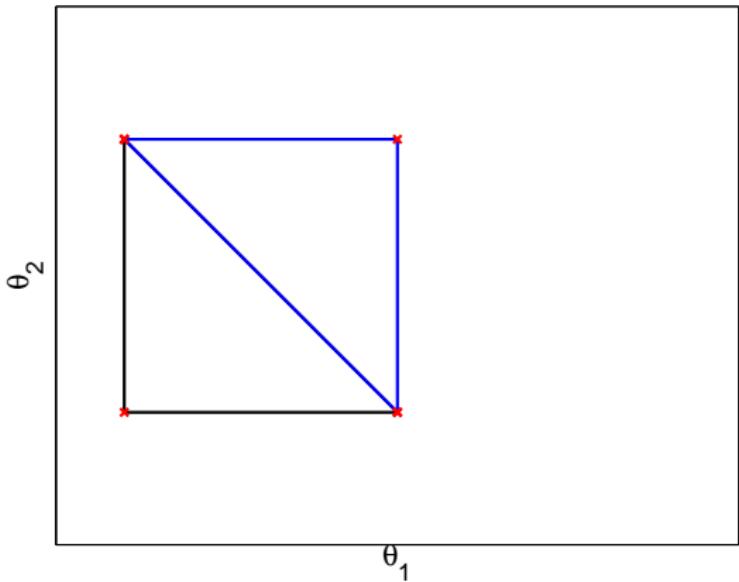
- ◊ Evaluate  $f$  on the  $n + 1$  vertices of the simplex  $x_k + \Delta_k \mathcal{S}^{(k)}$
- ◊ Reflect worst vertex about the best face
- ◊ Shrink, contract, or expand  $\Delta_k \mathcal{S}^{(k)}$



# The Nelder-Mead Method [1965]

Basic iteration ( $k \geq 0$ ):

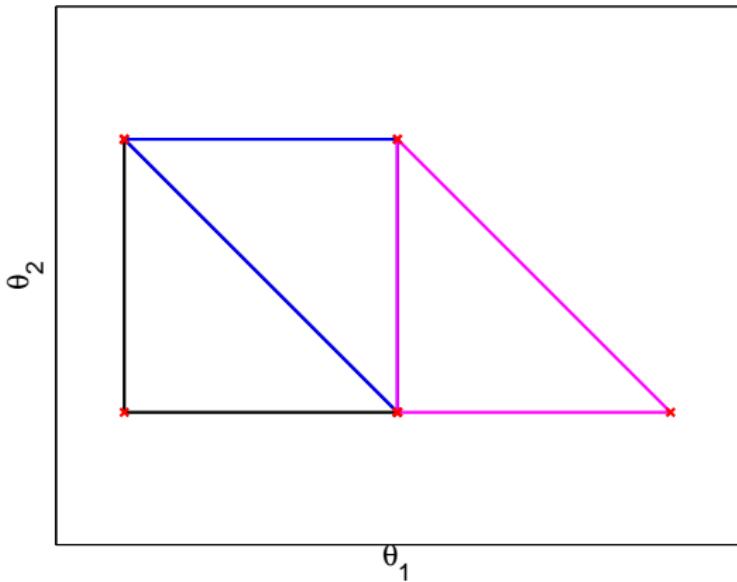
- ◊ Evaluate  $f$  on the  $n + 1$  vertices of the simplex  $x_k + \Delta_k \mathcal{S}^{(k)}$
- ◊ Reflect worst vertex about the best face
- ◊ Shrink, contract, or expand  $\Delta_k \mathcal{S}^{(k)}$



# The Nelder-Mead Method [1965]

Basic iteration ( $k \geq 0$ ):

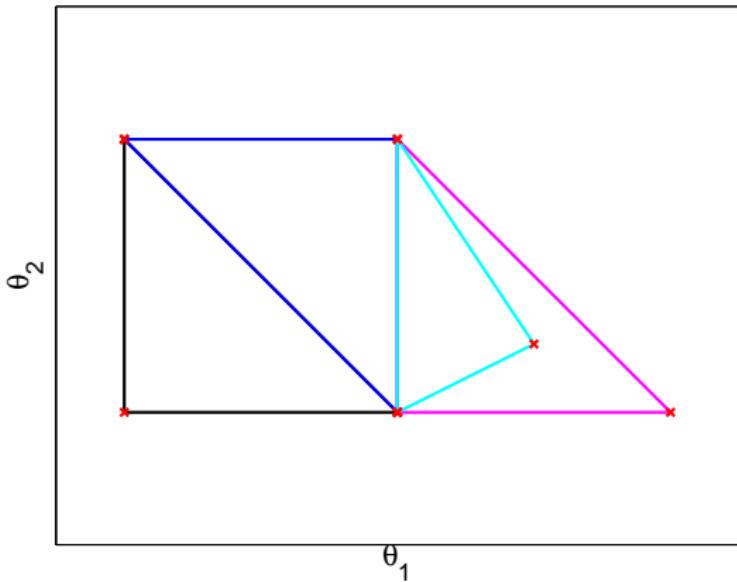
- ◊ Evaluate  $f$  on the  $n + 1$  vertices of the simplex  $x_k + \Delta_k \mathcal{S}^{(k)}$
- ◊ Reflect worst vertex about the best face
- ◊ Shrink, contract, or expand  $\Delta_k \mathcal{S}^{(k)}$



# The Nelder-Mead Method [1965]

Basic iteration ( $k \geq 0$ ):

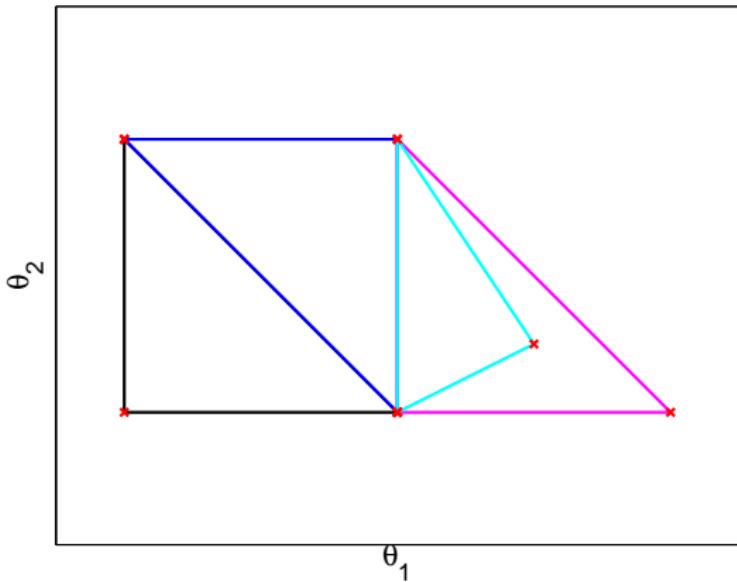
- ◊ Evaluate  $f$  on the  $n + 1$  vertices of the simplex  $x_k + \Delta_k \mathcal{S}^{(k)}$
- ◊ Reflect worst vertex about the best face
- ◊ Shrink, contract, or expand  $\Delta_k \mathcal{S}^{(k)}$



# The Nelder-Mead Method [1965]

Basic iteration ( $k \geq 0$ ):

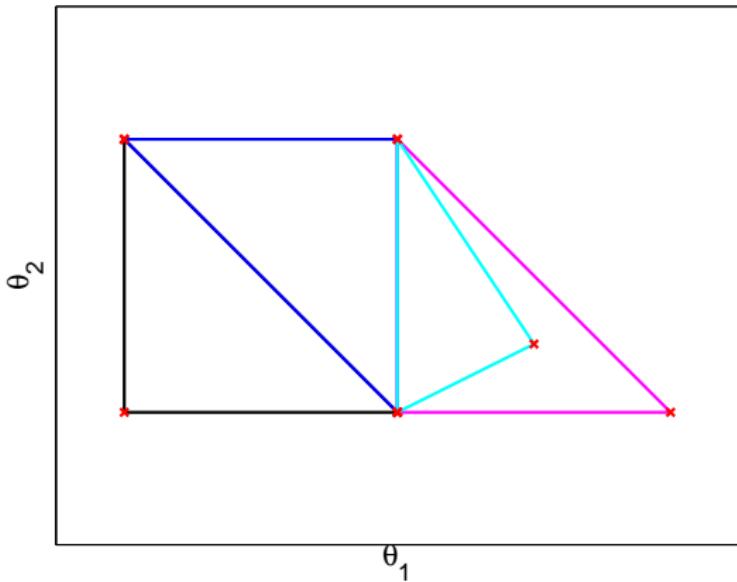
- ◊ Evaluate  $f$  on the  $n + 1$  vertices of the simplex  $x_k + \Delta_k \mathcal{S}^{(k)}$
- ◊ Reflect worst vertex about the best face
- ◊ Shrink, contract, or expand  $\Delta_k \mathcal{S}^{(k)}$



# The Nelder-Mead Method [1965]

Basic iteration ( $k \geq 0$ ):

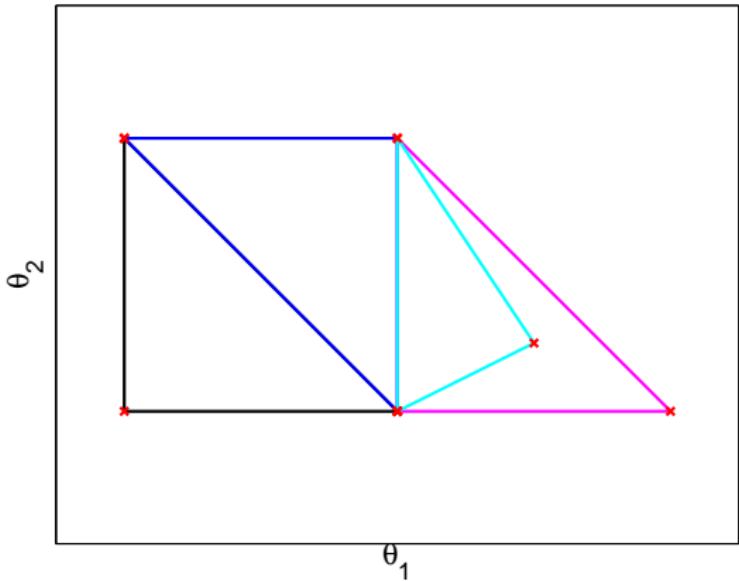
- ◊ Evaluate  $f$  on the  $n + 1$  vertices of the simplex  $x_k + \Delta_k \mathcal{S}^{(k)}$
- ◊ Reflect worst vertex about the best face
- ◊ Shrink, contract, or expand  $\Delta_k \mathcal{S}^{(k)}$



# The Nelder-Mead Method [1965]

Basic iteration ( $k \geq 0$ ):

- ◊ Evaluate  $f$  on the  $n + 1$  vertices of the simplex  $x_k + \Delta_k \mathcal{S}^{(k)}$
- ◊ Reflect worst vertex about the best face
- ◊ Shrink, contract, or expand  $\Delta_k \mathcal{S}^{(k)}$



Only the order of the function values matter:

$f(\hat{x}) = 1, f(\tilde{x}) = 1.0001$  is the same as  $f(\hat{x}) = 1, f(\tilde{x}) = 10000$ .

→ A very popular (due to “Numerical Recipes”), robust first choice

# What Are We Missing?

These methods will (eventually) find a local solution

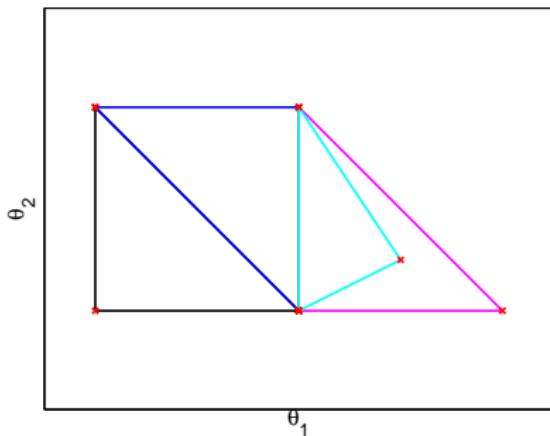
Overview: → [Kolda, Lewis, Torczon, SIREV 2003]

Each evaluation of  $f$  is expensive (valuable)



N-M:

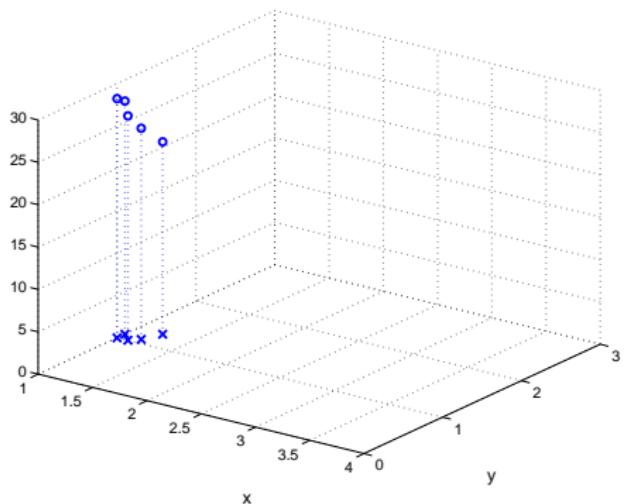
1. Only remembers the last  $n + 1$  evaluations
2. Neglects the magnitudes of the function values (order only)
3. Doesn't take into account the special (LS) problem structure



→ This is the reason many direct search methods use a search phase on top of the usual poll phase

## Making the Most of Little Information on $f$

- ◊  $f$  is expensive  $\Rightarrow$  can afford to make better use of points
- ◊ Overhead of the optimization routine is negligible relative to the cost of evaluating the simulation.



Bank of data,  $\{x_i, f(x_i)\}_{i=1}^k$ :

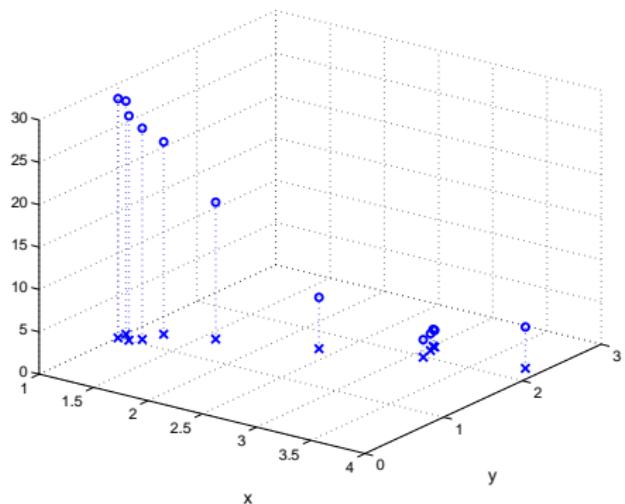
- Points (& function values) evaluated so far
- Everything known about  $f$

Goal:

- ◊ Make use of growing Bank as optimization progresses

## Making the Most of Little Information on $f$

- ◊  $f$  is expensive  $\Rightarrow$  can afford to make better use of points
- ◊ Overhead of the optimization routine is negligible relative to the cost of evaluating the simulation.



Bank of data,  $\{x_i, f(x_i)\}_{i=1}^k$ :

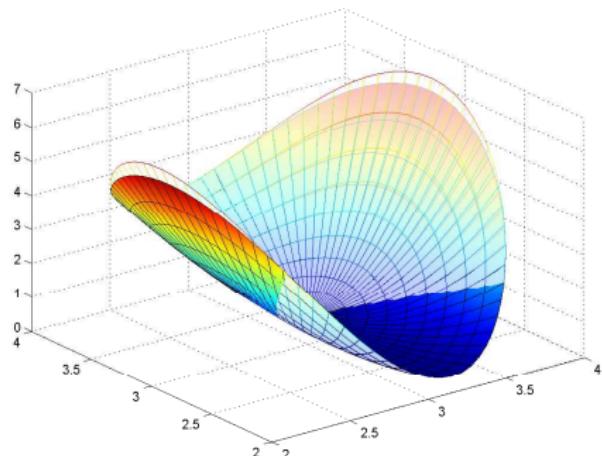
- = Points (& function values) evaluated so far
- = Everything known about  $f$

Goal:

- ◊ Make use of growing Bank as optimization progresses

## Trust-Region Methods Use Models Instead of $f$

To reduce the number of expensive  $f$  evaluations  
→ Replace difficult optimization problem  $\min f(x)$   
with a much simpler one  $\min \{m(x) : x \in \mathcal{B}\}$

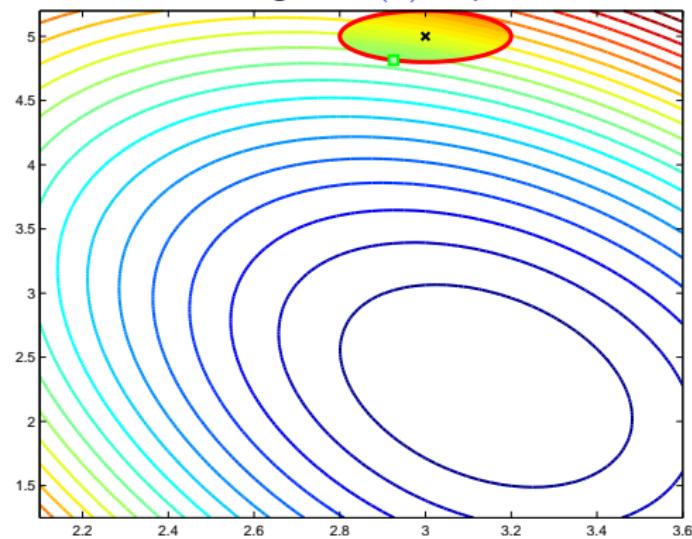


Classic NLP Technique:

- $f$  Original function: computationally expensive, no derivatives
- $m$  Surrogate model: computationally attractive, analytic derivatives

## Basic Trust-Region Idea

Use a surrogate  $m(x)$  in place of the unwieldy  $f(x)$



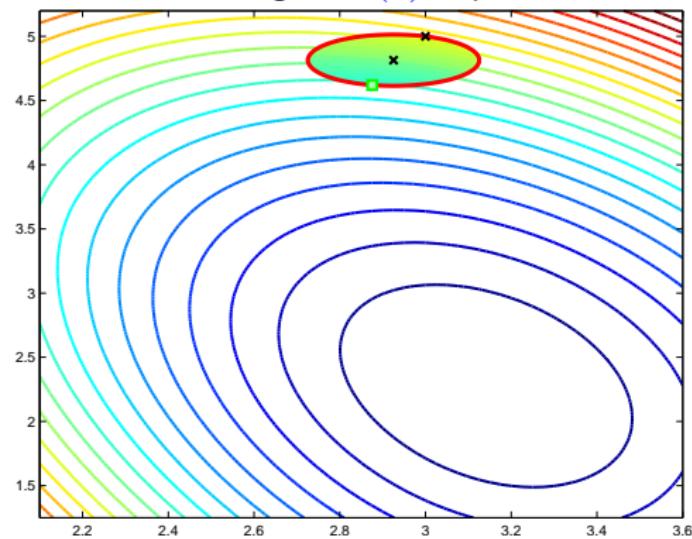
Optimize over  $m$  to avoid expense of  $f$ :

- ◊ Trust  $m$  to approximate  $f$  within  $\mathcal{B} = \{x \in \mathbb{R}^n : \|x - x_k\| \leq \Delta_k\}$ ,
- ◊ Obtain next point from  $\min \{m(x) : x \in \mathcal{B}\}$ ,
- ◊ Evaluate function and update  $(x_k, \Delta_k)$  based on how good the model's prediction was.

[Trust-Region Methods; Conn, Gould, Toint; SIAM, 2000]

## Basic Trust-Region Idea

Use a surrogate  $m(x)$  in place of the unwieldy  $f(x)$



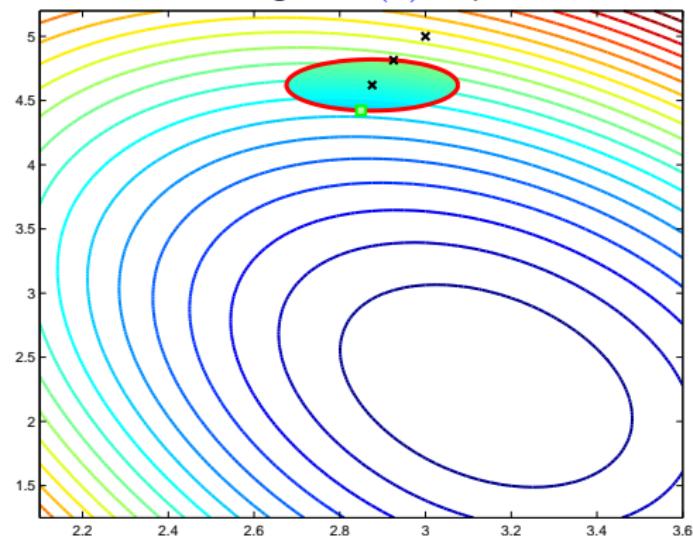
Optimize over  $m$  to avoid expense of  $f$ :

- ◊ Trust  $m$  to approximate  $f$  within  $\mathcal{B} = \{x \in \mathbb{R}^n : \|x - x_k\| \leq \Delta_k\}$ ,
- ◊ Obtain next point from  $\min \{m(x) : x \in \mathcal{B}\}$ ,
- ◊ Evaluate function and update  $(x_k, \Delta_k)$  based on how good the model's prediction was.

[Trust-Region Methods; Conn, Gould, Toint; SIAM, 2000]

## Basic Trust-Region Idea

Use a surrogate  $m(x)$  in place of the unwieldy  $f(x)$



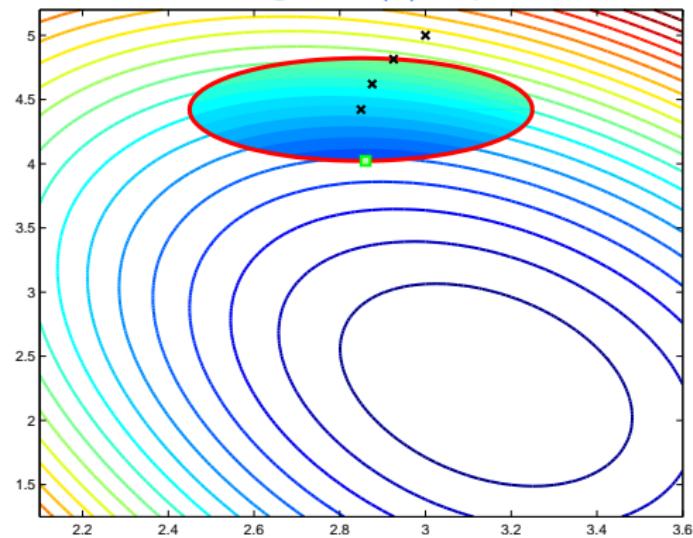
Optimize over  $m$  to avoid expense of  $f$ :

- ◊ Trust  $m$  to approximate  $f$  within  $\mathcal{B} = \{x \in \mathbb{R}^n : \|x - x_k\| \leq \Delta_k\}$ ,
- ◊ Obtain next point from  $\min \{m(x) : x \in \mathcal{B}\}$ ,
- ◊ Evaluate function and update  $(x_k, \Delta_k)$  based on how good the model's prediction was.

[Trust-Region Methods; Conn, Gould, Toint; SIAM, 2000]

## Basic Trust-Region Idea

Use a surrogate  $m(x)$  in place of the unwieldy  $f(x)$



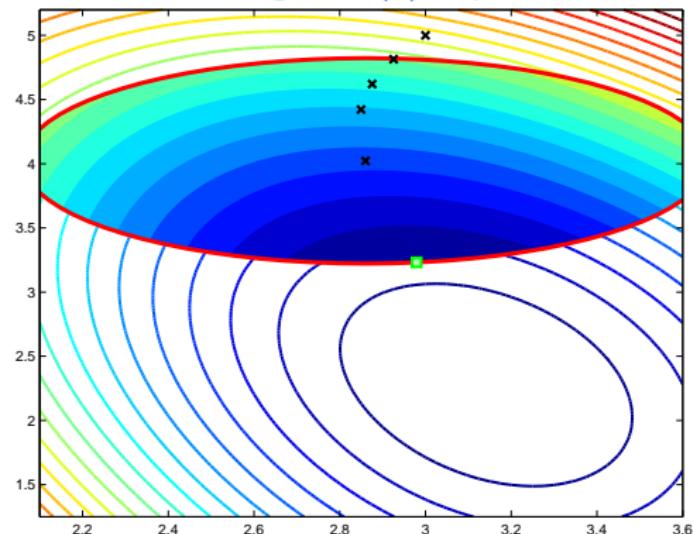
Optimize over  $m$  to avoid expense of  $f$ :

- ◊ Trust  $m$  to approximate  $f$  within  $\mathcal{B} = \{x \in \mathbb{R}^n : \|x - x_k\| \leq \Delta_k\}$ ,
- ◊ Obtain next point from  $\min \{m(x) : x \in \mathcal{B}\}$ ,
- ◊ Evaluate function and update  $(x_k, \Delta_k)$  based on how good the model's prediction was.

[Trust-Region Methods; Conn, Gould, Toint; SIAM, 2000]

## Basic Trust-Region Idea

Use a surrogate  $m(x)$  in place of the unwieldy  $f(x)$



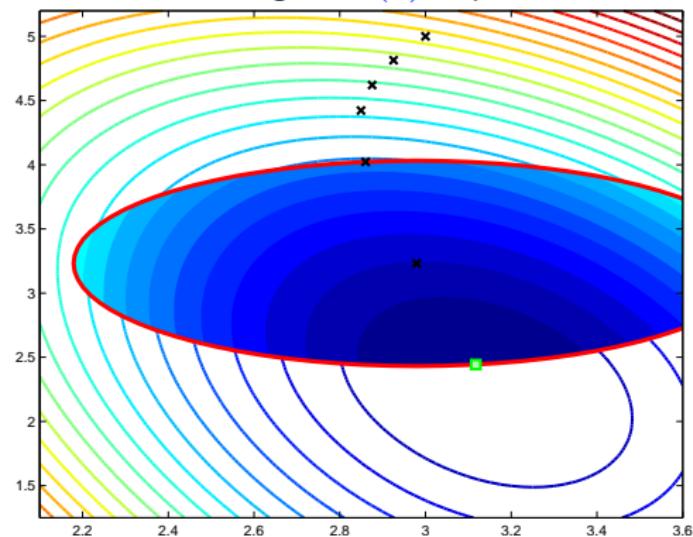
Optimize over  $m$  to avoid expense of  $f$ :

- ◊ Trust  $m$  to approximate  $f$  within  $\mathcal{B} = \{x \in \mathbb{R}^n : \|x - x_k\| \leq \Delta_k\}$ ,
- ◊ Obtain next point from  $\min \{m(x) : x \in \mathcal{B}\}$ ,
- ◊ Evaluate function and update  $(x_k, \Delta_k)$  based on how good the model's prediction was.

[Trust-Region Methods; Conn, Gould, Toint; SIAM, 2000]

## Basic Trust-Region Idea

Use a surrogate  $m(x)$  in place of the unwieldy  $f(x)$



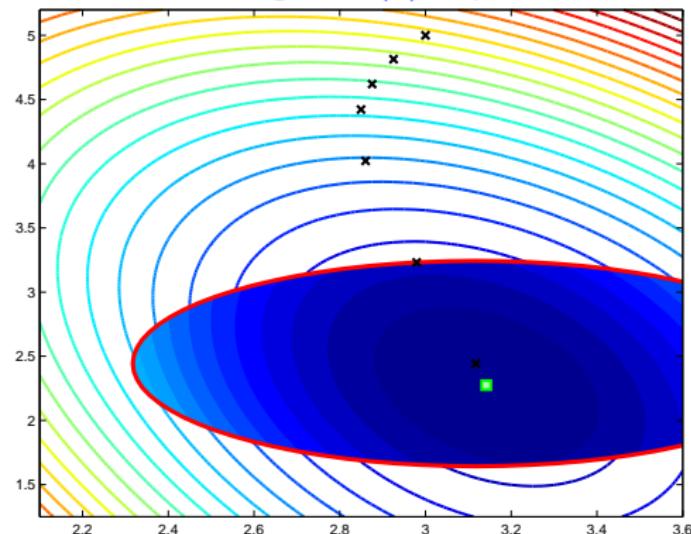
Optimize over  $m$  to avoid expense of  $f$ :

- ◊ Trust  $m$  to approximate  $f$  within  $\mathcal{B} = \{x \in \mathbb{R}^n : \|x - x_k\| \leq \Delta_k\}$ ,
- ◊ Obtain next point from  $\min \{m(x) : x \in \mathcal{B}\}$ ,
- ◊ Evaluate function and update  $(x_k, \Delta_k)$  based on how good the model's prediction was.

[Trust-Region Methods; Conn, Gould, Toint; SIAM, 2000]

## Basic Trust-Region Idea

Use a surrogate  $m(x)$  in place of the unwieldy  $f(x)$



Optimize over  $m$  to avoid expense of  $f$ :

- ◊ Trust  $m$  to approximate  $f$  within  $\mathcal{B} = \{x \in \mathbb{R}^n : \|x - x_k\| \leq \Delta_k\}$ ,
- ◊ Obtain next point from  $\min \{m(x) : x \in \mathcal{B}\}$ ,
- ◊ Evaluate function and update  $(x_k, \Delta_k)$  based on how good the model's prediction was.

[Trust-Region Methods; Conn, Gould, Toint; SIAM, 2000]

## Where Does the Model Come From?

When derivatives are available:

Taylor-based model  $m(x_k + s) = c + g_k^T s + \frac{1}{2} s^T H_k s$

- ◊  $g_k = \nabla_x f(x_k)$
- ◊  $H_k \approx \nabla_{x,x}^2 f(x_k)$



# Where Does the Model Come From?

When derivatives are available:

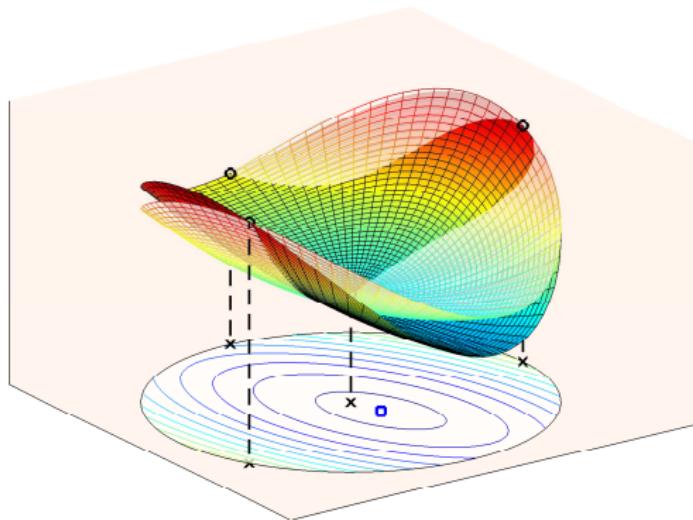
Taylor-based model  $m(x_k + s) = c + g_k^T s + \frac{1}{2} s^T H_k s$

- ◊  $g_k = \nabla_x f(x_k)$
- ◊  $H_k \approx \nabla_{x,x}^2 f(x_k)$

Without derivatives

- ◊ Interpolation-based models
- ◊ Regression-based models
- ◊ Stochastic/randomized models

# Interpolation-Based Quadratic Models



An interpolating quadratic in  $\mathbb{R}^2$

$$m(x_k + s) = c + g^T s + \frac{1}{2} s^T H s:$$

Get the model parameters  
 $c, g, H = H^T$  by demanding  
interpolation:

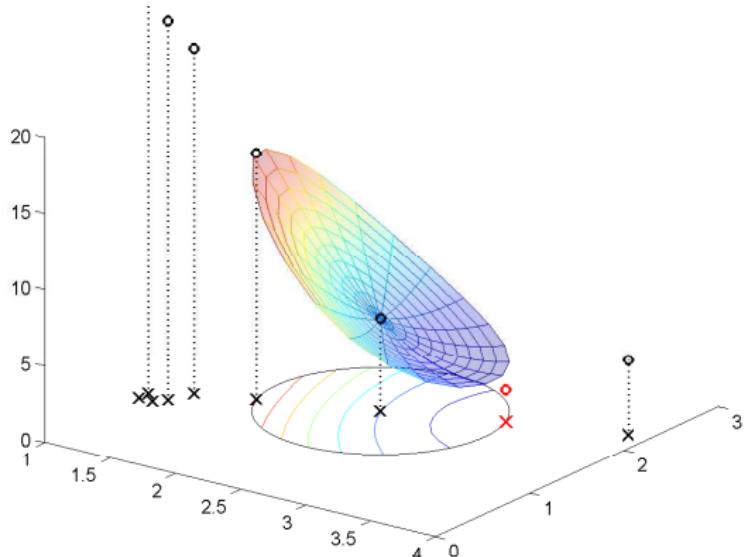
$$m(x_k + y_i) = f(x_k + y_i)$$

for all  $y_i \in \mathcal{Y}$  = interpolation set

Main difficulty is  $\mathcal{Y}$ :

- ◊ Use prior function evaluations,
- ◊  $m$  well-defined and approximates  $f$  locally.

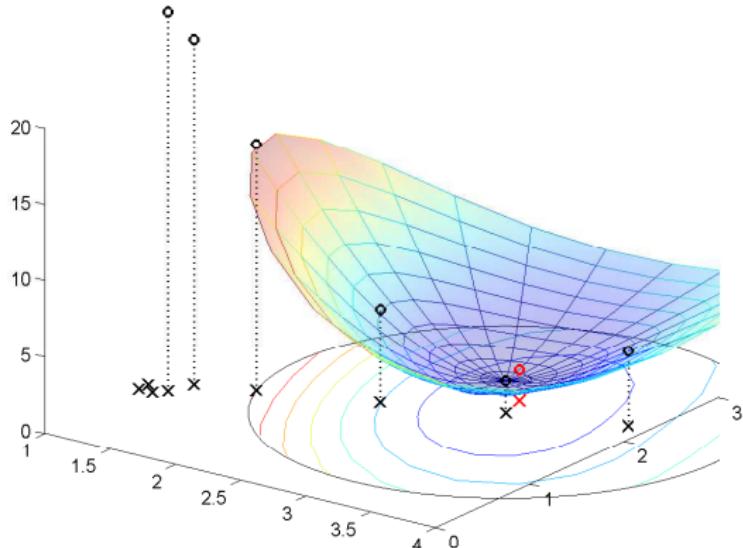
# Interpolation-Based Trust-Region Methods



Iteration  $k$ :

- ◊ Build a model  $m_k$  interpolating  $f$  on  $\mathcal{Y}$
- ◊ Trust  $m_k$  within region  $\mathcal{B}_k$
- ◊ Minimize  $m_k$  within  $\mathcal{B}_k$  to obtain next point for evaluation
- ◊ Do expensive evaluation
- ◊ Update  $m_k$  and  $\mathcal{B}_k$  based on how good model prediction was

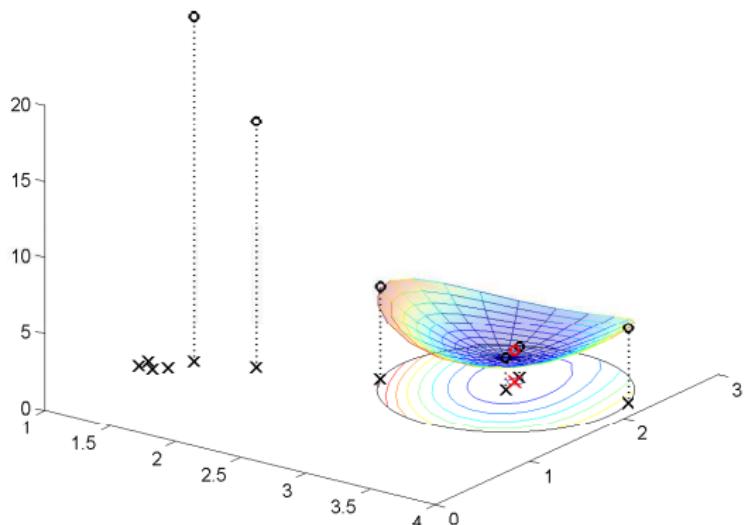
# Interpolation-Based Trust-Region Methods



Iteration  $k$ :

- ◊ Build a model  $m_k$  interpolating  $f$  on  $\mathcal{Y}$
- ◊ Trust  $m_k$  within region  $\mathcal{B}_k$
- ◊ Minimize  $m_k$  within  $\mathcal{B}_k$  to obtain next point for evaluation
- ◊ Do expensive evaluation
- ◊ Update  $m_k$  and  $\mathcal{B}_k$  based on how good model prediction was

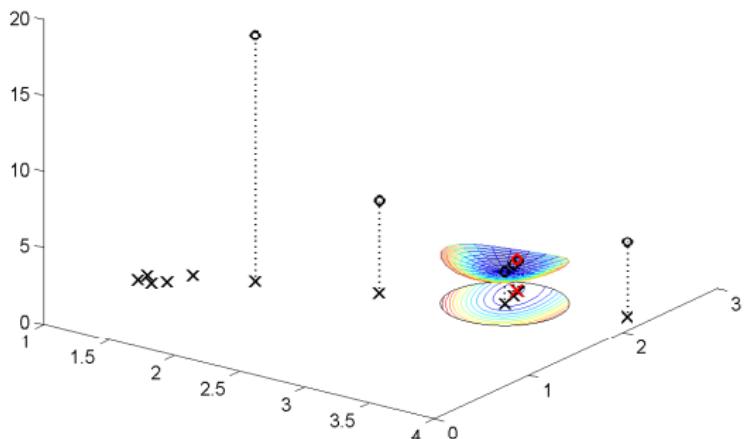
# Interpolation-Based Trust-Region Methods



Iteration  $k$ :

- ◊ Build a model  $m_k$  interpolating  $f$  on  $\mathcal{Y}$
- ◊ Trust  $m_k$  within region  $\mathcal{B}_k$
- ◊ Minimize  $m_k$  within  $\mathcal{B}_k$  to obtain next point for evaluation
- ◊ Do expensive evaluation
- ◊ Update  $m_k$  and  $\mathcal{B}_k$  based on how good model prediction was

## Iteration $k$ :

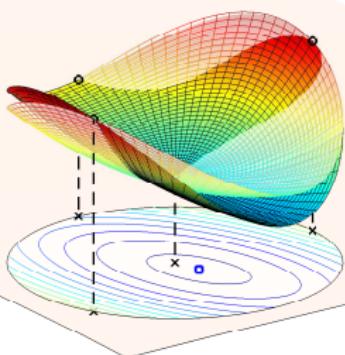


- ◊ Build a model  $m_k$  interpolating  $f$  on  $\mathcal{Y}$
- ◊ Trust  $m_k$  within region  $\mathcal{B}_k$
- ◊ Minimize  $m_k$  within  $\mathcal{B}_k$  to obtain next point for evaluation
- ◊ Do expensive evaluation
- ◊ Update  $m_k$  and  $\mathcal{B}_k$  based on how good model prediction was

# Building Models Without Derivatives

Given  $(\mathcal{Y}_k, f(\mathcal{Y}_k))$ , "solve"

$$\Phi(\mathcal{Y}_k)z = \begin{bmatrix} \Phi_c & \Phi_g & \Phi_H \end{bmatrix} \begin{bmatrix} z_c \\ z_g \\ z_H \end{bmatrix} = \underline{\mathbf{f}} = f(\mathcal{Y}_k)$$



$n = 2, |\mathcal{Y}_k| = 4$

Full quadratics,  $|\mathcal{Y}_k| = \frac{(n+1)(n+2)}{2}$

- ◊ Geometric conditions on points in  $\mathcal{Y}_k$

Undetermined interpolation,  
 $|\mathcal{Y}_k| < \frac{(n+1)(n+2)}{2}$

- ◊ Use (Powell) Hessian updates

$$\begin{array}{ll} \min_{g_k, H_k} & \|H_k - H_{k-1}\|_F^2 \\ \text{s.t.} & q_k = \underline{\mathbf{f}} \text{ on } \mathcal{Y}_k \end{array}$$

Regression,  $|\mathcal{Y}_k| > \frac{(n+1)(n+2)}{2}$

- ◊ Solve  $\min_z \|\Phi z - \underline{\mathbf{f}}\|$

# Multivariate Interpolation is a Different Kind of Animal

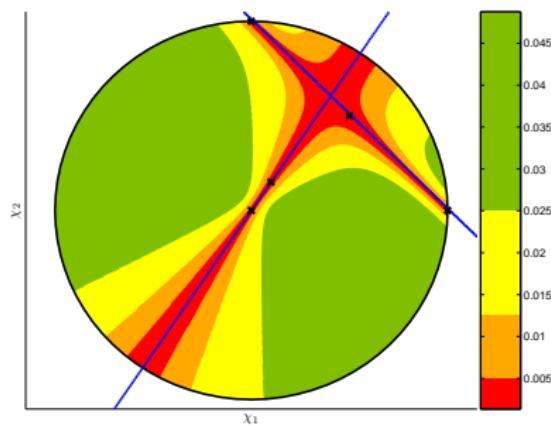
$$m(x_k + y_i) = f(x_k + y_i) \quad \forall y_i \in \mathcal{Y}$$

$n = 1$  Given  $p$  distinct points, can find a unique degree  $p - 1$  polynomial  $m$

$n > 1$  Not true! (see Mairhuber-Curtis Theorem)

For quadratic models in  $\mathbb{R}^n$ :

- ◊  $\frac{(n+1)(n+2)}{2}$  coefficients
- ◊ Unique interpolant may not exist, even when  $|\mathcal{Y}| = \frac{(n+1)(n+2)}{2}$
- ◊ Locations of the points in  $\mathcal{Y}$  must satisfy additional geometric conditions (has nothing to do with  $f$  values)



→ [Scattered Data Approximation; Wendland; Cambridge University Press, 2010]

## “Taylor-like” Error Bounds

1. Assuming underlying  $f$  is sufficiently smooth  
= derivatives of  $f$  exist but are unavailable
2. A model  $m_k$  is locally **fully linear** if:

For all  $x \in \mathcal{B}_k = \{x \in \Omega : \|x - x_k\| \leq \Delta_k\}$

- ♦  $|m_k(x) - f(x)| \leq \kappa_1 \Delta_k^2$
- ♦  $\|\nabla m_k(x) - \nabla f(x)\| \leq \kappa_2 \Delta_k$

for constants  $\kappa_i$  independent of  $x$  and  $\Delta_k$ .

→ [Intro. to DFO; Conn, Scheinberg, Vicente; SIAM 2009]

## “Taylor-like” Error Bounds

1. Assuming underlying  $f$  is sufficiently smooth
2. A model  $m_k$  is locally **fully quadratic** if:

For all  $x \in \mathcal{B}_k = \{x \in \Omega : \|x - x_k\| \leq \Delta_k\}$

- ♦  $|m_k(x) - f(x)| \leq \kappa_1 \Delta_k^3$
- ♦  $\|\nabla m_k(x) - \nabla f(x)\| \leq \kappa_2 \Delta_k^2$
- ♦  $\|\nabla^2 m_k(x) - \nabla^2 f(x)\| \leq \kappa_3 \Delta_k$

for constants  $\kappa_i$  independent of  $x$  and  $\Delta_k$ .

→ [Intro. to DFO; Conn, Scheinberg, Vicente; SIAM 2009]

## Ingredients for Convergence to Stationary Points

Assuming underlying  $f$  is sufficiently smooth and regular (e.g., has bounded level sets)

$\lim_{k \rightarrow \infty} \nabla f(x_k) = 0$  provided:

1. Control  $\mathcal{B}_k$  based on model quality
2. (Occasional) approximation within  $\mathcal{B}_k$   
Our quadratics satisfy
  - ♦  $|q_k(x) - f(x)| \leq \kappa_1(\gamma_f + \|H_k\|)\Delta_k^2, \quad x \in \mathcal{B}_k$
  - ♦  $\|g_k + H_k(x - x_k) - \nabla f(x)\| \leq \kappa_2(\gamma_f + \|H_k\|)\Delta_k, \quad x \in \mathcal{B}_k$
3. Sufficient decrease

*At least model gradient should be good*

Radial Basis Function (RBF) models also fit in this framework

General introduction

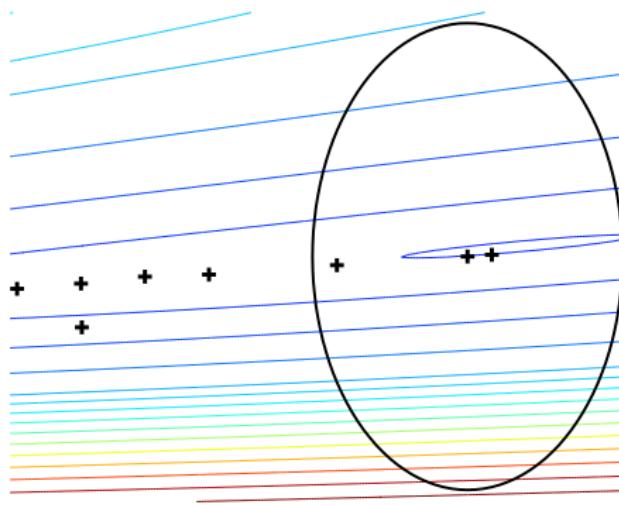
→ [W. & Shoemaker, SIREV 2013]

→ [[Intro. to DFO](#); Conn, Scheinberg, Vicente; SIAM 2009]

## Greed. Alone. Can. Hurt.

Model-improvement may be needed when:

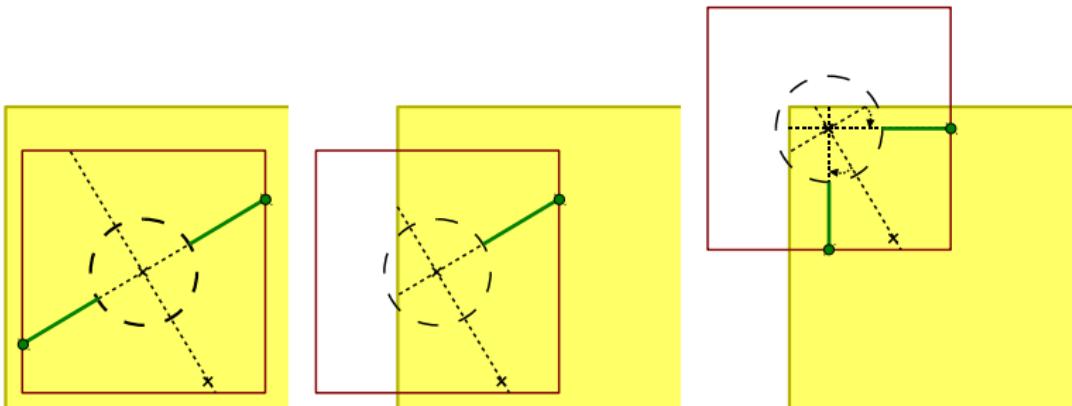
- ◊ Nearby points line up
- ◊ May not have enough points to ensure model quality in all directions



→ May need  $n$  additional evaluations

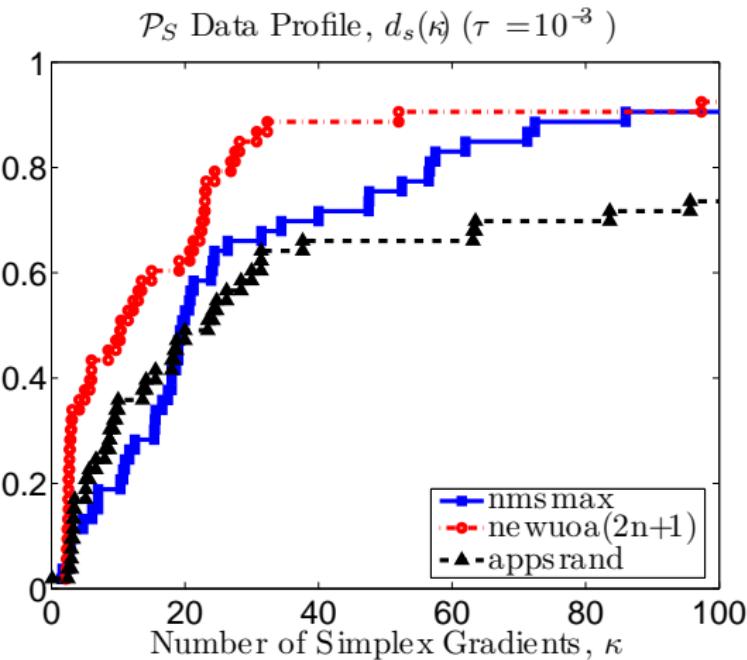
# Constraints and Model Quality

Constraints complicate matters  
... if one does not allow evaluation of infeasible points



→ May need directions normal to nearby constraints

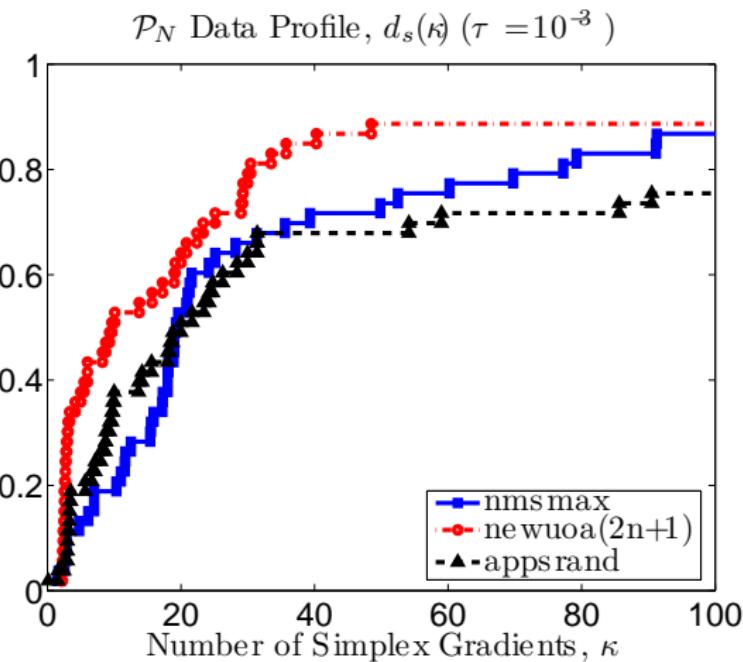
## Performance Comparisons on Test Functions



- When evaluations are sequential,  
model-based methods  
(**NEWUOA**) regularly  
outperform direct  
search methods  
without a search  
phase  
(**nmsmax,apprsrand**)

→ [Moré & W., SIAM 2009]

## Performance Comparisons on Test Functions



Noisy problems

- When evaluations are sequential,  
model-based methods  
(**NEWUOA**) regularly  
outperform direct  
search methods  
without a search  
phase  
(**nmsmax,apprsrand**)

→ [Moré & W., SIOP 2009]

## Many Practical Details In Implementations

- ◊ Choice of interpolation points  $\mathcal{Y}_k$
- ◊ Updating of trust region  $\mathcal{B}_k$
- ◊ Improvement of models

## Many Practical Details In Implementations

- ◊ Choice of interpolation points  $\mathcal{Y}_k$
- ◊ Updating of trust region  $\mathcal{B}_k$
- ◊ Improvement of models

**BOBYQA** [Powell], **DFO** [Scheinberg], **POUNDer** [W.]

Initialization  $p = |\mathcal{Y}_k|$  structured evaluations

Based on input,  $\approx 2n + 1$

Based on input, no more than  $n + 1$

Interpolation Set  $p = |\mathcal{Y}_k|, \forall k$

Bootstrap to  $|\mathcal{Y}_k| = \frac{(n+1)(n+2)}{2}$ , then fixed

Varies in  $\{n + 1, \dots, \frac{(n+1)(n+2)}{2}\}$  based on available points

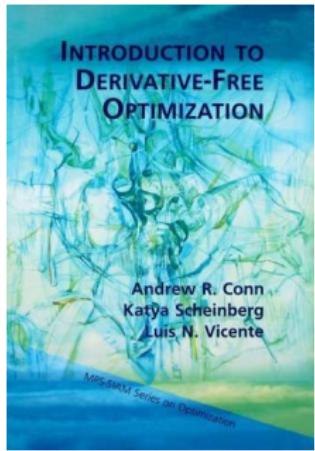
Linear Algebra If  $p = \mathcal{O}(n)$ , model formation costs only  $\mathcal{O}(n^2)$

Expensive

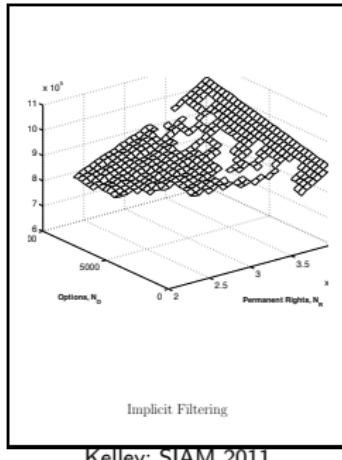
Expensive

# Growing Body of Tools and Resources for Local DFO

? What to use on problems with characteristics X, Y, and Z ?



Conn, Scheinberg,  
Vicente; SIAM 2009



Kelley; SIAM 2011

## Many solvers

Sample considered by Rios & Sahinidis, 2010:

ASA,  
CMA-ES,  
DAKOTA/\*,  
FMINSEARCH,  
HOPSPACK,  
MCS,  
NOMAD,  
SID-PSM,

BOBYQA,  
DFO,  
TOMLAB/\*,  
GLOBAL,  
IMFIL,  
NEWUOA,  
PSWARM,  
SNOBFIT

A quick sketch of a **multistart** methods and some practical details

- ◊ useful in derivative-based and derivative-free cases
- ◊ obtain a list of distinct minimizers (for post-processing, etc.)
- ◊ simple to get started
- ! simple to abuse/misuse ("I found all minimizers")

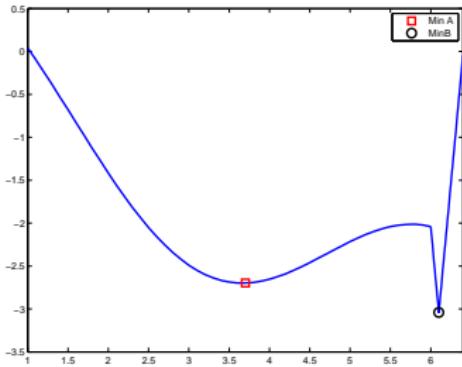
# Why Multistart?

Multiple local minima are often of interest in practice:

**Design:** Multiple objectives (or even constraints) might later be of interest

**Simulation Errors:** Could have spurious local minima from anomalies in the simulator

**Uncertainty:** Some minima are more sensitive to perturbations than others (gentle valleys versus steep cliffs)



## Two phase iterative method

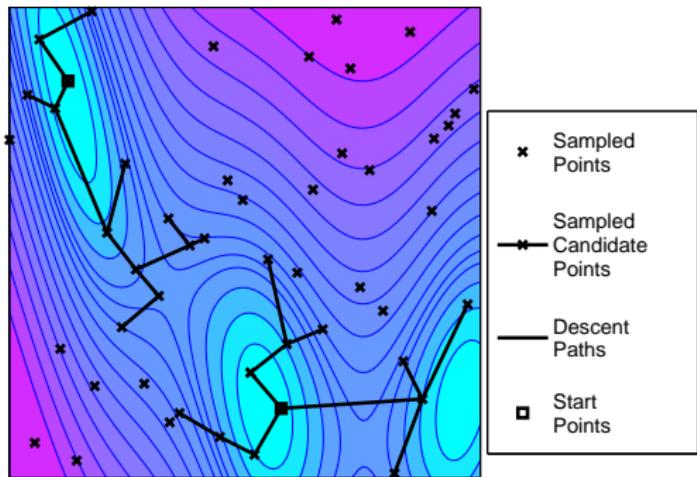
**Global Exploration:** Sample  $N$  points in  $\mathcal{D}$ .  $\leftarrow$  Guarantees convergence

**Local Refinement:** Start a local minimization algorithm  $\mathcal{A}$  from some promising subset of the sample points.

Want to find many (good) local minima while avoiding repeatedly finding the same local minima.

# Multi Level Single Linkage (MLSL) Clustering Procedure

Where to start  $\mathcal{A}$  in  $k$ th iteration [Rinnooy Kan & Timmer, Math. Programming 1987]



Start  $\mathcal{A}$  at each sample point  $x_i$  provided:

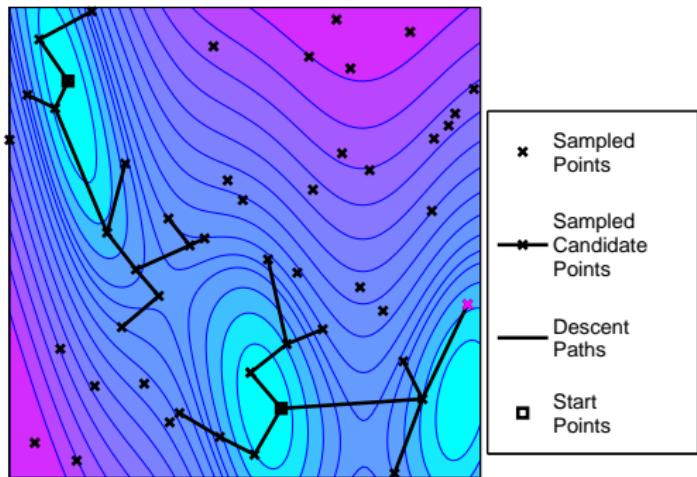
- ◊  $\mathcal{A}$  has not been started from  $x_i$ , and
- ◊ no other sample point  $x_j$  with  $f(x_j) < f(x_i)$  is within a distance

$$r_k = \frac{1}{\sqrt{\pi}} \sqrt[n]{\text{vol}(\mathcal{D}) \frac{5\Gamma\left(1 + \frac{n}{2}\right) \log(kN)}{kN}},$$

Ex.: It. 1 Exploration

# Multi Level Single Linkage (MLSL) Clustering Procedure

Where to start  $\mathcal{A}$  in  $k$ th iteration [Rinnooy Kan & Timmer, Math. Programming 1987]



Start  $\mathcal{A}$  at each sample point  $x_i$  provided:

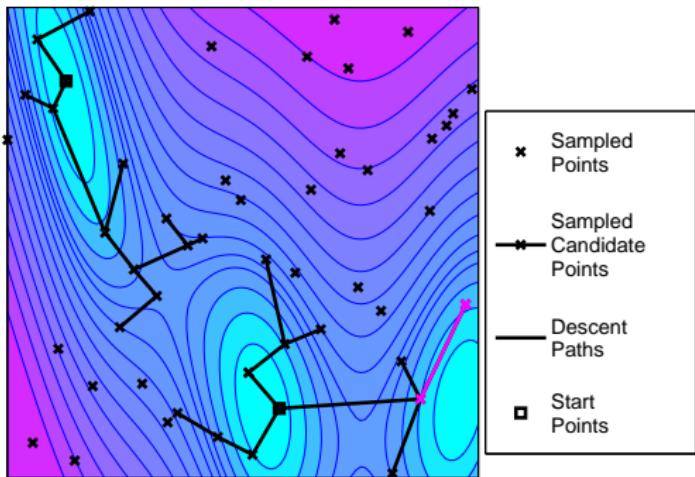
- ◊  $\mathcal{A}$  has not been started from  $x_i$ , and
- ◊ no other sample point  $x_j$  with  $f(x_j) < f(x_i)$  is within a distance

$$r_k = \frac{1}{\sqrt{\pi}} \sqrt[n]{\text{vol}(\mathcal{D}) \frac{5\Gamma\left(1 + \frac{n}{2}\right) \log(kN)}{kN}},$$

Ex.: It. 1 Exploration

# Multi Level Single Linkage (MLSL) Clustering Procedure

Where to start  $\mathcal{A}$  in  $k$ th iteration [Rinnooy Kan & Timmer, Math. Programming 1987]



Start  $\mathcal{A}$  at each sample point  $x_i$  provided:

- ◊  $\mathcal{A}$  has not been started from  $x_i$ , and
- ◊ no other sample point  $x_j$  with  $f(x_j) < f(x_i)$  is within a distance

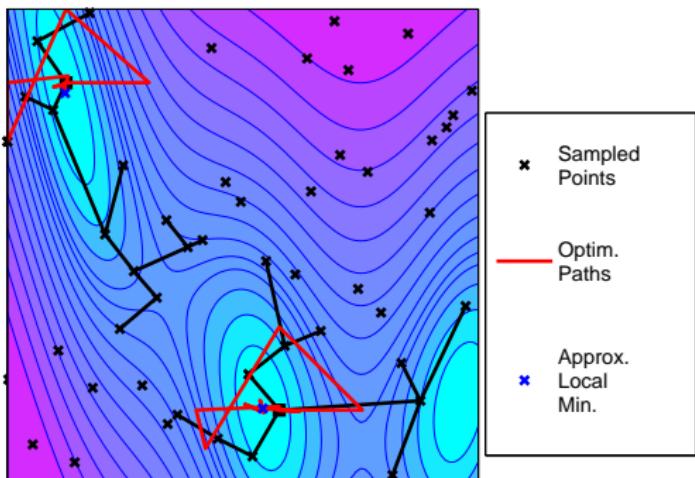
$$r_k = \frac{1}{\sqrt{\pi}} \sqrt[n]{\text{vol}(\mathcal{D}) \frac{5\Gamma\left(1 + \frac{n}{2}\right) \log(kN)}{kN}},$$

Ex.: It. 1 Exploration

Thm [RK-T]- Will start finitely many local runs with probability 1.

# Multi Level Single Linkage (MLSL) Clustering Procedure

Where to start  $\mathcal{A}$  in  $k$ th iteration [Rinnooy Kan & Timmer, Math. Programming 1987]



Start  $\mathcal{A}$  at each sample point  $x_i$  provided:

- ◊  $\mathcal{A}$  has not been started from  $x_i$ , and
- ◊ no other sample point  $x_j$  with  $f(x_j) < f(x_i)$  is within a distance

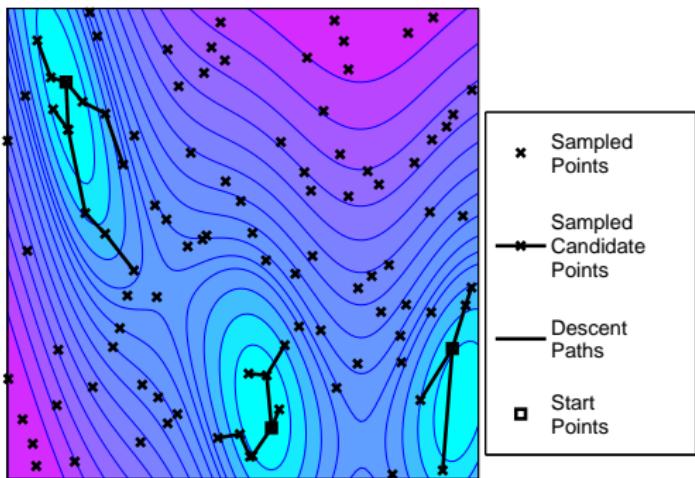
$$r_k = \frac{1}{\sqrt{\pi}} \sqrt[n]{\text{vol}(\mathcal{D}) \frac{5\Gamma\left(1 + \frac{n}{2}\right) \log(kN)}{kN}},$$

Ex.: It. 1 Refinement

Thm [RK-T]- Will start finitely many local runs with probability 1.

# Multi Level Single Linkage (MLSL) Clustering Procedure

Where to start  $\mathcal{A}$  in  $k$ th iteration [Rinnooy Kan & Timmer, Math. Programming 1987]



Start  $\mathcal{A}$  at each sample point  $x_i$  provided:

- ◊  $\mathcal{A}$  has not been started from  $x_i$ , and
- ◊ no other sample point  $x_j$  with  $f(x_j) < f(x_i)$  is within a distance

$$r_k = \frac{1}{\sqrt{\pi}} \sqrt[n]{\text{vol}(\mathcal{D}) \frac{5\Gamma\left(1 + \frac{n}{2}\right) \log(kN)}{kN}},$$

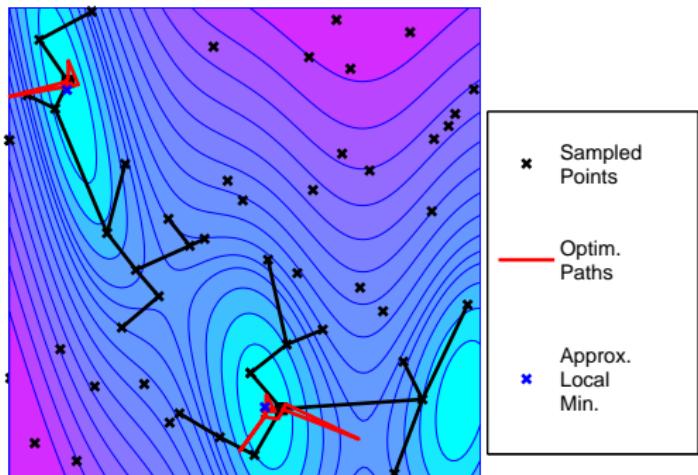
Ex.: It. 2 Exploration

Thm [RK-T]- Will start finitely many local runs with probability 1.

# How To Do Even Better

MIPE: Maximum Information from Previous Evaluations:

Take advantage of the information gained from running  $\mathcal{A}$



Start  $\mathcal{A}$  at each sample point  $x_i$  provided:

- ◊  $\mathcal{A}$  has not been started from  $x_i$ , and
- ◊ no other sample point  $x_j$  with  $f(x_j) < f(x_i)$  is within a distance  $r_k$
- ◊ no optimization point  $x^{\mathcal{A}}$  with  $f(x^{\mathcal{A}}) < f(x_i)$  is within a distance  $r_k$

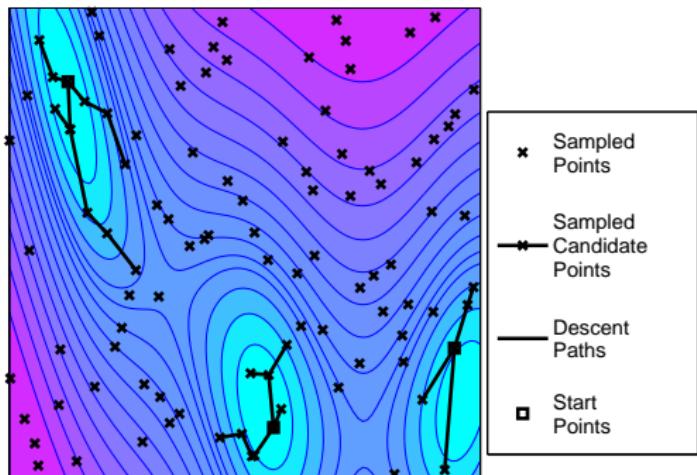
Ex.: It. 1 Refinement

MIPE does at least as few local runs as MLSL.

# How To Do Even Better

MIPE: Maximum Information from Previous Evaluations:

Take advantage of the information gained from running  $\mathcal{A}$



Start  $\mathcal{A}$  at each sample point  $x_i$  provided:

- ◊  $\mathcal{A}$  has not been started from  $x_i$ , and
- ◊ no other sample point  $x_j$  with  $f(x_j) < f(x_i)$  is within a distance  $r_k$
- ◊ no optimization point  $x^{\mathcal{A}}$  with  $f(x^{\mathcal{A}}) < f(x_i)$  is within a distance  $r_k$

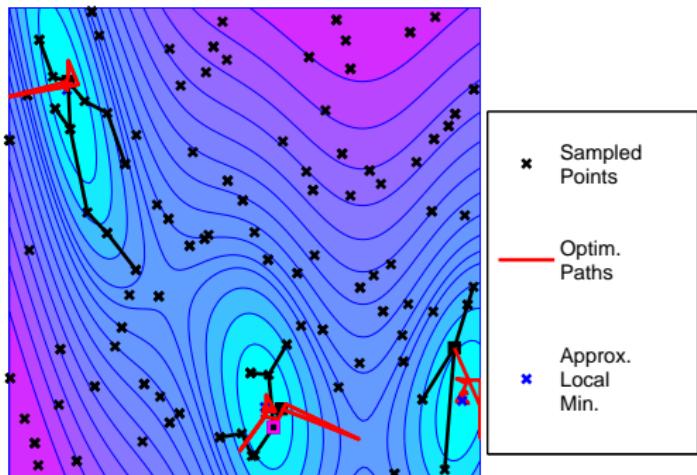
Ex.: It. 2 Exploration

MIPE does at least as few local runs as MLSL.

# How To Do Even Better

MIPE: Maximum Information from Previous Evaluations:

Take advantage of the information gained from running  $\mathcal{A}$



Start  $\mathcal{A}$  at each sample point  $x_i$  provided:

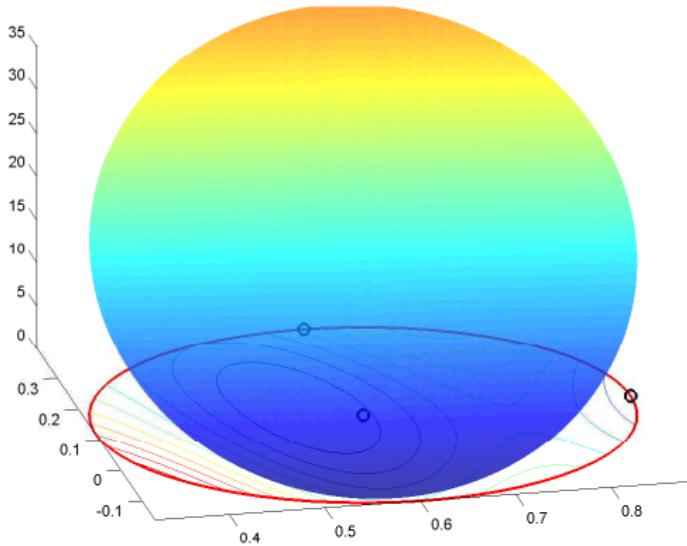
- ◊  $\mathcal{A}$  has not been started from  $x_i$ , and
- ◊ no other sample point  $x_j$  with  $f(x_j) < f(x_i)$  is within a distance  $r_k$
- ◊ no optimization point  $x^{\mathcal{A}}$  with  $f(x^{\mathcal{A}}) < f(x_i)$  is within a distance  $r_k$

Ex.: It. 2 Refinement

MIPE does at least as few local runs as MLSL.

## Using External Points To Form Better Local Models

Initial model interpolating  $n + 1 = 3$  points



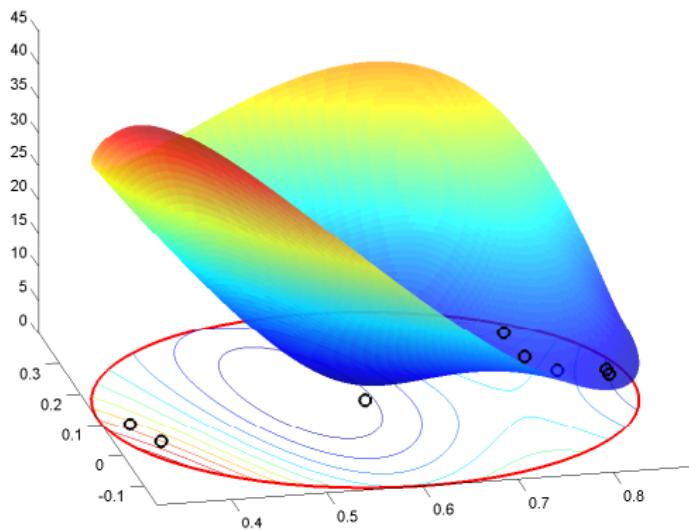
Take advantage of global history:

- ◊ from current minimization,
- ◊ from previous minimizations, and
- ◊ from the global sampling.

More information means rapid progress.

# Using External Points To Form Better Local Models

Initial model interpolating 8 sample points



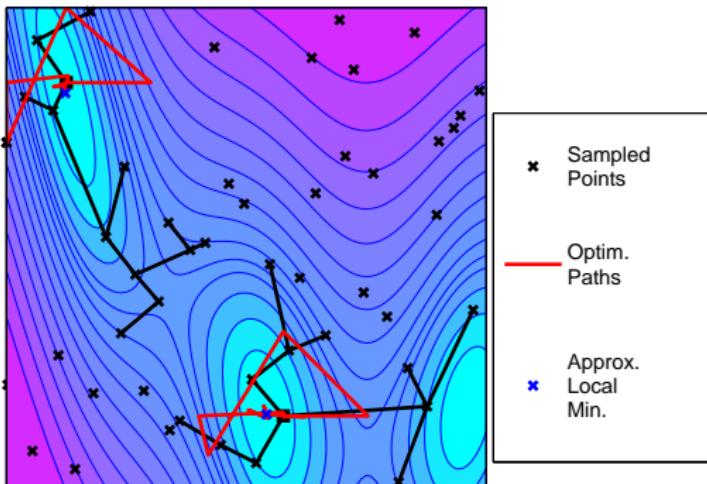
Take advantage of global history:

- ◊ from current minimization,
- ◊ from previous minimizations, and
- ◊ from the global sampling.

More information means rapid progress.

# Using External Points To Form Better Local Models

Optimization paths without external points



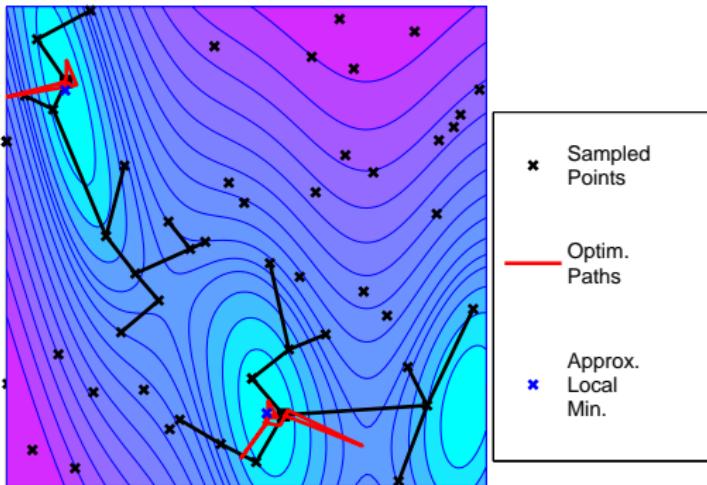
Take advantage of global history:

- ◊ from current minimization,
- ◊ from previous minimizations, and
- ◊ from the global sampling.

More information means rapider progress.

# Using External Points To Form Better Local Models

Optimization paths with external points

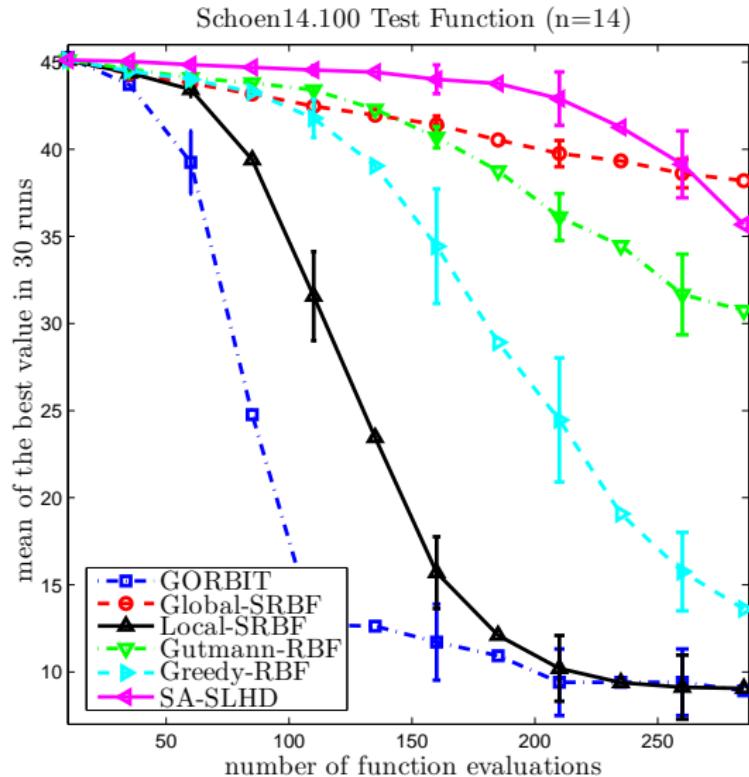


Take advantage of global history:

- ◊ from current minimization,
- ◊ from previous minimizations, and
- ◊ from the global sampling.

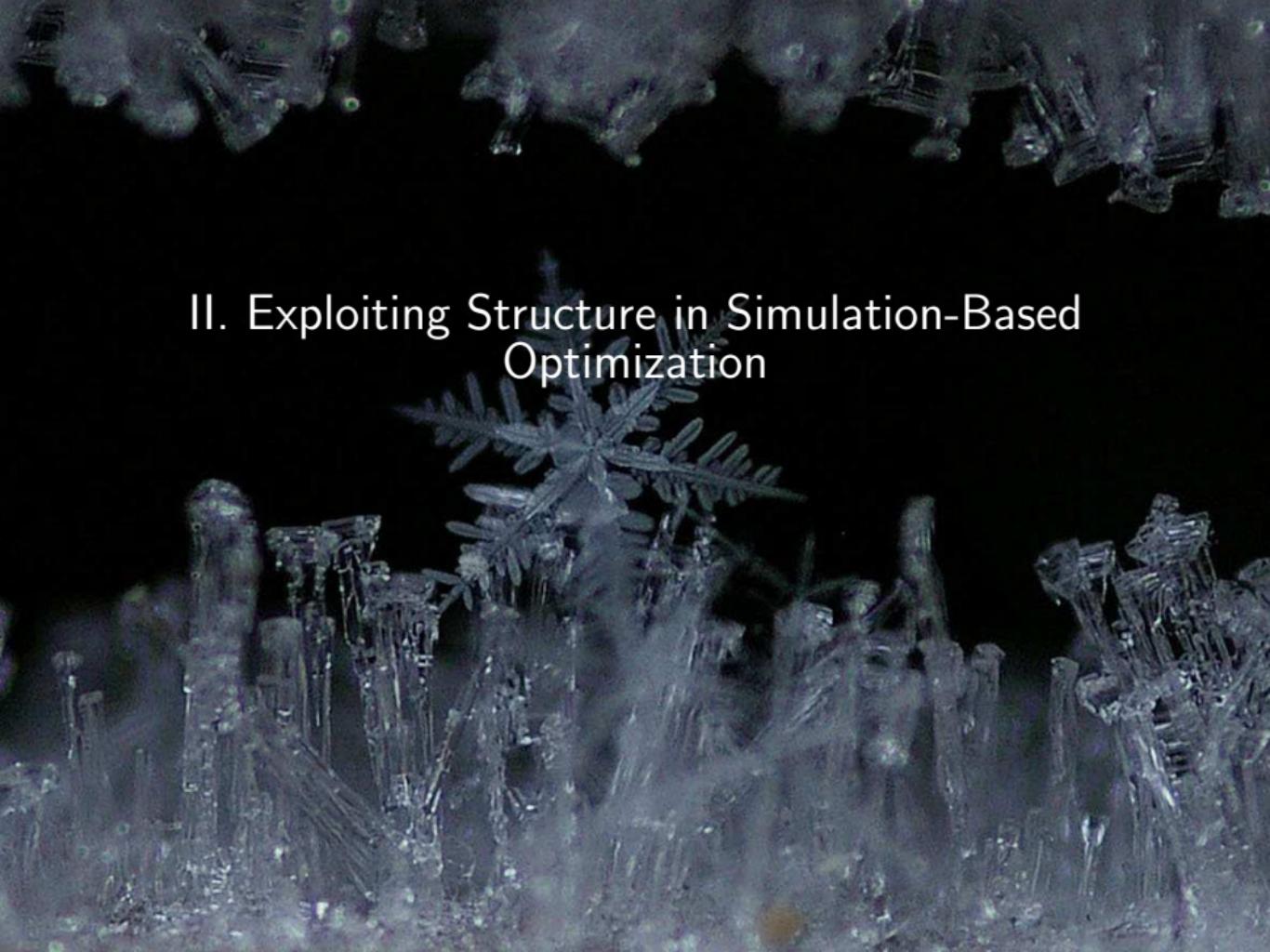
More information means rapider progress.

## Performance Comparisons on Test Functions



- ◊ **GORBIT** is multistart with RBF model-based method
- ◊ **SA-SLHD** is a heuristic (simulated annealing with a symmetric Latin hypercube design as initialization)

→ [W., Cornell University, 2009]



## II. Exploiting Structure in Simulation-Based Optimization

## Beyond the Black Box

$$\min f(x) = F[S(x)]$$

So far,  $f = S$

## Beyond the Black Box

$$\min f(x) = F[S(x)]$$

So far,  $f = S$

Your problems are not black-box problems

## Beyond the Black Box

$$\min f(x) = F[S(x)]$$

So far,  $f = S$

Your problems are not black-box problems

You formulated the problem

⇒ You know more than nothing

# Structure in Simulation-Based Optimization, $\min f(x) = F[S(x)]$

$f$  is often not a black box  $S$

- ◊ Nonlinear least squares

$$f(x) = \frac{1}{2} \sum_i (\textcolor{red}{S}_i(\textcolor{red}{x}) - d_i)^2$$

- ◊ Not all variables enter simulation

$$f(x) = g(x_I, x_J) + h(\textcolor{red}{S}(\textcolor{red}{x}_J))$$

- ◊ Only some constraints depend on simulation

$$\min\{f(x) : c_1(x) = 0, c_{\textcolor{red}{S}}(x) = 0\}$$

- ◊ Slack variables

$$\Omega_S = \{(x_I, x_J) : \textcolor{red}{S}(\textcolor{red}{x}_J) + x_I = 0, x_I \geq 0\}$$

Model-based methods can be a great way to exploit this structure

## Ex. 1- Least Squares $f(x) = \frac{1}{2} \sum_i F_i(x)^2$

Obtain a vector of output  $F_1(x), \dots, F_p(x)$

- ◊ Model each  $F_i$

$$F_i(x) \approx q_k^{(i)}(x) = F_i(x_k) + (x - x_k)^\top g_k^{(i)} + \frac{1}{2}(x - x_k)^\top H_k^{(i)}(x - x_k)$$

- ◊ Approximate:

$$\nabla f(x) = \sum_i \nabla \mathbf{F}_i(\mathbf{x}) F_i(x) \rightarrow \sum_i \nabla q_k^{(i)}(x) F_i(x)$$

$$\begin{aligned} \nabla^2 f(x) &= \sum_i \nabla \mathbf{F}_i(\mathbf{x}) \nabla \mathbf{F}_i(\mathbf{x})^T + \sum_i F_i(x) \nabla^2 \mathbf{F}_i(\mathbf{x}) \\ &\rightarrow \sum_i \nabla q_k^{(i)}(x) \nabla q_k^{(i)}(x)^T + \sum_i F_i(x) \nabla^2 q_k^{(i)}(x) \end{aligned}$$

- ◊ Model  $f$  via Gauss-Newton or similar

→ [DFLS; Zhang, Conn, Scheinberg] regularized Hessians  
→ [POUNDERS; W., More] uses full Newton

## Ex. 1- Consequences for $f(x) = \frac{1}{2} \sum_i F_i(x)^2$

- ◊ Save linear algebra using interpolation set  $\mathcal{Y}_k$  common to all models
  - ◆ Single system solve, multiple right hand sides
$$\Phi(\mathcal{Y}_k) \begin{bmatrix} z^{(1)} & \dots & z^{(p)} \end{bmatrix} = \begin{bmatrix} \underline{E}_1 & \dots & \underline{E}_p \end{bmatrix}$$
  - ◆ fully linear  $q^{(0)}$   $\Rightarrow$  all  $q^{(i)}$  fully linear
- ◊ (nearly) exact gradients for  $F_i$  (nearly) linear
- ◊ No longer interpolate function at data points

$$\begin{aligned} m(x_k + \delta) = & f(x_k) + \delta^T \sum_i g_k^{(i)} F_i(x_k) \\ & + \frac{1}{2} \delta^T \sum_i \left( g_k^{(i)} (g_k^{(i)})^T + F_i(x_k) H_k^{(i)} \right) \delta \\ & + \text{missing h.o. terms} \end{aligned}$$

## Ex. 1- Calibrating Energy Density Functionals

$$\min_x \left\{ f(x) = \sum_{i=1}^p \left( \frac{s_i(x) - d_i}{w_i} \right)^2 \right\}$$

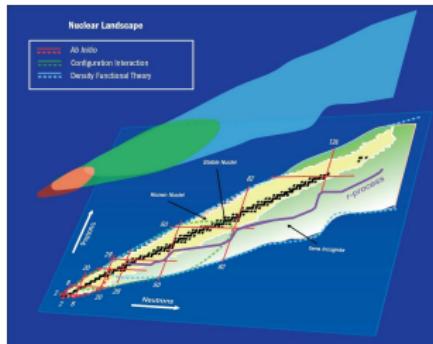
$s_i(x)$  Simulated (DFT) nucleus property

$d_i$  Nucleus  $i$  experimental data

$w_i$  Weight for data type  $i$

$p$  Parallel simulations

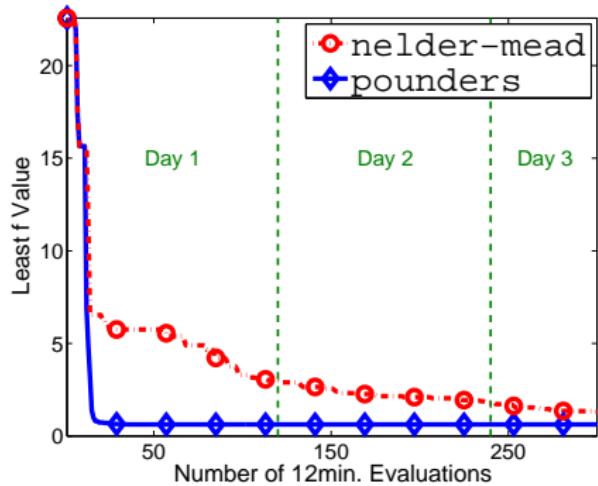
- ◊ Engineering starting point  $x_0$
- ◊ Bound constraints scaled to unit cube
- ◊ 12 CPU minutes per evaluation



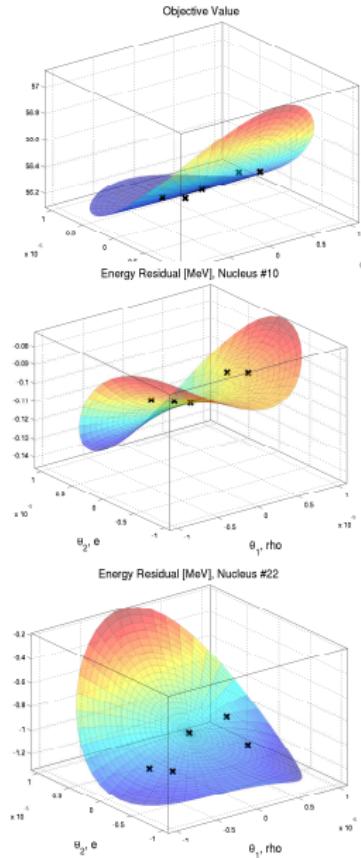
SciDAC physics project UNEDF

## Ex. 1- POUNDERS in Practice

- ◊ Models for each residual  $\frac{d_i - s_{t,i}(x)}{\sigma_i}$ .  
 $\{m^i(x)\}_{i=1}^{90}$
- ◊ Further reduces # of expensive evaluations



→ [Kortelainen et al., PhysRevC '10]



## Ex. 2- Some Known Partial

$x = (x_I, x_J)$ ; have  $\frac{\partial f}{\partial x_I}$  but not  $\frac{\partial f}{\partial x_J}$

“Solve”

$$\Phi z = \underline{f}$$

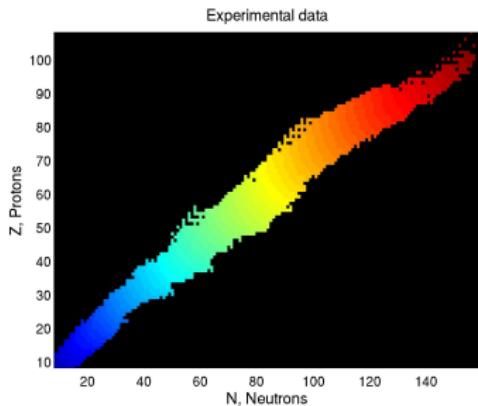
with known  $z_{g,I}, z_{H,I}$

$$\begin{bmatrix} \Phi_c & \Phi_{g,J} & \Phi_{H,J} \end{bmatrix} \begin{bmatrix} z_c \\ z_{g,J} \\ z_{H,J} \end{bmatrix} = \underline{f} - \Phi_{g,I} z_{g,I} - \Phi_{H,I} z_{H,I}$$

- ◊ Effectively lowers dimension to  $|J| = n - |I|$  for
  - ◆ approximation
  - ◆ model-improving evaluations
  - ◆ linear algebra
- ◊ Still have interpolation where required

## Ex. 2- Multi-level Optimization Structure

$$\min_x \left\{ f(x) = \sum_{i=1}^p (s_i(x) - d_i)^2 \right\}$$



[Bertolli, Papenbrock, W., PhysRevC '12]

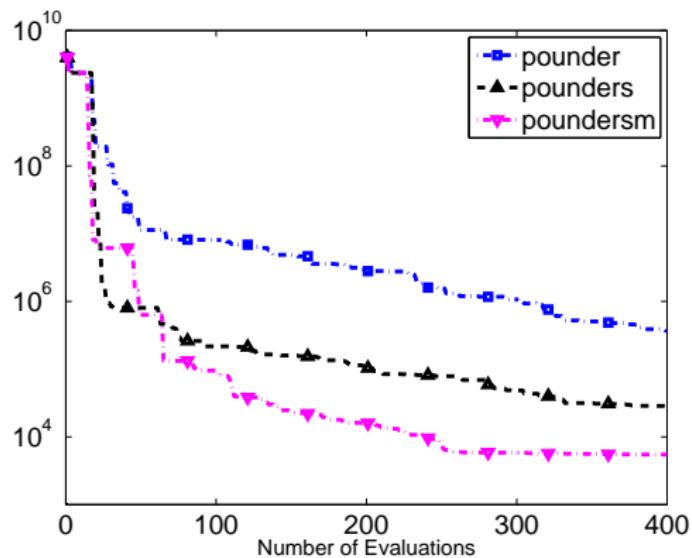
$s_i(x)$  solution to lower level problem

$$\begin{aligned} s_i(x) &= \tilde{g}_i(x) + \min_y \{h_i(x_J; y) : y \in \mathcal{D}_i\} \\ &= \tilde{g}_i(x) + h_i(x_J; y_{i,*}[x_J]) \\ \nabla_{x_J} s_i(x) &\approx \nabla_{x_J} \tilde{g}_i(x) + \nabla_{x_J} q^{(i)}(x_J) \end{aligned}$$

For  $x = (x_I, x_J)$

- ◊  $\nabla_{x_I} s_i(x_I, x_J)$  available
- ◊  $s_i(x)$  continuous and smooth in  $x_I$
- ◊  $\tilde{g}_i(x)$  cheap to compute!
- ◊ No noise/errors introduced in  $\tilde{g}_i(x)$

## Ex. 2- Numerical Results With Some Partials



**s** exploit least squares  
**m** use  $\nabla_{x_I}$  derivatives

- $\diamond n = 16, |I| = 3$
- $\diamond$  5-10 seconds/evaluation on 8 cores

“Same algorithmic framework”,  
performance advantages from exploiting structure

## Ex. 2- Convergence With Known Partial $\nabla_{x_I}$

$\lim_{k \rightarrow \infty} \nabla f(x_k) = 0$  as before

Approximation bounds

- ◊  $|q_k(x) - f(x)| \leq \kappa_1(\gamma_f + \|H_k\|)\Delta_k^2, \quad x \in \mathcal{B}_k$
- ◊  $\|g_k + H_k(x - x_k) - \nabla f(x)\| \leq \kappa_2(\gamma_f + \|H_k\|)\Delta_k, \quad x \in \mathcal{B}_k$

with most constants now a function of  $n - |I|$ .

Guaranteed

- ◊ strict descent
- ◊ optimality

in some directions.

## Ex. 3- General Constraints

$$\min\{f(x) : c_1(x) = 0, c_S(x) = 0\}$$

- ◊ Approximate Lagrangian:

$$\begin{aligned}\nabla L &= \nabla f + \lambda_1^T \nabla c_1 + \lambda_2^T \nabla \mathbf{c}_S \\ &\rightarrow \nabla f + \lambda_1^T \nabla c_1 + \lambda_2^T \nabla \mathbf{m}\end{aligned}$$

- ◊ Use favorite method: filters, augmented Lagrangian, ...
- ◊ Slack variables
  - ◆ Do not increase effective dimension
  - ◆ Subproblems can treat separately
  - ◆ Know derivatives

→ [Diniz-Ehrhardt, Martínez, Pedroso, C&A Math. 2011]: modified AL methods

## Ex. 3- What Constraint Derivatives Buy You

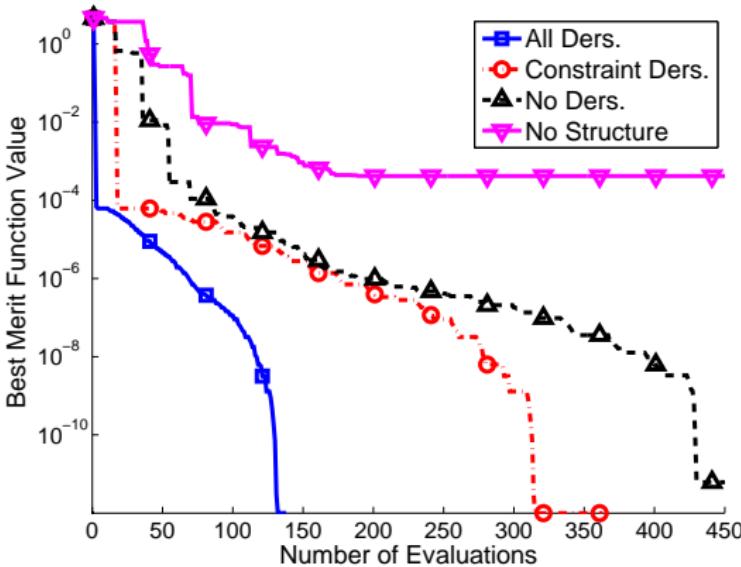
Augmented Lagrangian methods,  $L_A(x, \lambda; \mu) = f(x) - \lambda^T c(x) + \frac{1}{\mu} \|c(x)\|^2$

$$\min_x \{f(x) : c(x) = 0\}$$

Four choices:

1. Penalize constraints
2. Treat  $c$  and  $f$  both as (separate) black boxes
3. Work with  $f$  and  $\nabla_x c$
4. Have both  $\nabla_x f$  and  $\nabla_x c$

→ With Slava Kungurtsev (UCSD)



$n = 15, 11$  constraints

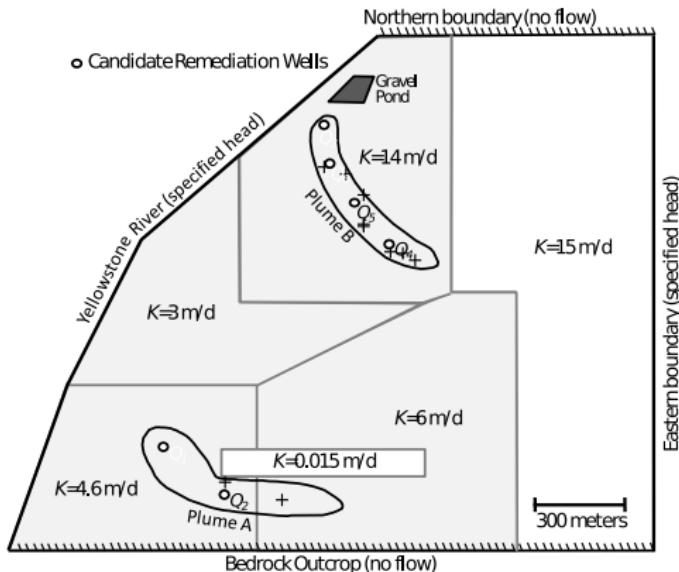
## Ex. 4- Remediation of Chlorinated Solvents

Determine extraction rates  
for 6 installed wells

### Structure

- Minimize operating cost (linear)
- Plume flux constraints (expensive simulation)

$$c_P(x) = \sum_{i=1}^p |S_{P,i}(x)|$$



Lockwood Solvent Ground Water Plume Site (LSGPS)

## Ex. 4- Exploiting Nonsmoothness

What if derivatives of  $f(x)$  do not always exist?

$$f(x) = g(x, S(x))$$

- ◊  $g$  **Nonsmooth, known**
- ◊  $S$  **Smooth, black-box**

Examples:

- ◆  $f(x) = \sum_{i=1}^p |F_i[S(x)]|$
- ◆  $f(x) = \max(S_1(x), S_2(x))$

→ With Aswin Kannan

## Ex. 4- Targeting Nonsmoothness in $f(x) = \sum_{i=1}^p |F_i(x)|$

### Model-based Approaches:

pounder Ignore structure, model  $f$  as usual

pounders-sqrt  $f = \sum_{i=1}^p \sqrt{|F_i|}^2,$

model  $\sqrt{|F_i|}$  by  $Q_i$

subproblem  $\min \sum_{i=1}^p \tilde{Q}_i(x)^2$

poundera-abs  $f = \sum_{i=1}^p |F_i|,$

model  $|F_i|$  by  $Q_i$

subproblem  $\min \sum_{i=1}^p Q_i(x)$

poundera-nsm  $f = \sum_{i=1}^p |F_i|,$

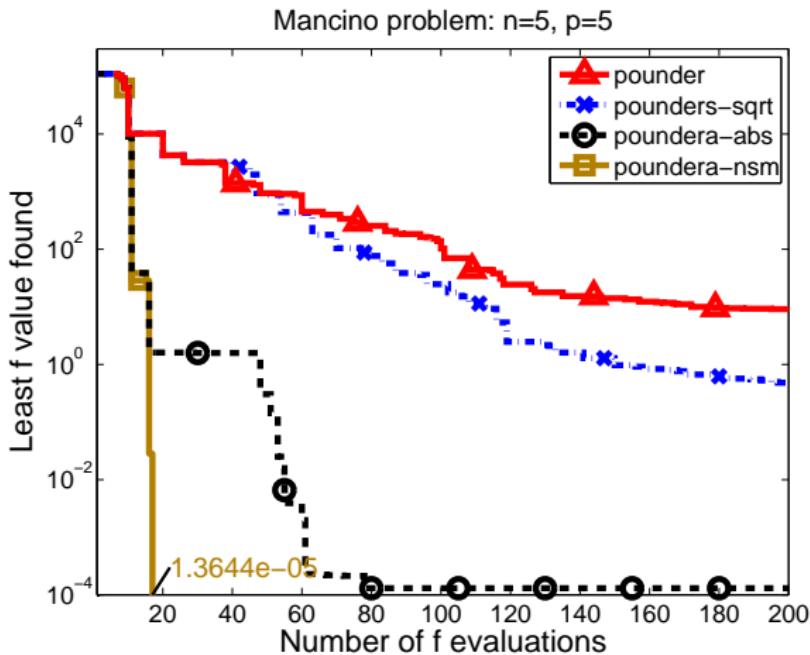
model  $F_i$  by  $Q_i$

subproblem  $\min \sum_{i=1}^p |Q_i(x)|$



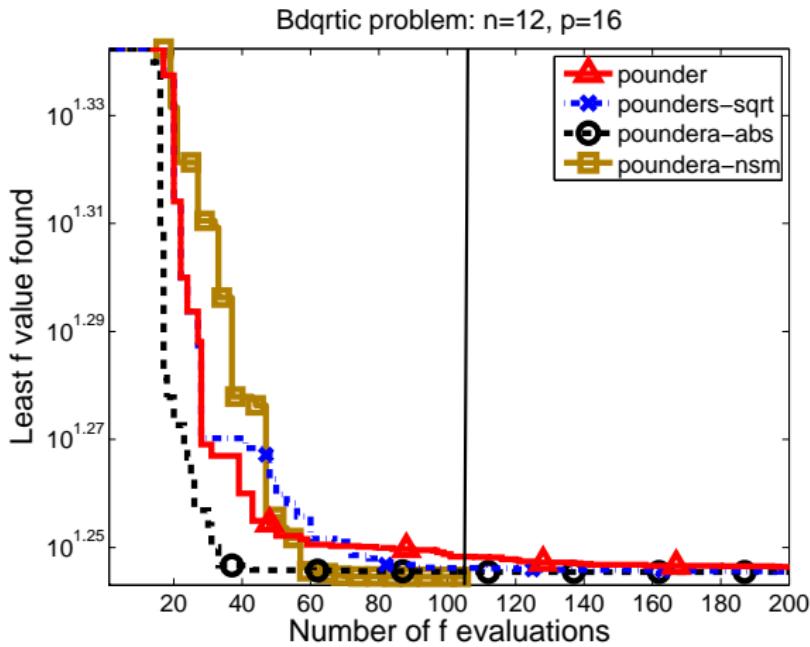
## Ex. 4- Preliminary Results, $\min \sum_{i=1}^p |F_i(x)|$

pounder black-box  
pounders-sqrt  $\sum_{i=1}^p \tilde{Q}_i(x)^2$   
poundera-abs  $\sum_{i=1}^p Q_i(x)$   
poundera-nsm  $\sum_{i=1}^p |Q_i(x)|$



## Ex. 4- Preliminary Results, $\min \sum_{i=1}^p |F_i(x)|$

pounder black-box  
pounders-sqrt  $\sum_{i=1}^p \tilde{Q}_i(x)^2$   
poundera-abs  $\sum_{i=1}^p Q_i(x)$   
poundera-nsm  $\sum_{i=1}^p |Q_i(x)|$



# So You Want To Solve A Hard Optimization Problem?

Mathematically unwrap problems to expose the deepest black boxes!

- ◊ It is easy to get started with derivative-free methods
- ◊ You should strive to obtain derivatives & apply methods from Todd's talk
  
- ◊ Model-based methods can make use of expensive function values
- ◊ Structure is everywhere, even in "black-box" / legacy code-driven optimization problems
- ◊ By exploiting structure, optimization can solve grand-challenge problems in `<insert your field here>`:
  - ◆ Model residuals  $\{r_i(x)\}_i$ , not  $\|r(x)\|$
  - ◆ Model constraints  $\{c_i(x)\}_i$ , not a penalty  $P(c(x))$
  - ◆ Explicitly handle nonsmoothness (and noise)



### III. Computational Noise

### III. Computational Noise

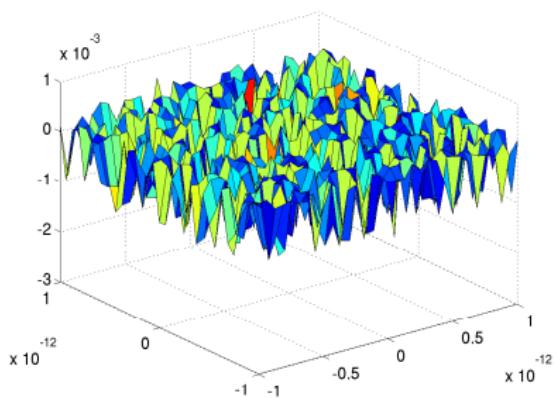
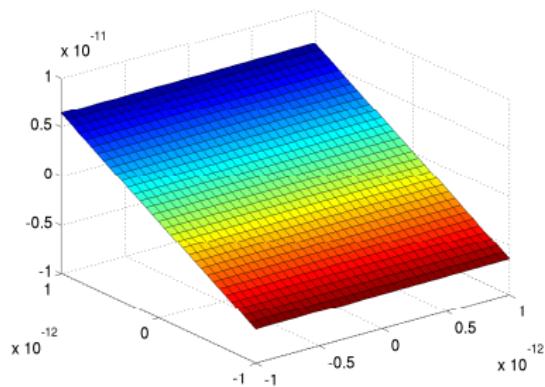
- ◊ What is computational noise?
- ◊ How can noise be estimated efficiently?
- ◊ How does noise affect numerical differentiation?
- ◊ How accurate are near-optimal finite-difference estimates?

## Two Questions To Ask Yourself

1. Do you know how accurate your derivatives are?
2. What do you do with this information?

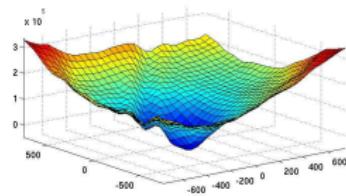
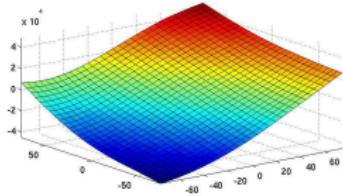
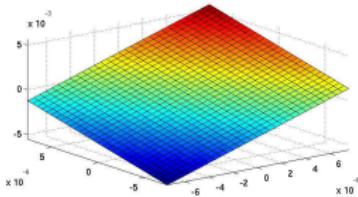
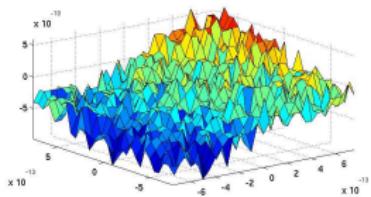
# Noise May Hurt You, Or It May Not

These are the same problem:



# Noise May Hurt You, Or It May Not

So are these:



# Computational Noise is not a Newcomer

From Hamming's 1971 Introduction to Numerical Analysis:

Where does this noise come from? . . . *infinite processes in mathematics which of necessity must be approximated by finite processes.*

Truncation vs. roundoff *Finite number length leads to roundoff. Finite processes lead to truncation.*



Competing errors *Smaller steps usually reduce truncation error and may increase roundoff error.*

Deterministic *In practice, the same input, barring machine failures, gives the same result.*

# Computational Noise is not a Newcomer

From Hamming's 1971 Introduction to Numerical Analysis:

Where does this noise come from? . . . *infinite processes in mathematics which of necessity must be approximated by finite processes.*

Truncation vs. roundoff *Finite number length leads to roundoff. Finite processes lead to truncation.*



Competing errors *Smaller steps usually reduce truncation error and may increase roundoff error.*

Deterministic *In practice, the same input, barring machine failures, gives the same result. ← changing!*

## Floating Point Arithmetic

Commutative:

$$A + B = B + A \quad \text{and} \quad A * B = B * A$$

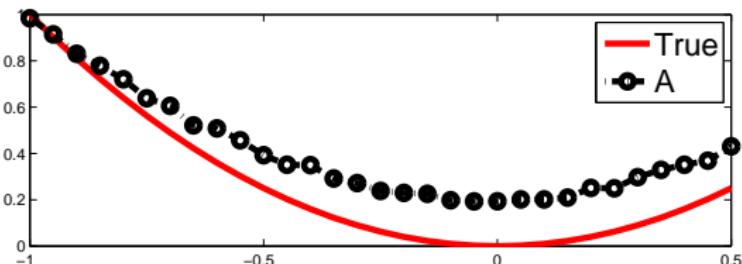
Non-associative:

$$A + (B + C) \neq (A + B) + C$$

- ◊ This is likely to affect the reproducibility of your calculations in the future (for performance reasons)

Many details → [What Every Computer Scientist Should Know About Floating-Point Arithmetic, Goldberg, 1991]

# The Effects of Computational Noise



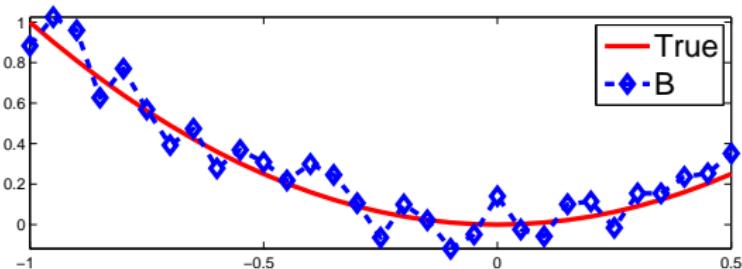
Noise is not **truncation error**

$$R_{m+1}(x) = f_a(x) - \sum_{i=0}^m P_i(x)$$

and is not **roundoff error**

$$f_\infty(x) - f(x)$$

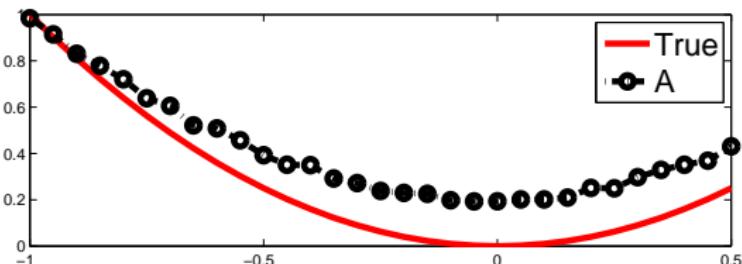
Which do you prefer?



A less noise, more error

B less error, more noise

# The Effects of Computational Noise



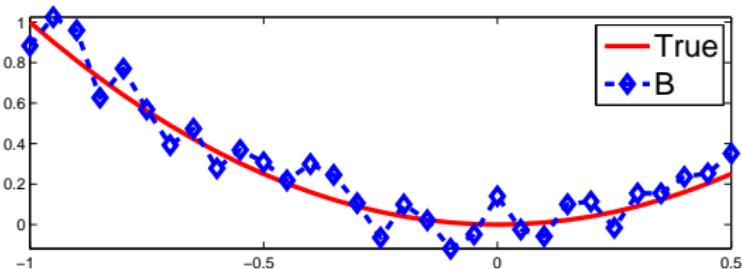
Noise is not **truncation error**

$$R_{m+1}(x) = f_a(x) - \sum_{i=0}^m P_i(x)$$

and is not **roundoff error**

$$f_\infty(x) - f(x)$$

Which do you prefer?



It matters how noisy your simulation is!

A less noise, more error

- Optimization
- Sensitivity Analysis

B less error, more noise

- Physics

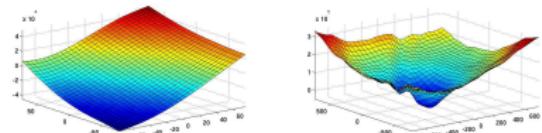
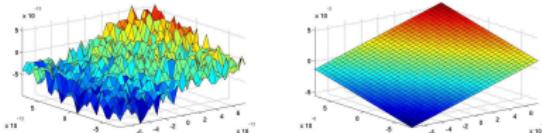
# Computational Noise in Deterministic Simulations

$$\text{Difference } |f(x) - f(x + Zw)|,$$

Finite precision + finite processes

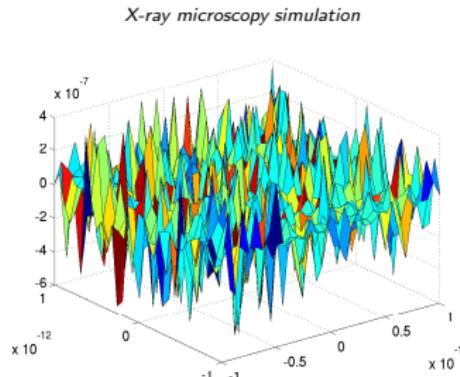
- ◊ Iteratively solving systems of PDEs or estimating eigenvalues
- ◊ Adaptively computing integrals
- ◊ Discretizations/meshes

destroy underlying smoothness



Goal: estimate the “variation” in  $f(x)$

- ◊ a few  $f$  evaluations
- ◊ deterministic and stochastic noise

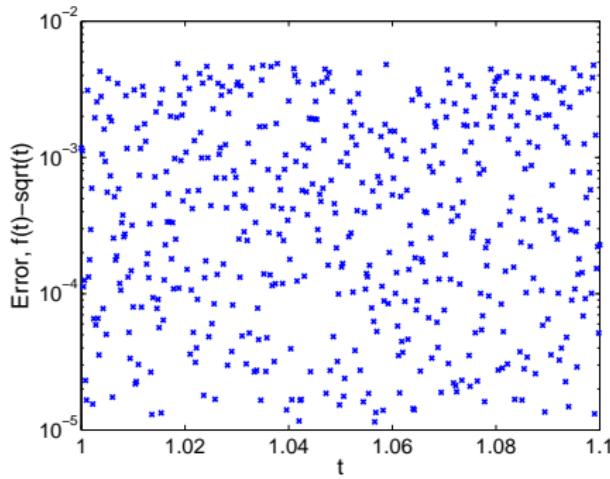
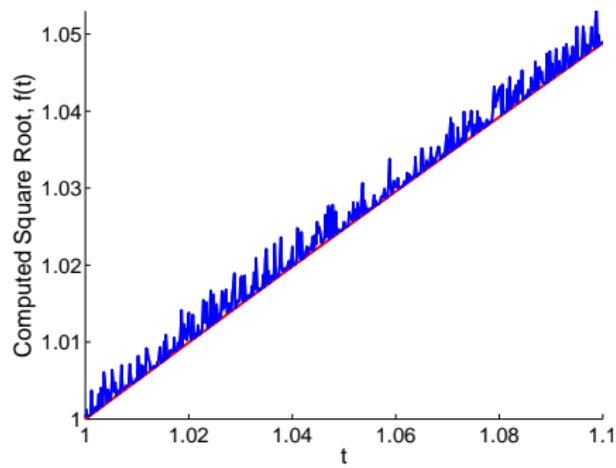


Sparse linear large-scale system

## Ex.- Square Roots with Newton's Method

Compute  $f(t) = \sqrt{t}$  near  $t = 1$

- ◊ Random starting point in  $[0,1]$
- ◊ Stop when  $|f(t)^2 - t| \leq \tau = 10^{-2}$



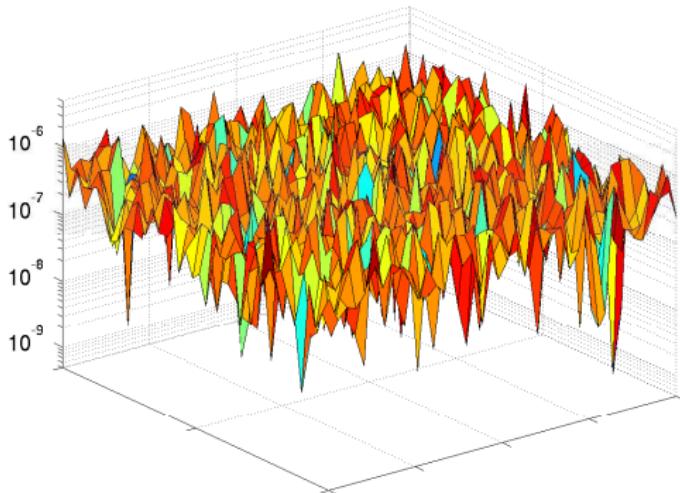
## Higher Dimensional Example: Eigenvalues

### Partial trace problem

$$f(x) = \sum_{i=1}^k \lambda_i(A + \text{diag}(x))$$

`A =delsq(numgrid('L',20))`  
243 × 243 symmetric matrix  
 $\kappa(A) \approx 10^2$  well-conditioned

$\lambda(\cdot)$  sorted eigenvalues  
compute with `eigs`,  $\tau = 10^{-3}$   
 $k = 5$   
 $x = \text{rand}(n, 1)$



Difference  $|f(x) - f(x + Z\omega)|$ ,  
 $Z$  a 2d random slice,  $\|\omega\|_\infty \leq 10^{-12}$

# The Noise Level $\epsilon_f$

Simple model for the noise

$$f(t) = f_s(t) + \varepsilon(t), \quad t \in \mathcal{I}$$

$f$  the computed function

$f_s$  a smooth, deterministic function

$\varepsilon$  is the noise with  $\{\varepsilon(t) : t \in \mathcal{I}\}$  iid

← only assumption

The noise level of  $f$  is  $\varepsilon_f = (\text{Var} \{\varepsilon(t)\})^{1/2}$

(*independent of  $t$* )

## The $k$ -th Order Difference $\Delta^k f(t)$

$$\Delta^{k+1} f(t) = \Delta^k f(t+h) - \Delta^k f(t), \quad \Delta^0 f(t) = f(t)$$

$$\Delta^k f(t) = \Delta^k f_s(t) + \Delta^k \varepsilon(t)$$

1. Differences of smooth  $f_s$  tend to zero rapidly
2. Differences of noise are bounded away from zero
  - ◆ If  $h$  is sufficiently small,

$$\Delta^k f(t) \approx \Delta^k \varepsilon(t)$$

- ◆ If  $f_s$  is  $k$ -times differentiable,

$$\Delta^k f(t) = f_s^{(k)}(\xi_k)h^k + \Delta^k \varepsilon(t), \quad \xi_k \in (t, t+kh)$$

Goal: make  $h$  small enough to remove smooth component

# Theory Underlying the ECNoise Algorithm

For  $\{\varepsilon(t + ih) : i = 0, \dots, m\}$  iid and  $k \leq m$ :

1.  $E\left\{\Delta^k \varepsilon(t)\right\} = 0$

2.  $\gamma_k E\left\{[\Delta^k \varepsilon(t)]^2\right\} = \varepsilon_f^2 \quad \gamma_k = \frac{(k!)^2}{(2k)!}$

3. If  $f_s$  is continuous at  $t$ , then

$$\lim_{h \rightarrow 0} \gamma_k E\left\{\left[\Delta^k f(t)\right]^2\right\} = \varepsilon_f^2$$

4. If  $f_s$  is  $k$ -times continuously differentiable at  $t$ , then

$$\lim_{h \rightarrow 0} \frac{\gamma_k E\left\{[\Delta^k f(t)]^2\right\} - \varepsilon_f^2}{h^{2k}} = \gamma_k \left[f_s^{(k)}(t)\right]^2$$

$$\Rightarrow \varepsilon_f^2 \approx \gamma_k E\left\{[\Delta^k f(t)]^2\right\},$$

when the sampling distance  $h$  is sufficiently small

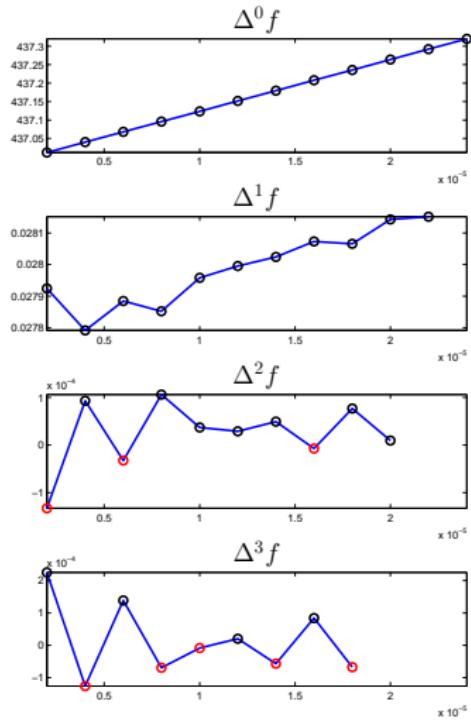


# The ECNoise Algorithm

Uses  $\sigma_k = \left( \frac{\gamma_k}{m+1-k} \sum_{i=0}^{m-k} [\Delta^k f(t + ih)]^2 \right)^{1/2}$

1. Chooses  $k$
  2. Verifies  $h$  is small enough
- ◇ Works for deterministic  $f$

[Estimating Computational Noise. Moré & W., SISC 2011]



$$\text{ECNoise Estimator } \sigma_k = \left( \frac{\gamma_k}{m+1-k} \sum_{i=0}^{m-k} [\Delta^k f(t_i)]^2 \right)^{1/2}$$

For  $f(t) = \cos(t) + \sin(t) + 10^{-3}U_{[0, 2\sqrt{3}]}$  ( $m = 6, t_i = \frac{i}{100}$ )

$f(t_i)$	$\Delta f(t_i)$	$\Delta^2 f(t_i)$	$\Delta^3 f(t_i)$	$\Delta^4 f(t_i)$	$\Delta^5 f(t_i)$	$\Delta^6 f(t_i)$
1.003	7.54e-3	2.15e-3	1.87e-4	-5.87e-3	1.46e-2	-2.49e-2
1.011	9.69e-3	2.33e-3	-5.68e-3	8.73e-3	-1.03e-2	
1.021	1.20e-2	-3.35e-3	3.05e-3	-1.61e-3		
1.033	8.67e-3	-2.96e-4	1.44e-3			
1.041	8.38e-3	1.14e-3				
1.050	9.52e-3					
1.059						
$\sigma_k$	6.78e-3	8.96e-4	9.02e-4	9.93e-4	1.10e-3	1.14e-3

## Extension to Multivariate $g : \mathbb{R}^n \mapsto \mathbb{R}$

Given base point  $x_b \in \mathbb{R}^n$ , unit direction  $p \in \mathbb{R}^n$ , consider

$$f_p(t) = g(x_b + tp), \quad t \geq 0$$

Apply univariate theory

- ◊ Directional differences, directional derivatives
- ◊  $\varepsilon_f$  may now depend on a direction  $p \in \mathbb{R}^n$
- ◊ **ECnoise** uses  $T_{i,0} = f(x_b + ihp)$  with random unit direction  $p \in \mathbb{R}^n$

# Computational Experience with Stochastic Noise

Validate ECnoise and empirical properties of

$$\sigma_k^2 = \frac{\gamma_k}{m+1-k} \sum_{i=0}^{m-k} T_{i,k}^2$$

under known conditions:

- ◊ Known noise level  $\varepsilon_f$
- ◊ Theory directly applies

Target: every estimate within a factor  $\eta = 4$  of the mean

Noisy Quadratic,  $f(x) = (x^T x)(1 + R)$ ,  $x \in \mathbb{R}^{10}$

Estimate relative noise

$$\frac{\sigma_k}{f(x_b)} \approx \sqrt{\text{Var}\{R\}} = 10^{-3}$$

$x_b$  random base point

$p$  10000 random unit  
directions

$m$  evaluations

Noisy Quadratic,  $f(x) = (x^T x)(1 + R)$ ,  $x \in \mathbb{R}^{10}$

$$R \sim \text{Uniform}[-\sqrt{3} \cdot 10^{-3}, \sqrt{3} \cdot 10^{-3}]$$

Estimate relative noise

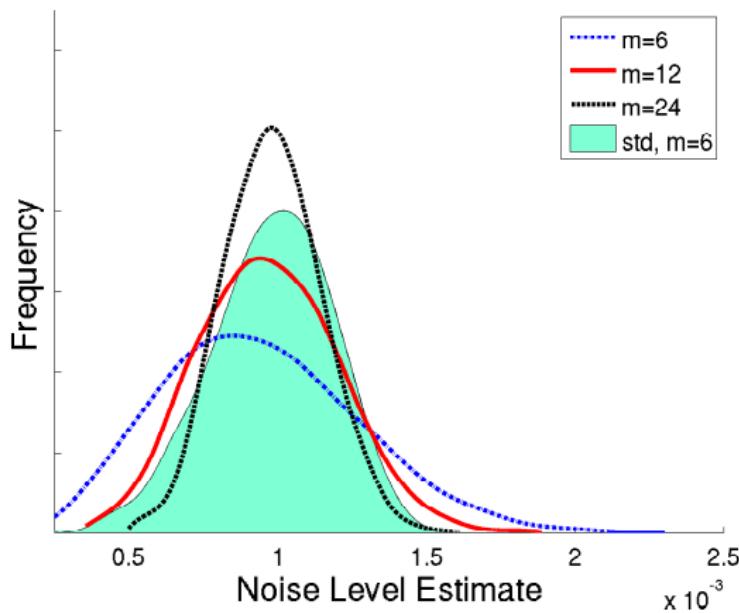
$$\frac{\sigma_k}{f(x_b)} \approx \sqrt{\text{Var}\{R\}} = 10^{-3}$$

$x_b$  random base point

$p$  10000 random unit  
directions

$m$  evaluations

99.2% within a factor  $\eta = 4$  for  
 $m = 6$



Noisy Quadratic,  $f(x) = (x^T x)(1 + R)$ ,  $x \in \mathbb{R}^{10}$

$$R \sim \text{Normal}(0, 10^{-6})$$

Estimate relative noise

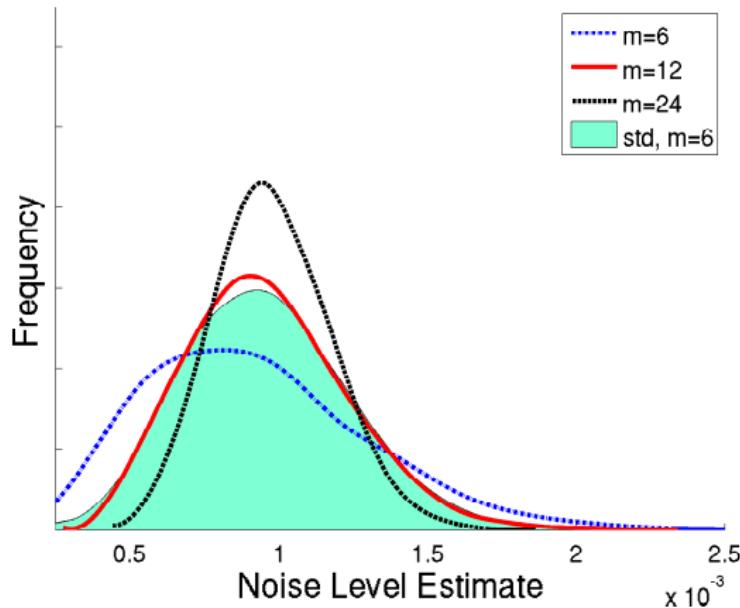
$$\frac{\sigma_k}{f(x_b)} \approx \sqrt{\text{Var}\{R\}} = 10^{-3}$$

$x_b$  random base point

$p$  10000 random unit directions

$m$  evaluations

98.9% within a factor  $\eta = 4$   
for  $m = 6$



# MC Finance Example with Higher Order Derivatives

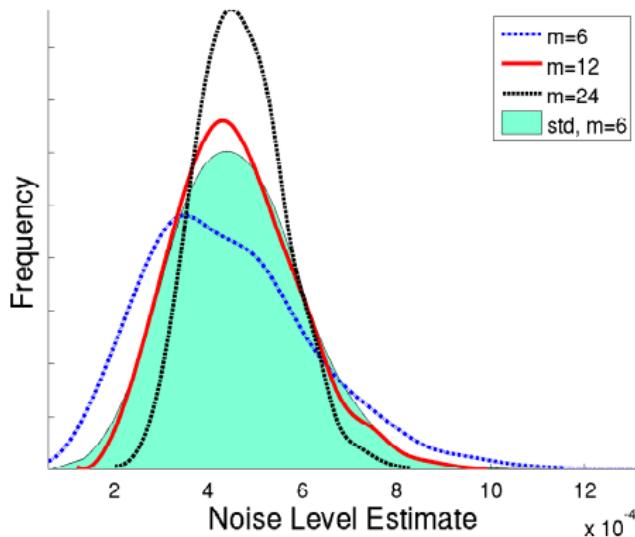
Today's value of a \$1 payment  $n$  years from now rates [Cafisch]:

$$f(x) = \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} \prod_{i=0}^n \frac{e^{-\frac{\|u\|^2}{2}}}{1+r_i(u,x)} du, \quad r_i(u,x) = \begin{cases} \frac{1}{10} & i=0 \\ r_{i-1}(u,x) e^{x_i u_i - x_i^2/2} & i \geq 1 \end{cases}$$

10000 MC integrations  
(directions  $p$ ) with

- ◊  $n = 3$  years,  
 $x_b = [.1, .1, .1]$
- ◊  $tol = 5000$  standard  
normal random  
variables
- ◊ no variance reduction

99.6% within a factor 4 for  $m = 6$



## Statistical Consistency of the Estimator $\sigma_k^2$ (with J. Gómez)

Entries in the difference table  $T_{i,k}, T_{j,k}$  are dependent!

Further assume  $f_s^{(k)}$  bounded and  $E\{\varepsilon(t)^4\} < \infty$

$$\lim_{m \rightarrow \infty} \text{Var}\{\sigma_k^2\} = 0$$

→ Can estimate stochastic noise without having to repeatedly sample at a single  $x$

- ◊ Choose  $k$  and  $h$  to control the bias,  $E\{\sigma_k^2 - \varepsilon_f^2\} = \gamma_k h^{2k} f_s^{(k)}(\xi_m)^2$
- ◊ Choose  $m \geq k$  to control the variance and ensure noise captured

Noisy Quadratic,  $f(x) = (x^T x)(1 + R)$ ,  $x \in \mathbb{R}^{10}$

$$R \sim \text{PearsonV}(3, 2 \cdot 10^{-3})$$

Estimate relative noise

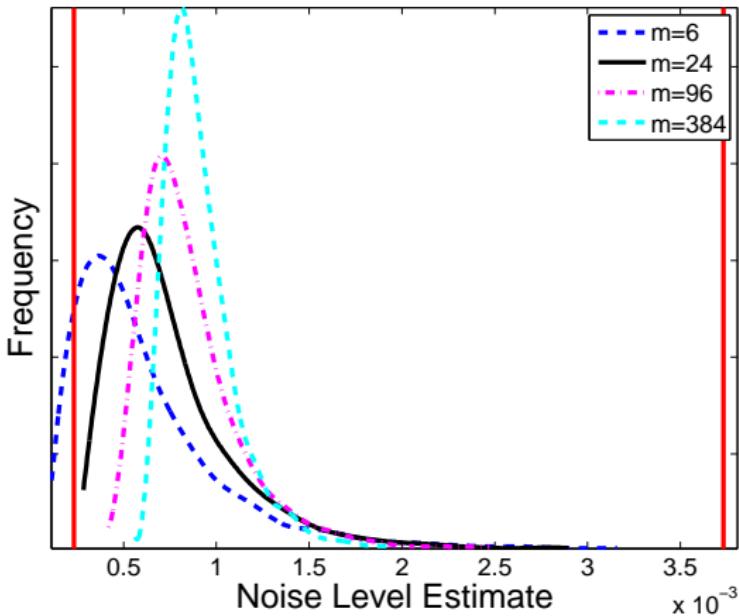
$$\frac{\sigma_k}{f(x_b)} \approx \sqrt{\text{Var}\{R\}} = 10^{-3}$$

$x_b$  base point

$p$  10000 random unit  
directions  $\mathcal{B}_1$

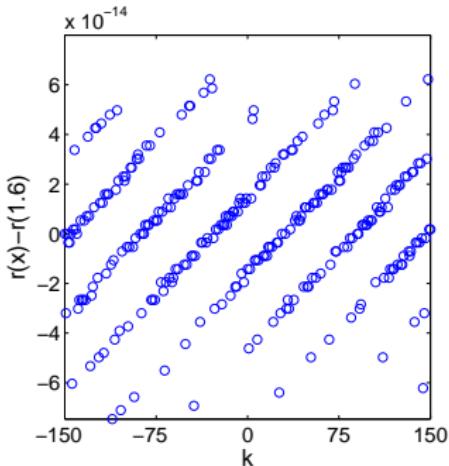
$m$  evaluations

89.3% (99.3%) within a factor  
 $\eta = 4$  for  $m = 6$  ( $m = 24$ )

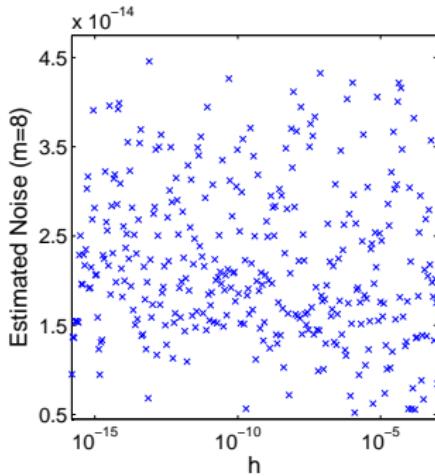


## Transition to Non-IID & Deterministic Noise

Kahan's  $r(x) = \frac{622-x(751-x(324-x(59-4x)))}{112-x(151-x(72-x(14-x)))}$  violates iid assumption



Kahan's  $r$  at  $x = 1.6 + 2^{-52}k$



Noise estimates for  $r(1.6)$

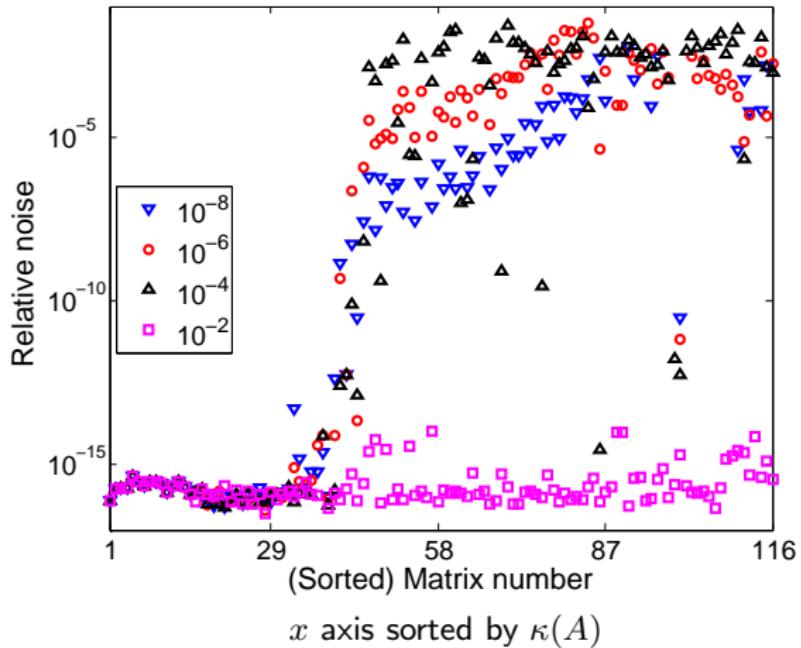
- ◊ All noise estimates within factor 4 of  $2 \cdot 10^{-14}$
- ◊ (Unlikely)  $m+1$  points solely on one line  $\Rightarrow \epsilon_f \approx 2 \cdot 10^{-15}$

$$f_\tau(t) = \|y_\tau(x_b + tp)\|_2^2,$$

$y_\tau$  from iterative solver for  $Ay_\tau(x) = x$   
with tolerance  $\tau > 0$

- ◊  $A = 116$  spd UF matrices ( $n < 10^4$ ), scaled by diagonal
- ◊ 28 with  $\kappa(A) \leq 10$ , 10 with  $\kappa(A) \geq 10^{10}$
- ◊ random direction  $p \in \mathbb{R}^n$
- ◊ variety of tolerances  $\tau$  in  $[10^{-8}, 10^{-2}]$
- ◊ only  $m = 8$  additional evaluations
- ◊ tested `bicgstab`, `gmres`, `idr(s)`, `pcg`, `minres`, `minresqlp`, `symmlq`

## ECNoise on bicgstab Functions $f_\tau$

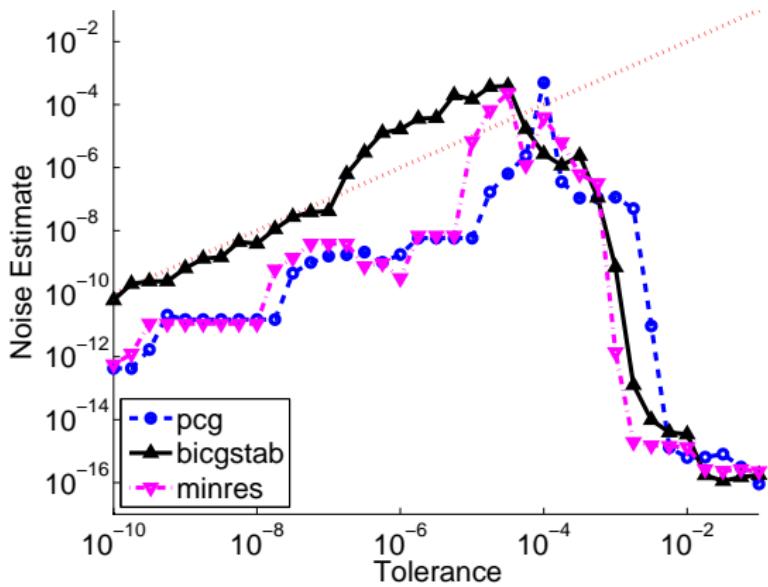


## Effect of the Tolerance $\tau$ for $Ay_\tau(x) = x$

$$f_\tau(t) = \|y_\tau(x_b + tp)\|_2^2,$$

$y_\tau$  from iterative solver for  $Ay_\tau(x) = x$  with tolerance  $\tau > 0$

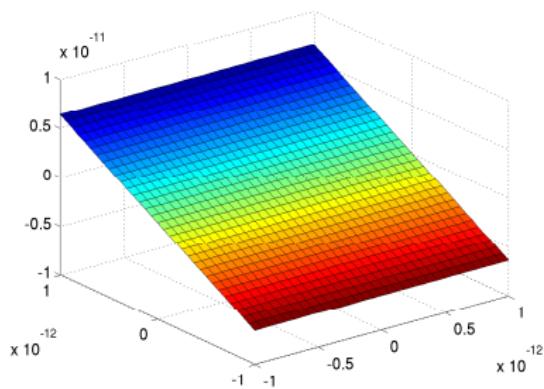
- ◊ Reliable estimates,  $m = 8$  additional evaluations
- ◊ Non-monotone relationship between the relative noise and tolerance  $\tau$
- ◊ random direction  $p$
- ◊ spd UF matrix  $A$



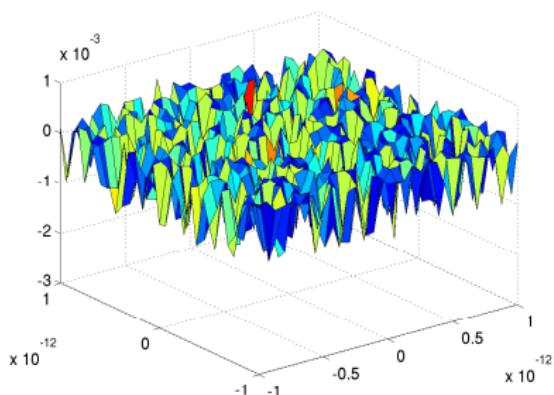
# How does bcsstk02 Noise Change with the Tolerance?

2d slice of  $f$  for bcsstk02 ( $n = 66, \kappa(A) = 1833$ )

bicgstab



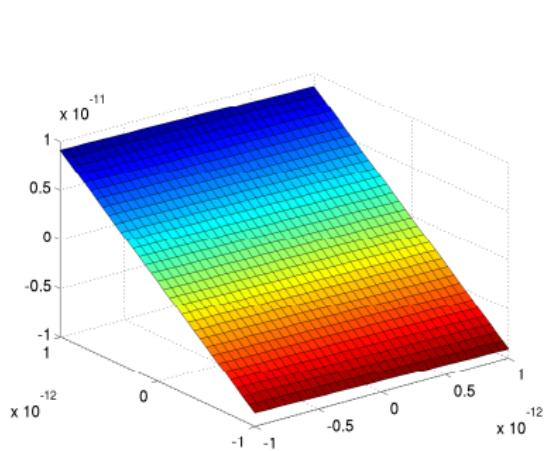
$$\tau = 10^{-2}, \text{ std} = 3.86\text{e-}12$$



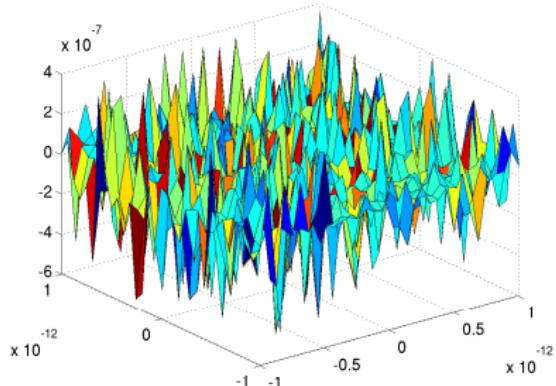
$$\tau = 10^{-5}, \text{ std} = 4.98\text{e-}04$$

## How does bcsstk02 Noise Change with the Tolerance?

2d slice of  $f$  for bcsstk02 ( $n = 66, \kappa(A) = 1833$ )

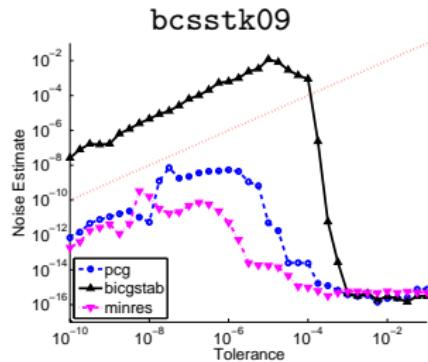
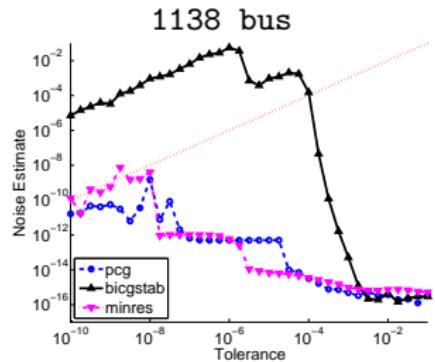
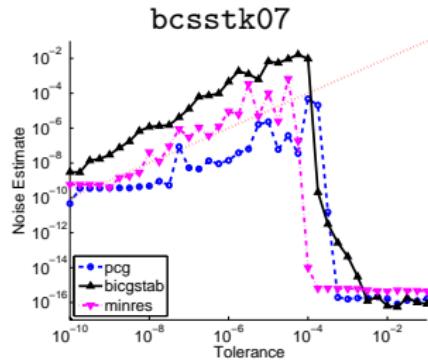
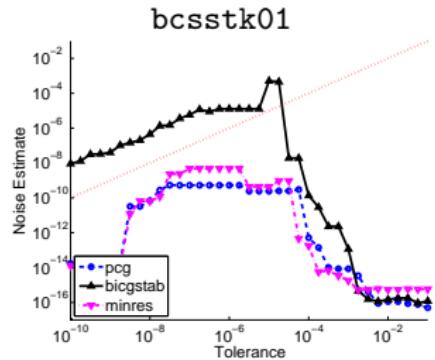


$$\tau = 10^{-2}, \text{ std} = 5.29\text{e-}12$$



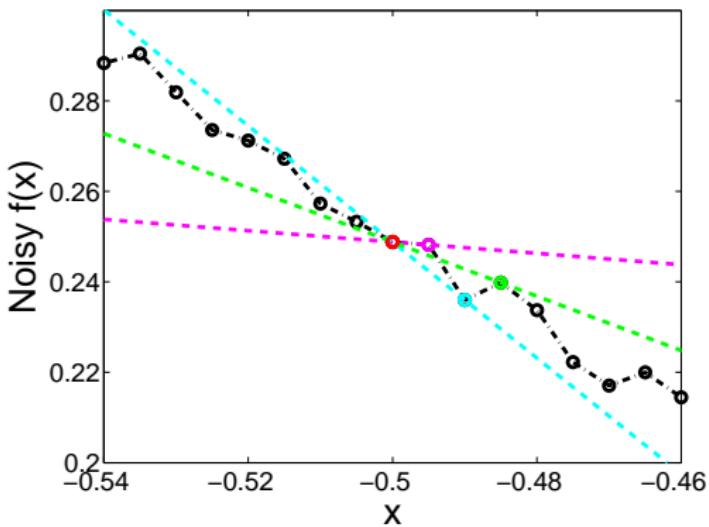
$$\tau = 10^{-5}, \text{ std} = 1.90\text{e-}07$$

# Noise Estimates for Different Tolerances



## Finite Differences Sensitive to Choice of $h$

$$\frac{f(t_0 + h) - f(t_0)}{h} \approx f'_s(t_0)$$



## Noisy Forward Differences

Minimize       $E \{ \mathcal{E}(h) \} = E \left\{ \left( \frac{f(t_0+h) - f(t_0)}{h} - f'_s(t_0) \right)^2 \right\}$

Our  $h$  will depend on

- ◊ Loose estimate of noise
- ◊ Loose estimate of  $|f''|$  !
- ◊ Stochastic theory:
  1.  $f(t) = f_s(t) + \epsilon$  on  $I = \{t_0 + h : 0 \leq h \leq h_0\}$
  2.  $f_s$  twice differentiable
  3.  $\mu_L \leq |f''_s| \leq \mu_M$  on  $I$

[Estimating Noisy Derivatives. Moré & W., TOMS 2012]]

# Optimal Forward Difference Parameter $h$

$$\frac{1}{4}\mu_L^2 h^2 + 2\frac{\varepsilon_f^2}{h^2} \leq E\{\mathcal{E}(h)\} \leq \frac{1}{4}\mu_M^2 h^2 + 2\frac{\varepsilon_f^2}{h^2}$$

$h \downarrow$  Variance (noise) dominates

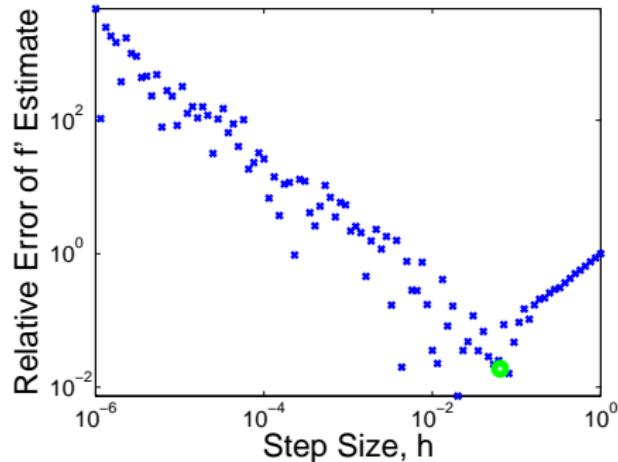
$h \uparrow$  Bias ( $f''$ ) dominates

For  $h_0$  sufficiently large

1. Upper bound minimized by

$$h_M = 8^{1/4} \left( \frac{\varepsilon_f}{\mu_M} \right)^{1/2}$$

2. When  $\mu_L > 0$ ,  $h_M$  is near-optimal:



$$E\{\mathcal{E}(h_M)\} = \sqrt{2}\mu_M\varepsilon_f \leq \left( \frac{\mu_M}{\mu_L} \right) \min_{0 \leq h \leq h_0} E\{\mathcal{E}(h)\}.$$

## Alternative FD Step Sizes

[Gill, Murray, Saunders, Wright; 1983]

Given uniform bound on roundoff error,

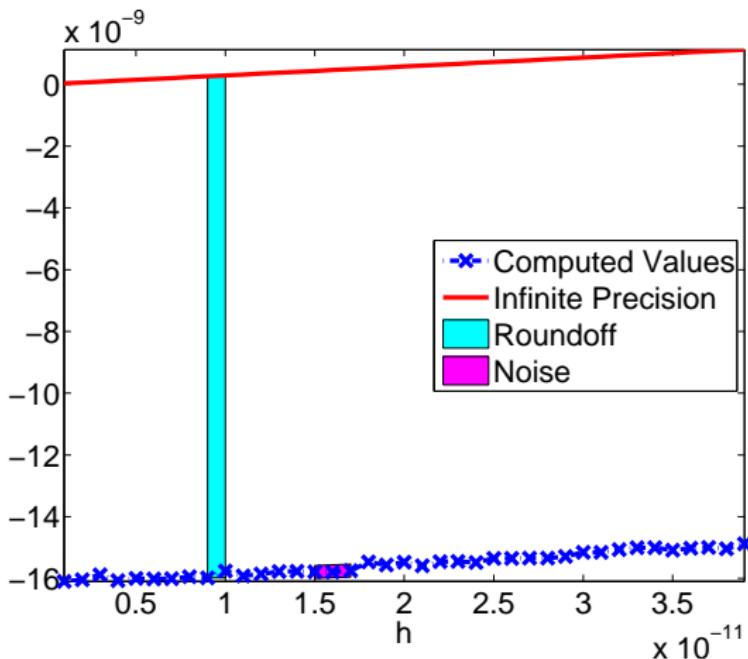
$$|f(t) - f_\infty(t)| \leq \varepsilon_A \quad t \in I,$$

Minimizer of (upper bound on)  
 $l_1$  error is

$$h_A = 2 \left( \frac{\varepsilon_A}{\mu_M} \right)^{1/2}$$

Assumes:

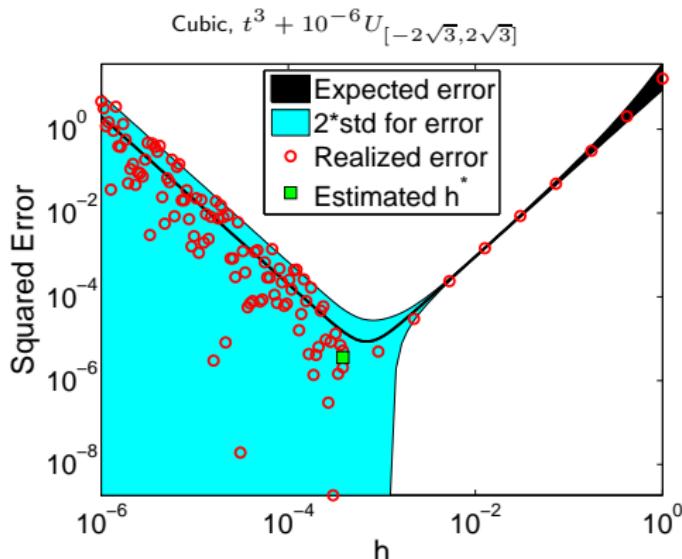
- ◊  $h_A \leq h_0$
- ◊ Estimate of  $\varepsilon_A$  available



## Stochastic Examples

Estimate  $f'_s(t) = E\{f(t)\}'$  at  $t = 1$

$(\varepsilon_f = 10^{-6})$



$$\text{Log-log realizations of } \mathcal{E}(h) = E \left\{ \left( \frac{f(t_0+h) - f(t_0)}{h} - f'_s(t_0) \right)^2 \right\}$$

Expected error and uncertainty regions predicted by the theory

## Extension: Central Differences

First derivatives,  $\frac{f(t_0+h)-f(t_0-h)}{2h}$

- ◊  $|h_M| = \gamma_5 \left( \frac{\varepsilon_f}{\mu_M} \right)^{1/3}, \quad \gamma_5 = 3^{1/3} \approx 1.44$
- ◊  $\mu_L \leq |f_s^{(3)}| \leq \mu_M$
- ◊  $E\{\mathcal{E}_c(h_M)\} \leq \left( \frac{\mu_M}{\mu_L} \right)^{2/3} \min_{|h| \leq h_0} E\{\mathcal{E}_c(h)\}$

Second derivatives,  $\frac{f(t_0+h)-2f(t_0)+f(t_0-h)}{h^2}$

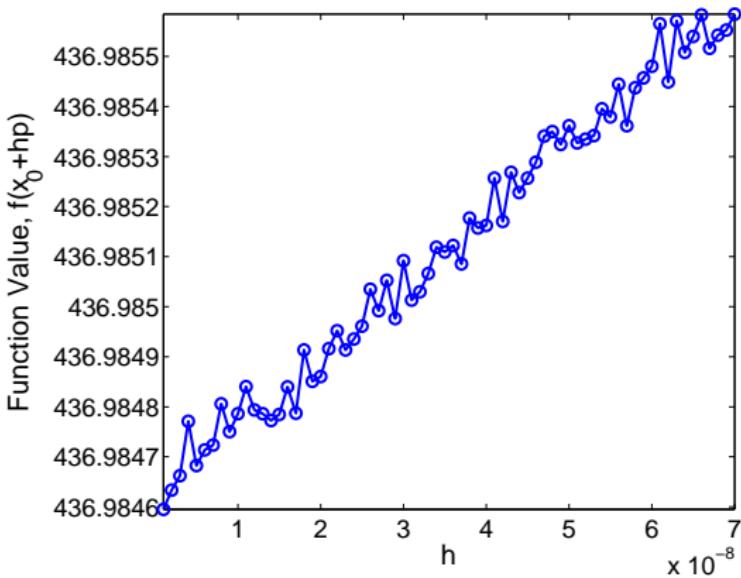
- ◊  $|h_M| = \gamma_7 \left( \frac{\varepsilon_f}{\mu_M} \right)^{1/4}, \quad \gamma_7 = 2^{5/8} 3^{1/8} \approx 2.33$
- ◊  $\mu_L \leq |f_s^{(4)}| \leq \mu_M$
- ◊  $E\{\mathcal{E}_2(h_M)\} \leq \left( \frac{\mu_M}{\mu_L} \right) \min_{|h| \leq h_0} E\{\mathcal{E}_2(h)\}$ 
  - ◆ use to obtain rough estimate of  $|f_s''|$  for forward-difference  $h$



## Ex.- Noisy Deterministic Functions (bicgstab, $\tau = 10^{-3}$ )

Subset of 100 UF matrices

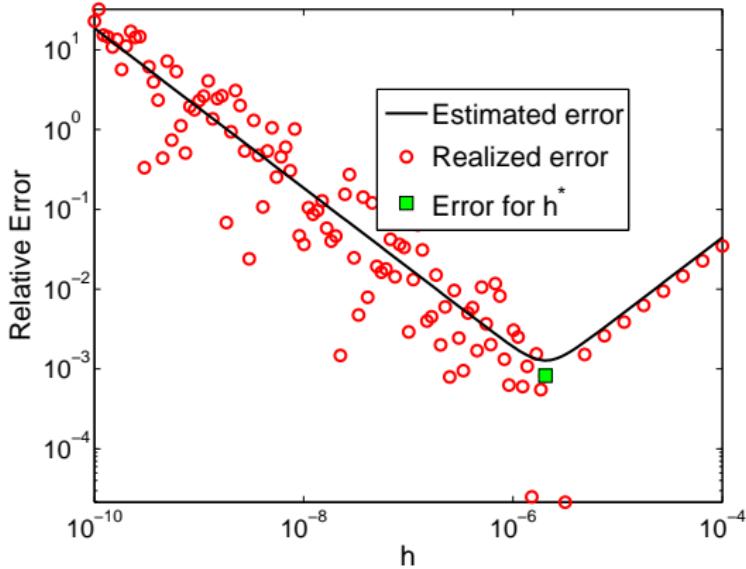
- ◊ FD sensitive to noise



## Ex.- Noisy Deterministic Functions (bicgstab, $\tau = 10^{-3}$ )

Subset of 100 UF matrices

- ◊ FD sensitive to noise
- ◊ Exhibits behavior similar to stochastic FD

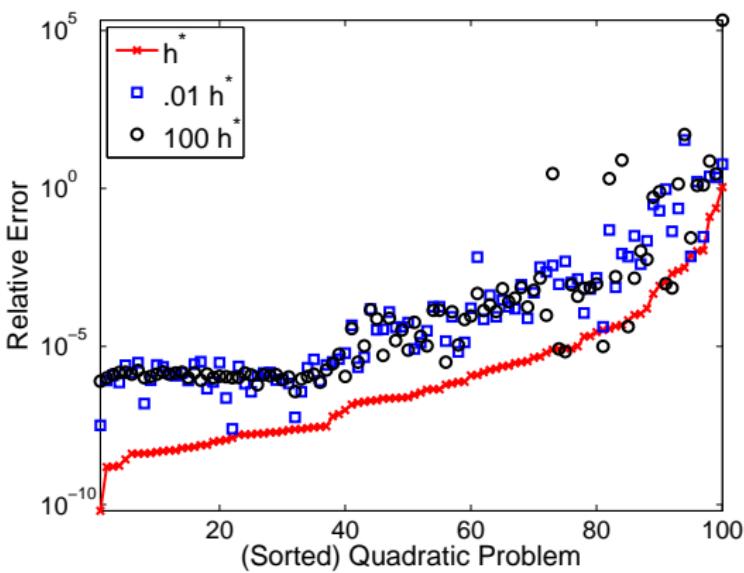


*Compared with AD (INTLAB) derivative*

## Ex.- Noisy Deterministic Functions (`bicgstab`, $\tau = 10^{-3}$ )

Subset of 100 UF matrices

- ◊ FD sensitive to noise
- ◊ Exhibits behavior similar to stochastic FD
- ◊  $h_M$  obtains 2 more correct digits than  $10^{\pm 2} h_M$
- ◊  $h_M$  significantly better than  $\sqrt{\epsilon_{\text{mach}}}$

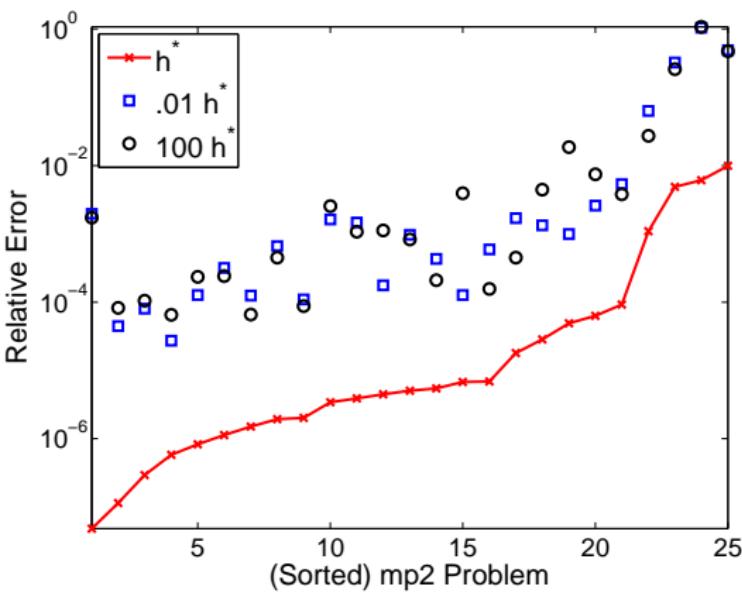


*Compared with AD (INTLAB) derivative*

## Ex.- Highly Nonlinear MINPACK-2 Problems

25 problems,  $n \leq 64 \cdot 10^4$

- ◊ Accurate estimates obtained even when  $f''$  not constant



*Compared with hand-coded derivative*

## Sensitivity to $\varepsilon_f$ and $\mu_M$ Estimates

If  $h^* = \alpha h_M$ , then

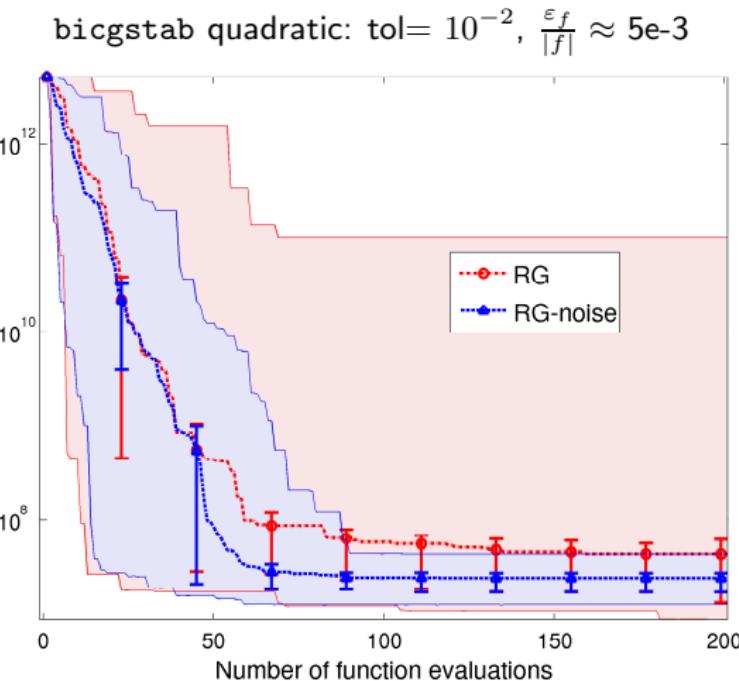
$$E\{\mathcal{E}(h^*)\}^{1/2} \leq \left(\frac{1}{\alpha} + \alpha\right) E\{\mathcal{E}(h_M)\}^{1/2}$$

- ◊ mis-estimating  $\varepsilon_f/\mu_M$  by factor  $10^d$  ( $\alpha = 10^{\pm d/2}$ ) gives  
 $E\{\mathcal{E}(h^*)\}^{1/2} \leq 10^{d/2} E\{\mathcal{E}(h_M)\}^{1/2}$
- ◊ rms error grows slowly with ratio  $\alpha = h^*/h_M$

# Using the Noise in Nesterov's Random Gradient Method

## General RG iteration

1. Generate direction  $d_k$
2. Evaluate gradient-free oracle  $g(x_k; h_k) = \frac{f(x_k + h_k d_k) - f(x_k)}{h} d_k$
3. Compute  $x_{k+1} = x_k - \delta_k g(x_k; h_k)$ , evaluate  $f(x_{k+1})$



## Summary: How Loud Are Your Simulations?

- ◊ Computational noise complicates analysis of simulation-based functions, worst-case bounds overly pessimistic (see Baudouin talk)
- ◊ With a few (6-8) additional evaluations, ECNoise reliably estimates the noise
- ◊ Stochastic theory for near-optimal difference parameters
- ◊ Coarse estimates of  $|f''|$  (2-4 evaluations) yield more accurate directional derivatives
- ◊ Both work on deterministic functions in practice



Some refs <http://mcs.anl.gov/~wild>:

[*Estimating Computation Noise*, SISC 2011]

[*Estimating Derivatives of Noisy Simulations*, TOMS 2012]

[*Do You Trust Derivatives or Differences?* Preprint, 2013]

[*Obtaining Quadratic Models of Noisy Functions*, Preprint, 2013]

Computing <http://mcs.anl.gov/~wild/cnoise>

Merci!