

On the empirical scaling of running time of d4 for solving random instances with $\mu = 1.3$

Empirical Scaling Analyzer

7th February 2022

1 Introduction

This is the automatically generated report on the empirical scaling of the running time of d4 for solving random instances with $\mu = 1.3$.

2 Methodology

For our scaling analysis, we considered the following parametric models:

- $Exp[a, b](n) = a \times b^n$ (2-parameter Exp)
- $Poly[a, b](n) = a \times n^b$ (2-parameter Poly)

Note that the approach could be easily extended to other scaling models. For fitting parametric scaling models to observed data, we used iteratively re-weighted linear least squares to perform quantile regression. Since this method works best for linear models, we used log transformations to convert non-linear models into linear models. This transformation biases the fitted models to more heavily favour the smaller training instance sizes, so we used a heuristic error correction term to compensate. Preliminary studies using this fitting method when applied to simulated running time datasets with known scaling properties show that using the heuristic error correction term improves the quality of the fitted models and allows the procedure to fit scaling models consistent with the true, underlying scaling of the data.

The fitted models correspond to performance predictions for the empirical scaling of the median of the distribution of running times. To assess the fit of a given scaling model to observed data, we used the mean absolute error as a loss function.

Closely following [1], we computed 95.0% bootstrap confidence intervals for the performance predictions obtained from our scaling models, based on 1001 bootstrap samples per instance set and 1001 automatically fitted variants of each scaling model. To extend this idea, we calculated training and challenge losses for each of the fitted models' predictions and the corresponding bootstrap samples of the observed data. We used these bootstrap sample losses to calculate median and 95% confidence intervals of the support and challenge losses for each model.

We calculated the observed point estimates for the medians of the data by fitting a linear model to local data with Gaussian weights, and then recording the observed statistic as the prediction from the linear model as the mid-point of the local data. In the following, we say that a scaling model prediction is in-consistent with observed data if the bootstrap confidence interval for the observed data is disjoint from the bootstrap confidence interval for the predicted median running times; we say that a scaling model prediction is weakly consistent with the observed data if the bootstrap confidence interval for

n	14	15	16
# instances	122	201	365
mean	1.329	1.808	2.772
Q(0.1)	0.6403	0.78	1.191
Q(0.25)	0.8504	1.1	1.64
median	1.17	1.521	2.39
Q(0.75)	1.65	2.15	3.44
Q(0.9)	2.09	2.96	4.98

n	17	18
# instances	441	508
mean	4.231	6.256
Q(0.1)	1.64	2.56
Q(0.25)	2.42	3.6
median	3.54	5.24
Q(0.75)	5.36	7.83
Q(0.9)	7.57	11.27

Table 1: Details of the running time dataset used as support data for model fitting. The “# of instances” is the sum of the weights of the instances used to calculate these statistics.

Figure 1: Fitted models of the medians of the running times. The models correspond to predictions for the medians of the running times of d4 solving the set of random instances with $\mu = 1.3$ with $11 \leq n \leq 18$ variables, and are challenged by the medians of the running times of $18 < n \leq 27$ variables.

the prediction overlaps with the bootstrap confidence interval for the observed data; and, we say that a scaling model is strongly consistent with observed data, if the bootstrap confidence interval for the observed median is fully contained within the bootstrap confidence interval for predicted running times. Also, we define the residue of a model at a given size as the observed point estimate less the predicted value using the fitted running time scaling model (fitted to the set of training data).

3 Dataset Description

The dataset contains running times of the d4 algorithm solving 5064 instances of different sizes. We split the running times into two categories, *support* or *training* instances ($n \leq 18$) and *challenge* or *test* instances ($n > 18$) with 1706 and 3358 instances, respectively. The details of the dataset can be found in Tables 1 and 2.

4 Empirical Scaling of Solver Performance

We first fitted our parametric scaling models to the medians of the running times of d4, as described in Section 2. The models were fitted using the training instance data and later challenged with the test instance data. This resulted in the models, shown along with losses on support and challenge data, shown in Table 3. In addition, we illustrate the fitted models of d4 in Figure 1, and the residues for the models in Figure 2.

n	19	20	21	22
# instances	557	520	588	621
mean	9.053	13.48	21.03	32.47
Q(0.1)	3.922	5.1	8.639	13.56
Q(0.25)	5.31	7.81	12.51	18.39
median	7.75	11.75	17.68	27.52
Q(0.75)	10.87	16.51	26.12	39.84
Q(0.9)	15.42	23.46	36.85	56.09

n	23	24	25
# instances	544	344	149
mean	43.59	65.06	76.91
Q(0.1)	18.57	27.9	38.83
Q(0.25)	25.57	39	53.03
median	36.79	56.8	70.29
Q(0.75)	51.85	81.1	93.62
Q(0.9)	80.98	114.7	126.4

Table 2: Details of the running time dataset used as challenge data for model fitting. The “# of instances” is the sum of the weights of the instances used to calculate these statistics.

	Model	Support loss	Challenge loss
d4	Exp. Model	$0.004660024 \times 1.477099^n$	1435.7
	Poly. Model	$4.69294 \times 10^{-8} \times n^{6.407237}$	25702

Table 3: Fitted models of the medians of the running times and loss values (in CPU sec). The models yielding the most accurate predictions (as per losses on challenge data) are shown in boldface.

But how much confidence should we have in these models? Are the losses small enough that we should accept them? To answer this question, we assessed the fitted models using the bootstrap approach outlined in Section 2. Table 4 shows the bootstrap intervals of the model parameters, Table 5 shows the bootstrap intervals of the model prediction losses, and Table 6 contains the bootstrap intervals for the support data. Challenging the models with extrapolation, as shown in Table 7, it is concluded that the Exp model tends to under-estimate the data, and the Poly model tends to under-estimate the data (as also illustrated in Figure 1). We base these statements on an analysis of the fraction of predicted bootstrap intervals that are strongly consistent, weakly consistent and disjoint from the observed bootstrap intervals for the challenge data. To provide stronger emphasis for the largest instance sizes, we also consider these fractions for the largest half of the challenge instance sizes. To be precise, ; and we say a model tends to under-estimate the data if $> 10\%$ of the confidence intervals for predictions on challenge instance sizes are disjoint from the confidence intervals for observed running time data and $\geq 90\%$ of the predicted intervals are below or are consistent with the observed intervals.

Figure 2: Residues of the fitted models of the medians of the running times.

Solver	Model	Confidence interval of a	Confidence interval of b
d4	Exp.	[0.002979635, 0.006767255]	[1.444915, 1.517095]
	Poly.	$[1.547496 \times 10^{-8}, 1.664548 \times 10^{-7}]$	[5.954207, 6.806386]

Table 4: 95% bootstrap intervals of model parameters for the medians of the running times

Solver	Model	Support Loss	Challenge Loss
d4	Exp.	[1352.3, 1518.4]	[22526, 24578]
	Poly.	[1353.9, 1520.5]	[24390, 28481]

Table 5: 95% bootstrap confidence intervals of model prediction losses for the medians of the running times.

Solver	n	Predicted confidence intervals	Observed median run-time	
		Exp. model	Point estimates	Confidence intervals
d4	14	[1.003, 1.17]	1.17	[1.08, 1.25]
	n	Predicted confidence intervals	Observed median run-time	
		Poly. model	Point estimates	Confidence intervals
	14	[0.9653, 1.121]	1.17	[1.08, 1.25]

Table 6: 95% bootstrap confidence intervals for the medians of the running time predictions and observed running times on random instances with $\mu = 1.3$. The instance sizes shown here are those used for fitting the models. Bootstrap intervals on predictions that are weakly consistent with the observed point estimates are shown in boldface and those that are strongly consistent are marked by asterisks (*).

5 Conclusion

In this report, we presented an empirical analysis of the scaling behaviour of d4 on random instances with $\mu = 1.3$. We found the Exp model tends to under-estimate the data, and the Poly model tends to under-estimate the data.

References

- [1] Yasha Pushak and Holger H. Hoos. Advanced Statistical Analysis of Empirical Performance Scaling. in *Proceedings of the 22nd Genetic and Evolutionary Computation Conference*, (GECCO 2020), pages 236–244, 2020.
- [2] Yasha Pushak, Zongxu Mu and Holger H. Hoos. Empirical scaling analyzer: An automated system for empirical analysis of performance scaling. *AI Communications* – to appear, 2020.
- [3] Jérémie Dubois-Lacoste, Holger H. Hoos, and Thomas Stützle. On the empirical scaling behaviour of state-of-the-art local search algorithms for the Euclidean TSP. In *Proceedings of the 17th Genetic and Evolutionary Computation Conference*, (GECCO 2015), pages 377–384, 2015.
- [4] Holger H. Hoos. A bootstrap approach to analysing the scaling of empirical run-time data with problem size. Technical report, Technical Report TR-2009-16, Department of Computer Science, University of British Columbia, 2009.

Solver	n	Predicted confidence intervals	Observed median run-time	
		Exp. model	Point estimates	Confidence intervals
Solver	n	Predicted confidence intervals	Observed median run-time	
		Poly. model	Point estimates	Confidence intervals

Table 7: 95% bootstrap confidence intervals for the medians of the running time predictions and observed running times on random instances with $\mu = 1.3$. The instance sizes shown here are larger than those used for fitting the models. Bootstrap intervals on predictions that are weakly consistent with the observed data are shown in boldface and those that are strongly consistent are marked by asterisks (*).

- [5] Holger H. Hoos and Thomas Stützle. On the empirical scaling of run-time for finding optimal solutions to the travelling salesman problem. *European Journal of Operational Research*, 238(1):87–94, 2014.