PSO facing non-separable and ill-conditioned problems

Anne Auger Nikolaus Hansen Nikolas Mauny Marc Schoenauer

TAO Team, INRIA Futurs, FRANCE
 http://tao.lri.fr
 {Anne.Auger@inria.fr}

Séminaire OEP, Avril 2007

Continuous optimization and stochastic search

- Non-separable problems
- Ill-Conditioned Problems
- PSO and CMA-ES

PSO facing ill-conditioned and non-separable problems

- Test problems
- Results



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

$$f: \mathcal{S} \subseteq \mathbb{R}^n \to \mathbb{R}, \qquad \mathbf{x} \mapsto f(\mathbf{x})$$

• Black Box scenario (direct search)

$$x \longrightarrow f(x)$$

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

$$f: \mathcal{S} \subseteq \mathbb{R}^n \to \mathbb{R}, \qquad \mathbf{x} \mapsto f(\mathbf{x})$$

• Black Box scenario (direct search)

$$x \longrightarrow f(x)$$

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

$$f: \mathcal{S} \subseteq \mathbb{R}^n \to \mathbb{R}, \qquad \mathbf{x} \mapsto f(\mathbf{x})$$

• Black Box scenario (direct search)



- gradients are not available or not useful
- problem domain specific knowledge is used only within the black box
- search costs: number of function evaluations
- goal (1): find a solution *x* with a function value as small as possible with as small as possible search costs
- goal (2): convergence to an optimum
- goal (3): convergence to the global optimum

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

$$f: \mathcal{S} \subseteq \mathbb{R}^n \to \mathbb{R}, \qquad \mathbf{x} \mapsto f(\mathbf{x})$$

• Black Box scenario (direct search)



- Typical Examples:
 - shape optimization (e.g. using CFD)
 - model calibration
 - parameter calibration

curve fitting, airfoils

biological, physical

controller, plants, images

What makes a problem hard? Fitness function properties - Why stochastic search?

ruggedness

non-smooth, discontinuous, multimodal, and/or noisy function

 dimensionality (considerably) larger than three



cut from 3-D example solvable with an ES

- non-separability
- ill-conditioning

What makes a problem hard? Fitness function properties - Why stochastic search?

ruggedness

non-smooth, discontinuous, multimodal, and/or noisy function

 dimensionality (considerably) larger than three



cut from 3-D example solvable with an ES

- non-separability
- ill-conditioning

What makes a problem hard? Fitness function properties - Why stochastic search?



Non-separable problems

Separable Problems

Definition (Separable Problem)

A function f is separable if

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

 \Rightarrow it follows that f can be optimized in a sequence of n independent 1-D optimization processes

Non-separable problems

Separable Problems

Definition (Separable Problem)

A function f is separable if

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

 \Rightarrow it follows that f can be optimized in a sequence of n independent 1-D optimization processes

Example: Additively decomposable functions

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

Rastrigin function

-33	-2	-1	0	1	2	3
	0		0	0	0	0
-2	0			0		
- 1	0	0	0	0	0	O
_1			20	0	20	e.
0						
	0	0	\odot	\odot	0	0
1						
-	0	0	0	0	0	0
			0	0	O	
Ŭ						

Summary and conclusion

Non-separable problems

Non-Separable Problems Building a non-separable problem from a separable one



¹Hansen, Ostermeier, Gawelczyk (1995). On the adaptation of arbitrary normal mutation distributions in evolution strategies: The generating set adaptation. Sixth ICGA, pp. 57-64, Morgan Kaufmann

²Salomon (1996). "Reevaluating Genetic Algorithm Performance under Coordinate Rotation of Benchmark Functions; A survey of some theoretical and practical aspects of genetic algorithms." BioSystems, 39(3):263-278 **III-Conditioned Problems**

III-Conditioned Problems

If *f* is quadratic, $f : x \mapsto x^T H x$, ill-conditioned means a high condition number of Hessian Matrix H

III-Conditioned Problems

III-Conditioned Problems

If *f* is quadratic, $f : x \mapsto x^T H x$, ill-conditioned means a high condition number of Hessian Matrix H

ill-conditioned means "squeezed" lines of equal function value



consider the curvature of iso-fitness lines

Black Box Stochastic Optimization Algorithms

Template to minimize $f : \mathbb{R}^n \to \mathbb{R}$

Initialize distribution parameters $\theta,$ set population size $\lambda \in \mathbb{N}$ While not terminate

- **1** Sample distribution $P(\mathbf{x}|\boldsymbol{\theta}) \rightarrow \mathbf{x}_1, \ldots, \mathbf{x}_{\lambda} \in \mathbb{R}^n$
- 2 Evaluate x_1, \ldots, x_{λ} on f
- **3** Update parameters $\theta \leftarrow F_{\theta}(\theta, \mathbf{x}_1, \dots, \mathbf{x}_{\lambda}, f(\mathbf{x}_1), \dots, f(\mathbf{x}_{\lambda}))$

Examples



Black Box Stochastic Optimization Algorithms

Template to minimize $f : \mathbb{R}^n \to \mathbb{R}$

Initialize distribution parameters $\theta,$ set population size $\lambda \in \mathbb{N}$ While not terminate

- **1** Sample distribution $P(\mathbf{x}|\boldsymbol{\theta}) \rightarrow \mathbf{x}_1, \ldots, \mathbf{x}_{\lambda} \in \mathbb{R}^n$
- 2 Evaluate x_1, \ldots, x_{λ} on f
- **3** Update parameters $\theta \leftarrow F_{\theta}(\theta, \mathbf{x}_1, \dots, \mathbf{x}_{\lambda}, f(\mathbf{x}_1), \dots, f(\mathbf{x}_{\lambda}))$

Examples



Particle Swarm Optimization (PSO)

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

Sample new position For every particle index $i \in \{1, ..., \lambda\}$ and every coordinate $j \in \{1, ..., n\}$

$$x_{i}^{j}(t+1) = x_{i}^{j}(t) + w \times \underbrace{\left(x_{i}^{j}(t) - x_{i}^{j}(t-1)\right)}_{i}$$

previous step, momentum

+
$$\underbrace{c_1 \, \mathcal{U}_i^j(0,1)(p_i^j - x_i^j(t))}_{i}$$
 + $\underbrace{c_2 \, \tilde{\mathcal{U}}_i^j(0,1)(g_i^j - x_i^j(t))}_{i}$

approach the "previous" best approach the "global" best

inertia weight $w \approx 0.7$, and coefficients $c_1 = c_2 \approx 1.2$

- 2 Evaluate $x_1(t+1), \ldots, x_{\lambda}(t+1)$
- Opdate parameter of distribution

$$\begin{aligned} & p_i = x_i(t+1) \text{ if } f(x_i(t+1)) < f(p_i) \\ & g_i = x_*(t+1) \text{ where } f(x_*(t+1)) = \min \{ f(x_k(t+1)), k \in \text{neighborhood of } i \} \end{aligned}$$

Summary and conclusion

PSO and CMA-ES



Summary and conclusion

PSO and CMA-ES



Summary and conclusion

PSO and CMA-ES



Summary and conclusion

PSO and CMA-ES



Summary and conclusion

PSO and CMA-ES



Summary and conclusion

PSO and CMA-ES





Summary and conclusion

PSO and CMA-ES





Summary and conclusion

PSO and CMA-ES





Evolution Strategies (ES)

New search points are sampled normally distributed

$$\boldsymbol{x}_i \sim \mathcal{N}_i(\boldsymbol{m}, \sigma^2 \boldsymbol{C})$$
 for $i = 1, \dots, \lambda$

where $x_i, m \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, and $C \in \mathbb{R}^{n \times n}$





Adaptation of the sampling distribution

one-fifth success rule

rotational invariant isotropic distribution Rechenberg, 60's



adaptive search distribution along coordinate axis not rotational invariant

Rechenberg, Schwefel, 70's, 80's

- CMA-ES
 - adaptation of full covariance matrix
 - invariance to linear transformation of the search space

adaptive *and* rotational invariant Hansen, Ostermeier, 90's

Performed best for CEC 05 challenge on parameter optimization



PSO facing ill-conditioned and non-separable problems

- Test problems
- Results

3 Summary and conclusion

Experimental setting

Test problems

- Ellipsoid function
- Rosenbrock function
- Rastrigin function

separable / non-separable different condition number

non-separable different condition number

separable / non-separable

2 Algorithms

CMA-ES

PSO (Standard PSO 2006)

code available at

http://www.particleswarm.info/Standard_PSO_2006.c

code and references available at http://www.bionik.tu-berlin.de/user/niko/

21 runs in each case and default paramters.

Stochastic Optimization

Numerical experiments

000000

Summary and conclusion

Test problems

Ellipsoid

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

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

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

R random rotation
 $H_{\text{elli}}^{\text{rot}} = \mathbf{R}^T H_{\text{elli}} \mathbf{R}$
non-separable

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

$$\label{eq:alpha} \begin{split} \alpha = 6 \equiv \text{axis ratio of } 10^3 \text{, typical for real-world} \\ \text{problem} \end{split}$$

1

Results

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

PSO vs CMA-ES





Summary and conclusion

 $\alpha = 1, \ldots, 4$

Results

Rosenbrock (Banana) function

$$f_{\text{rosen}}(x) = \sum_{i=1}^{n-1} \left[(1-x_i)^2 + 10^{\alpha} (x_{i+1} - x_i^2)^2 \right]$$

- non-separable
- $\alpha = 2$, classical Rosenbrock function



Stochastic	0	ptim	izati	on
00000000	0			

Results

Rosenbrock functions $f_{\text{rosen}}(x) = \sum_{i=1}^{n-1} \left[(1-x_i)^2 + 10^{\alpha} (x_{i+1} - x_i^2)^2 \right]$



Results

Rastrigin function



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

- **R** random rotation
- on non-separable
- multimodal



Results

Rastrigin function PSO exploits the separability

PSO & Rastrigin: 21 simulations, dimension 10, tol 0,0001, taille 300, eval max 10000000



PSD & Rastrigin (rot): 21 simulations, dimension 10, tol 0.0001, taille 300, eval max 1000000



Results

Rastrigin function PSO vs CMA-ES



PSO & Rastrigin: 21 simulations, dimension 10, tol 0,0001, taille 300, eval max 10000000

10be+00 1e+05 2e+05 3e+05 4e+05 5e+05 6e+05 7e+05 8e+05 9e+05 1e+06







- PSO facing ill-conditioned and non-separable problems
- 3 Summary and conclusion

Summary and conclusion

• PSO exploits the separability

ellipsoid vs rotated ellipsoid rastrigin vs rotated rastrigin

• faster than CMA-ES for separable ellipsoid and cond > 3

```
up to 6 times faster
```

- more than 100 times slower on rotated ellipsoid than non-rotated (cond = 10^4)
- PSO virtually unable to solve rotated ellipsoid for cond > 10⁴
- PSO 70 times slower than CMA-ES on classical Rosenbrock function
- PSO non-invariant

Can we built an PSO that performs well on functions with dependencies between variables?