



| | |
|-------------------------------------|---|
| Title | Guidelines for defining benchmark problems in Genetic Programming |
| Authors(s) | Nicolau, Miguel, Agapitos, Alexandros, O'Neill, Michael, Brabazon, Anthony |
| Publication date | 2015-05-28 |
| Publication information | Nicolau, Miguel, Alexandros Agapitos, Michael O'Neill, and Anthony Brabazon. "Guidelines for Defining Benchmark Problems in Genetic Programming." IEEE, 2015. |
| Conference details | IEEE Congress on Evolutionary Computation, CEC 2015, Sendai, Japan, May 25-28, 2015, Proceedings, Sendai, Japan, May, 2015 |
| Publisher | IEEE |
| Item record/more information | http://hdl.handle.net/10197/8249 |
| Publisher's statement | © 2015 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works. |
| Publisher's version (DOI) | 10.1109/CEC.2015.7257019 |

Downloaded 2024-04-20 01:22:22

The UCD community has made this article openly available. Please share how this access benefits you. Your story matters! (@ucd_oa)



© Some rights reserved. For more information

Guidelines for defining benchmark problems in Genetic Programming

Miguel Nicolau and Alexandros Agapitos and Michael O’Neill and Anthony Brabazon

Natural Computing Research & Applications Group

Complex & Adaptive Systems Laboratory

University College Dublin, Ireland

Email: {miguel.nicolau, alexandros.agapitos, m.oneill, anthony.brabazon}@ucd.ie

Abstract—The field of Genetic Programming has recently seen a surge of attention to the fact that benchmarking and comparison of approaches is often done in non-standard ways, using poorly designed comparison problems. We raise some issues concerning the design of benchmarks, within the domain of symbolic regression, through experimental evidence. A set of guidelines is provided, aiming towards careful definition and use of artificial functions as symbolic regression benchmarks.

I. INTRODUCTION

The Genetic Programming research community has recently recognised the need for defining a suite of benchmark problems. So far, the effort has focussed on surveying the literature and choosing a set of problems depending on their popularity and difficulty [1], [2]. Clearly, the definition of a benchmark suite is an ongoing process – new problems should gain traction and will be added to sets of benchmarks. The aim of the present paper is to suggest guidelines for generating synthetic problems for symbolic regression tasks. Our work therefore has implications for the broader project of defining a detailed framework for benchmark creation.

In a symbolic regression problem, one is given a set of N training examples $\{(x_i, y_i)\}_1^N$, where $y \in \mathbb{R}$ is the response variable and $\mathbf{x} \in \mathbb{R}^d$ is a vector of explanatory variables. The goal is to find a function $F^* : \mathbb{R}^d \rightarrow \mathbb{R}$, such that over the joint distribution $P(\mathbf{x}, y)$ the expected value $\mathbb{E}_{x,y}$ of some specified loss function $L(y, F(\mathbf{x}))$ is minimised:

$$F^*(\mathbf{x}) = \arg \min_{F(\mathbf{x})} \mathbb{E}_{x,y} [L(y, F(\mathbf{x}))] \quad (1)$$

A typical loss function for such problems is *squared error*, where $L(y, F(\mathbf{x})) = (y - F(\mathbf{x}))^2$.

The motivation for working towards a framework for synthetic problem generation comes from the following observations of typical common practice:

- 1) The input range of synthetic functions is often arbitrarily chosen.
- 2) The training sample size N is often arbitrarily defined.
- 3) Function-set elements are often tailored for specific problems.
- 4) Lack of predefined and available training/test sets, which often hinders the reproducibility of results.

- 5) The lack of baselines for contrasting performance of evolved models.
- 6) The lack of noisy versions of the problems (real-world problems are inherently noisy).
- 7) In general, the constraint imposed by most learning algorithms can be described as *smoothness* restrictions of one kind or another. This essentially requires a regular behaviour of the underlying target function in small neighbourhoods of the input space. That is, for all input points sufficiently close to each other in some metric, the target function exhibits some behaviour that may be approximated using a constant, linear or low-order polynomial fit.

In this paper, we ran a set of experiments, to empirically test the effect of these observations in the performance of two typical Genetic Programming (GP) systems: Koza-like tree-based GP [3] and Grammatical Evolution (GE) [4]. The objective of these experiments is not to accurately measure and compare both systems, but rather to highlight the effect of neglecting some of the observations made above, when designing benchmark tests.

The next section presents an analysis of steps required when creating benchmarks. Section III presents the experimental setup of the problems and algorithms used, which are analysed in Section IV; finally, Section V draws conclusions and future work directions.

II. CREATING BENCHMARKS

A. Generating Datasets

When designing a benchmark problem using a known function, a dataset must be generated. This usually involves designing a function, choosing the range of the inputs, and generating a dataset.

Not all functions are suitable to be used as benchmarks, however. As an experiment, we generated datasets from the functions shown in Table I. Most of these were extracted from seminal works [5], [6], [7], [8], [9], suggested in the GP benchmarks effort [2]; functions $F14$, $F16$ and $F20$ were developed for this study. All functions were uniformly sampled from the input range $[-5 \dots 5]$ ¹.

¹This was done to deliberately highlight the problem of arbitrarily choosing input ranges; when designing benchmarks, input range sampling techniques (such as Latin Hypercube Sampling [10]) should be used.

TABLE I. SYMBOLIC REGRESSION PROBLEMS

| | |
|----------|---|
| F_1 | $f(x_1, x_2) = \frac{e^{-(x_1-1)^2}}{1.2+(x_2-2.5)^2}$ |
| F_2 | $f(x_1, x_2) = e^{-x_1} x_1^3 \cos(x_1) \sin(x_1) (\cos(x_1) \sin^2 x_1 - 1) (x_2 - 5)$ |
| F_3 | $f(x_1, x_2, x_3, x_4, x_5) = \frac{10}{5 + \sum_{i=1}^5 (x_i - 3)^2}$ |
| F_4 | $f(x_1, x_2, x_3) = 30 \frac{(x_1-1)(x_3-1)}{x_2^2(x_1-10)}$ |
| F_5 | $f(x_1, x_2) = 6 \sin(x_1) \cos(x_2)$ |
| F_6 | $f(x_1, x_2) = (x_1 - 3)(x_2 - 3) + 2 \sin((x_1 - 4)(x_2 - 4))$ |
| F_7 | $f(x_1, x_2) = \frac{(x_1-3)^4 + (x_2-3)^3 - (x_2-3)}{(x_2-2)^4 + 10}$ |
| F_8 | $f(x_1, x_2) = \frac{1}{1+x_1^{-4}} + \frac{1}{1+x_2^{-4}}$ |
| F_9 | $f(x_1, x_2) = x_1^4 - x_1^3 + x_2^2/2 - x_2$ |
| F_{10} | $f(x_1, x_2) = \frac{8}{2+x_1^2+x_2^2}$ |
| F_{11} | $f(x_1, x_2) = x_1^3/5 + x_2^3/2 - x_2 - x_1$ |
| F_{12} | $f(x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9, x_{10}) = x_1 x_2 + x_3 x_4 + x_5 x_6 + x_1 x_7 x_9 + x_3 x_6 x_{10}$ |
| F_{13} | $f(x_1, x_2, x_3, x_4, x_5) = -5.41 + 4.9 \frac{x_4 - x_1 + x_2/x_5}{3x_4}$ |
| F_{14} | $f(x_1, x_2, x_3, x_4, x_5, x_6) = (x_5 x_6) / (\frac{x_1}{x_2} \frac{x_3}{x_4})$ |
| F_{15} | $f(x_1, x_2, x_3, x_4, x_5) = 0.81 + 24.3 \frac{2x_2+3x_3^2}{4x_4^3+5x_5^4}$ |
| F_{16} | $f(x_1, x_2, x_3, x_4, x_5) = 32 - 3 \frac{\tan(x_1) \tan(x_3)}{\tan(x_2) \tan(x_4)}$ |
| F_{17} | $f(x_1, x_2, x_3, x_4, x_5) = 22 - 4.2(\cos(x_1) - \tan(x_2)) (\frac{\tanh(x_3)}{\sin(x_4)})$ |
| F_{18} | $f(x_1, x_2, x_3, x_4, x_5) = x_1 x_2 x_3 x_4 x_5$ |
| F_{19} | $f(x_1, x_2, x_3, x_4, x_5) = 12 - 6 \frac{\tan(x_1)}{e^{x_2}} (x_3 - \tan(x_4))$ |
| F_{20} | $f(x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9, x_{10}) = \sum_{i=1}^5 1/x_i$ |
| F_{21} | $f(x_1, x_2, x_3, x_4, x_5) = 2 - 2.1 \cos(9.8x_1) \sin(1.3x_5)$ |

Although using inputs sampled from the same range, the response variable for each function can have a completely different distribution. Fig. 1 shows an histogram of the response variable for each function defined in Table I, with 50000 samples drawn randomly from the input range specified. The responses were normalised to the range $[0 \dots 1]$, using the formula:

$$r_i = \frac{v_i - v_{\min}}{v_{\max} - v_{\min}}$$

where v_i is a response value, v_{\min} is the minimum response observed, v_{\max} is the maximum response observed, and r_i is the normalised response value.

The figure shows how varied the distribution of response values can be. Some functions are evenly distributed over the response domain, such as F_5 or F_{21} , while others present more challenging skewed distributions, such as F_8 or F_9 . Other functions present very challenging distributions, such as F_{18} , with a tail of outlier values spread upwards and downwards from the median, or even functions such as F_{16} , with only a handful of outliers, very distant and unevenly distributed from the median.

The latter two types of response distributions are very challenging as benchmarks, as they are very erratic, and unlikely to be found in real-world problems. To model such responses, a much larger sample set would be required (50000 samples were used to plot each function). Furthermore, the use of MSE as an error measure (as in the majority of evolutionary systems for symbolic regression) will result in very large errors, due to the extreme values found in the response variable. For example, the 50000 samples drawn from function F_{16} have a median of 32.0, with an interquartile range of 6.07, but minimum and maximum observed samples of $-6.527053E6$ and $2.50911394E8$ respectively.

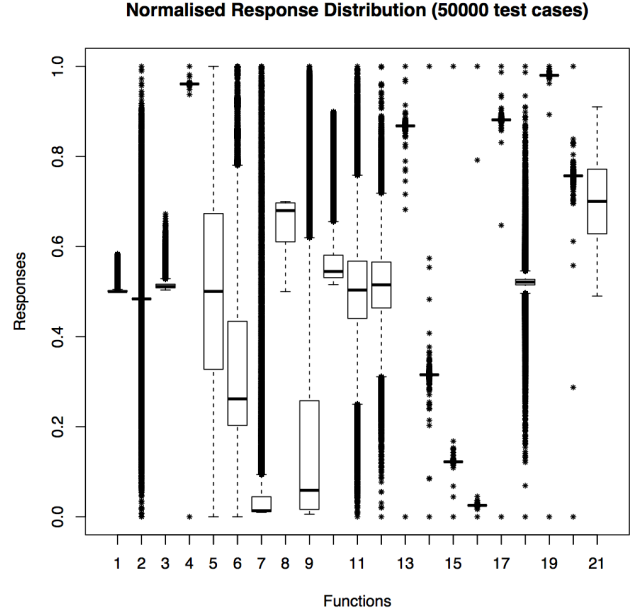


Fig. 1. Normalised response variable distribution over 50000 samples, for each function.

B. Sample Size

The distributions shown in Fig. 1 highlight how unevenly distributed a response variable can be, over the chosen input range. The size of the training and test samples must reflect this, and be adjusted accordingly.

Fig. 2 shows data generated in the same way as in Fig. 1, but with sample sizes of 50 (top) and 1000 (bottom), for training purposes. These figures show how the sample size can drastically affect the distribution of responses. Some functions like F_5 and F_{21} present essentially the same distribution across the response range. Functions such as F_{11} and F_{12} show essentially the same distribution with different sample sizes, but with a much higher density of outliers, making them very challenging. Finally, functions such as F_{13} to F_{20} show very radical changes in the distribution of responses, an indication both of ruggedness and insuitable small sample sizes.

C. Different Sample Sizes

The previous sections also highlight how the size of a drawn sample can paint a different picture of the behaviour of the underlying function. As such, when designing benchmarks, different sample sizes can be used, to control the difficulty of the problem. In this study, training sets of 50, 100, 500 and 1000 samples were used.

D. Using the same samples

The previous section also shows how different input sample sets can have very different response distributions. Depending on the ruggedness of the function, this can have the effect of substantially altering the performance of the learning method applied to a benchmark.

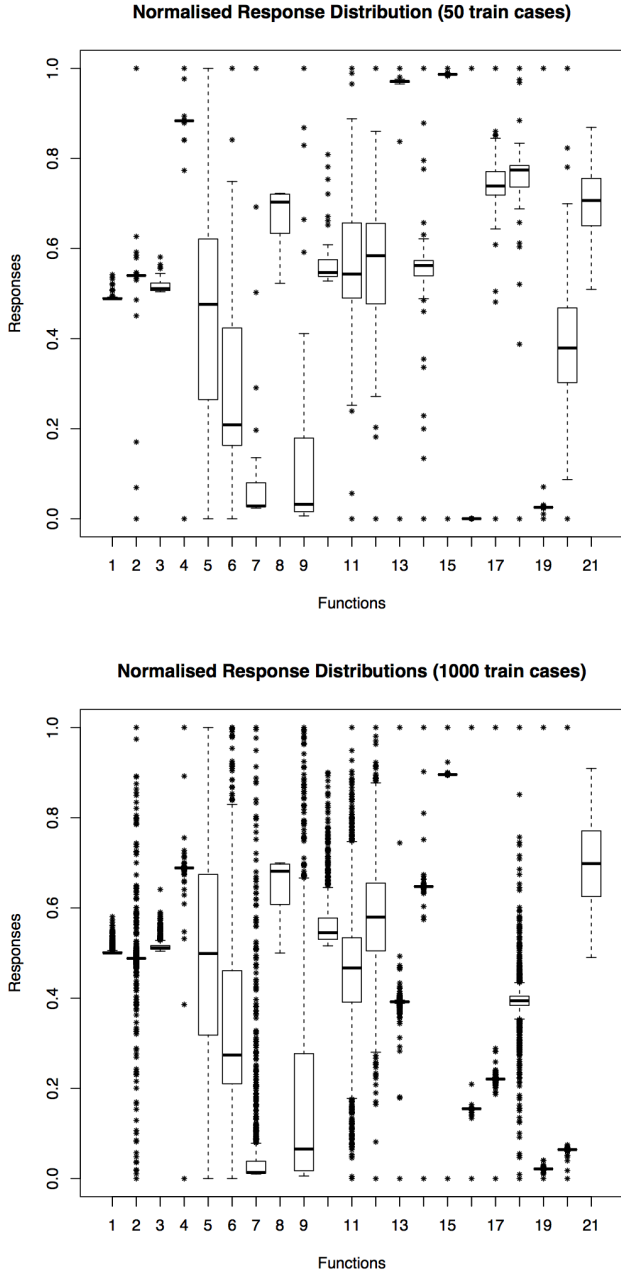


Fig. 2. Normalised response variable distribution over 50 (top) and 1000 (bottom) samples, for each function.

To exemplify this, 1000 sets of 50 samples each were created for both $F5$ and $F16$. Fig. 3 shows boxplots of the standard deviations of all 1000 sets, for both functions. This figure highlights the range of standard deviations that different samples can have; while most samples drawn from $F5$ exhibit a similar range, samples for $F16$ show an extreme variance in their range.

Optimally large sample sizes for each function can reduce this effect, but can also be computationally infeasible. Therefore different algorithms should be compared using the exact same samples, both for training and for test performance

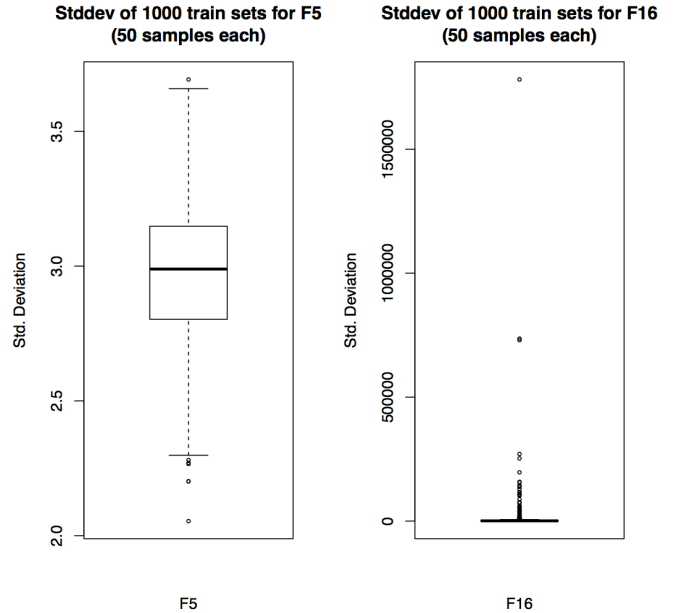


Fig. 3. Boxplots of distribution of standard deviation of response variable, for 1000 sets of 50 random samples each.

measurement (data can be easily made available online).

E. Large Test Set

This analysis of the variability of the response within small samples also raises questions about how to consistently measure the performance of an algorithm on unseen data. Based both on the range of the inputs and the variability of the response, a suitably large set of unseen samples should be used; in this study, a test set of 50000 unseen samples was used. As mentioned in the previous section, this sample should also always be the same, when comparing different algorithms.

F. Artificial Noise

The objective behind benchmarking a modelling algorithm is to simulate its application to real world data, which is inherently noisy. As such, benchmarks should provide datasets with noise added.

Given a synthetic function $F(x)$, noisy datasets are generated according to $y_i = F(x_i) + \epsilon_i$. We propose the addition of normally-distributed errors. In this case, the error ϵ_i is generated from a normal distribution with zero mean and variance adjusted [11] so that

$$\mathbb{E}|\epsilon| = \mathbb{E}_x|F(x) - \text{median}_x F(x)| \quad (2)$$

giving a 1/1 signal-to-noise ratio. For that, we set the standard deviation of our Gaussian noise σ as follows:

$$\sigma = \mathbb{E}_x|F(x) - \text{median}_x F(x)| / \sqrt{2.0/\pi} \quad (3)$$

Other types of error can be used, such as slash-distributed errors [11]. In this case, $\epsilon_i = s \cdot (u/v)$, where $u \sim N(0, 1)$ (normal distribution with zero mean and unit standard deviation) and $v \sim U[0, 1]$ (uniform distribution in $[0.0, 1.0]$). The scale factor is adjusted to give a 1/1 signal-to-noise ratio. This is performed as follows. We first generate a number of N ratios $\{r_1, \dots, r_N\}$, where $r_i = u/v$. We then calculate the standard deviation σ_r of the sample. Each ϵ_i is then set as:

$$\epsilon_i = r_i \cdot (\mathbb{E}_x |F(x) - \text{median}_x F(x)| / \sigma_r) \cdot \sqrt{\pi/2} \quad (4)$$

The slash distribution has very thick tails and is often used as extreme to test robustness.

G. Baselines

It is common practice in Evolutionary Computation to compare results using a baseline system. In many publications, this is usually one of the earliest systems developed: a simple GA [12] when using Genetic Algorithms, or Koza’s original system [3] when using Genetic Programming.

This however is not enough. GP is known to sometimes severely underperform, even in symbolic regression problems, and baselines such as a constant can sometimes outperform it [13]. Therefore a significant gain in precision when comparing to the original GA or GP systems is not necessarily evidence of good performance.

In this study, two baselines were used. One is a constant, which is just the average response observed in the training set; another is a linear regression model, applied to the training set. The results section shows how these two baselines can sometimes provide as good or better performance as GP.

III. EXPERIMENTS

A. Benchmark setup

Each of the functions listed in Table I was used, with eight versions: 4 training set sizes (50, 100, 500, 1000), and their equivalent with added noise, as discussed previously.

B. Run Parameters

We applied standard GP and GE systems to all benchmarks, along with constant and linear regression baselines. The experimental setup for GP and GE was very similar; both are shown in Table II. These are typical setups as seen in literature, consisting of the four arithmetic operators and five unary functions. Sub-tree crossover was used with GP, along with sub-tree mutation (sub-trees of depth d are replaced with a random sub-tree of depth between 1 and d); with GE, linear variable 1-point crossover was used, along with integer mutation (with l being the length of the individual). The constants used with GP were $[-0.9..-0.1] \cup [0.1..0.9]$, in steps of 0.1 (18 constants total); in the case of GE, digit concatenation [14] was used, such as in the following grammar (for a three input function):

```
<e> ::= + <e> <e> | - <e> <e>
      | * <e> <e> | / <e> <e>
```

TABLE II. EXPERIMENTAL SETUP

| Parameter | GP | GE |
|-----------------------------|---|-----------|
| Number of independent runs | | 100 |
| Total number of generations | | 50 |
| Population size | | 500 |
| Initialisation tree depth | | 5 |
| Maximum tree depth | 10 | unlimited |
| Tournament Size | 4 | 5 |
| Crossover rate | .9 | .5 |
| Mutation rate | .1 | 1/l |
| Number of elites | 1 | 50 |
| Function set | +,*,J,e ^x ,ln(x),√x,sin(x),tanh(x) | |

```
| exp <e> | ln <e> | sqrt <e>
| sin <e> | tanh <e>
| x\[0\] | x\[1\] | x\[2\]
| <nums>.<nums>
<nums> ::= <nums>0 | <nums>1
| <nums>2 | <nums>3
| <nums>4 | <nums>5
| <nums>6 | <nums>7
| <nums>8 | <nums>9
| 0 | 1 | 2 | 3 | 4
| 5 | 6 | 7 | 8 | 9
```

The initial populations were initialised using well-known techniques: ramped half-and-half for GP [3], and sensible initialisation for GE [15]. No wrapping operator was used with GE.

Some operators used protected versions. Division returned 1 for x/y if $y < 1e-5$; natural logarithm returned x if $x \leq 0$; and squared root returned x if $x < 0$.

IV. RESULTS AND ANALYSIS

A. Measuring test performance

GP systems tend to model input data quite well. This sometimes leads to overfitting of the training data. For this reason, training performance should never be reported as a means to measure system performance ².

In order to avoid overfitting the training data, some approaches use various approaches, such as validation sets. These should be drawn from the training data provided for the benchmark, and never from the test data.

B. Expected performance

Measuring the test performance of stochastic algorithms requires specific measurements. The GP community has moved on from reporting the performance of a single run of the system, and nowadays the mean best performance of a minimum of 30 independent runs is usually reported. This is done in order to measure the expected performance of an algorithm, and to test its statistical significance with other approaches.

Many publications report the performance of the model achieved by the “best run”, however (e.g. to compare its performance against other reported results). This is a mistake, as to choose a best run requires a comparison between all runs, making them no longer independent. If running GP to

²It can sometimes still be useful to report, however, to highlight the degree of overfitting (or lack of) in the system.

achieve a good model requires n runs, then a minimum of $n \times 30$ runs should be done, for statistical purposes.

For this reason, in this study, the median performance of all independent runs is reported, for GP and GE, in the (previously unseen) test set. This performance is then compared against the constant and linear regression models (which, being deterministic, require only one run). As a measure of variability of test performance between runs, the inter-quartile range (iqr) of the set of 30 runs is also reported, for GP and GE.

C. Results

Tables III to V show the results obtained. For each benchmark setup, the performance of GP, GE, a constant and Linear Regression are reported. The performance of GP and GE is reported in several ways: the median test performance of the models from all 100 runs, along with the iqr; the percentage of such models returning an infinity value in at least one test case; the median number of descents per run (i.e. number of times the training performance of the best model improved), along with the iqr; and the amount of training error decrease by the best model of a run, when compared to its initial population.

These results correlate quite well to the difficulty of the benchmarks, as seen in Fig. 1. Functions $F1$, $F3$, and $F13$ to $F20$ are extremely hard to model, with the given input range and training set sizes, and the test performance of all approaches varies from $1.0E4$ all the way to over $2.0E13$.

For the other functions, the GP approach generally provides a better performing median model when compared to the other methods, along with a smaller iqr than GE.

An interesting exception is $F21$ (also known as korns-12 [9]), which although well behaved in the response range, proved to be a very hard function to model. The best performance was achieved by both a constant and GE; an analysis of the GE model revealed that a constant had been evolved. Another challenging function was $F3$, for which GP and GE achieved a relatively low test error, but a linear regression model was a better choice³.

Linear regression and constants are in fact occasionally as good or better than the evolutionary models. Fig. 4 plots the performance improvement of the median GP model over the constant predictor, versus twice the standard deviation of the response variable in the test set, for all 21 functions. There is a clear relationship between variance of response variable and GP performance improvement; a notable exception is $F2$, a notoriously hard to model function.

In terms of training sample size, a larger sample almost always leads to better test performance of the median model. Training with noisy data also tends to produce less precise models, as expected. Notable exceptions are the performance of GP in $F1$, $F5$ and $F8$, where training with noisy data and a large sample size produces better performing median models; a probable explanation is that noise helps in preventing overfitting the training data.

³Note that both GP and GE achieved better performance than linear regression for this problem in several runs, but not in their median run.

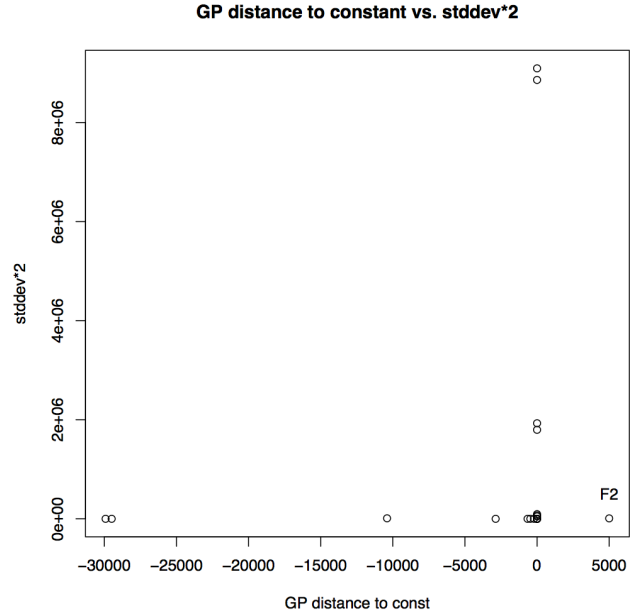


Fig. 4. GP performance difference to constant predictor versus double standard deviation of response variable in test set, for all 21 functions.

V. CONCLUSIONS & FUTURE WORK

This paper presented an experimental analysis of the difficulties in designing good benchmarks for Genetic Programming and similar systems. Several points were identified as crucial in designing artificial datasets from known functions, in the domain of symbolic regression. Crucially, the “smoothness” of synthetic regression problems is seldom studied in GP benchmarks, and the results obtained highlight this, with some results worse than those obtained by simple constant or linear scaling models. We do not imply that only smooth functions should be used, but if a synthetic problem is known to be highly non-smooth, it should be avoided.

This study is preliminary work, highlighting many of the issues in designing regression benchmarks, adding to those already raised previously [1]. Future work will continue this effort, addressing the raised issues and designing benchmark datasets, which will be available to the community. The issue of which metrics to report in the context of benchmarking will also be further investigated.

REFERENCES

- [1] J. McDermott, D. R. White, S. Luke, L. Manzoni, M. Castelli, L. Vanneschi, W. Jaskowski, K. Krawiec, R. Harper, K. D. Jong, and U.-M. O’Reilly, “Genetic programming needs better benchmarks,” in *Genetic and Evolutionary Computation - GECCO 2012, Genetic and Evolutionary Computation Conference, Philadelphia, USA, July 7-11, 2012, Proceedings*, T. S. et al., Ed. ACM, 2012, pp. 791–798.
- [2] D. R. White, J. McDermott, M. Castelli, L. Manzoni, B. W. Goldman, G. Kronberger, W. Jaskowski, U.-M. O’Reilly, and S. Luke, “Better gp benchmarks: community survey results and proposals,” *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*, vol. 14, no. 1, pp. 3–29, 2013.
- [3] J. R. Koza, *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. Cambridge, MA: MIT Press, 1992.

TABLE V. RESULTS - PART 3

| Data setup | GENETIC PROGRAMMING | | | | GRAMMATICAL EVOLUTION | | | | CONST. | L.R. |
|--------------------|----------------------|----------|------------------|------------------|-----------------------|----------|------------------|------------------|---------|---------|
| | Test MSE | ∞ | Descents per run | Decrease per run | Test MSE | ∞ | Descents per run | Decrease per run | | |
| Function 17 | | | | | | | | | | |
| T1000 | 5.808386E8 (1.31E3) | 4% | 17 (8) | 19.94% (2.62) | 5.8084E8 (4.51E2) | 0% | 22 (9) | 1.70% (2.66) | 5.81E8 | 5.81E8 |
| T500 | 5.808392E8 (7.10E3) | 7% | 17 (10) | 23.06% (3.74) | 5.8084E8 (9.55E2) | 0% | 22 (5) | 1.20% (5.36) | 5.81E8 | 5.81E8 |
| T100 | 5.808426E8 (2.34E6) | 20% | 17 (9) | 34.12% (7.74) | 5.8084E8 (1.19E4) | 2% | 22 (8) | 10.54% (14.31) | 5.81E8 | 5.81E8 |
| T50 | 5.809396E8 (1.80E8) | 48% | 20 (9) | 55.71% (13.24) | 5.8086E8 (5.94E5) | 5% | 25 (7) | 32.20% (23.02) | 5.81E8 | 5.81E8 |
| T1000-N | 5.808385E8 (1.30E3) | 3% | 16 (7) | 15.24% (2.89) | 5.8084E8 (8.20E2) | 0% | 22 (6) | 1.95% (1.47) | 5.81E8 | 5.81E8 |
| T500-N | 5.808393E8 (4.56E3) | 5% | 17 (8) | 18.49% (2.86) | 5.8084E8 (2.55E3) | 0% | 23 (7) | 3.26% (3.97) | 5.81E8 | 5.81E8 |
| T100-N | 5.808417E8 (1.26E5) | 8% | 19 (8) | 27.82% (5.68) | 5.8084E8 (6.11E4) | 0% | 22 (6) | 9.30% (13.93) | 5.81E8 | 5.81E8 |
| T50-N | 5.849402E8 (7.45E5) | 25% | 22 (7) | 45.35% (13.08) | 5.8093E8 (1.37E6) | 6% | 25 (6) | 28.39% (21.57) | 5.81E8 | 5.81E8 |
| Function 18 | | | | | | | | | | |
| T1000 | 1.086989E4 (3.69E4) | 1% | 24 (8) | 78.80% (89.13) | 4.1258E4 (4.62E4) | 0% | 20 (8) | 3.99% (82.03) | 4.08E4 | 4.08E4 |
| T500 | 2.459615E4 (4.93E4) | 5% | 26 (9) | 56.89% (86.52) | 4.2307E4 (3.49E4) | 0% | 20 (7) | 5.40% (4.15) | 4.08E4 | 4.10E4 |
| T100 | 6.713410E4 (4.55E5) | 3% | 26 (8) | 43.30% (49.03) | 4.8814E4 (1.65E5) | 1% | 18 (8) | 27.90% (15.63) | 4.08E4 | 4.25E4 |
| T50 | 5.094047E4 (1.62E6) | 3% | 24 (8) | 55.28% (48.85) | 4.9665E4 (3.76E5) | 0% | 20 (8) | 42.94% (23.00) | 4.09E4 | 4.13E4 |
| T1000-N | 1.373869E4 (4.01E4) | 1% | 25 (7) | 46.65% (56.20) | 4.1610E4 (5.55E4) | 0% | 18 (9) | 3.28% (62.06) | 4.08E4 | 4.08E4 |
| T500-N | 1.225567E4 (4.08E4) | 2% | 27 (8) | 57.70% (50.54) | 4.3329E4 (6.26E4) | 0% | 23 (6) | 4.69% (2.52) | 4.09E4 | 4.11E4 |
| T100-N | 6.779343E4 (4.16E5) | 4% | 26 (8) | 32.99% (13.18) | 5.1310E4 (8.09E4) | 0% | 19 (7) | 23.76% (10.07) | 4.08E4 | 4.30E4 |
| T50-N | 2.539613E5 (5.84E6) | 7% | 24 (8) | 43.21% (12.03) | 4.7245E4 (3.21E5) | 1% | 19 (7) | 33.42% (17.93) | 4.08E4 | 4.21E4 |
| Function 19 | | | | | | | | | | |
| T1000 | 8.072103E11 (3.25E5) | 5% | 22 (7) | 4.46% (2.16) | 8.0721E11 (1.62E5) | 0% | 24 (6) | 3.04% (2.26) | 8.07E11 | 8.07E11 |
| T500 | 8.072103E11 (4.49E5) | 5% | 23 (7) | 6.97% (1.85) | 8.0721E11 (2.66E5) | 0% | 21 (7) | 5.27% (2.06) | 8.07E11 | 8.07E11 |
| T100 | 8.072270E11 (2.43E8) | 8% | 23 (7) | 23.64% (7.40) | 8.0721E11 (1.24E8) | 2% | 19 (7) | 15.75% (9.42) | 8.07E11 | 8.07E11 |
| T50 | 8.072174E11 (1.34E8) | 11% | 23 (6) | 54.89% (9.31) | 8.0721E11 (8.59E6) | 4% | 20 (8) | 37.11% (22.43) | 8.07E11 | 8.07E11 |
| T1000-N | 8.072102E11 (2.15E5) | 4% | 23 (6) | 4.24% (1.64) | 8.0721E11 (1.82E5) | 0% | 24 (5) | 2.22% (1.45) | 8.07E11 | 8.07E11 |
| T500-N | 8.072102E11 (3.72E5) | 3% | 23 (7) | 5.48% (1.53) | 8.0721E11 (1.36E5) | 0% | 22 (7) | 3.69% (1.51) | 8.07E11 | 8.07E11 |
| T100-N | 8.072162E11 (2.31E9) | 9% | 26 (7) | 20.36% (5.84) | 8.0721E11 (1.25E7) | 7% | 21 (10) | 10.80% (7.64) | 8.07E11 | 8.07E11 |
| T50-N | 8.072163E11 (2.71E7) | 9% | 25 (7) | 44.05% (9.32) | 8.0721E11 (6.99E6) | 3% | 20 (7) | 26.73% (16.54) | 8.07E11 | 8.07E11 |
| Function 20 | | | | | | | | | | |
| T1000 | 2.680063E7 (3.13E4) | 5% | 19 (11) | 16.85% (20.19) | 2.6810E7 (8.30E3) | 0% | 19 (13) | 17.31% (12.98) | 2.68E7 | 2.68E7 |
| T500 | 2.680764E7 (5.65E3) | 2% | 18 (9) | 23.91% (13.47) | 2.6807E7 (3.11E3) | 0% | 18 (7) | 16.70% (10.12) | 2.68E7 | 2.68E7 |
| T100 | 2.681777E7 (1.93E5) | 7% | 19 (8) | 28.78% (12.93) | 2.6814E7 (2.50E4) | 0% | 17 (8) | 29.41% (13.02) | 2.68E7 | 2.68E7 |
| T50 | 2.682065E7 (8.78E4) | 6% | 19 (10) | 47.87% (16.37) | 2.6813E7 (1.08E5) | 0% | 14 (8) | 44.38% (16.41) | 2.68E7 | 2.68E7 |
| T1000-N | 2.680719E7 (2.66E4) | 5% | 19 (10) | 17.58% (13.21) | 2.6807E7 (3.81E3) | 0% | 20 (10) | 10.96% (7.64) | 2.68E7 | 2.68E7 |
| T500-N | 2.680746E7 (1.73E3) | 2% | 18 (8) | 18.90% (7.76) | 2.6807E7 (1.47E3) | 0% | 19 (7) | 13.13% (7.69) | 2.68E7 | 2.68E7 |
| T100-N | 2.682727E7 (8.35E5) | 8% | 19 (7) | 26.73% (13.07) | 2.6815E7 (4.69E4) | 0% | 17 (6) | 26.83% (12.05) | 2.68E7 | 2.68E7 |
| T50-N | 2.683119E7 (7.36E4) | 7% | 19 (5) | 40.66% (18.47) | 2.6843E7 (1.32E5) | 0% | 14 (7) | 39.72% (23.96) | 2.68E7 | 2.68E7 |
| Function 21 | | | | | | | | | | |
| T1000 | 1.074 (0.007) | 0% | 10 (5) | 5.45% (13.34) | 1.064 (0.000) | 0% | 14 (12) | 0.06% (0.05) | 1.064 | 1.071 |
| T500 | 1.088 (0.024) | 0% | 11 (5) | 4.63% (7.78) | 1.066 (0.000) | 0% | 21 (7) | 0.14% (0.34) | 1.066 | 1.085 |
| T100 | 1.204 (0.082) | 1% | 14 (5) | 22.44% (10.15) | 1.065 (0.000) | 0% | 15 (13) | 0.08% (0.45) | 1.065 | 1.099 |
| T50 | 1.364 (0.155) | 1% | 16 (5) | 39.92% (11.74) | 1.110 (0.404) | 4% | 16 (10) | 7.50% (22.33) | 1.071 | 1.273 |
| T1000-N | 1.078 (0.012) | 1% | 11 (6) | 2.86% (4.10) | 1.064 (0.000) | 0% | 18 (12) | 0.02% (0.03) | 1.064 | 1.074 |
| T500-N | 1.111 (0.036) | 1% | 11 (6) | 4.43% (4.12) | 1.067 (0.000) | 0% | 20 (10) | 0.03% (0.10) | 1.067 | 1.084 |
| T100-N | 1.305 (0.192) | 4% | 14 (5) | 17.64% (6.89) | 1.064 (0.054) | 0% | 14 (12) | 0.02% (3.89) | 1.064 | 1.080 |
| T50-N | 1.509 (0.317) | 8% | 14 (5) | 29.66% (5.82) | 1.064 (0.513) | 2% | 20 (7) | 0.36% (18.24) | 1.064 | 1.320 |

- [4] M. O'Neill and C. Ryan, *Grammatical Evolution - Evolutionary Automatic Programming in an Arbitrary Language*, ser. Genetic Programming. Kluwer Academic, 2003, vol. 4.
- [5] E. J. Vladislavleva, G. F. Smits, and D. den Hertog, "Order of non-linearity as a complexity measure for models generated by symbolic regression via pareto genetic programming," *IEEE Transactions on Evolutionary Computation*, vol. 13, no. 2, pp. 333–349, 2009.
- [6] L. Pagie and P. Hogeweg, "Evolutionary consequences of coevolving targets," *Evolutionary Computation*, vol. 5, no. 4, pp. 401–418, 1997.
- [7] M. Keijzer, "Improving symbolic regression with interval arithmetic and linear scaling," in *Genetic Programming, 6th European Conference, EuroGP 2003, Essex, UK, April 14-16, 2003, Proceedings*, ser. Lecture Notes in Computer Science, C. Ryan, T. Soule, M. Keijzer, E. Tsang, R. Poli, and E. Costa, Eds., vol. 2610. Springer, 2003, pp. 70–82.
- [8] R. Poli, "A simple but theoretically-motivated method to control bloat in genetic programming," in *Genetic Programming, 6th European Conference, EuroGP 2003, Essex, UK, April 14-16, 2003, Proceedings*, ser. Lecture Notes in Computer Science, C. Ryan, T. Soule, M. Keijzer, E. Tsang, R. Poli, and E. Costa, Eds., vol. 2610. Springer, 2003, pp. 204–217.
- [9] M. F. Korns, "Accuracy in symbolic regression," in *Genetic Programming Theory and Practice IX*, ser. Genetic and Evolutionary Computation, R. Riolo, E. Vladislavleva, and J. H. Moore, Eds. New York: Springer, 2011, pp. 129–151.
- [10] M. D. McKay, "Latin hypercube sampling as a tool in uncertainty analysis of computer models," in *Winter Simulation, 24th Conference, WSC 1992, Arlington, VA, USA, Proceedings*, J. J. Swain, D. Goldsman, R. C. Crain, and J. R. Wilson, Eds. ACM, 1992, pp. 557–564.
- [11] J. H. Friedman, "Greedy function approximation: A gradient boosting machine," *Annals of Statistics*, vol. 29, no. 5, pp. 1189–1232, 2000.
- [12] D. E. Goldberg, *Genetic Algorithms in Search, Optimization and Machine Learning*. Addison Wesley, 1989.
- [13] J. McDermott. (2014, March) Genetic programming needs better baselines. [Online]. Available: <http://jmmcd.net/2013/12/19/gp-needs-better-baselines.html>
- [14] M. Nicolau and I. Dempsey, "Introducing grammar based extensions for grammatical evolution," in *IEEE Congress on Evolutionary Computation, CEC 2006, Vancouver, BC, Canada, July 16-21, 2006, Proceedings*. IEEE Press, 2006, pp. 2663–2670.
- [15] C. Ryan and A. Azad, "Sensible initialisation in grammatical evolution," in *Genetic and Evolutionary Computation - GECCO 2003, Genetic and Evolutionary Computation Conference, Chicago, IL, USA, July 12-16, 2004, Workshops, Proceedings*, E. Cantú-Paz, J. A. Foster, K. Deb, L. Davis, R. Roy, U.-M. O'Reilly, H.-G. Beyer, R. K. Standish, G. Kendall, S. W. Wilson, M. Harman, J. Wegener, D. Dasgupta, M. A. Potter, A. C. Schultz, K. A. Dowsland, N. Jonoska, and J. F. Miller, Eds. AAAI, 2003.