Derivative-Free Optimization Benchmarking and Performance Assessment

December 13, 2024 M2 Optimization, Université Paris-Saclay





Dimo Brockhoff Inria Saclay – Ile-de-France



Practical Group Project

What do we want to achieve?

- understand some (new or not so new) algorithms
- and how they behave subject to noise/outliers
- learn how to conduct scientific testing, how to assess performance, and how to draw scientific conclusions from benchmarking data

How does it work?

- Today:
 - introduction to benchmarking
 - installation of COCO on your machines
 - first exercise to interpret (existing) benchmarking data
- In groups of 3
 - please use this to organize yourself: http://tinyurl.com/cocoproject2024

Practical Group Project

- We propose some possible algorithms (max. 1 group/algo)
 - but you can take others (e.g. your own)
 - benchmark this algorithm with the help of COCO on noisy functions
- 24th of January:
 - report (PDF) sent by email, 8 pages (based on COCO LaTeX template?)
- 31st of January:
 - 12 minutes presentation per group
 - + 10 minutes questions
 - about the algorithm
 - about its performance (compared to other algos)
 - order of the talks to be decided early next year
- 10th of January:
 - send a short progress statement by email (<10 lines)

The Group Project's Context

A real science project

- Teaching how to do research
- Results will be new
- "Expect the unexpected"
- Anne and I are open for discussions along the project

We are very curious about your results

Scientific Context: Noisy optimization / outliers

- How is the performance of algorithms "perturbed" if the objective function gets "perturbed"?
- We provide a wrapper around deterministic test functions (from the bbob suite of the COCO platform) and an experiment script to do the benchmarking for varying levels of noise/outlier probability

Outliers + Noise Level

The noise model is the following:

- deterministic noise ("outliers"):
 - if re-evaluated, a search point x will always result in the same f(x)
 - the points *x* which are "outliers" are fixed as well
- noise-level $\in \{0, 0.02, 0.05, 0.1, 0.2, 0.3, 0.4\}$ (probability that a search point x is not the underlying objective function)
- if a point x is "noisy" or an "outlier", we add or substract a value to the original objective function
 - *x* becomes worse more often than it becomes better
 - the difference between observed and original function value is Cauchy distributed
 - more details can be found in the <u>code</u>

Excerpt from the Code

[...]

go

Benchmarking Optimization Algorithms

or: critical performance assessment







challenging optimization problems appear in many scientific, technological and industrial domains







Practical (Numerical) Blackbox Optimization



derivatives not available or not useful

Not clear:

which of the many algorithms should I use on my problem?

visualizing the raw data (single runs)

© Anne Auger and Dimo Brockhoff, Inria 2024

Displaying 3 runs (3 trials)



Displaying 3 runs (3 trials)



Displaying 3 runs (3 trials)





Even better like this: substract minimum over all runs

Displaying 51 runs



Displaying 51 runs



Which Statistics?



Which Statistics?



More Problems with Averages/Expectations

- to reliably estimate an expectation (from the *average*) we need to make *assumptions* on the tail of the underlying distribution
 - these can not be implied from the observed data
 - AKA: the average is well-known to be (highly) sensitive to outliers (extreme events)
- rare events can only be analyzed by collecting a large enough number of data

from Hansen GECCO 2019 Experimentation tutorial

Which Statistics?



- unique for uneven number of data
- independent of log-scale, offset...

median(log(data))=log(median(data))

same when taken over x- or y-direction

Implications

• use the median as summary datum

unless there are good reasons for a different statistics out of practicality: use an odd number of repetitions

• more general: use quantiles as summary data

for example out of 15 data: 2nd, 8th, and 14th value represent the 10%, 50%, and 90%-tile

from Hansen GECCO 2019 Experimentation tutorial

Benchmarking =

evaluate the performance of optimization algorithms compare the performance of different algorithms

understand strengths and weaknesses of algorithms help in design of new algorithms

How Do We Measure Performance?

Meaningful quantitative measure

- quantitative on the ratio scale (highest possible)
 "algo A is two *times* better than algo B" is a meaningful statement
- assume a wide range of values
- meaningful (interpretable) with regard to the real world

possible to transfer from benchmarking to real world

How Do We Measure Performance?

Meaningful quantitative measure

- quantitative on the ratio scale (highest possible)
 "algo A is two *times* better than algo B" is a meaningful statement
- assume a wide range of values
- meaningful (interpretable) with regard to the real world

possible to transfer from benchmarking to real world

CPU timing not a good candidate

 \rightarrow depends on implementation/language/machine/...

time is spent on code optimization instead of science J.N. Hooker: Testing heuristics, we have it all wrong, 1995, J. of Heuristics

Measuring Performance Empirically

convergence graphs is all we have to start with...



How Do We Measure Performance?

Two objectives:

- Find solution with small(est possible) function/indicator value
- With the least possible search costs (number of function evaluations)

For measuring performance: fix one and measure the other

Measuring Performance Empirically

convergence graphs is all we have to start with...



Runtime is the prime candidate

Why?

- runtimes are easier to interpret than f-values
- runtimes are interpretable in a meaningful way
- runtimes have a natural zero (hence: ratio scale)

Another reason:

- in both cases, we have missing values
- in the fixed budget view, we miss values of "too good" algorithms
- in the fixed target view, we miss values of "bad/slow" algorithms

Collect for a given target (several target), the number of function evaluations needed to reach a target

Repeat several times:

if algorithms are stochastic, never draw a conclusion from a single run

if deterministic algorithm, repeat by changing (randomly) the initial conditions

ECDF:

Empirical Cumulative Distribution Function of the Runtime [aka data profile]

Cumulative Distribution Function (CDF)

Given a random variable T, the cumulative distribution function (CDF) is defined as

 $CDF_T(t) = Pr(T \le t)$ for all $t \in \mathbb{R}$

It characterizes the probability distribution of T

If two random variables have the same CDF, they have the same probability distribution

Cumulative Distribution Function (CDF)

Given a random variable T, the cumulative distribution function (CDF) is defined as

 $CDF_T(t) = Pr(T \le t)$ for all $t \in \mathbb{R}$

It characterizes the probability distribution of T



Empirical Cumulative Distribution Function

Given a collection of data $T_1, T_2, ..., T_k$ (e.g. an empirical sample of a random variable) the *empirical* cumulative distribution function (ECDF) is a step function that jumps by 1/k at each value in the data.



It is an estimate of the CDF that generated the points in the sample.

Empirical Cumulative Distribution Function

$$ECDF_{(T_1,...,T_k)}(t) = \frac{\text{number of } T_i \le t}{k} = \frac{1}{k} \sum_{i=1}^k \mathbb{1}_{\{T_i \le t\}}$$

For $\{T_i: i \ge 1\}$ i.i.d. realization of a random variable T, by the LLN



A Convergence Graph



First Hitting Time is Monotonous


15 Runs



15 Runs ≤ 15 Runtime Data Points



Empirical Cumulative Distribution



- the ECDF of run lengths to reach the target
- has for each data point a vertical step of constant size
- displays for each x-value (budget) the count of observations to the left (first hitting times)

Empirical Cumulative Distribution



- interpretations possible:
- 80% of the runs reached the target
- e.g. 60% of the runs need between 2000 and 4000 evaluations





50 equally spaced targets







the empirical **CDF** makes a step for each star, is monotonous and displays for each budget the fraction of targets achieved within the budget







15 runs



15 runs50 targets



15 runs 50 targets



15 runs 50 targets ECDF with 750 steps



50 targets from 15 runs ...integrated in a single graph

Interpretation



area over the ECDF curve = average log runtime (or geometric avg. runtime) over all targets (difficult and easy) and all runs

Fixed-target: Measuring Runtime



Fixed-target: Measuring Runtime

Idea: Simulated restarts by bootstrapping from measured runtimes until we see a success



Fixed-target: Measuring Runtime

Expected running time of the restarted algorithm:

$$E[RT^{r}] = \frac{1 - p_{s}}{p_{s}} E[RT_{unsuccessful}] + E[RT_{successful}]$$

Estimator average running time (aRT/ERT/Enes/SP2*, w/o proof):

$$\widehat{p_s} = \frac{\# \text{successes}}{\# \text{runs}}$$

 $\widehat{RT_{unsucc}}$ = Average evals of unsuccessful runs

 $\widehat{RT_{succ}}$ = Average evals of successful runs

$$aRT = \frac{\text{total #evals}}{\text{#successes}}$$

* The concept is so essential that it has been discovered/proposed multiple times under different names in the past.

© Anne Auger and Dimo Brockhoff, Inria 2024

ECDFs with Simulated Restarts

What we typically plot are ECDFs of the simulated restarted algorithms:



Measuring Performance

On

- real world problems
 - expensive
 - comparison typically limited to certain domains
 - · experts have limited interest to publish
- "artificial" benchmark functions
 - cheap
 - controlled
 - data acquisition is comparatively easy
 - problem of representativeness

Test Functions

• define the "scientific question"

the relevance can hardly be overestimated

- should represent "reality"
- are often too simple?

remind separability

- a number of testbeds are around
- account for invariance properties

prediction of performance is based on "similarity", ideally equivalence classes of functions

What to Benchmark?

Furious activity is no substitute for understanding (H.H. Williams)

Taking all possible functions from a repository?

Copyright: A. Auger and N. Hansen

What to Benchmark?

Furious activity is no substitute for understanding (H.H. Williams)

- Taking all possible functions from a repository?
- Bad idea if
 - function difficulties are unbalanced too many small dimensional problems, convex problems...
 - and performance are aggregated
- Leads to bias in the performance assessment

Copyright: A. Auger and N. Hansen

What to Benchmark?

- test functions should be representative of difficulties we want to test therefore NFL has no relevance as assumption of being closed under permutation has no relevance wrt real world problems
- related to real-word difficulties

for performance to be generalizable to RW

scalable

dimension plays a big role in performance curse of dimensionality

comprehensible but not too easy

BB optimization does not mean BB benchmarking

• we should still hide properties from the solver (hide optimum, ...) solvers should not be able to exploit the benchmark intentionally or not

Copyright: A. Auger and N. Hansen

automated benchmarking: COCO

Comparing Continuous Optimizers Platform https://coco-platform.org

COCO implements a reasonable, well-founded, and well-documented pre-chosen methodology







Experiments are simple (after pip install coco-experiment)

```
import cocoex
import scipy.optimize
### input
suite name = "bbob" # or "bbob-biobj" or ...
output folder = "scipy-optimize-fmin"
### prepare
suite = cocoex.Suite(suite name, "", "")
observer = cocoex.Observer(suite name,
                           "result_folder: " + output_folder)
### qo
for problem in suite: # this loop will take several minutes
   problem.observe with(observer) # generates the data for
                                    # cocopp post-processing
    scipy.optimize.fmin(problem, problem.initial solution)
```





https://coco-platform.org



data from 300+ algorithms can be accessed directly through their name

(see https://coco-platform.org/data-archive/)

How to benchmark algorithms with COCO? [in python]
https://coco-platform.org/getting-started/

installation I: experiments

As easy (in python) as:

pip install coco-experiment

installation II: postprocessing

pip install cocopp

for other languages, things might be slightly more involved...

Simplified Example Experiment in Python

```
import cocoex
import scipy.optimize
                                     coupling algo + COCO
### input
suite name = "bbob"
output folder = "scipy-optimize-fmin"
fmin = scipy.optimize.fmin
### prepare
suite = cocoex.Suite(suite name, "", "")
observer = cocoex.Observer(suite name,
                           "result_folder: " + output folder)
### ao
for problem in suite: # this loop will take several minutes
   problem.observe with(observer) # generates the data for
                                    # cocopp post-processing
    fmin(problem, problem.initial solution)
Note: the actual example experiment complete.py contains
more advanced things like restarts, batch experiments, other
algorithms (e.g. CMA-ES), etc.
```

https://coco-platform.org/getting-started/

running the experiment:

python example_experiments2.py 3

budget_multiplier: experiments will be run dimension times this number many fevals

tip:

start with small #funevals (until bugs fixed ©) then increase budget to get a feeling how long a "long run" will take

https://coco-platform.org/getting-started/

postprocessing

python -m cocopp YOURALGORITHMFOLDER ALG1 ALG2

ALG1, ALG2, ... can be chosen from any of the 300+ algorithms directly through its name (see https://coco-platform.org/data-archive/)



Automatically Generated Results

 $\overline{\bigcirc}$ index $\leftarrow \rightarrow C$ × +

0

A https://numbbo.github.io/ppdata-archive/bbob/2009/index.html

Home

Runtime profiles (dynamic navigation) Runtime profiles per function Runtime profiles summary and function groups Scaling with dimension Tables for selected targets

Runtime profiles over all targets



© Anne Auger and Dimo Brockhoff, Inria 2024

Available Test Suites in COCO

bbob 24 noiseless fcts 250+ algo data sets bbob-noisy 30 noisy fcts 40+ algo data sets bbob-biobj 55 bi-objective fcts 30+ algo data sets bbob-largescale 24 noiseless fcts 16 algo data sets bbob-mixint 24 mixed-int. fcts 5 algo data sets bbob-constrained 54 fcts w/ varying constr. 9 algo data sets bbob-boxed 3 algo data sets 24 box-constr. fcts

Worth to Note: ECDFs in COCO

In COCO, ECDF graphs

- never aggregate over dimension
 - but often over targets and functions
- can show data of more than 1 algorithm at a time



The single-objective BBOB functions

https://numbbo.github.io/gforge/downloads/download16.00/bbobdocfunctions.pdf

The bbob Testbed

24 functions in 5 groups:

1 Separable Functions		4 1	4 Multi-modal functions with adequate global structure	
f1	Sphere Function	f15	Rastrigin Function	
f2	Ellipsoidal Function	f16	Weierstrass Function	
f3	Rastrigin Function	f17	Schaffers F7 Function	
f4	Büche-Rastrigin Function	f18	Schaffers F7 Functions, moderately ill-conditioned	
f5	Linear Slope	f19	Composite Griewank-Rosenbrock Function F8F2	
2 Functions with low or moderate conditioning		5 Multi-modal functions with weak global structure		
f6	OAttractive Sector Function	f20	Schwefel Function	
f7	Step Ellipsoidal Function	f21	Gallagher's Gaussian 101-me Peaks Function	
f8	Rosenbrock Function, original	f22	Gallagher's Gaussian 21-hi Peaks Function	
f9	Rosenbrock Function, rotated	f23	GKatsuura Function	
3 F	unctions with high conditioning and unimodal	f24	SLunacek bi-Rastrigin Function	
f10	Ellipsoidal Function			
f11	ODiscus Function			
f12	Bent Cigar Function			
f13	Sharp Ridge Function			
f14	Different Powers Function			

6 dimensions: 2, 3, 5, 10, 20, (40 optional)

Notion of Instances

- All COCO problems come in form of instances
 - e.g. as translated/rotated versions of the same function
- Prescribed instances typically change from year to year
 - avoid overfitting
 - 5 instances are always kept the same

Notion of Instances



© Anne Auger and Dimo Brockhoff, Inria 2024

Notion of Instances

Plus:

 the bbob functions are locally perturbed by nonlinear transformations



Exercise

Objectives:

- investigate the performance of these 6 algorithms:
 - CMA-ES ("IPOP-CMA-ES" version)
 - CMA-ES ("BIPOP-CMA-ES" version)
 - Nelder-Mead simplex (use "NelderDoerr" version here)
 - BFGS quasi-Newton
 - Genetic Algorithm: discretization of cont. variables ("GA")
 - ONEFIFTH: (1+1)-ES with 1/5 rule
- pip install cocopp --pre
- python -m cocopp 2010/ipop! bipop! nelder! bfgs! 2009/GA ONEFIFTH!
- postprocessed data available here: <u>http://www.cmap.polytechnique.fr/~dimo.brock</u> <u>hoff/advancedOptSaclay/2024/</u>
- so now: investigate the data!

Exercise

Objective:

investigate the data:

- a) which algorithms are the best ones?
- b) does this depend on the dimension? Or on other things?
- c) look at single graphs: can we say something about the algorithms' invariances, e.g. wrt. rotations of the search space?
- d) what's the impact of covariance-matrix-adaptation?
- e) what do you think: are the displayed algorithms well-suited for problems with larger dimension?