



# A note on performance profiles for benchmarking software

NIM Gould, J Scott,

July 2015

Submitted for publication in ACM Transactions on Mathematical Software

RAL Library  
STFC Rutherford Appleton Laboratory  
R61  
Harwell Oxford  
Didcot  
OX11 0QX

Tel: +44(0)1235 445384  
Fax: +44(0)1235 446403  
email: [libraryral@stfc.ac.uk](mailto:libraryral@stfc.ac.uk)

Science and Technology Facilities Council preprints are available online  
at: <http://epubs.stfc.ac.uk>

**ISSN 1361- 4762**

Neither the Council nor the Laboratory accept any responsibility for loss or damage arising from the use of information contained in any of their reports or in any communication about their tests or investigations.

# A note on performance profiles for benchmarking software

Nicholas Gould<sup>1</sup> and Jennifer Scott<sup>1</sup>

## ABSTRACT

In recent years, performance profiles have become a popular and widely used tool for benchmarking and evaluating the performance of several solvers when run on a large test set. Here we use data from a real application as well as a simple artificial example to illustrate that caution should be exercised when trying to interpret performance profiles to assess the relative performance of the solvers.

**Keywords:** performance profiles, benchmarking software.

---

<sup>1</sup> Scientific Computing Department, STFC Rutherford Appleton Laboratory, Harwell Oxford, Oxfordshire, OX11 0QX, UK.  
nick.gould@stfc.ac.uk and jennifer.scott@stfc.ac.uk  
Supported by EPSRC grant EP/I013067/1.

June 24, 2015

# 1 Introduction to performance profiles

The quantities of data that results from benchmarking mathematical software (such as optimization packages or sparse linear solvers) with large problem sets have naturally led to researchers developing tools to analyse the data. A popular and widely used tool is the performance profile, which was proposed by Dolan and Moré in 2002 [1] as a means of providing objective information when benchmarking optimization software. Since their introduction, performance profiles have been used in many studies; as of June 2015, there were more than 1500 citations of the original paper [1] listed on Google Scholar.

Benchmark results are generated by running a solver on a set  $\mathcal{T}$  of problems and recording the information of interest (which might include, for example, the computation time, the number of function evaluations, the number of iterations or the memory used). Let  $\mathcal{S}$  represent the set of solvers that are to be compared. Suppose that a given solver  $i \in \mathcal{S}$  reports a statistic  $s_{ij} \geq 0$  when run on example  $j$  from the test set  $\mathcal{T}$ , and that the smaller this statistic the better the solver is considered to be. For  $j \in \mathcal{T}$ , let  $\hat{s}_j = \min\{s_{ij}; i \in \mathcal{S}\}$  and define  $r_{ij} = s_{ij}/\hat{s}_j$  to be the *performance ratio*<sup>1</sup>. Then for  $\omega \geq 1$  and each  $i \in \mathcal{S}$  define

$$k(r_{ij}, \omega) = \begin{cases} 1 & \text{if } r_{ij} \leq \omega \\ 0 & \text{otherwise.} \end{cases}$$

The *performance profile* of solver  $i$  is given by the function

$$p_i(\omega) = \frac{\sum_{j \in \mathcal{T}} k(r_{ij}, \omega)}{|\mathcal{T}|}, \quad \omega \geq 1.$$

Thus  $p_i(\omega)$  is the probability for solver  $i \in \mathcal{S}$  that a performance ratio  $r_{ij}$  is within a factor  $\omega$  of the best possible ratio. In particular,  $p_i(1)$  gives the fraction of the examples for which solver  $i$  is the winner (that is, the best according to the statistic  $s_{ij}$ ), while  $p_i^* := \lim_{\omega \rightarrow \infty} p_i(\omega)$  gives the fraction for which solver  $i$  is successful. If we are just interested in the number of wins, we need only compare the values of  $p_i(1)$  for all the solvers but, if we are interested in solvers with a high probability of success, we should choose those for which  $p_i^*$  is largest.

As many researchers have found, for a selected test set, performance profiles provide a very useful and convenient means of assessing the performance of a solver relative to the best solver on that set. When commenting on a performance profile presented in their paper, Dolan and Moré state that it “gives a clear indication of the relative performance of each solver” and they go on to say that “performance profiles provide an estimate of the expected performance difference between solvers”. Data from a practical study of solvers applied to a large test set and a simple artificial example will show that using performance profiles to assess the relative performance of the solvers should be undertaken with a degree of caution.

## 2 Example

We recently carried out a study to assess the performance of a number of sparse solvers (here denoted as diag, mi35 and ma97) on a set of 207 linear least squares problems; details may be found in [2]. One of the time performance profiles we obtained during the preliminary stages of our study is given in Figure 2.1(a). From this figure, it is clear that while it has the most failures (22 failures compared to 5 failures for solvers mi35 and diag), the solver ma97 has the highest number of wins (it is the fastest on 59% of the problems) and over our chosen range of  $\omega$  it dominates the other solvers, while the solver diag wins a respectable 34% of the time. Solver mi35 has the lowest number of wins and, if we are only interested in solvers that are within a factor 5 of the best, then it is tempting to conclude that, as the curve for solver mi35 lies below the other curves for  $\omega \in [1, 5]$ , it is the worst solver (see, for example, [3] where a similar conclusion is drawn). However, if we remove solver ma97 and redraw the performance profiles, we obtain Figure 2.1(b). We see that solver mi35 is the better of the remaining two solvers.

---

<sup>1</sup>If a solver  $i \in \mathcal{S}$  fails to solve problem  $j$ ,  $r_{ij} = \infty$ .

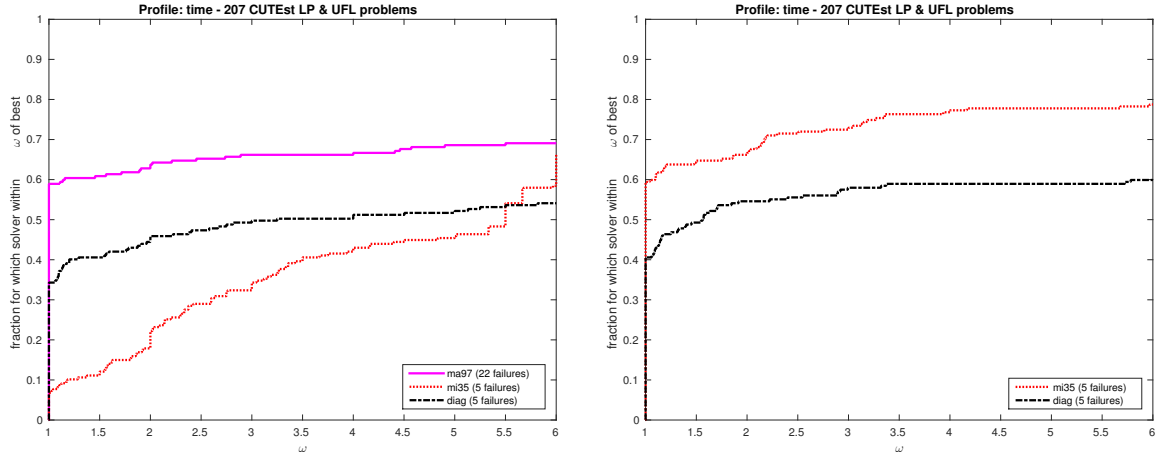


Figure 2.1: Time performance profiles for a real test case (a) with and (b) without the solver ma97.

This apparent change in fortunes can be seen clearly using the sample data given in Table 2.1 and corresponding performance profiles given in Figure 2.2. While Solver 1 is the best on 80% of the problems, Solver 2 is not the winner on any and, if we are interested in having a solver that can solve at least 60% of the problems with the greatest efficiency, then Solver 1 or 3 should be chosen. However, once Solver 1 is removed, Solver 2, which was second best on 60% of the problems, is now the best solver. In Figure 2.1(a) it is not apparent that solver mi35 is the second best solver.

Problem	Solver 1	Solver 2	Solver 3
1	2	1.5	1
2	1	1.2	2
3	1	4	2
4	1	5	20
5	2	5	20

Table 2.1: Performance of three solvers on a set of 5 problems; here, the smaller the statistic, the better the solver performance.

### 3 Conclusions

When comparing two solvers on a given test set, performance profiles give a clear measure of which is the better solver for a selected range of  $\omega$ . But as the examples above illustrate, if performance profiles are used to compare more than two solvers (and Dolan and Moré state that “performance profiles are most useful in comparing several solvers”), we can determine which solver has the highest probability  $p_i(\omega)$  of being within a factor  $\omega$  of the best solver for  $\omega$  in a chosen interval, but we cannot necessarily assess the performance of one solver relative to another that is not the best. In some situations, being able to rank (or partially rank) the solvers may be important. For example, a user may not have access to the best solver and so may want to know which is second (or perhaps third) best. To rank the solvers for a chosen range  $[1, \omega]$ , an obvious approach is to produce a series of performance profiles, excluding the best solver over the range from successive profiles until only two remain. We illustrate this in Figure 3.1, again using real data from our least squares study but now with a larger number of solvers and a larger range for  $\omega$ . Notice how that,

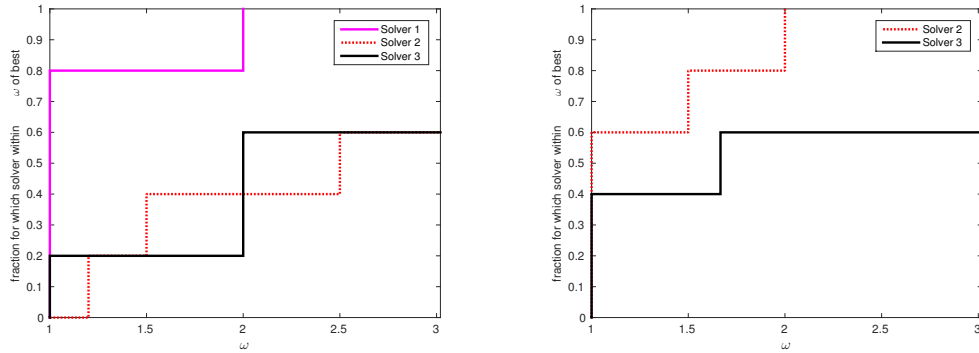


Figure 2.2: Performance profiles for small test case with and without Solver 1.

as before, removing the “leading” solver ma97 exposes mi35 as the runner up, and further removals illustrate that solver rif is higher in the performance hierarchy than the initial profiles might suggest.

A switch in the expected ordering may indicate the test set contains a large number of problems for which each solver performs in a consistent way and further examination of the test set and how it was selected may be advisable. However, our experience has been that, even without such a subset apparently present within the test set, switches can occur. We conclude that, while performance profiles are a powerful tool for benchmarking a solver relative to the best solver, as Dolan and Moré point out, “performance profiles must be used with care”.

## Acknowledgements

Many thanks to our colleagues in the Numerical Analysis Group at the Rutherford Appleton Laboratory for discussions on our findings and commenting on a draft of this note.

## References

- [1] E. D. DOLAN AND J. J. MORÉ, *Benchmarking optimization software with performance profiles*, Mathematical Programming, 91 (2002), pp. 201–213.
- [2] N. I. M. GOULD AND J. A. SCOTT, *The state-of-the-art of preconditioners for sparse linear least squares problems*, Technical Report, Rutherford Appleton Laboratory, 2015. In preparation.
- [3] N. J. HIGHAM, *The scaling and squaring method for the matrix exponential revisited*, SIAM Review, 51 (2009), pp. 747–767.

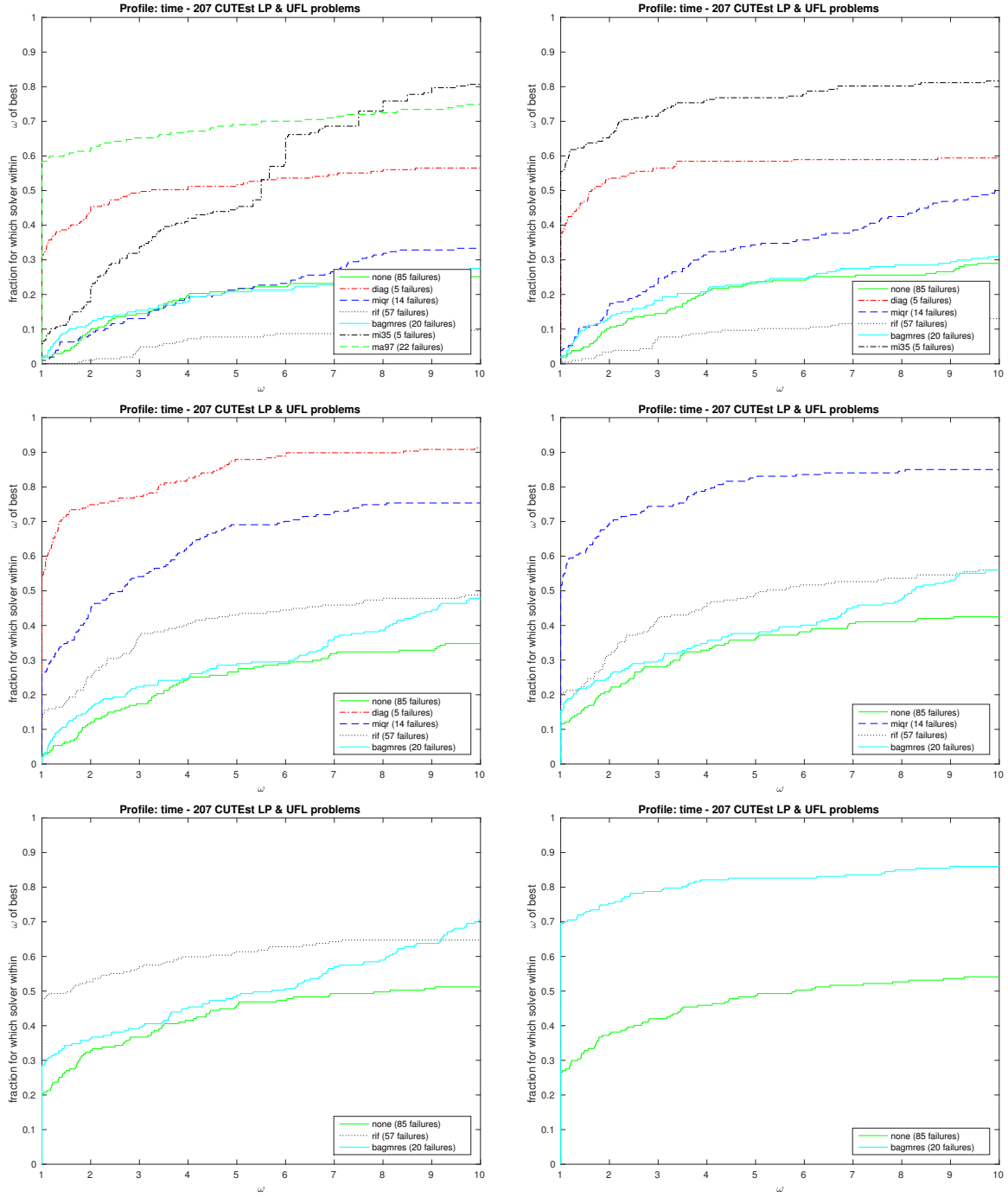


Figure 3.1: A sequence of time performance profiles for the real test case from Section 2 in which the “best” solver is removed until only two remain.