Anytime Performance Assessment: Runtime Distributions Beyond Data Profiles

Nikolaus Hansen Inria & École polytechnique, France

June 2024









Check for updates

COCO: a platform for comparing continuous optimizers in a black-box setting

Nikolaus Hansen^{a,b}, Anne Auger^{a,b}, Raymond Ros^c, Olaf Mersmann^d, Tea Tušar^e and Dimo Brockhoff^{a,b}

^aRandopt Team, Inria, Palaiseau, France; ^bCMAP, CNRS, Ecole Polytechnique, Institut Polytechnique de Paris, Palaiseau, France; ^cLRI, Université Paris-Sud, Orsay, France; ^dComputational Statistics, TU Dortmund University, Dortmund, Germany; ^eJožef Stefan Institute, Ljubljana, Slovenia

ABSTRACT

We introduce COCO, an open-source platform for Comparing Continuous Optimizers in a black-box setting. COCO aims at automatizing the tedious and repetitive task of benchmarking numerical optimization algorithms to the greatest possible extent. The platform and the underlying methodology allow to benchmark in the same framework deterministic and stochastic solvers for both single and multiobjective optimization. We present the rationals behind the (decade-long) development of the platform as a general proposition for guidelines towards better benchmarking. We detail underlying fundamental concepts of COCO such as the definition of a problem as a function instance, the underlying idea of instances, the use of target values, and runtime defined by the number of function calls as the central performance measure. Finally, we give a quick overview of the basic code structure and the currently available test suites.

ARTICLE HISTORY

Received 17 June 2019 Accepted 8 August 2020

KEYWORDS

Numerical optimization; black-box optimization; derivative-free optimization; benchmarking; performance assessment; test functions; runtime distributions: software

1. Introduction

We consider the continuous black-box optimization or search problem to minimize

$$f: X \subset \mathbb{R}^n \to \mathbb{R}^m \quad n, m \ge 1, \tag{1}$$

where the search domain X is typically a bounded hypercube or the entire continuous space.¹ More specifically, we aim to find, as quickly as possible, one or several solutions **x** in the search space X with *small* value(s) of $f(\mathbf{x}) \in \mathbb{R}^m$.

A continuous optimization algorithm, denoted as *solver*, addresses the above problem. In this paper, we only consider zero-order black-box optimization [19,57,58]: while the search domain $X \subset \mathbb{R}^n$ and its boundaries are accessible, no other prior knowledge about f is available to the solver.² That is, f is considered as a black-box, also known as an oracle, and the only way the solver can acquire information on f is by querying the value f(x) of a solution $x \in X$. Zero-order black-box optimization is thus a derivative-free optimization setting.³ We generally consider 'time' to be the number of calls to the function f and will define 'runtime' correspondingly.

Anytime Performance Assessment in Blackbox **Optimization Benchmarking**

Nikolaus Hansen, Anne Auger, Dimo Brockhoff¹⁰, and Tea Tušar¹⁰

Abstract—We present concepts and recipes for the anytime performance assessment when benchmarking optimization algorithms in a blackbox scenario. We consider runtime-oftentimes measured in the number of blackbox evaluations needed to reach a target quality-to be a universally measurable cost for solving a problem. Starting from the graph that depicts the solution quality versus runtime, we argue that runtime is the only performance measure with a generic, meaningful, and quantitative interpretation. Hence, our assessment is solely based on runtime measurements. We discuss proper choices for solution quality indicators in single- and multi-objective optimization, as well as in the presence of noise and constraints. We also discuss the choice of the target values, budget-based targets, and the aggregation of runtimes by using simulated restarts, averages, and empirical cumulative distributions which generalize convergence graphs of single runs. The presented performance assessment is to a large extent implemented in the comparing continuous optimizers (COCO) platform freely available at https://github.com/numbbo/coco.

Index Terms—Anytime optimization, benchmarking, blackbox optimization, performance assessment, quality indicator.

I. INTRODUCTION

WE PRESENT practical concepts and ideas for the performance assessment of optimization algorithms when benchmarked in a blackbox and anytime scenario. Going beyond a simple ranking of algorithms, we aim to provide a quantitative and meaningful performance assessment, which allows for conclusions like algorithm A is seven times faster than algorithm B in solving a given problem or in solving problems with certain characteristics. To achieve this end in a comparative and timeless manner, we argue that we should measure the number of blackbox evaluations to reach a predefined *quality indicator* value (a target). More generally, we argue to measure a cost that is defined on a ratio scale and is comparable across publications. We call this measure

Nikolaus Hansen, Anne Auger, and Dimo Brockhoff are with Inria and Ecole Polytechnique, Institut Polytechnique de Paris, 91128 Palaiseau, France (e-mail: dimo.brockhoff@inria.fr).

Tea Tušar is with the Department of Intelligent Systems, Jožef Stefan Institute, 1000 Ljubljana, Slovenia.

Color versions of one or more figures in this article are available at https://doi.org/10.1109/TEVC.2022.3210897.

the *runtime* of the algorithm to reach a given target. Yet, our assessment methodology does not depend on any specific cost measure, as long as the costs are quantitative and comparable.¹

In this article, we formalize the optimization goal by a so-called quality indicator. Its definition may heavily depend on the optimization scenario, e.g., the number of objectives or constraints. Broadly speaking, a quality indicator is based on the sequence of all so-far visited solutions. In the simplest case, it is the objective function value of the last visited solution.

Runtimes represent the cost of optimization. Compared to the quality indicator, the definition of costs depends to a lesser extent on the specific optimization scenario. To sustain reproducibility and comparability across publications, we recommend against CPU or wall-clock time as cost measure² (see also Hooker [22] for a further discussion on the unwanted consequences of benchmarking based on CPU time).

Benchmarking is usually computationally expensive and benchmarking for a *single* budget seems vastly inefficient by 1) addressing only one of many possible budget scenarios (scenarios heavily depend on the software and hardware environment) and 2) throwing away most of the data generated during the experiment. An anytime approach to benchmarking prevents these drawbacks. To allow for a budget-free performance assessment even for non-anytime algorithms that have a maximum or timeout budget as decisive or mandatory *input* parameter (decided by the user), we collect data with an any-budget experimental procedure that runs repeated experiments with increasing input budget [31].³ Non-anytime algorithms that do not take a maximum budget as an input parameter can be accurately assessed only by the time of their final solution proposal.

In this article, we advocate to routinely use (anytime) empirical runtime distributions to assess the performance of optimization algorithms. We demonstrate how to directly compare the runtime distributions of algorithms that have

¹We are grateful to the anonymous reviewer pointing this out to us.

²An exploratory CPU timing experiment to get an estimate of the internal time complexity of the algorithm is still advisable, like it is prescribed in the comparing continuous optimizer (COCO) platform [20].

³We can stop the procedure when the last budget was not fully exhausted. Increasing the budget each time by a factor of r > 1 adds to the overall computational costs for the experimentation less than $\sum_{i=0}^{\infty} 1/r^i = r/(r-1)$ times the last consumed budget. For the performance assessment, always the data from the smallest eligible budget is used. The performance assessment will be too optimistic by tacitly assuming that the budget can be set properly without additional costs. On the other hand, runtimes may be overestimated (by less than a factor of r). A code example is provided in

1293

Manuscript received 18 September 2021; revised 16 February 2022 and 16 June 2022; accepted 14 September 2022. Date of publication 29 September 2022; date of current version 1 December 2022. This work was supported in part by the French National Research Agency (NumBBO) under Grant ANR-12-MONU-0009, and in part by the Slovenian Research Agency under Grant P2-0209 and Grant N2-0254. (Corresponding author: Dimo Brockhoff.)

Overview

- The COCO platform
- Data profiles are runtime distributions
- What to aggregate?
- Target *f*-values
- Integrating out (low) success rates
- Runtime distributions versus convergence graphs \bullet

Nikolaus Hansen, Inria

З



Benchmarking is a tedious and repetitive task

with surprisingly many traps and pitfalls

That's where COCO comes into play: semi-automatized benchmarking

Black-Box Optimization Benchmarking Template for the Comparison of More than Two Algorithms on the Noiseless Testbed' Draft version

ABSTRACT

CCS CONCEPTS

KEYWORDS nchmarking, Black-box opti

CPU TIMING ler to evaluate the CPU timing of the a

It becomes easy to give a comprehensive quantitative *context* of the presented work by using previous data...

https://dl.acm.org/doi/pdf/10.1145/1830761.1830790 https://github.com/numbbo/coco



Figure 2: Bootstrapped empirical cumulative distribution of the number of objective function evaluations divided by dimen sion (FEvals/DIM) for 51 targets with target precision in $10^{[-8.,2]}$ for all functions and subgroups in 5-D. As reference algorithm, the best algorithm from BBOB 2009 is shown as light thick line with diamond markers.



COCO/bbob



Figure by Tea Tušar in Hansen et al (2021), COCO: A platform for comparing continuous optimizers in a black-box setting. Optimization Methods and Software, 36(1), 114-144. 5

Important Characteristics of Benchmarking with COCO/bbob

- Functions are instantiated multiple times with
 - different locations of the optimum in x- and f-space
 - different search space rotations
 - several monotonous local transformations to obfuscate simple regularities \bullet making inadvertent exploitation of artificial regularities (aka overfitting) much less likely
- Functions are (generally) scalable with dimension
- Performance evaluation via (generalized) data profiles
- Provision of ~300 data sets to directly compare with

first and foremost



Important Characteristics of Benchmarking with COCO/bbob

- Functions are instantiated multiple times with
 - different locations of the optimum in x- and f-space
 - different search space rotations
 - several monotonous local transformations to obfuscate simple regularities making inadvertent exploitation of artificial regularities (aka overfitting) much less likely
- Functions are (generally) scalable with dimension
- Performance evaluation via (generalized) data profiles
- Provision of ~300 data sets to directly compare with

first and foremost



COCO (COmparing Continuous Optimizers) Impact Numbers

Table 1. Visibility of COCO.

Data sets online bbob suite bbob-noisy suite bbob-biobj suite bbob-largescale suite bbob-mixint suite BBOB workshop papers using COCO Unique authors on the workshop papers Papers in Google Scholar found with the search phrase OR 'black-box optimization benchmarking (BBOB)' Citations to the COCO documentation including [16,17,28 Citations to Hansen et al 2021. COCO: A platform for c

can be loaded and processed.

Hansen et al 2021. OMS...

	227
	45
	32
	11
	4
	143
	109 from 28 countries
'comparing continuous optimizers'	559
8,35–37,39–42]	1455
comparing OMS 36,1.	491

Any `cocopp.archiving.create() `-ed data sets that are provided under an URL



"data profiles [...] have been designed to provide [...] the percentage of problems that can be solved (for a given tolerance τ) with a given number of function evaluations [μ_f]."

- 1. Create a set of recorded sequences f_t , $t = 1, 2, ..., \mu_f$
- 2. Define for each problem a *target f*-value based on the data (aka convergence test)

Data Profiles

Moré and Wild 2009. SIAM J. OPTIM. 20,1.

for different algorithms and different problems/functions

based on the smallest observed *f*-value f_L and a tolerance parameter $\tau \in \{10^{-1}, 10^{-3}, 10^{-5}, 10^{-7}\}$: $f_{\text{target}} = f_L + \tau \cdot (f(x_0) - f_L)$

3. Find the "runtimes" t, as number of f-evaluations, when, for the first time, $f_t \leq f_{\text{target}}$ was met

for all problems and each algorithm and each tolerance τ

4. Plot the empirical cumulative distribution of $\frac{t}{1}$ over all problems for each algorithm and each τ n + 1different tolerances τ are usually shown in different plots





FIG. 5.1. Data profiles $d_s(\kappa)$ for the smooth problems \mathcal{P}_S show the percentage of problems solved as a function of a computational budget of simplex gradients.



Remarks/Discussion

problems.

the values are aggregated in the same graph, hence they should have the same unit

• The target value f_{target} is based on (the same) empirical data

- Aggregation is
 - done over different problems and different dimensions
 - not done over different algorithms or different tolerances τ

The measured values t can be any cost value that is comparable between all

even for the same problem, performance results are not directly comparable between publications



Remarks/Discussion

problems.

the values are aggregated in the same graph, hence they should have the same unit

• The target value f_{target} is based on (the same) empirical data

- Aggregation is
 - done over different problems and different dimensions
 - not done over different algorithms or different tolerances au

The measured values t can be any cost value that is comparable between all

even for the same problem, performance results are not directly comparable between publications



Our Insights

- 1. Aggregation over a wide range of dimensions is not conducive dimension can (and should) be used as algorithm selection decision parameter!
- 3. Restricting the budget of the *benchmarked* algorithms has no methodological advantage

2. Aggregation over various tolerance values is possible and seems useful

arguably, a problem is only defined by a function and a target then we still just aggregate over all problems

at least none I am aware of our measured data are runtimes (budgets)



Generalizing the setting of f_{target}

- We can separate the set of algorithms that determine $f_{\text{target}} = f_L + \tau \cdot (f(x_0) f_L)$ (by determining f_I) from the set of benchmarked algorithms. Thereby, because we can keep f_L constant, data and graphs become comparable across publications
- Then, instead of setting and varying the f-tolerance τ , we can vary the budget μ_f to get different target values (AKA budget-based or runlength-based targets)

- In COCO, we (usually) set f_{target} -values based on the known optimal f-value
- We use a budget-free experimental setup

increasing the "time-out" budget does not affect results for smaller budgets, hence results remain compatible larger budgets allow for a better quantification of performance losses

the meaning of budgets is somewhat easier to understand we don't need to change the budget of the benchmarked algorithms

thereby defining the objective extrinsically/absolute rather than relative



Restarted Algorithm and Simulated Runtimes

• The runtime (RT) of a (randomized) restarted algorithm is

• The expected RT is (where t is a random variable) • An estimator from the data is $\hat{E}[\text{RT}] = \frac{\sum t^{\text{unsuc}} + \sum t^{\text{succ}}}{N_{\text{succ}}} = \frac{\text{overall costs}}{\# \text{successes}}$ aka ERT, Enes, SP2, aRT.

Implications:

- negotiate success vs speed

Price 1997, Auger & Hansen 2005 Nunsuc $RT = \sum_{i=1}^{n} t_i^{\text{unsuc}} + t^{\text{succ}}$ i=1 $E[\mathsf{RT}] = \frac{1 - p_{\mathsf{succ}}}{E[t^{\mathsf{unsuc}}]} + E[t^{\mathsf{succ}}]$ $p_{\sf succ}$ $(1 - p_{succ})E[t^{unsuc}] + p_{succ}E[t^{succ}]$ *p*succ

• if we have at least one successful run, we can simulate runtimes, estimate the expected runtime, and don't need to

• if we have no successful run, we can estimate a lower bound for the expected runtime $\hat{E}[RT] > \sum_{t} t^{t}$



Restarted Algorithm and Simulated Runtimes

• The runtime (RT) of a (randomized) restarted algorithm is

• The expected RT is (where t is a random variable)

• An estimator from the data is

 $\hat{E}[\text{RT}] = \frac{\sum t^{\text{unsuc}} + \sum t^{\text{succ}}}{N_{\text{succ}}} = \frac{\text{overall costs}}{\#_{\text{successes}}}$

aka ERT, Enes, SP2, aRT.

Implications:

- negotiate success vs speed

Price 1997, Auger & Hansen 2005 Nunsuc $RT = \sum_{i=1}^{n} t_i^{\text{unsuc}} + t^{\text{succ}}$ i=1 $E[\mathsf{RT}] = \frac{1 - p_{\mathsf{succ}}}{E[t^{\mathsf{unsuc}}]} + E[t^{\mathsf{succ}}]$ $p_{\sf succ}$ $- \frac{(1 - p_{succ})E[t^{unsuc}] + p_{succ}E[t^{succ}]}{E[t^{succ}]}$ *p*succ

• if we have at least one successful run, we can simulate runtimes, estimate the expected runtime, and don't need to

• if we have no successful run, we can estimate a lower bound for the expected runtime $\hat{E}[RT] > \sum_{t} t^{t}$



Restarted Algorithm and Simulated Runtimes

• The runtime (RT) of a (randomized) restarted algorithm is

• The expected RT is (where t is a random variable)

• An estimator from the data is

 $\hat{E}[\text{RT}] = \frac{\sum t^{\text{unsue}} + \sum t^{\text{suce}}}{N_{\text{outor}}} = \frac{\text{overall costs}}{\#\text{successes}}$

aka ERT, Enes, SP2, aRT.

Implications:

- negotiate success vs speed

Price 1997, Auger & Hansen 2005 Nunsuc $RT = \sum_{i=1}^{n} t_i^{\text{unsuc}} + t^{\text{succ}}$ i=1 $E[\mathsf{RT}] = \frac{1 - p_{\mathsf{succ}}}{E[t^{\mathsf{unsuc}}]} + E[t^{\mathsf{succ}}]$ $p_{\sf succ}$ $= \frac{(1 - p_{\text{succ}})E[t^{\text{unsuc}}] + p_{\text{succ}}E[t^{\text{succ}}]}{1 - p_{\text{succ}}E[t^{\text{succ}}]}$ *p*succ

• if we have at least one successful run, we can simulate runtimes, estimate the expected runtime, and don't need to

• if we have no successful run, we can estimate a lower bound for the expected runtime $\hat{E}[RT] > \sum t^{unsuc}$





Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).





Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

equidistance "target" values





Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

Nikolaus Hansen, Inria

20





Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).





Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. IEEE Trans. on EC, 26(6).





23

Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. IEEE Trans. on EC, 26(6).

for the remaining construction, we could use any runtimes, for example, from different runs or different functions





Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. IEEE Trans. on EC, 26(6).





Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. IEEE Trans. on EC, 26(6).





optimization benchmarking. IEEE Trans. on EC, 26(6).





Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. IEEE Trans. on EC, 26(6).





optimization benchmarking. IEEE Trans. on EC, 26(6).



when we maximize (instead of minimize), the graph can be considered as an empirical runtime distribution as is



Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. IEEE Trans. on EC, 26(6).



the area above the curve represent a (truncated) average runtime

When the x-axis is in logscale, the area is the (truncated) geometric average



optimization benchmarking. IEEE Trans. on EC, 26(6).



the area above the curve represent a (truncated) average runtime

When the x-axis is in logscale, the area is the (truncated) geometric average



optimization benchmarking. IEEE Trans. on EC, 26(6).



Aggregated Runtime Distributions



























Stop writing "statistically significant"

- recommending that the calculation and use of continuous p-values be *quantities* (e.g., p = 0.08)."
- result refutes or supports a scientific hypothesis."
- replication crisis."

• Wasserstein et al. 2019: "We conclude, based on our review of the articles in this special issue and the broader literature, that it is time to stop using the term "statistically significant" entirely. Nor should variants such as "significantly different," "p < 0.05," and "nonsignificant" survive, [...] however, we are not discontinued. Where *p-values* are used, they should be reported as continuous

Moving to a World Beyond "p < 0.05". The American Statistician, 73, S1.

• Amrhein et al. + 800 signatories, 2019: "We agree, and call for the entire concept of statistical significance to be abandoned. [...] we are calling for a stop to the use of P values in the conventional, dichotomous way — to decide whether a

Retire statistical significance. Scientists rise up against statistical significance. *Nature*, 567(7748).

• Cockburn et al. 2020: "misuse of statistical significance as the standard of evidence for experimental success has been identified as a key contributor in the

Threats of a replication crisis in empirical computer science. Communications of the ACM, 63(8).







How to use a *p*-value wisely

An observed *p*-value indicates by how much we should *update* our confidence in H_0 (not: how confident we should be in H_0)

posterior odds

If we do not provide an estimate for the prior odds, we have no argument to reject H_0 (and that's perfectly fine too)

If we improved a well-established state-of-the-art algorithm or invented cold fusion or find a room-temperature superconductor at atmospheric pressures, the prior odds of H_0 are usually high, say, at least 10^4 .

> the higher the prior odds for H_0 , the more exceptional or surprising is the scientific result to accept $\neg H_0$ with the same confidence we had in H_0 before, we need $p \approx P(\neg H_0)^2$

$Odds(H_0 \mid D) \approx Odds(H_0) \times 2p$

prior odds

a small p stands on its own merits: we can conclude that the odds for H_0 have decreased by a factor of about 2p



Summary

- Target definitions (the convergence test condition) can be separated from benchmarked algorithms
- We aggregate over different target values but not over dimension
- (simulated) restarts can integrate out success rates
- We use a budget-free experimental setup
- We read data profiles preferably as horizontal data (runtimes) rather then vertically data (success rates)
- Data profiles do not obscure the problem difficulty (as performance profiles do)

thereby, results can become comparable across publications



Summary

- algorithms
- We aggregate over different target values but not over dimension
- (simulated) restarts can integrate out success rates lacksquare
- We use a budget-free experimental setup.
- (success rates)
- Data profiles do not obscure the problem difficulty (as performance profiles do)

• Target definitions (the convergence test condition) can be separated from benchmarked

thereby, results can become comparable across publications



• We read data profiles preferably as horizontal data (runtimes) rather then vertically data

