

Experimental Comparisons of Derivative Free Optimization Algorithms

Anne Auger Nikolaus Hansen J. M. Perez Zerpa
Raymond Ros **Marc Schoenauer**

TAO Project-Team, INRIA Saclay – Île-de-France, and
Microsoft-INRIA Joint Centre, Orsay, FRANCE

<http://tao.lri.fr>
First.Last@inria.fr

SEA'09, 4 juin 2009

Problem Statement

Continuous Domain Search/Optimization

The problem

- Minimize a **objective function** (*fitness function, loss function*) in continuous domain

$$f : \mathcal{S} \subseteq \mathbb{R}^n \rightarrow \mathbb{R},$$

- in the **Black Box scenario** (direct search)



Hypotheses

- domain specific knowledge only used within the black box
- gradients are not available

Problem Statement

Continuous Domain Search/Optimization

The problem

- Minimize a **objective function** (*fitness function, loss function*) in continuous domain

$$f : \mathcal{S} \subseteq \mathbb{R}^n \rightarrow \mathbb{R},$$

- in the **Black Box scenario** (direct search)



Typical Examples

- shape optimization (e.g. using CFD) curve fitting, airfoils
- model calibration biological, physical
- parameter identification controller, plants, images

The practitioner's point of view

Issues

- How to choose the best algorithm?
 - ▶ For a given objective function set of functions
 - ▶ Without theoretical support
- Empirical comparisons on extensive test suites
 - ▶ what performance measures?
 - ▶ what test functions? representative of real-world

Some proposals

- Expected Running Time + Empirical Cumulative Distributions
- An artificial testbed, with controlled typical difficulties
- A (partial) case study, involving 2 deterministic and 3 bio-inspired algorithms
- in back-box scenario without specific intensive parameter tuning

1 Performance Measures and Experimental Comparisons

- How to empirically compare algorithms?

2 Problem difficulties and algorithm invariances

3 Derivative-Free Optimization Algorithms

4 Experiments and Results

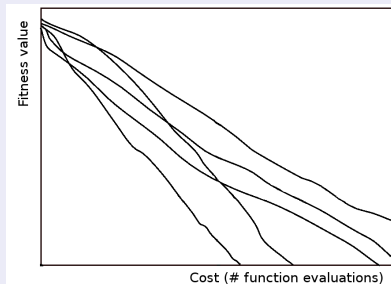
5 Conclusion

Context

Optimization goal(s)

- Find best possible objective value \neq get close to the global optimum
- At minimal cost \approx number of function evaluations

Available data



- Objective value vs Computational cost
- What to record/analyze/summarize?

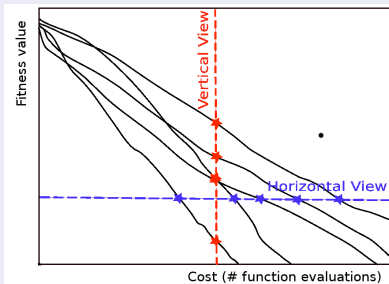
Minimization

Context

Optimization goal(s)

- Find best possible objective value \neq get close to the global optimum
- At minimal cost \approx number of function evaluations

Two points of view

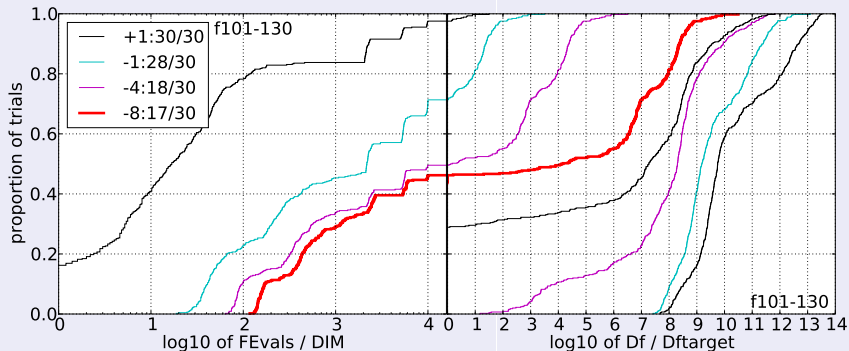


- **Vertical:** Value reached for a given effort
- **Horizontal:** Effort required to reach a given objective value

Both Views: Empirical Cumulative Distributions Fns

Horizontal

Vertical



- ... of normalized running times,
up to an upper bound ($10^4 n$)

- for different target values
 10^{-8} , 10^{-4} , 10^{-1} , 10^1

- ... of precision,
w.r.t. given target value (10^{-8})

- for different max. running times
 n , $10n$, $100n$, $1000n$

Discussion

Vertical vs horizontal

- **Vertical**: Value reached for a given effort
 - ▶ **Fixed budget** scenario
 - ▶ **Qualitative** comparisons *Algo. A reaches better value than Algo. B*
- **Horizontal**: Effort required to reach a given objective value
 - ▶ **Baseline** requirement *e.g. beat the opponent!*
 - ▶ **Absolute** comparisons: *Algo. A is X times faster than Algo. B*
 - ▶ Monotonous-invariant criterion

Statistics

- Difficult to summarize multiple viewpoints into a single measure
- ... and to find a sound estimator for it
compute its variance, perform statistical tests, ...

Performance measures

ECDFs

- Require arbitrary
 - ▶ maximal target precision,
 - ▶ maximal run length
- Can be used for sets of benchmark functions
- Need to be sub-sampled for comparisons

[previous slide](#)

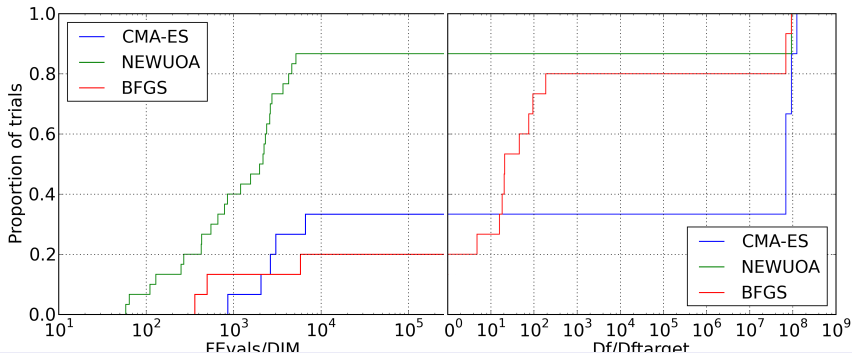
Horizontal performance measures

- Fix a *target objective value*,
- compute Expected Running Times Distribution,
- measure average effort to success

from Empirical Cumulative Distribution Functions

Horizontal

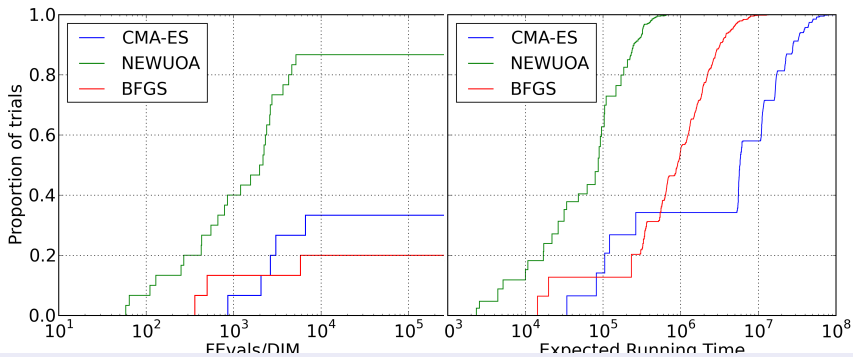
Vertical



from Empirical Cumulative Distribution Functions to Expected Running Time

Runtime ECDF

Estimated Runtime Distribution



using *resampling* technique

Expected Running Time

Experiments and notations

- Fixed number of runs
- Arbitrary target objective value f_{target}
- Arbitrary bound on # evaluations
- p_{succ} : proportion of successful runs that reached f_{target}
- \widehat{RT}_{succ} (resp. \widehat{RT}_{fail}): empirical average number of evaluations of successful (resp. unsuccessful) runs

Expected Running Time Measures

$$SP1(f_{target}) = \frac{\widehat{RT}_{succ}}{p_{succ}}$$

$$SP2(f_{target}) = \frac{p_{succ}\widehat{RT}_{succ} + (1 - p_{succ})\widehat{RT}_{fail}}{p_{succ}}$$

Expected Running Time (2)

Discussion

- Both measures
 - ▶ reflect some average effort to reach f_{target}
 - ▶ are equivalent in case of 100% success
 - ▶ are unreliable estimators in case of small p_{succ}
 - ▶ **can be used to easily compare algorithms** on sets of functions
by **normalizing w.r.t. best algorithm on each function**
- SP1 insensitive to the running length of unsuccessful runs
- SP2 very sensitive to the **stopping criterion** and the **restart strategy**, that are part of the algorithm fine tuning ...

History

- CEC'05 Challenge on Continuous Optimization used SP1
- GECCO'09 Workshop on Black-Box Optimization Benchmarking uses SP2

1 Performance Measures and Experimental Comparisons

2 Problem difficulties and algorithm invariances

- What makes a continuous optimization problem hard?

3 Derivative-Free Optimization Algorithms

4 Experiments and Results

5 Conclusion

Problem Difficulties and Algorithm Invariances

What makes a problem hard?

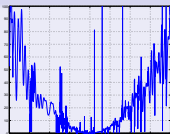
- Non-convexity

invalidates most of deterministic theory

- ▶ Ruggedness

- ▶ Multimodality

non-smooth, discontinuous, noisy
presence of local optima



- Dimensionality

line search is 'trivial'

The magnificence of high dimensionality ...

- Ill-conditioning

Very different scalings along different directions

- Non-separability

Correlated variables

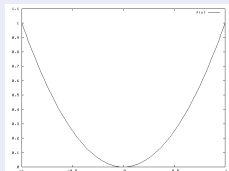
The benefits of invariance

- Some difficulties become harmless

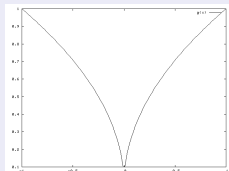
- More robust parameter setting

Ruggedness and Monotonous Invariance

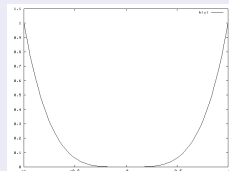
Monotonous transformations of the objective function



$$\sum x_i^2$$



$$\sqrt{\sqrt{\sum x_i^2}}$$



$$(\sum x_i^2)^2$$

Monotonous Invariance

- Invariance w.r.t. monotonous transformations
- A guarantee against ill-scaled objective functions
- Comparison-based algorithms are monotonous-invariant

Multimodality

Presence of multiple local optima

Restart strategies

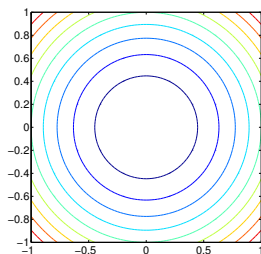
- For local optimizers, starting point is crucial on multimodal functions
- Multiple restarts are mandatory
 - ▶ from uniformly distributed points
 - ▶ from the final point of some previous run after some parameter reset
- Also efficient with any optimization algorithm

global restart
local restart

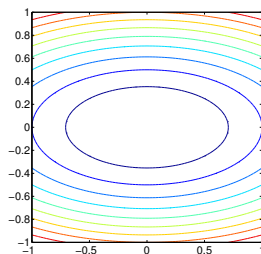
III-Conditionning

- The **Condition Number (CN)** of a positive-definite matrix \mathbf{H} is the ratio of its largest and smallest eigenvalues
- If f is **quadratic**, $f(\mathbf{x}) = \mathbf{x}^T \mathbf{H} \mathbf{x}$, the CN of f is that of its Hessian \mathbf{H}
- More generally, the CN of f is that of its Hessian wherever it is defined.

Graphically, ill-conditioned means “squeezed” lines of equal function value



Increased
→
condition
number



Issue: The gradient does not point toward the minimum ...

Separability

Definition (Separable Problem)

A function f is separable if

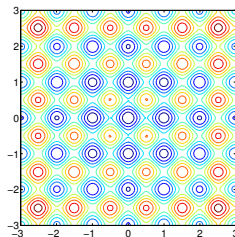
$$\arg \min_{(x_1, \dots, x_n)} f(x_1, \dots, x_n) = \left(\arg \min_{x_1} f(x_1, \dots), \dots, \arg \min_{x_n} f(\dots, x_n) \right)$$

solve n independent 1D optimization problems

Example: Additively decomposable functions

$$f(x_1, \dots, x_n) = \sum_{i=1}^n f_i(x_i)$$

e.g. Rastrigin function



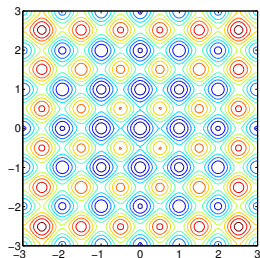
Designing Non-Separable Problems

Rotating the coordinate system

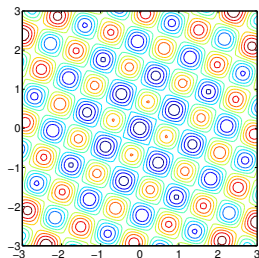
- $f : \mathbf{x} \mapsto f(\mathbf{x})$ separable
- $f : \mathbf{x} \mapsto f(\mathbf{R}\mathbf{x})$ non-separable

\mathbf{R} rotation matrix

Hansen, Ostermeier, & Gawelczyk, 95; Salomon, 96



\mathbf{R}
→



- 1 Performance Measures and Experimental Comparisons
- 2 Problem difficulties and algorithm invariances
- 3 Derivative-Free Optimization Algorithms**
 - Deterministic and Bio-Inspired Algorithms
- 4 Experiments and Results
- 5 Conclusion

Optimization techniques

Numerical methods

- Applied Mathematicians Long history
- Classical methods based on First Principles gradient-based
require regularity numerical gradient amenable to numerical pitfalls
- Recent DFO methods quadratic interpolation of the objective function

Bio-inspired algorithms

- Computer Scientists mostly from AI field
- Many recent (and trendy) methods

but

- Computationally heavy
- No convergence proof almost ...

BFGS

Broyden, Fletcher, Goldfarb, & Shanno, 1970

Gradient-based methods

$$\begin{cases} x_{t+1} = x_t - \rho_t d_t \\ \rho_t = \mathit{Argmin}_{\rho} \{ \mathcal{F}(x_t - \rho d_t) \} \end{cases} \text{ Line search}$$

Choice of d_t , the **descent direction**?

BFGS: a Quasi-Newton method

- Maintain an approximation \hat{H}_t of the Hessian of f
- Solve for d_t

$$\hat{H}_t d_t = -\nabla f(x_t)$$

- Compute x_{t+1} and update $\hat{H}_t \rightarrow \hat{H}_{t+1}$
- Converges if **quadratic approximation** of \mathcal{F} holds

around the optimum

- Reliable and robust

on quadratic functions!

BFGS: A priori Discussion

Properties

- Gradient-based algorithms are
 - ▶ **local** optimizers
 - ▶ **not** monotonous invariant
 - ▶ **independent** of the coordinate system
- But numerical gradient is not rotational invariant

and can suffer from ill-conditioning

Implementation

- Matlab built-in `fminunc`
- using numerical gradient
- Function improvement threshold set to 10^{-25}
- Multiple restarts mandatory

widely blindly used

local or global

NEWUOA

Powell, 2006

A Derivative-Free Optimization Algorithm

- Builds a **quadratic interpolation** of the objective function
- Maintains a **trust region** size ρ where interpolation is accurate
- Alternates
 - ▶ **trust region iterations**:
one conjugate gradient step on the surrogate model
 - ▶ **alternative iterations**:
new trust region and surrogate model

Parameters

- Number of interpolation points
from $n + 6$ to $\frac{n(n+1)}{2}$ $2n + 1$ recommended
- Initial (ρ_{init}) and final (ρ_{end}) sizes of trust region

NEWUOA: A priori Discussion

Properties

- a **global** optimizer (re)start with large trust region size
- **not** monotonous invariant
- Full model ($\frac{n(n+1)}{2}$ points) is independent of coordinate system
Most efficient model ($2n + 1$) is **not**

Implementation

- Matthieu Guibert's implementation
<http://www.inrialpes.fr/bipop/people/guilbert/newuoa/newuoa.html>
- $\rho_{init} = 100$ and $\rho_{end} = 10^{-15}$ after some preliminary experiments
- Local restarts by resetting ρ to ρ_{init}

Stochastic Search

A unified point of view

A black box search template to minimize $f : \mathbb{R}^n \rightarrow \mathbb{R}$

Initialize distribution parameters θ , set sample size $\lambda \in \mathbb{N}$

While not terminate

- 1 **Sample** distribution $P(\mathbf{x}|\theta) \rightarrow \mathbf{x}_1, \dots, \mathbf{x}_\lambda \in \mathbb{R}^n$
- 2 **Evaluate** $\mathbf{x}_1, \dots, \mathbf{x}_\lambda$ on f
- 3 **Update** parameters $\theta \leftarrow F_\theta(\theta, \mathbf{x}_1, \dots, \mathbf{x}_\lambda, f(\mathbf{x}_1), \dots, f(\mathbf{x}_\lambda))$

Covers

- Deterministic algorithms including BFGS and NEWUOA
- Evolutionary Algorithms, PSO, DE
 P implicitly defined by the variation operators (crossover/mutation)
- Estimation of Distribution Algorithms

Particle Swarm Optimization (PSO)

Eberhart & Kennedy, 1995

The basic algorithm

Let $\mathbf{x}_1, \dots, \mathbf{x}_\lambda \in \mathbb{R}^n$ be a set of particles

- 1 **Sample** new positions

$$\begin{aligned}
 x_i^j(t+1) = x_i^j(t) + w \underbrace{(x_i^j(t) - x_i^j(t-1))}_{\text{aka velocity}} \\
 + \underbrace{c_1 \mathcal{U}_i^j(0,1)(p_i^j - x_i^j(t))}_{\text{approach the "previous" best}} + \underbrace{c_2 \tilde{\mathcal{U}}_i^j(0,1)(g_i^j - x_i^j(t))}_{\text{approach the "global" best}}
 \end{aligned}$$

- 2 **Evaluate** $\mathbf{x}_1(t+1), \dots, \mathbf{x}_\lambda(t+1)$
- 3 **Update** distribution parameters

$$p_i = \mathbf{x}_i(t+1) \text{ if } f(\mathbf{x}_i(t+1)) < f(p_i)$$

$$g_i = \mathbf{x}_*(t+1) \text{ where } f(\mathbf{x}_*(t+1)) = \min \{f(\mathbf{x}_k(t+1)), k \in \text{neighbor}(i)\}$$

PSO: A priori Discussion

Properties

- Comparison-based Monotonous invariance
- **not** rotational invariant sampled distribution is an hyperrectangle

Implementation

- Standard PSO 2006, C code,
≠ Matlab code! from *PSO Central* at
<http://www.particleswarm.info/>
- Default settings:
 - ▶ inertia weight $w \approx 0.7$,
 - ▶ $c_1 = c_2 \approx 1.2$
- Swarm size = $\lambda = 10 + \text{floor}(2\sqrt{n})$

Differential Evolution (DE)

Rainer & Storn, 1995

Basic algorithm

- Generate NP individuals uniformly in bounds
- Until *stopping criterion*
 - ▶ For each individual x in the population
 - ★ Perturbation using an intra-population **difference** vector

$$\hat{x} = x^\alpha + F(x^\beta - x^\gamma)$$

- ★ (Uniform) crossover with probability $1 - CR$ Keep \hat{x} if $CR = 1$

$$y_i = \begin{cases} \hat{x}_i & \text{if } U(0, 1) < CR \\ x_i & \text{otherwise} \end{cases}$$

- ★ Keep best from x and y

DE: A priori Discussion

Properties

- Comparison-based Monotonous invariance
- **Crossover** is not rotational invariant exchange coordinates

Implementation

- C code from DE home page
- **No** default setting!
- Extensive DOE: On 10-D ellipsoid function
 - ▶ Strategy 2
 - ▶ $F = 0.8$,
 - ▶ $CR = 1.0$ Rotational invariance
- Population size: $NP = 10 * n$

The (μ, λ) –CMA–Evolution Strategy

Minimize $f : \mathbb{R}^n \rightarrow \mathbb{R}$

Initialize distribution parameters θ , set population size $\lambda \in \mathbb{N}$

While not terminate

- 1 **Sample** distribution $P(\mathbf{x}|\theta) \rightarrow \mathbf{x}_1, \dots, \mathbf{x}_\lambda \in \mathbb{R}^n$
- 2 **Evaluate** offspring $\mathbf{x}_1, \dots, \mathbf{x}_\lambda$ on f
- 3 **Update** parameters $\theta \leftarrow F_\theta(\theta, \mathbf{x}_1, \dots, \mathbf{x}_\lambda, f(\mathbf{x}_1), \dots, f(\mathbf{x}_\lambda))$

- P is a **multi-variate normal** distribution

$$\mathcal{N}(\mathbf{m}, \sigma^2 \mathbf{C}) \sim \boxed{\mathbf{m} + \sigma \mathcal{N}(\mathbf{0}, \mathbf{C})}$$

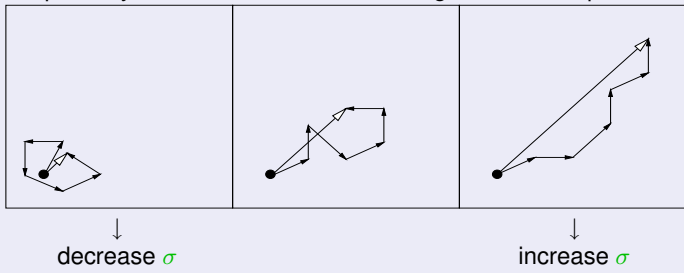
- $\theta = \{\mathbf{m}, \mathbf{C}, \sigma\} \in (\mathbb{R}^n \times \mathbb{R}^{n \times n} \times \mathbb{R}_+)$
- $F_\theta(\theta, \mathbf{x}_{1:\lambda}, \dots, \mathbf{x}_{\mu:\lambda})$ only depends on the $\mu \leq \lambda$ best offspring

Cumulative Step-Size Adaptation (CSA)

$$\begin{aligned} \mathbf{x}_i &= \mathbf{m} + \sigma \mathbf{z}_i \\ \mathbf{m} &\leftarrow \mathbf{m} + \sigma \langle \mathbf{z} \rangle_{\text{sel}} \end{aligned}$$

Measure the length of the *evolution path*

the pathway of the mean vector \mathbf{m} in the generation sequence



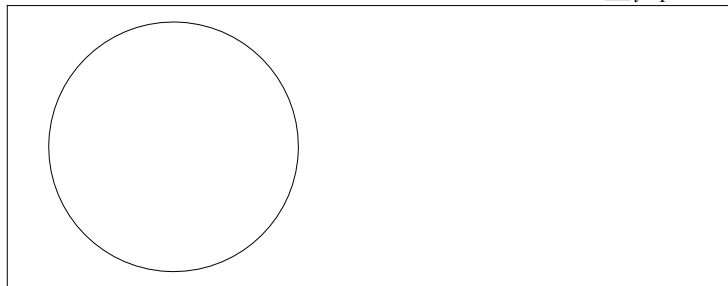
loosely speaking steps are

- perpendicular under random selection (in expectation)
- perpendicular in the desired situation (to be most efficient)

Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \langle \mathbf{z} \rangle_{\text{sel}}, \langle \mathbf{z} \rangle_{\text{sel}} = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda}, \mathbf{z}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



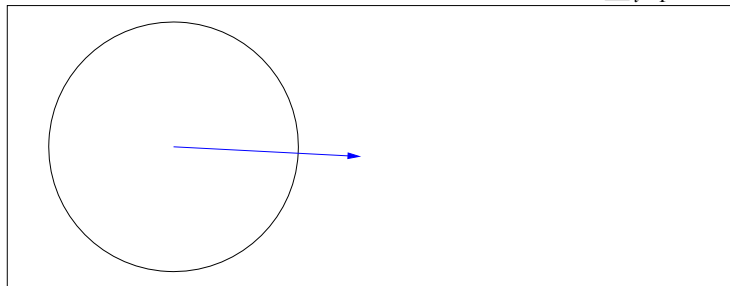
initial distribution, $\mathbf{C} = \mathbf{I}$

- new distribution: $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \langle \mathbf{z} \rangle_{\text{sel}} \langle \mathbf{z} \rangle_{\text{sel}}^T$
- ruling principle: the adaptation increases the probability of successful steps, $\langle \mathbf{z} \rangle_{\text{sel}}$, to appear again

Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \langle \mathbf{z} \rangle_{\text{sel}}, \langle \mathbf{z} \rangle_{\text{sel}} = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda}, \mathbf{z}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



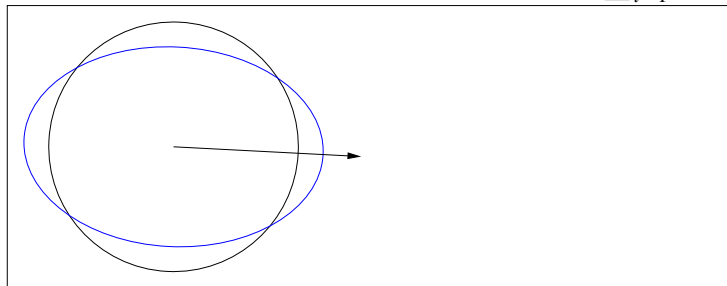
$\langle \mathbf{z} \rangle_{\text{sel}}$, movement of the population mean \mathbf{m} (disregarding σ)

- new distribution: $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \langle \mathbf{z} \rangle_{\text{sel}} \langle \mathbf{z} \rangle_{\text{sel}}^T$
- ruling principle: the adaptation increases the probability of successful steps, $\langle \mathbf{z} \rangle_{\text{sel}}$, to appear again

Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \langle \mathbf{z} \rangle_{\text{sel}}, \langle \mathbf{z} \rangle_{\text{sel}} = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda}, \mathbf{z}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



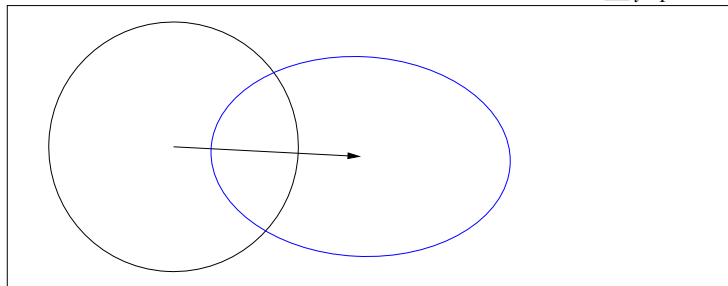
mixture of \mathbf{C} and step $\langle \mathbf{z} \rangle_{\text{sel}}$, $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \langle \mathbf{z} \rangle_{\text{sel}} \langle \mathbf{z} \rangle_{\text{sel}}^T$

- new distribution: $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \langle \mathbf{z} \rangle_{\text{sel}} \langle \mathbf{z} \rangle_{\text{sel}}^T$
- ruling principle: the adaptation increases the probability of successful steps, $\langle \mathbf{z} \rangle_{\text{sel}}$, to appear again

Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \langle \mathbf{z} \rangle_{\text{sel}}, \langle \mathbf{z} \rangle_{\text{sel}} = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda}, \mathbf{z}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



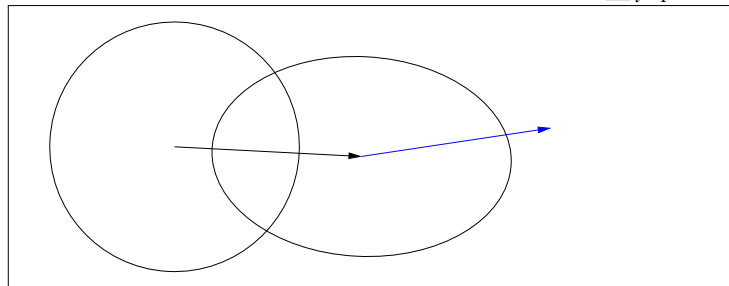
new distribution (disregarding σ)

- new distribution: $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \langle \mathbf{z} \rangle_{\text{sel}} \langle \mathbf{z} \rangle_{\text{sel}}^T$
- ruling principle: the adaptation increases the probability of successful steps, $\langle \mathbf{z} \rangle_{\text{sel}}$, to appear again

Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \langle \mathbf{z} \rangle_{\text{sel}}, \langle \mathbf{z} \rangle_{\text{sel}} = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda}, \mathbf{z}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



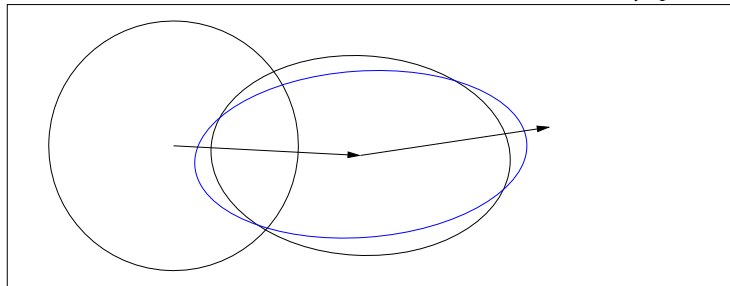
movement of the population mean \mathbf{m}

- new distribution: $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \langle \mathbf{z} \rangle_{\text{sel}} \langle \mathbf{z} \rangle_{\text{sel}}^T$
- ruling principle: the adaptation increases the probability of successful steps, $\langle \mathbf{z} \rangle_{\text{sel}}$, to appear again

Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \langle \mathbf{z} \rangle_{\text{sel}}, \langle \mathbf{z} \rangle_{\text{sel}} = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda}, \mathbf{z}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



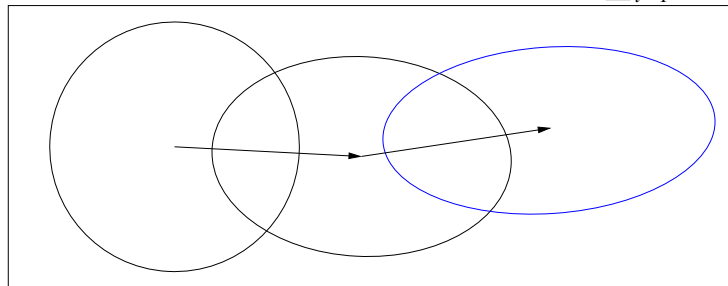
mixture of \mathbf{C} and step $\langle \mathbf{z} \rangle_{\text{sel}}$, $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \langle \mathbf{z} \rangle_{\text{sel}} \langle \mathbf{z} \rangle_{\text{sel}}^T$

- new distribution: $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \langle \mathbf{z} \rangle_{\text{sel}} \langle \mathbf{z} \rangle_{\text{sel}}^T$
- ruling principle: the adaptation increases the probability of successful steps, $\langle \mathbf{z} \rangle_{\text{sel}}$, to appear again

Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \langle \mathbf{z} \rangle_{\text{sel}}, \langle \mathbf{z} \rangle_{\text{sel}} = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda}, \mathbf{z}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



- new distribution: $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \langle \mathbf{z} \rangle_{\text{sel}} \langle \mathbf{z} \rangle_{\text{sel}}^T$
- ruling principle: the adaptation increases the probability of successful steps, $\langle \mathbf{z} \rangle_{\text{sel}}$, to appear again

CMA-ES: A priori Discussion

Properties

- Comparison-based Monotonous invariance
- Rotational invariant Adapts the coordinate system

Implementation

- **CMA-ES**: Matlab code from author's home page
<http://www.bionik.tu-berlin.de/user/niko/>
- Using rank- μ update faster adaptation Hansen et al., 2003
- Default settings
- Population size: $\lambda = 4 + \text{floor}(3 \ln n)$

- 1 Performance Measures and Experimental Comparisons
- 2 Problem difficulties and algorithm invariances
- 3 Derivative-Free Optimization Algorithms
- 4 Experiments and Results**
 - Empirical Comparisons
- 5 Conclusion

The parameters

Common Parameters

- Default parameters
- Upper bound on run-time = 10^7 evaluations
- $f_{target} = 10^{-9}$
- Domain: $[-20, 80]^d$
- 21 runs in each case

except for DE

Optimum not at the center
except BFGS with little success

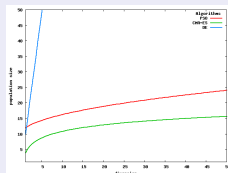
Population size

Standard values:

for $n = 10, 20, 40$

- PSO: $10 + \text{floor}(2\sqrt{n})$ 16, 18, 22
- CMA-ES: $4 + \text{floor}(3 \ln n)$ 10, 12, 15
- DE: $10 * n$ 100, 200, 400

To be increased for multi-modal functions



e.g. Rastrigin

Ellipsoid

- $f_{\text{elli}}(x) = \sum_{i=1}^n 10^{\alpha \frac{i-1}{n-1}} x_i^2 = x^T H_{\text{elli}} x$

$$H_{\text{elli}} = \begin{pmatrix} 1 & 0 & \dots \\ & \ddots & \\ \dots & 0 & 10^\alpha \end{pmatrix}$$

convex, separable

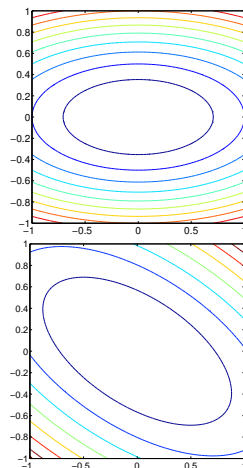
- $f_{\text{elli}}^{\text{rot}}(x) = f_{\text{elli}}(\mathbf{R}x) = x^T H_{\text{elli}}^{\text{rot}} x$

\mathbf{R} random rotation

$$H_{\text{elli}}^{\text{rot}} = \mathbf{R}^T H_{\text{elli}} \mathbf{R}$$

convex, non-separable

- $\text{cond}(H_{\text{elli}}) = \text{cond}(H_{\text{elli}}^{\text{rot}}) = 10^\alpha$



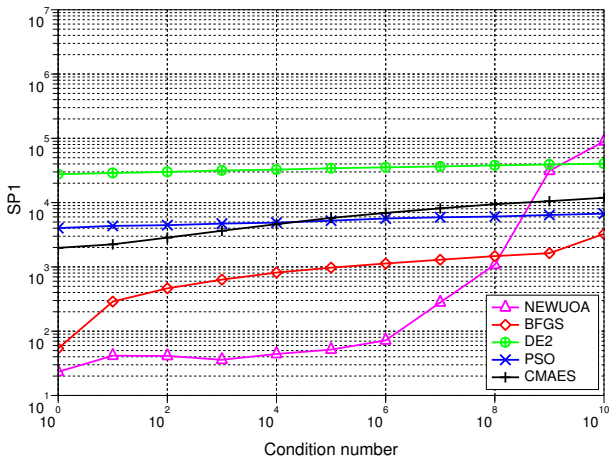
$\alpha = 1, \dots, 10$

$\alpha = 6 \equiv$ axis ratio of 10^3 , typical for real-world problem

Separable Ellipsoid function

PSO, DE2, CMA-ES, NEWUOA, and BFGS - Dimension 10

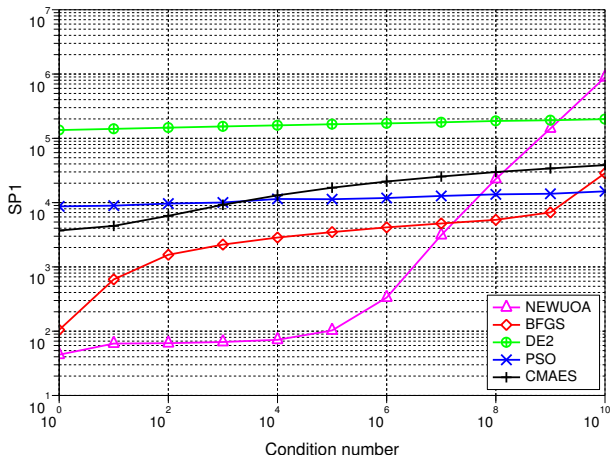
Ellipsoid dimension 10, 21 trials, tolerance $1e-09$, eval max $1e+07$



Separable Ellipsoid function

PSO, DE2, CMA-ES, NEWUOA, and BFGS - Dimension 20

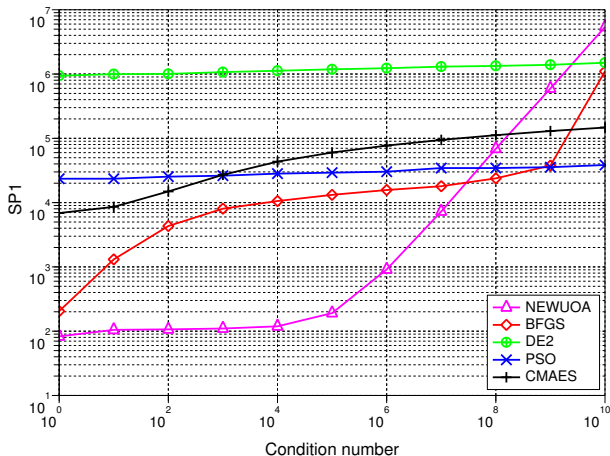
Ellipsoid dimension 20, 21 trials, tolerance $1e-09$, eval max $1e+07$



Separable Ellipsoid function

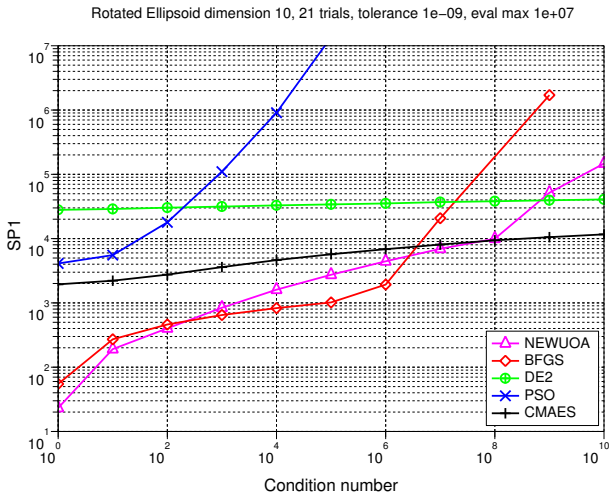
PSO, DE2, CMA-ES, NEWUOA, and BFGS - Dimension 40

Ellipsoid dimension 40, 21 trials, tolerance $1e-09$, eval max $1e+07$



Rotated Ellipsoid function

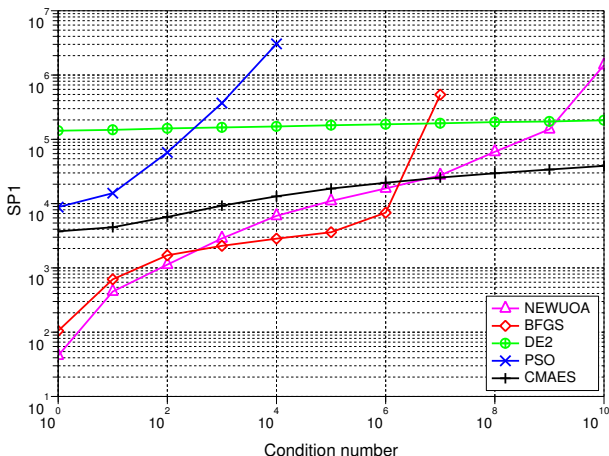
PSO, DE2, CMA-ES, NEWUOA, and BFGS - Dimension 10



Rotated Ellipsoid function

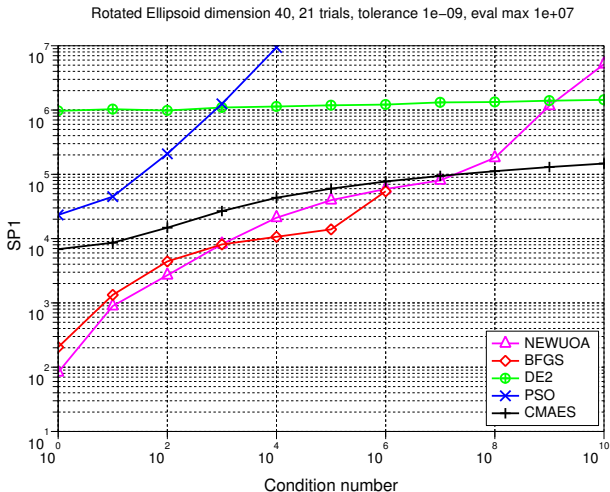
PSO, DE2, CMA-ES, NEWUOA, and BFGS - Dimension 20

Rotated Ellipsoid dimension 20, 21 trials, tolerance $1e-09$, eval max $1e+07$



Rotated Ellipsoid function

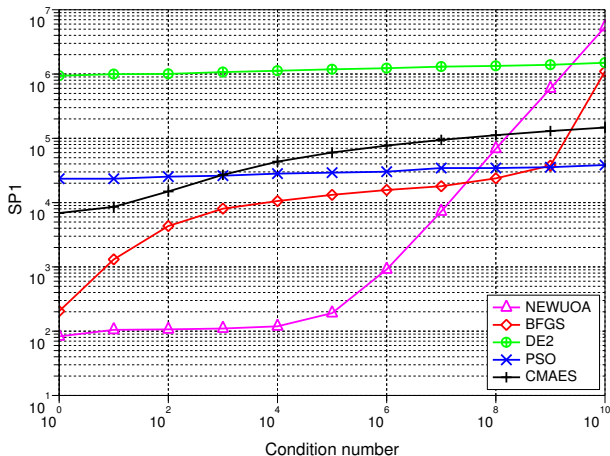
PSO, DE2, CMA-ES, NEWUOA, and BFGS - Dimension 40



Separable Ellipsoid function

PSO, DE2, CMA-ES, NEWUOA, and BFGS - Dimension 40

Ellipsoid dimension 40, 21 trials, tolerance $1e-09$, eval max $1e+07$



Ellipsoid: Discussion

Bio-inspired algorithms

- Separable case: PSO and DE insensitive to conditioning
- ... but PSO rapidly fails to solve the rotated version
- ... while CMA-ES and DE ($CR = 1$) are rotation invariant
- DE scales poorly with dimension
 $d^{2.5}$ compared to $d^{1.5}$ for PSO and CMA-ES

and BFGS

... vs deterministic

- BFGS fails to solve ill-conditioned cases Matlab "Roundoff error"
- CMA-ES only 7 times slower on pure quadratic functions!
- NEWUOA 70 times better than CMA-ES on separable cases
- performs worse for very high conditioning and rotated cases

Away from “quadraticity”

Monotonous invariance

- Comparison-based algorithms are insensitive to monotonous transformations
True for DE, PSO and all ESs
- BFGS and NEWUOA are not
theory behind BFGS depends on
convexity

Another test function

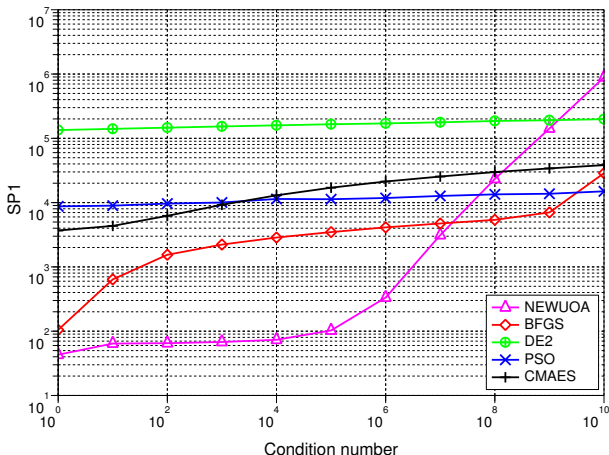
- Simple transformation of ellipsoid

$$f_{\text{SSE}}(x) = \sqrt{\sqrt{f_{\text{elli}}(x)}}$$

Separable Ellipsoid function

PSO, DE2, CMA-ES, NEWUOA, and BFGS - Dimension 20

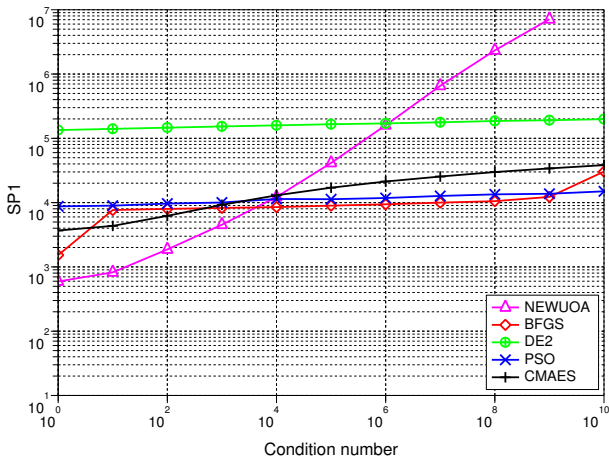
Ellipsoid dimension 20, 21 trials, tolerance $1e-09$, eval max $1e+07$



Separable Ellipsoid $^{\frac{1}{4}}$ function

PSO, DE2, CMA-ES, NEWUOA, and BFGS - Dimension 20

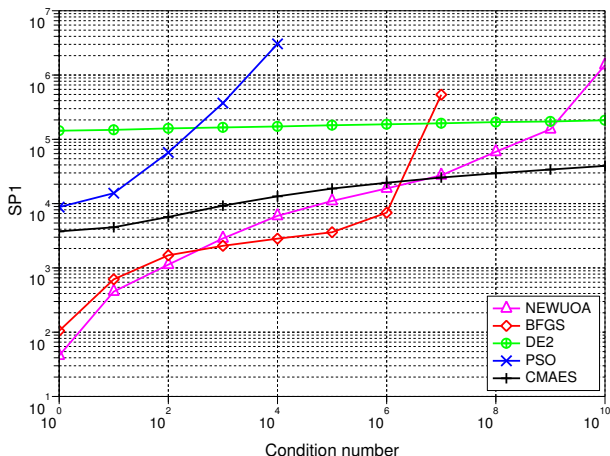
Sqrt of sqrt of ellipsoid dimension 20, 21 trials, tolerance 1e-09, eval max 1e+07



Rotated Ellipsoid function

PSO, DE2, CMA-ES, NEWUOA, and BFGS - Dimension 20

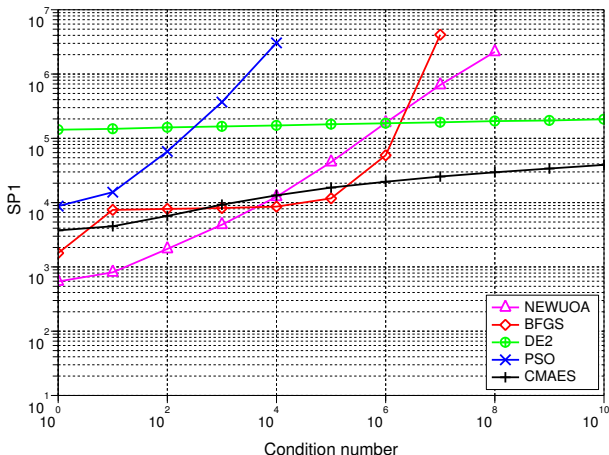
Rotated Ellipsoid dimension 20, 21 trials, tolerance $1e-09$, eval max $1e+07$



Rotated Ellipsoid $^{\frac{1}{4}}$ function

PSO, DE2, CMA-ES, NEWUOA, and BFGS - Dimension 20

Sqrt of sqrt of rotated ellipsoid dimension 20, 21 trials, tolerance $1e-09$, eval max $1e+07$



Non-quadratic: Discussion

Bio-inspired algorithms

- Invariant

as expected!

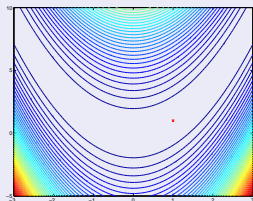
NEWUOA and BFGS

- Worse on $\sqrt{\sqrt{\text{Ellipsoid}}}$ than on Ellipsoid
- Premature numerical convergence for high CN for BFGS ... fixed by the 'local restart' strategy
- NEWUOA suffers from conjunction of rotation, high CN and non-quadraticity

Rosenbrock function (Banana)

$$f_{\text{rosen}}(x) = \sum_{i=1}^{n-1} [(1 - x_i)^2 + \beta(x_{i+1} - x_i^2)^2]$$

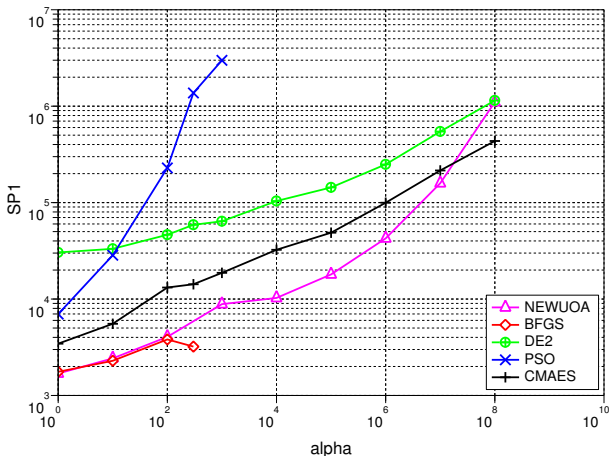
- Non-separable, but ...
also ran rotated version
- $\beta = 100$, classical Rosenbrock function
 $\beta = 1, \dots, 10^8$
- Multi-modal for dimension > 3



Rosenbrock functions

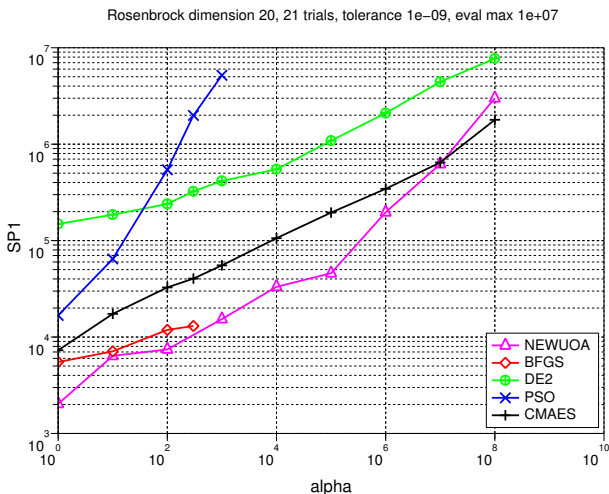
PSO, DE2, CMA-ES, NEWUOA, and BFGS - Dimension 10

Rosenbrock dimension 10, 21 trials, tolerance $1e-09$, eval max $1e+07$



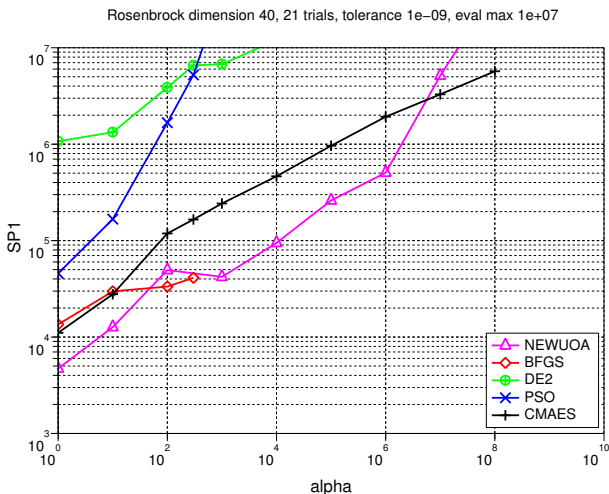
Rosenbrock functions

PSO, DE2, CMA-ES, NEWUOA, and BFGS - Dimension 20



Rosenbrock functions

PSO, DE2, CMA-ES, NEWUOA, and BFGS - Dimension 40



Rosenbrock: Discussion

Bio-inspired algorithms

- PSO sensitive to non-separability
- DE still scales badly with dimension

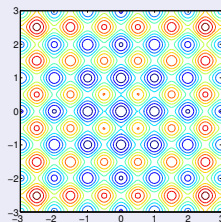
... vs BFGS

- Numerical premature convergence on ill-condition problems
- Both local and global restarts improve the results

Rastrigin function

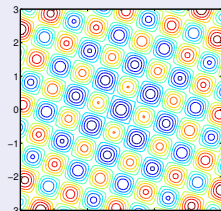
$$f_{\text{rast}}(x) = 10n + \sum_{i=1}^n x_i^2 - 10 \cos(2\pi x_i)$$

- separable
- multi-modal



$$f_{\text{rast}}^{\text{rot}}(x) = f_{\text{rast}}(\mathbf{R}x)$$

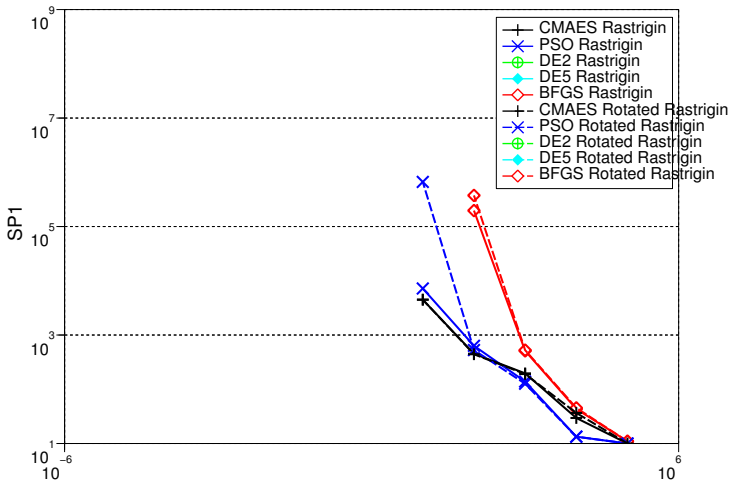
- \mathbf{R} random rotation
- non-separable
- multimodal



Rastrigin function - SP1 vs objective value

PSO, DE2, DE5, CMA-ES, and BFGS - PopSize 10

Rastrigin: dimension 10, NP = 10



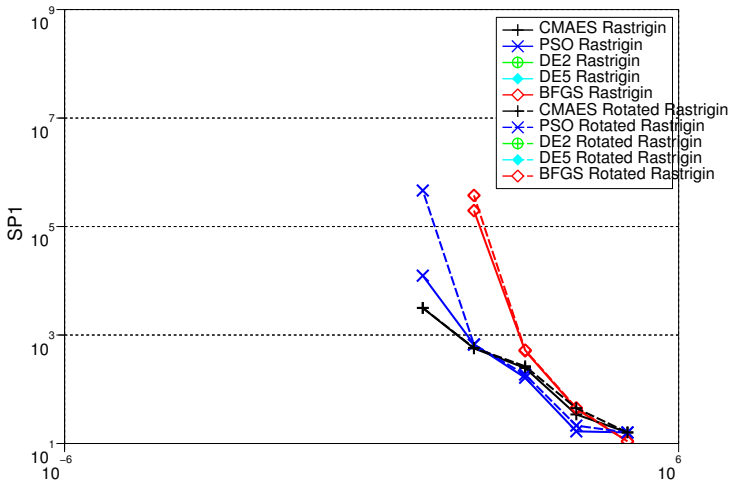
Function Value To Reach



Rastrigin function - SP1 vs objective value

PSO, DE2, DE5, CMA-ES, and BFGS - PopSize 16

Rastrigin: dimension 10, NP = 16

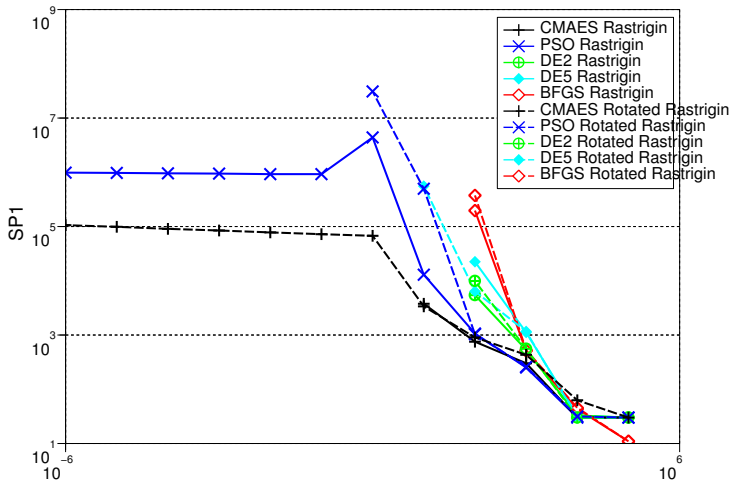


Function Value To Reach

Rastrigin function - SP1 vs objective value

PSO, DE2, DE5, CMA-ES, and BFGS - PopSize 30

Rastrigin: dimension 10, NP = 30



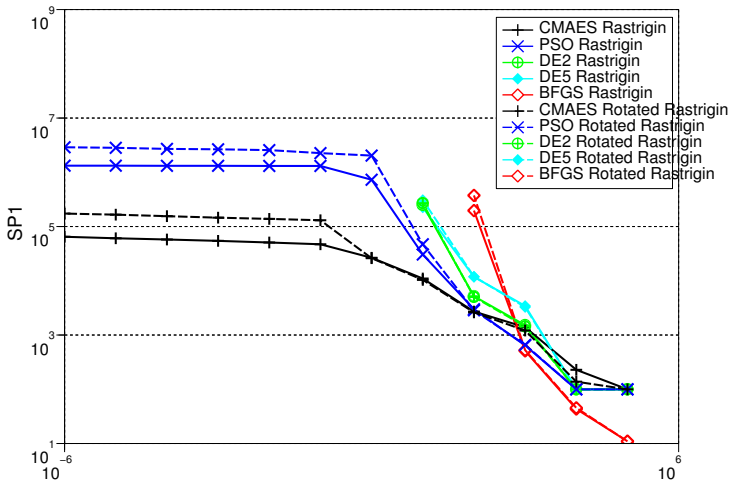
Function Value To Reach



Rastrigin function - SP1 vs objective value

PSO, DE2, DE5, CMA-ES, and BFGS - PopSize 100

Rastrigin: dimension 10, NP = 100

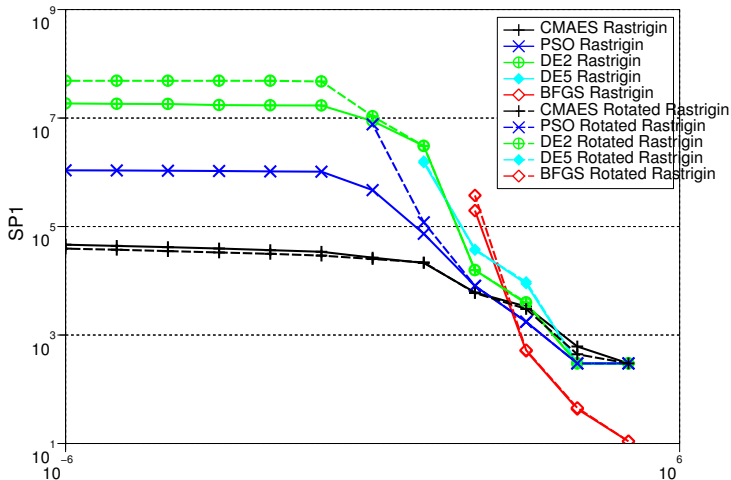


Function Value To Reach

Rastrigin function - SP1 vs objective value

PSO, DE2, DE5, CMA-ES, and BFGS - PopSize 300

Rastrigin: dimension 10, NP = 300

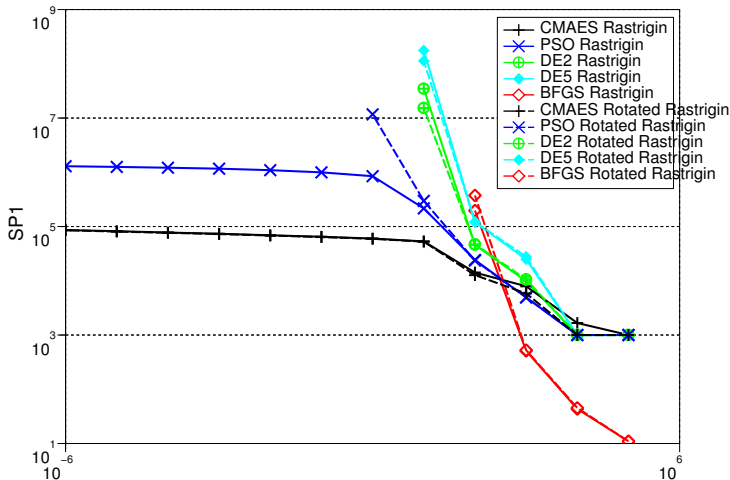


Function Value To Reach

Rastrigin function - SP1 vs objective value

PSO, DE2, DE5, CMA-ES, and BFGS - PopSize 1000

Rastrigin: dimension 10, NP = 1000



Function Value To Reach



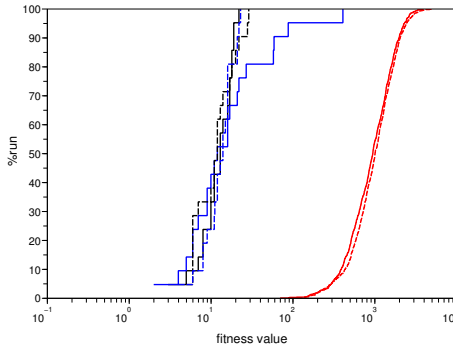
Rastrigin function - Cumulative distributions

PSO, **DE2**, **DE5**, **CMA-ES**, and **BFGS** - PopSize 10

astrigin : 21 trials, dimension 10, tol 1.000E-09, alpha 10, default size , eval max 10000000

Rastrigin : 21 trials, dimension 10, tol 1.000E-09, alpha 10, default size , eval max 10000000

evaluation number



% success
vs
eval to reach success threshold (= 10^{-9})

% success
vs
objective value reached before max eval (= 10^7)

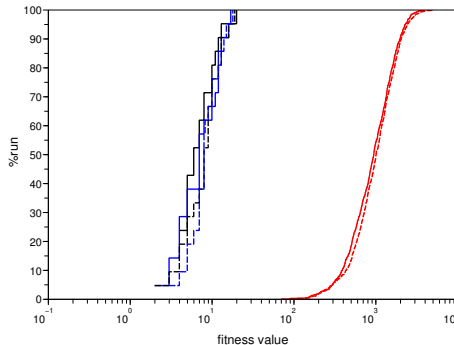
Rastrigin function - Cumulative distributions

PSO, **DE2**, **DE5**, **CMA-ES**, and **BFGS** - PopSize 16

astrigin : 21 trials, dimension 10, tol 1.000E-09, alpha 16, default size , eval max 10000000

Rastrigin : 21 trials, dimension 10, tol 1.000E-09, alpha 16, default size , eval max 10000000

evaluation number



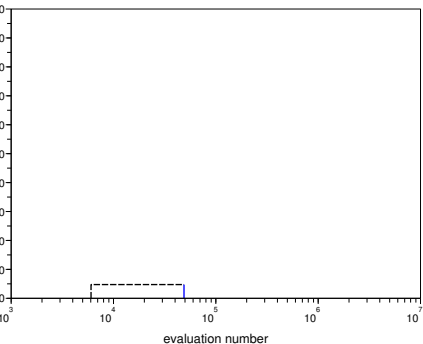
% success
vs
eval to reach success threshold (= 10^{-9})

% success
vs
objective value reached before max eval (= 10^7)

Rastrigin function - Cumulative distributions

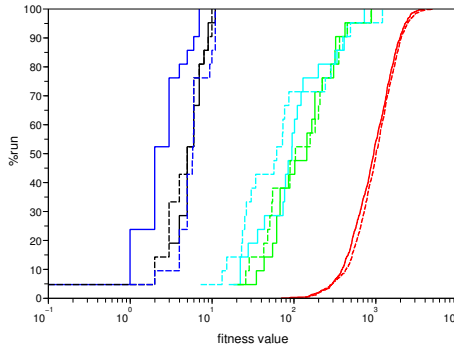
PSO, **DE2**, **DE5**, **CMA-ES**, and **BFGS** - PopSize 30

Rastrigin : 21 trials, dimension 10, tol 1.000E-09, alpha 30, default size , eval max 10000000



% success
vs
eval to reach success threshold (= 10^{-9})

Rastrigin : 21 trials, dimension 10, tol 1.000E-09, alpha 30, default size , eval max 10000000

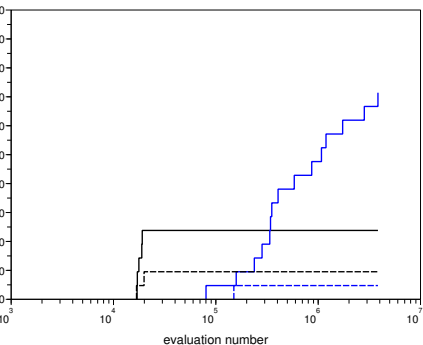


% success
vs
objective value reached before max eval (= 10^7)

Rastrigin function - Cumulative distributions

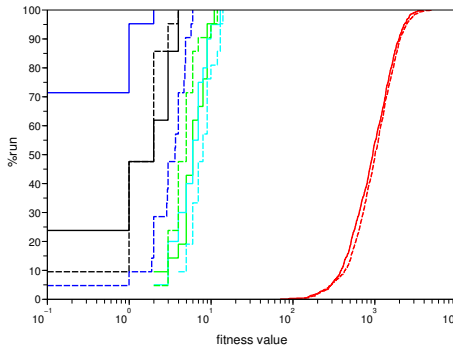
PSO, **DE2**, **DE5**, **CMA-ES**, and **BFGS** - PopSize 100

Rastrigin : 21 trials, dimension 10, tol 1.000E-09, alpha 100, default size , eval max 10000000



% success
vs
eval to reach success threshold (= 10^{-9})

Rastrigin : 21 trials, dimension 10, tol 1.000E-09, alpha 100, default size , eval max 10000000

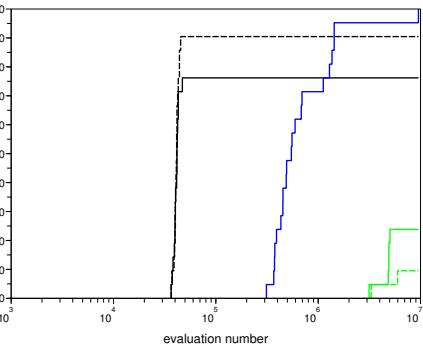


% success
vs
objective value reached before max eval (= 10^7)

Rastrigin function - Cumulative distributions

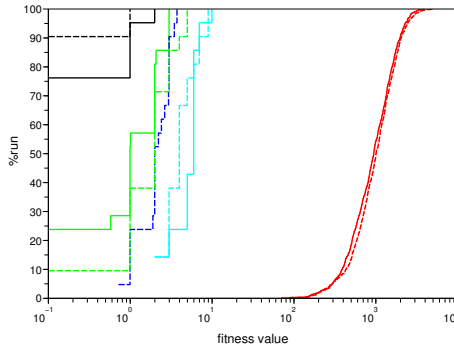
PSO, **DE2**, **DE5**, **CMA-ES**, and **BFGS** - PopSize 300

Rastrigin : 21 trials, dimension 10, tol 1.000E-09, alpha 300, default size , eval max 10000000



% success
vs
eval to reach success threshold (= 10^{-9})

Rastrigin : 21 trials, dimension 10, tol 1.000E-09, alpha 300, default size , eval max 10000000

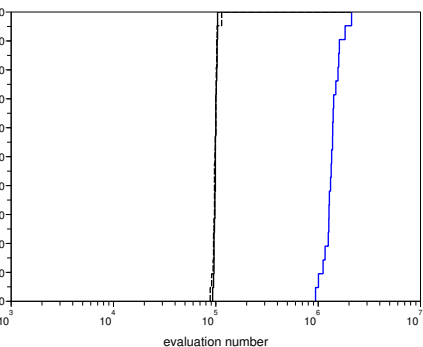


% success
vs
objective value reached before max eval (= 10^7)

Rastrigin function - Cumulative distributions

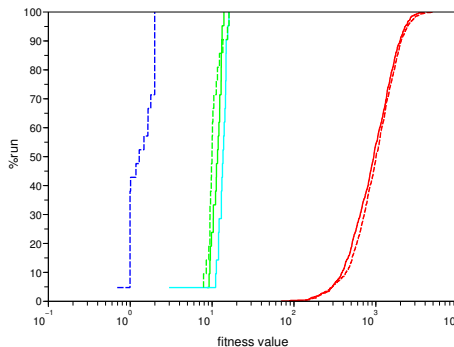
PSO, **DE2**, **DE5**, **CMA-ES**, and **BFGS** - PopSize 1000

Rastrigin : 21 trials, dimension 10, tol 1.000E-09, alpha 1000, default size , eval max 10000000



% success
vs
eval to reach success threshold (= 10^{-9})

Rastrigin : 21 trials, dimension 10, tol 1.000E-09, alpha 1000, default size , eval max 10000000



% success
vs
objective value reached before max eval (= 10^7)

Rastrigin: Discussion

Bio-inspired algorithms

- Increasing population size improves the results
Optimal size is algorithm-dependent
- CMA-ES and PSO solve separable case
PSO 100 times slower
- Only CMA-ES solves the rotated Rastrigin **reliably**
requires popSize ≥ 300

... vs BFGS

- Gets stuck in local optima
- Whatever the restart strategies
No numerical premature convergence

... and NEWUOA?

- solves it in 5 iterations!
identifies the global parabola

Some functions from GECCO'09 BBOB Workshop

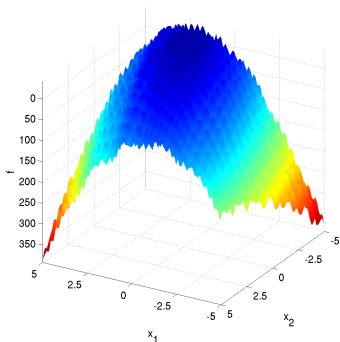
Properties

- No tunable difficulty e.g. single condition number
- but systematic non-linear transformations
- and asymetrisation of global and local structures

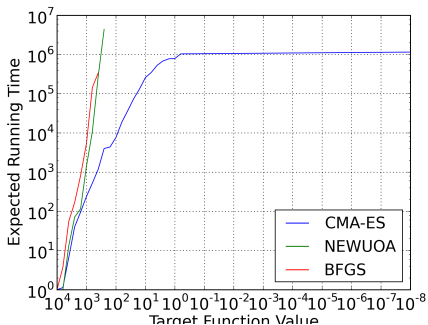
Issue

- Difficult to interpret or generalize
- Except on some families of functions e.g. with high conditioning

Asymmetric Rastrigin



2D plot

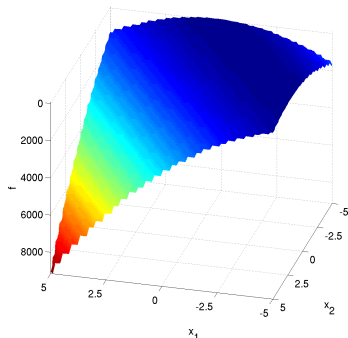


SP2 for $f_{target} = 10^{-8}$

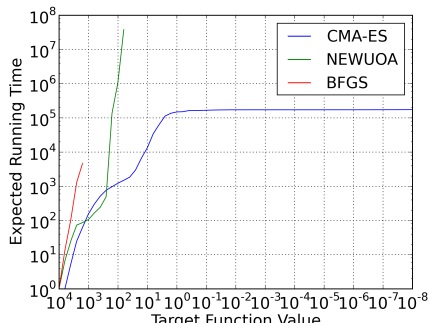
CMA-ES, NEWUOA, and BFGS

Step ellipsoid

Piecewise constant with global ellipsoid structure



2D plot

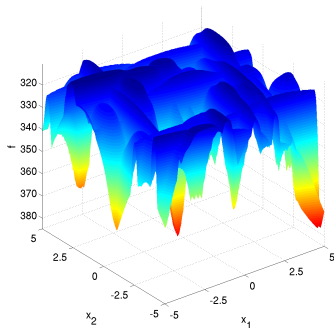


SP2 for $f_{target} = 10^{-8}$

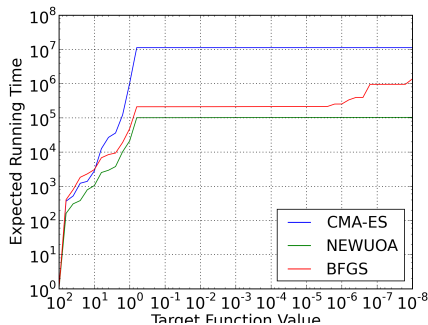
CMA-ES, NEWUOA, and BFGS

Gallagher's Gaussian Peaks

Very weak global structure



2D plot



SP2 for $f_{\text{target}} = 10^{-8}$

CMA-ES, NEWUOA, and BFGS

- 1 Performance Measures and Experimental Comparisons
- 2 Problem difficulties and algorithm invariances
- 3 Derivative-Free Optimization Algorithms
- 4 Experiments and Results
- 5 Conclusion**
 - and perspectives

Empirical Comparisons of DFO algorithms

An ill-defined problem

- Trade-off between precision and speed
- Task-dependent

Performance Measures

- **Empirical Cumulative Distributions** display all information
- Otherwise, need to chose one view-point e.g. horizontal
- → **Expected Running Times** for easy comparisons

Identified problem difficulties

- allow us to define test-suite for specific comparison purposes
- highlight the benefits of algorithm invariance properties

Bio-inspired Algorithms: The coming of age

Bio-inspired vs Deterministic

- CMA-ES only 7 (70) times slower than BFGS (DFO)

on (quasi-)quadratic functions.

but

- is less hindered by high conditioning,
- is monotonous-transformation invariant,
- is a global search method!

Moreover,

- Theoretical results are catching up
 - ▶ Linear convergence for SA-ES with bound on the CV speed
 - ▶ On-going work for CMA-ES

Auger, 05

Perspectives

Empirical comparisons

- GECCO'09 BBOB Workshop and Challenge
 - ▶ Noise, and/or **non-linear asymmetrizing** transformations
 - ▶ Compare wide set of algorithms **optimally tuned**
- Longer term
 - ▶ Constrained functions
 - ▶ Real-world functions **but which ones ???**

COCO, an **open** platform for **C**OMparison of **C**ontinuous **O**ptimizers

Toward automatic algorithm choice

- Need **descriptors** of problem difficulty **easy to compute**
- informative enough to guide algorithmic choice

Inspiration from SAT domain (**Hoost et al, 06**) and Statistical Physics?