





	Problem Statement	Black Box Optim	ization and Its Difficulties	
rc onti •	blem Statement nuous Domain Search/Optimization Task: minimize an objective fur function) in continuous domain	<mark>nction</mark> (<i>fitne</i> า	ess function, <i>loss</i>	
	$f:\mathcal{X}\subseteq\mathbb{R}^n$ –	$\rightarrow \mathbb{R}, x$	$\mapsto f(\boldsymbol{x})$	
Black Box scenario (direct search scenario)				
			f(x)	
•	 gradients are not available o problem domain specific kno box, e.g. within an appropria Search costs: number of funct 	r not useful wledge is us te encoding ion evaluati	sed only within the black	
·				
		4	<ロ>	



Problem Statement Black Box Optimization and Its Difficulties











Problem Statement Black Box Optimization and Its Difficulties

Curse of Dimensionality

The term *Curse of dimensionality* (Richard Bellman) refers to problems caused by the rapid increase in volume associated with adding extra dimensions to a (mathematical) space.

Example: Consider placing 20 points equally spaced onto the interval [0, 1]. Now consider the 10-dimensional space $[0, 1]^{10}$. To get similar coverage in terms of distance between adjacent points requires $20^{10} \approx 10^{13}$ points. 20 points appear now as isolated points in a vast empty space.

Remark: distance measures break down in higher dimensionalities (the central limit theorem kicks in)

Consequence: a search policy that is valuable in small dimensions might be useless in moderate or large dimensional search spaces. Example: exhaustive search.







What Makes a	a Function Difficult to	o Solve?			
The Problem	Possible Approac	hes			
Dimensionality	exploiting the problem stru separability, I	ucture ocality/neighborhood, encoding			
III-conditioning	second order approach ch	anges the neighborhood metric			
Ruggedness	non-local policy, large sampling width (step-size) as large as possible while preserving a reasonable convergence speed				
	population-based method	population-based method, stochastic, non-elitistic			
	recombination operator	serves as repair mechanism			
	restarts				
		metaphors			
	16	< 口 > < 图 > < 토 > < 토 > < 臣 > < 回 > < 이 오 아			

nt III-Conditioned Proble

Problem Statement		III-Conditioned Problems	
Metaphors			
Evolutionary Computation		Optimization/Nonlinear Programmin	
individual, offspring, parent	\longleftrightarrow	candidate solution decision variables design variables	
population fitness function	$\stackrel{\longleftrightarrow}{\longleftrightarrow}$	set of candidate solutions objective function loss function cost function	
generation	\longleftrightarrow	error function iteration	
		methods: ESs	
	1	< □ > < 億 > < 差 > < 差 > 差 - 釣んで	

Overview

Problem Statement Continuous Black-Box Optimization Typical Difficulties

2 Stochastic Black-Box Algorithms

General Template Invariance Comparisons of a few DFOs

Zoom on Evolution Strategies Step-size Adaptation

Displaying results and visualization Statistics

Average Runtime Empirical Cumulative Distribution Function (ECDF)

19

Problem Statement	
Landscape of Continuous Black-Box Optimization	
Deterministic algorithms Quasi-Newton with estimation of gradient (BFGS) [Broyden et al. 1970] Simplex downhill [Nelder & Mead 1965] Pattern search [Hooke and Jeeves 1961] Trust-region methods (NEWUOA, BOBYQA) [Powell 2006, 2009]	
Stochastic (randomized) search methods	
Evolutionary Algorithms (continuous domain)	
Differential Evolution [Storn & Price 1997]	
Particle Swarm Optimization [Kennedy & Eberhart 1995]	
Evolution Strategies, CMA-ES [Rechenberg 1965, Hansen & Ostermeier 2001	
Estimation of Distribution Algorithms (EDAs) [Larrañaga, Lozano, 2002]	
Cross Entropy Method (same as EDA) [Rubinstein, Kroese, 2004]	
Genetic Algorithms [Holland 1975, Goldberg 1989]	
Simulated annealing [Kirkpatrick et al. 1983]	
Simultaneous perturbation stochastic approximation (SPSA) [Spall 2000]	

Evolution Strategies (ES) A Search Template

Stochastic Search

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

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

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

Everything depends on the definition of P and F_{θ}

deterministic algorithms are covered as well

In many Evolutionary Algorithms the distribution *P* is implicitly defined via operators on a population, in particular, selection, recombination and mutation

Natural template for (incremental) Estimation of Distribution Algorithms





Evolution Strategies (ES) The Normal Distribution

The Multi-Variate (n-Dimensional) Normal Distribution

Any multi-variate normal distribution $\mathcal{N}(m, \mathbb{C})$ is uniquely determined by its mean value $m \in \mathbb{R}^n$ and its symmetric positive definite $n \times n$ covariance matrix \mathbb{C} .

The mean value m

- determines the displacement (translation)
- value with the largest density (modal value)
- the distribution is symmetric about the distribution mean



The covariance matrix C

- determines the shape
- geometrical interpretation: any covariance matrix can be uniquely identified with the iso-density ellipsoid { $x \in \mathbb{R}^n | (x m)^T \mathbb{C}^{-1} (x m) = 1$ }





The solution Strategies (ES) The Normal Distribution $\begin{aligned}
& \text{The } (\mu/\mu, \lambda) \text{-ES} \\
& \text{Non-elitist selection and intermediate (weighted) recombination} \\
& \text{Given the } i\text{-th solution point } x_i = m + \sigma \underbrace{\mathcal{N}_i(\mathbf{0}, \mathbf{C})}_{=:y_i} = m + \sigma y_i \\
& \text{Let } x_{i:\lambda} \text{ the } i\text{-th ranked solution point, such that } f(\mathbf{x}_{1:\lambda}) \leq \cdots \leq f(\mathbf{x}_{\lambda:\lambda}). \\
& \text{The new mean reads} \\
& m \leftarrow \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda} = m + \sigma \underbrace{\sum_{i=1}^{\mu} w_i y_{i:\lambda}}_{=:y_w} \\
& \text{where} \\
& w_1 \geq \cdots \geq w_\mu > 0, \quad \sum_{i=1}^{\mu} w_i = 1, \quad \frac{1}{\sum_{i=1}^{\mu} w_i^2} =: \mu_w \approx \frac{\lambda}{4} \\
& \text{The best } \mu \text{ points are selected from the new solutions (non-elitistic)} \\
& \text{and weighted intermediate recombination is applied.}
\end{aligned}$











ě



Overview

Problem Statement Continuous Black-Box Optimization Typical Difficulties

Stochastic Black-Box Algorithms General Template Invariance

Comparisons of a few DFOs

3 Zoom on Evolution Strategies

Step-size Adaptation Covariance Matrix Adaptation

O Evaluating Black-Box Algorithms

Displaying results and visualization Statistics Average Runtime Empirical Cumulative Distribution Function (ECDF)

35

Zoom On Evolution Strategies

Zoom on ESs: Objectives

Illustrate why and how sampling distribution is controlled

step-size control (overall standard deviation) allows to achieve linear convergence

covariance matrix control

allows to solve ill-conditioned problems













Overview
Problem Statement
Continuous Black-Box Optimization
l ypical Difficulties
Stochastic Black-Box Algorithms
General Template
Invariance
Comparisons of a few DFOs
O Zoom on Evolution Strategies
Step-size Adaptation
Covariance Matrix Adaptation
Evaluating Black-Box Algorithms
Displaying results and visualization
Statistics
Average Runtime
Empirical Cumulative Distribution Function (ECDF)
43























55

Covariance Matrix Adaptation (CMA) Covariance Matrix Bank-One Update **Covariance Matrix Adaptation** Rank-One Update Initialize $m \in \mathbb{R}^n$, and $\mathbb{C} = \mathbb{I}$, set $\sigma = 1$, learning rate $c_{cov} \approx 2/n^2$ While not terminate $\mathbf{x}_i = \mathbf{m} + \sigma \mathbf{y}_i, \quad \mathbf{y}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C}),$ $\boldsymbol{m} \leftarrow \boldsymbol{m} + \sigma \boldsymbol{y}_w$ where $\boldsymbol{y}_w = \sum_{i=1}^{\mu} w_i \boldsymbol{y}_{i:\lambda}$ $\mathbf{C} \leftarrow (1 - c_{\text{cov}})\mathbf{C} + c_{\text{cov}}\mu_{w} \underbrace{\mathbf{y}_{w}\mathbf{y}_{w}^{\mathrm{T}}}_{\mathbf{v}_{w}\mathbf{v}_{w}} \quad \text{where } \mu_{w} = \frac{1}{\sum_{i=1}^{\mu} w_{i}^{2}} \ge 1$ The rank-one update has been found independently in several domains^{6 7 8 9} ⁶Kiellström&Taxén 1981. Stochastic Optimization in System Design, IEEE TCS 7 Hansen&Ostermeier 1996. Adapting arbitrary normal mutation distributions in evolution strategies: The covariance matrix adaptation, ICEC ⁸Ljung 1999. System Identification: Theory for the User ⁹Haario et al 2001. An adaptive Metropolis algorithm, JSTOR ・ロト・個ト・モト・モト ヨー のへで





Overview

- Problem Statement Continuous Black-Box Optimization Typical Difficulties
- Stochastic Black-Box Algorithms General Template Invariance Comparisons of a few DFOs
- Zoom on Evolution Strategies Step-size Adaptation Covariance Matrix Adaptation

Output State St

Displaying results and visualization Statistics Average Runtime Empirical Cumulative Distribution Function (ECDF)

59



Performance Evaluation Subsection

Evaluation of Anytime Black-Box Optimizers Particularly Randomized Search Algorithms

Randomized optimization is mostly an empirical science

Hence it is crucial to properly conduct numerical experiments and assess performance

to not fool ourself on what our favorite algorithm is good/not good at

in order to not fool others ...

"The first principle is that you must not fool yourself and you are the easiest person to fool." Richard P. Feynman

Performance Evaluation Subsection

Evaluation of Anytime Black-Box Optimizers Particularly Randomized Search Algorithms

Evaluation of performance in a broad sense

not only be able to say "Algorithm A is better than B"

61

but

understand where and why algorithm work

quantify performance















(see next slide)













Performance Evaluation

Implication

Which Statistics?



Which Statistics?

Performance Evaluation

Performance Evaluation Statistical Assessment

Statistical Assessment

Apply rank-sum test (Wilcoxon, Mann-Whitney U) only assumption: no equal data values

hypothesis:

compares $sPr(x > y) \neq Pr(x < y) \neq 1/2$ inking

two-sided 1%-significance p-value needs only 2x5 data values

For the same p-value, fewer significant data are better using enough data, any difference can be made significant

Generally: non-parametric tests, Kolmogorov-Smirnov test for ECDFs, no need to use the t-test

77



Overview

Problem Statement
Continuous Black-Box Optimization
Typical Difficulties
Stochastic Black-Box Algorithms
General Template
Invariance
Comparisons of a few DFOs
O Zoom on Evolution Strategies
Step-size Adaptation
Covariance Matrix Adaptation
Output State St
Displaying results and visualization
Statistics
Average Runtime
Empirical Cumulative Distribution Function (ECDF)
78

Performance Evaluation On performance measure

Evaluation of Search Algorithms Behind the scene

a performance should be

quantitative on the ratio scale (highest possible) "algorithm A is two *times* better than algorithm B" is a meaningful statement

can assume a wide range of values

meaningful (interpretable) with regard to the real world possible to transfer from benchmarking to real world

runtime or first hitting time is the prime candidate (we don't have many choices anyway)









Overview	Em
Problem Statement	Sing
Continuous Black-Box Optimization	
Typical Difficulties	
Stochastic Black-Box Algorithms	
General Template	
Invariance	
Comparisons of a few DFOs	
October 2007 Strategies	
Step-size Adaptation	
Covariance Matrix Adaptation	
4 Evaluating Black-Box Algorithms	
Displaying results and visualization	
Statistics	
Average Runtime	
Empirical Cumulative Distribution Function (ECDF)	
85	









