

A Hitchhiker's Guide to Statistical Tests for Assessing Randomized Algorithms in Software Engineering¹

Andrea Arcuri and Lionel Briand

Simula Research Laboratory, P.O. Box 134, Lysaker, Norway.
Email: {arcuri,briand}@simula.no

Abstract

Randomized algorithms have been used to successfully address many different types of software engineering problems. This type of algorithms entail a significant degree of randomness as part of their logic. Randomized algorithms are useful to address difficult problems where a precise solution cannot be derived in a deterministic way within reasonable time. However, randomized algorithms can produce different results on every run when applied to the same problem instance. It is hence important to assess the effectiveness of randomized algorithms by collecting data from a large enough number of runs. The rigorous use of statistical tests is then essential to provide support to the conclusions derived by analyzing such data. In this paper, we provide a systematic review of the use of randomized algorithms in selected software engineering venues in 2009/2010. Its goal is not to perform a complete survey but to get a representative and up-to-date snapshot of current practice in software engineering research. We show that randomized algorithms are used in a significant percentage of papers but that, in most cases, randomness is not properly accounted for. This casts doubts on the validity of most empirical results assessing randomized algorithms for various applications. There are numerous statistical tests, based on different assumptions, and it is not always clear when and how to use these tests. We hence provide practical guidelines to support empirical research on randomized algorithms in software engineering.

Keyword: Statistical difference, effect size, parametric test, non-parametric test, confidence interval, Bonferroni adjustment, systematic review, survey.

1 Introduction

Many problems in software engineering can be alleviated through automated support. For example, automated techniques exist to generate test cases that satisfy some desired coverage criteria on the system under test, such as for example branch [48] and path coverage [42]. Because often these problems are undecidable, deterministic algorithms that are able to provide optimal solutions in reasonable time do not exist. The use of randomized algorithms [74] is hence necessary to address this type of problems.

The most well-known example of randomized algorithm in software engineering is perhaps *random testing* [26, 11]. Techniques that use random testing are of course randomized, as for example DART [42] (which combines random testing with symbolic execution). Furthermore, there is a large body of work on the application of *search algorithms* in software engineering [47], as for example Genetic Algorithms. Since all search algorithms are typically randomized and numerous software engineering problems can be addressed with search algorithms, randomized algorithms therefore play an increasingly important role. Applications of search algorithms include software testing [69], requirement engineering [16], project planning and cost estimation [2], bug fixing [12], automated maintenance [72], service-oriented software engineering [19], compiler optimisation [21] and quality assessment [56].

¹This paper is an extension of a conference paper [9] published in the International Conference on Software Engineering (ICSE), 2011.

A randomized algorithm may be strongly affected by chance. It may find an optimal solution in a very short time or may never converge towards an acceptable solution. Running a randomized algorithm twice on the same instance of a software engineering problem usually produces different results. Hence, researchers in software engineering that develop novel techniques based on randomized algorithms face the problem of how to properly evaluate the effectiveness of these techniques.

To analyze the cost and effectiveness of a randomized algorithm, it is important to study the *probability distribution* of its output and various performance metrics [74]. For example, a practitioner might want to know what is the execution time of those algorithms *on average*. But randomized algorithms can yield very complex and high variance probability distributions, and hence looking only at average values can be misleading, as we will discuss in more details in this paper.

The probability distribution of a randomized algorithm can be analyzed by running such an algorithm several times in an independent way, and then collecting appropriate data about its results and performance. For example, consider the case in which we want to find failures in software by using random testing (assuming that an automated oracle is provided). As a way to assess its cost and effectiveness, we can sample test cases at random until the first failure is detected. For example, in the first experiment, we might find a failure after sampling 24 test cases. We hence repeat this experiment a second time (if a pseudo-random generator is employed, we need to use a different seed for it) and then, for example, trigger the first failure when executing the second random test case. If in a third experiment we obtain the first failure after generating 274 test cases, the *mean* value of these three experiments would be 100. Using such a mean to characterize the performance of random testing on a set of programs would clearly be misleading given the extent of its variation.

Since such randomness might hinder the reliability of conclusions when performing the empirical analysis of randomized algorithms, researchers hence face two problems: (1) how many experiments should be run to obtain reliable results, and (2) how to assess in a rigorous way whether such results are indeed reliable. The answer to these questions lies in the use of *statistical tests* [86]. There are many books on various aspects of statistics (e.g., [86, 20, 60, 45, 102]), and that research field is still growing [102]. Notice that though statistical testing is used in most if not all scientific domains (e.g., medicine and behavioral science), each field has its own set of constraints to work with. Even within a field like software engineering the application context of statistical testing can vary significantly. When human resources and factors introduce randomness (e.g., [28, 52]) in the phenomena under study, the use of statistical tests is also required but the constraints we work with are quite different from those of randomized algorithms, such as for example the size of data samples and the types of distributions.

Because of the widely varying situations across domains and the overwhelming number of statistical tests, each one with its own characteristics and assumptions, many practical guidelines have been provided targeting different scientific domains, such as biology [77] and medicine [53]. There are also guidelines for running experiment with human subjects in software engineering [103]. In this paper, we intend to do the same for randomized algorithms in software engineering, as they entail specific issues and the application of statistical testing is far from easy, as we will see.

To assess whether the results obtained with randomized algorithms are properly analyzed in software engineering research, and therefore whether precise guidelines are required, we carried out a systematic review. We limited our analyses to the years 2009 and 2010, as our goal was not to perform an exhaustive systematic review but to obtain a recent, representative sample on which to draw conclusions about current practices. We focused on research venues that deal with all aspects of software engineering, such as IEEE Transactions of Software Engineering (TSE), IEEE/ACM International Conference on Software Engineering (ICSE) and International Symposium on Search Based Software Engineering (SSBSE). The former two are meant to get an estimate of the extent to which randomized algorithms are used in software engineering. The latter, more specialized venue provides us with additional insight into the way randomized algorithms are assessed in software engineering. The review shows that, in many cases, statistical analyses are either missing, inadequate, or incomplete. For example, though journal guidelines in medicine require a mandatory use of standardized *effect size* measurements [45] to quantify the effect of treatments, we have found only one case in which a standardized effect size was used to measure the relative effectiveness of a randomized algorithm [83]. Furthermore, in many of the surveyed empirical analyses, randomized algorithms were evaluated based on the results of only one run. Only few empirical studies reported the use of statistical analysis.

Given our survey's results, we hence found necessary to devise *practical* guidelines for the use of statistical

testing in assessing randomized algorithms in software engineering applications. Note that though guidelines have been provided for other scientific domains [77, 53] and for other types of empirical analyses in software engineering [28, 52], they are not necessarily applicable in the context of randomized algorithms. Our objective is therefore account for the specific properties of randomized algorithms in software engineering applications.

Notice that Ali *et al.* [3] have recently carried out a systematic review of search-based software testing which includes some limited guidelines on the use of statistical testing. This paper builds upon that work by: (1) analyzing software engineering as whole and not just software testing, (2) considering all types of randomized algorithms and not just search algorithms, and (3) giving precise, practical, and complete suggestions on many aspects related to statistical testing that were either not discussed or just briefly mentioned in [3].

The main contributions of this paper can be summarized as follows:

- We provide a systematic review of the current state of practice of the use of statistical testing to analyze randomized algorithms in software engineering. The review shows that randomness is not properly taken into account in the research literature.
- We provide practical guidelines on the use of statistical testing that are tailored to randomized algorithms in software engineering applications and the specific properties and constraints they entail.

The paper is organized as follows. Section 2 discusses a motivating example. The systematic review we carried out follows in Section 3. Section 4 presents the concept of statistical difference in the context of randomized algorithms. Section 5 compares two kinds of statistical tests and discussed their implications in our context. The problem of censored data and how it applies to randomized algorithms is discussed in Section 6. How to measure effect sizes and therefore the practical impact of randomized algorithms is presented in Section 7. Section 8 investigates the question of how many times randomized algorithms should be run. The problems associated with multiple tests is discussed in Section 9, whereas Section 10 deals with the choice of artifacts, which has usually a significant impact on results. Practical guidelines on how to use statistical tests in our context are summarized in Section 11. The threats to validity associated with our work are discussed in Section 12. Finally, Section 13 concludes the paper.

2 Motivating Example

In this section, we provide a motivating example to show why the use of statistical tests is a necessity in the analyses of randomized algorithms in software engineering. Assume that two techniques \mathcal{A} and \mathcal{B} are used in a type of experiment in which the output is binary: either *pass* or *fail*. For example, in the context of software testing, \mathcal{A} and \mathcal{B} could be testing techniques (e.g., random testing [26, 11]), and the experiment would determine whether they trigger or not any failure given a limited testing budget. The technique with highest *success rate*, that is failure rate in the testing example, would be considered to be superior. Further assume that both techniques are run n times, and a represents the times \mathcal{A} was successful, whereas b is the number of successes for \mathcal{B} . The *estimated* success rates of these two techniques are defined as a/n and b/n , respectively.

Now, consider that such experiment is repeated $n = 10$ times, and the results show that \mathcal{A} has a 70% estimated success rate, whereas \mathcal{B} has a 50% estimated success rate. Would it be safe to conclude that \mathcal{A} is better than \mathcal{B} ? Even if $n = 10$ and the difference in estimated success rates is quite large (i.e., 20%), it would actually be unsound to draw any conclusion about the respective performance of the two techniques. Because this might not be intuitive, we provide below the exact mathematical reasons to explain why that is the case.

A series of repeated n experiments with binary outcome can be described as a *binomial distribution* [31], where each experiment has probability p of success, and the mean value of the distribution (i.e., number of successes) is pn . In the case of \mathcal{A} , we would have an estimated success rate $p = a/n$ and an estimated number of successes $pn = a$. The probability mass function of a binomial distribution $B(n, p)$ with parameters n and p is:

$$P(B(n, p) = k) = \binom{n}{k} p^k (1 - p)^{n-k}.$$

$P(B(n, p) = k)$ represents the probability that a binomial distribution $B(n, p)$ would result in k successes. Exactly k runs would be successful (probability p^k) while the others $n - k$ would fail (probability $(1 - p)^{n-k}$).

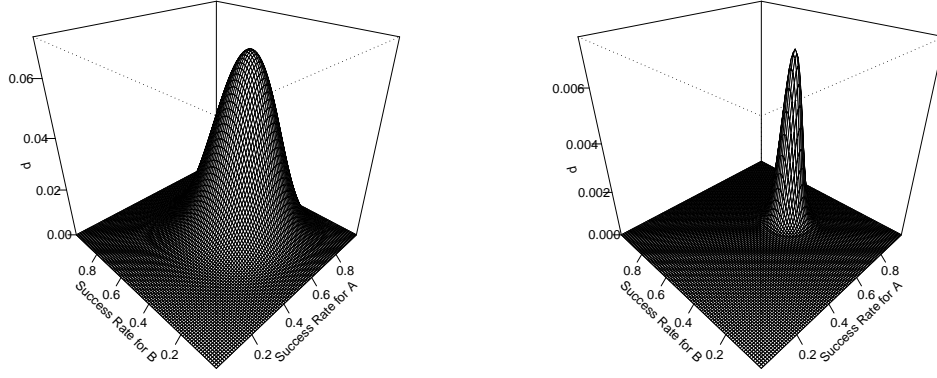


Figure 1: Probabilities to obtain $a = 0.7n$ and $b = 0.5n$ when $n = 10$ (left) and $n = 100$ (right) for different success rates of the algorithms \mathcal{A} and \mathcal{B} .

Since the order of successful experiments is not important, there are $\binom{n}{k}$ possible orders. Using this probability function, what is the probability that a equals the expected number of successes? Using our example, having a technique with an *actual* 70% success rate, what is the probability of having exactly 7 successes out of 10 experiments? This can be calculated with:

$$P(B(10,0.7) = 7) = \binom{10}{7} 0.7^7 (0.3)^3 = 0.26 .$$

This example shows there is only a 26% chance to have exactly $a = 7$ successes if the actual success rate is 70%! This shows a widespread misconception: expected values (e.g., successes) often have a relatively low probability of occurrence. Similarly, the probability that both techniques have a number of successes equal to their expected value would be even lower:

$$P(B(10,0.7) = 7) \times P(B(10,0.5) = 5) = 0.06 .$$

Reversely, even if we obtain $a = 7$ and $b = 5$, what would be the probability that both techniques have an equal actual success rate of 60%? We would have:

$$P(B(10,0.6) = 7) \times P(B(10,0.6) = 5) = 0.04 .$$

Though 0.04 seems a rather “low” probability, it is not much lower than 0.06, the probability of the observed number of successes to be actually equal to their expected values. Therefore, we cannot really say that the hypothesis of equal actual success rates (60%) is much more implausible than the one with 70% and 50% actual success rates. But what about the case where the two techniques have exactly the same actual success rate equal to 0.2? Or what about the cases in which \mathcal{B} would actually have a better actual success rate than \mathcal{A} ? What would be the probability for these situations to be true? Figure 1 shows all these probabilities, when $a = 0.7n$ and $b = 0.5n$, for two different numbers of runs: $n = 10$ and $n = 100$. For $n = 10$, there is a great deal of variance in the probability distribution of success rates. In particular, the cases in which \mathcal{B} has a higher actual success rate do not have a negligible probability. On the other hand, in the case of $n = 100$, the distribution variance has decreased significantly and high probabilities are all close to the expected average values (i.e., 70% for \mathcal{A} and 50% for \mathcal{B}). This clearly shows the importance of using sufficiently large samples, an issue we will get back to later in the paper.

In our example, with $n = 100$, the use of statistical tests (e.g., Fisher Exact test) would yield strong evidence to conclude that \mathcal{A} is better than \mathcal{B} . At an intuitive level, a statistical test would estimate the probability of mistakenly drawing the conclusion that \mathcal{A} is better than \mathcal{B} , under the form of a so-called p -value, as further discussed later in the paper. The resulting p -value would be quite small for $n = 100$ (i.e., 0.005), whereas for

Table 1: Number of publications grouped by venue, year and type.

Venue	Year	All	Regular	Randomized Algorithms
TSE	2009	48	48	3
	2010	48	48	12
ICSE	2009	70	50	4
	2010	111	54	10
SSBSE	2009	17	9	9
	2010	17	14	11
Total		311	223	49

$n = 10$ it would be much larger (i.e. 0.649), thus confirming and quantifying what is graphically visible in Figure 1. So even for what might appear to be large values of n , our capability to draw reliable conclusions could still be weak. Though some readers might find the above example rather basic, the fact of the matter is that many papers reporting on randomized algorithms overlook the principles and issues illustrated above.

3 Systematic Review

Systematic reviews are used to gather, in an unbiased and comprehensive way, published research on a specific subject and analyze it [54]. Systematic reviews are a useful tool to assess general trends in published research, and they are becoming increasingly common in software engineering [59, 28, 52].

In our review we want to analyze: (RQ1) how often randomized algorithms are used in software engineering, (RQ2) how many runs were used to collect data, and (RQ3) which types of statistical analyses were used to analyze these data.

To answer RQ1, we selected two of the main venues that deal with all aspects of software engineering: IEEE Transactions of Software Engineering (TSE) and IEEE/ACM International Conference on Software Engineering (ICSE). We also considered the International Symposium on Search-Based Software Engineering (SSBSE), which is a specialized venue devoted to search algorithms. Because our goal is not to perform an exhaustive survey of existing works, but simply to get an up-to-date snapshot of current practice regarding the application of randomized algorithms in software engineering research, we only considered 2009 and 2010 publications.

We only retained full length research papers and, as a result, 77 papers at ICSE and 11 at SSBSE were excluded. A total of 223 papers were considered: 96 in TSE, 104 in ICSE and 23 in SSBSE. These papers were manually checked to verify whether they made use of randomized algorithms, thus leading to a total of 49 papers. Table 1 summarizes the details of these publications divided by venue and year.

Notice that we excluded papers in which it was not clear whether randomized algorithms were used. For example, the techniques described in [50, 95] use external SAT solvers, and those might be based on randomized algorithms, though we cannot say for sure. Furthermore, we do not consider papers that involve *machine learning* algorithms that are randomized since they require different types of analysis [73]. On the other hand, if a paper focused on presenting a deterministic, novel technique, we included it when randomized algorithms were used for comparison purposes (e.g., fuzz testing [36]). Table 2 (for the year 2009) and Table 3 (for the year 2010) summarize the results of this systematic review for the final selection of 49 papers. The first clearly visible result is that randomized algorithms are widely used in software engineering (RQ1): we found them in 15% of the regular articles in TSE and ICSE, which are general-purpose and representative software engineering venues.

To answer RQ2, the data in Table 2 and Table 3 shows the number of times a technique was run to collect data regarding its performance on each artifact in the case study. Only 25 cases out of 49 show at least 10 runs. In many cases, data are collected from only one run of the randomized algorithms. Furthermore, notice that the case in which randomized algorithms are evaluated based on *only one run per case study artifact* is quite common in the literature. Even very influential papers such as DART [42] suffers of this problem, which poses a serious threat to the validity of those empirical analyses.

In the literature, there are empirical analyses in which randomized algorithms are run only once per case study artifact, but a large case study was generated at random (e.g., [78, 101]). The validity of such empirical analyses is questionable. However, the choice of a case study that is statistical relevant, and its relations with the needed number of runs for evaluating a randomized algorithm, needs proper care, and it will be discussed in more detail in Section 10.

Regarding RQ3, only 18 out of 49 articles include empirical analyses supported by some kind of statistical testing. More specifically, we can see t -tests, Welch and U-tests for when algorithms are compared in a pairwise fashion, whereas ANOVA and Kruskal-Wallis for multiple comparisons. Furthermore, in some cases linear regression is employed to build prediction models from a set of algorithm runs. However, in only one article [83] standardized *effect size* measures (see Section 7) are reported to quantify the relative effectiveness of algorithms.

Results in Table 2 and 3 clearly show that, when randomized algorithms are employed, empirical analyses in software engineering do not properly account for their random nature. Many of the novel proposed techniques may indeed be useful, but the results in Table 2 and 3 cast serious doubts on the validity of most existing results.

Notice that some of empirical analyses in Table 2 and 3 do not use statistical tests since they do not perform any comparison of the technique they propose with alternatives. For example, in the award winning paper at ICSE 2009, a search algorithm (i.e., Genetic Programming) was used and was run 100 times on each artifact in the case study [100]. However this algorithm was not compared against simpler alternatives or even random search. If we look more closely at the reported results in order to assess the implications of that lack of comparison, we see that the total number of fitness evaluations was 400 (a population size of 40 individuals that is evolved for 10 generations). This is an extremely low number (for example, for test data generation in branch coverage it is often the case of using 100,000 fitness evaluations for *each* branch [48]) and we can therefore conclude that there is very limited search taking place, which implies that a random search would have likely yielded similar results. This is directly confirmed in the reported results in [100], in which in half of the subject artifacts in the case study, the average number of fitness evaluations per run is at most 41, thus implying that, on average, appropriate patches are found in the random initialization of the first population before the actual evolutionary search even starts. This should not be surprising as the search operators were tailored to the specific, small set of bugs of the case study, which then led to an easy search problem. As discussed in [3], a search algorithm should always be compared against at least random search in order to check that the algorithm is not simply successful because the search problem is easy.

Since comparisons with simpler alternatives (at a very minimum random search) is a necessity when one proposes a novel randomized algorithm or addresses a new software engineering problem [3], statistical testing should be part of all publications reporting such empirical studies. In this paper we provide specific guidelines on how to use statistical tests to support comparisons among randomized algorithms.

4 Statistical Difference

When a novel randomized algorithm \mathcal{A} is developed to address a software engineering problem, it is common practice to compare it against existing techniques, in particular simpler alternatives. For simplicity, let us consider just one alternative randomized algorithm, and let us call it \mathcal{B} . For example, \mathcal{B} can be random testing, and \mathcal{A} can be a search algorithm such as Genetic Algorithms or an hybrid technique that combines symbolic execution with random testing (e.g., DART [42]).

To compare \mathcal{A} versus \mathcal{B} , we first need to decide which criteria are used in the comparisons. Many different measures (M), either attempting to capture the effectiveness or the cost of algorithms, can be selected depending on the problem at hand and contextual assumptions, e.g., source code coverage, execution time. Depending on our choice, we may want to either minimize or maximize M , for example maximize coverage and minimize execution time.

To enable statistical analysis, we should run both \mathcal{A} and \mathcal{B} a large enough number (n) of times, in an independent way, to collect information on the probability distribution of M for each algorithm. A *statistical test* should then be used to assess whether there is enough empirical evidence to claim, with a high level of confidence, that there is a difference between the two algorithms (e.g., \mathcal{A} is better than \mathcal{B}). A *null hypothesis* H_0 is typically defined to state that there is no difference between \mathcal{A} and \mathcal{B} . A statistical test is used to verify whether we should reject the null hypothesis H_0 . However, what aspect of the probability distribution of M

Table 2: Results of systematic review for the year 2009.

Reference	Venue	Repetitions	Statistical Tests
[1]	TSE	1/5	U-test
[68]	TSE	1	None
[78]	TSE	1	None
[71]	ICSE	100	t -test, U-test
[100]	ICSE	100	None
[36]	ICSE	1	None
[57]	ICSE	1	None
[7]	SSBSE	1000	Linear regression
[40]	SSBSE	30/500	None
[27]	SSBSE	100	U-test
[39]	SSBSE	50	None
[61]	SSBSE	10	Linear regression
[55]	SSBSE	10	None
[67]	SSBSE	1	None
[58]	SSBSE	1	None
[92]	SSBSE	1	None

is being compared depends on the used statistical test. For example, a t -test compares the mean values of two distributions whereas others tests focus on the median or proportions, as discussed in Section 5.

There are two possible types of error when performing statistical testing: (I) we reject the null hypothesis when it is true (we are claiming that there is a difference between two algorithms when actually there is none), and (II) we accept H_0 when it is false (there is a difference but we claim the two algorithms to be equivalent). The p -value of a statistical test denotes the probability of a Type I error. The *significant level* α of a test is the highest p-value we accept for rejecting H_0 . A typical value, inherited from widespread practice in natural and social sciences, is $\alpha = 0.05$.

Notice that the two types of error are conflicting; minimizing the probability of one of them necessarily tends to increase the probability of the other. But traditionally there is more emphasis on not committing a Type I error, a practice inherited from natural sciences where the goal is often to establish the existence of a natural phenomenon in a conservative manner. In our context we would only conclude that an algorithm \mathcal{A} is better than \mathcal{B} when the probability of a Type I error is below α . The price to pay for a small α value is that, when the data sample is small, the probability of a Type II error can be high. The concept of statistical *power* [20] refers to the probability of rejecting H_0 when it is false (i.e., the probability of claiming statistical difference when there is actually a difference).

Getting back to our comparison of techniques \mathcal{A} and \mathcal{B} , let us assume we obtain a p-value equal to 0.06. Even if one technique seems significantly better than the other in terms of effect size (Section 7), we would then conclude that there is no difference when using the traditional $\alpha = 0.05$ threshold. In software engineering, or in the context of *decision-making* in general, this type of reasoning can be counter-productive. The tradition of using $\alpha = 0.05$, discussed by Cowles [22], has been established in the early part of the last century, in the context of natural sciences, and is still applied by many across scientific fields. It has, however, an increasing number of detractors [43, 44] who believe that such thresholds are arbitrary, and that researchers should simply report p -values and let the readers decide in context what risks they are willing to take in their decision-making process.

When we need to make a choice between techniques \mathcal{A} and \mathcal{B} , we would like to use the one that is more likely to outperform the other. Whether we get a p-value lower than α bears little consequence from a practical standpoint, as in the end we *must* select an alternative, e.g., we must select a testing technique to verify the system. However, as we will show in Section 8, obtaining p-values lower than $\alpha = 0.05$ should not be a problem when experimenting with randomized algorithms. The focus of such experiments should rather be on whether a given technique brings any practically significant advantage, usually measured in terms of an estimated effect size and its confidence interval, an important concept addressed in Section 7.

Table 3: Results of systematic review for the year 2010.

Reference	Venue	Repetitions	Statistical Tests
[38]	TSE	1000	None
[108]	TSE	100	t -test
[48]	TSE	60	U-test
[83]	TSE	32	U-test, \hat{A}_{12}
[25]	TSE	30	Kruskal-Wallis, undefined pairwise
[94]	TSE	20	None
[18]	TSE	10	U-test, t -test, ANOVA
[29]	TSE	3	U-test
[6]	TSE	1	None
[14]	TSE	1	None
[17]	TSE	1	None
[101]	TSE	1	None
[62]	ICSE	100	None
[109]	ICSE	50	None
[41]	ICSE	5	None
[75]	ICSE	5	None
[35]	ICSE	1	None
[46]	ICSE	1	None
[51]	ICSE	1	None
[106]	ICSE	1	None
[80]	ICSE	1	None
[90]	ICSE	1	None
[23]	SSBSE	100	t -test
[24]	SSBSE	100	None
[66]	SSBSE	50	t -test
[70]	SSBSE	50	U-test
[105]	SSBSE	30	U-test
[107]	SSBSE	30	t -test
[63]	SSBSE	30	Welch
[98]	SSBSE	30	ANOVA
[15]	SSBSE	3/5	None
[65]	SSBSE	3	None
[110]	SSBSE	1	None

In practice, the selection of an algorithm would depend on the p-value of effectiveness comparisons, the effectiveness effect size, and the cost difference among algorithms (e.g., in terms of user-provided inputs or execution time). Given a context-specific decision model, the reader, using such information, could then decide which technique is more likely to maximize benefits and minimize risk. In the simplest case where compared techniques would have comparable costs, we would simply select the technique with the highest effectiveness regardless of the p-values of comparisons, even if as a result there is a non-negligible probability that it will bring no particular advantage.

When one has to carry out a statistical test, one must choose between *one-tailed* and a *two-tailed* test. Briefly, in a two-tailed test, one would reject H_0 if the performance of \mathcal{A} and \mathcal{B} are different no matter of which one is the best. On the other hand, in a one-tailed test, one is making assumptions about the relative performance of the algorithms. For example, one could expect that a new sophisticated algorithm \mathcal{A} is better than a naive algorithm \mathcal{B} used in the literature. In such a case, one would detect statistically significant difference when \mathcal{A} is indeed better than \mathcal{B} , but ignoring the “unlikely” case of \mathcal{B} being better than \mathcal{A} . An historical example in the literature of statistics is the test to check whether there is the right percent of gold (carats) in coins. One could expect that a dishonest coiner might produce coins with lower percent of gold than declared, and so a one-tailed test would be used rather than a two-tailed. Such a test could be used if one wants to verify whether the coiner is actually dishonest, whereas giving more gold than declared would be very unlikely. Using a one-tailed test has the advantage, compared to a two-tailed test, that the resulting p-value is lower (so it is easier to detect statistically significant differences).

Are there cases in which a one-tailed test could be advisable in the analysis of randomized algorithms in software engineering? As a rule of thumb, we say no: two-tailed tests should be used. One should use a one-tailed test only if (s)he has strong arguments to support such a decision. In fact, most of the time we cannot make any assumption on the relative performance of randomized algorithms. Even naive testing techniques such as random testing can be better than more sophisticated techniques on some classes of problems (e.g., see [91]). If one wants to lower the p-values, it is recommended to increase the number of runs (see Section 8) rather than using an arguable one-tailed test.

5 Parametric vs Non-Parametric Tests

In our context, the two most used statistical tests are the t -test and the Mann-Whitney U-test. These tests are in general used to compare two data samples (e.g., the results of running n times algorithm \mathcal{A} compared to \mathcal{B}). The t -test is *parametric*, whereas the U-test is *non-parametric*.

A parametric test makes assumptions on the underlining distribution of the data. For example, the t -test assumes normality and equal variance of the two data samples. A non-parametric test makes no assumption about the distribution of the data. *Why* is there the need for two different types of statistical tests? A simple answer is that, in general, non-parametric tests are less powerful than parametric ones when the latter’s assumptions are fulfilled. When, due to cost or time constraints, only small data samples can be collected, one would like to use the most powerful test available if its assumptions are satisfied.

There is a large body of work regarding which of the two types of tests should be used [30]. The assumptions of the t -test are in general not met. Considering that the variance of the two data samples is most of the time different, a Welch test should be used instead of a t -test. But the problem of the normality assumption remains.

An approach would be to use a statistical test to assess whether the data is normal, and, if the test is successful, then use a Welch test. This approach increases the probability of Type I error but is often not necessary. In fact, the Central Limit theorem tells us that, for large samples, the t -test and Welch test are robust even when there is strong departure from a normal distribution [86, 89]. But in general we cannot know how many data points (n) we need to reach reliable results. A rule of thumb is to have at least $n = 30$ for each data sample [86].

There are three main problems with such an approach: (1) if we need to have a large n for handling departures from normality, then it might be advisable to use a non-parametric test since, for a large n , it is likely to be powerful enough; (2) the rule of thumb $n = 30$ stems from analyses in behavioral science, and, to the best of our knowledge, there is no supporting evidence of its efficacy for randomized algorithms in software engineering; (3) the Central Limit theorem has its own set of assumptions, which are too often ignored. We now discuss points (2) and (3) in more details by accounting for the specific properties of the application of

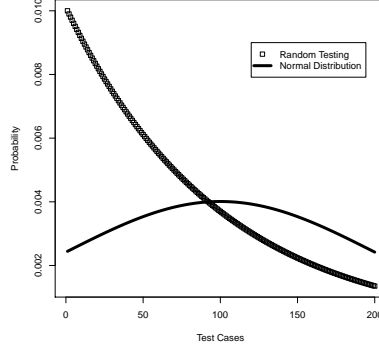


Figure 2: Density functions of random testing and normal distribution given same mean $\mu = 1/\theta$ and variance $\sigma^2 = (1 - \theta)/\theta^2$, where $\theta = 0.01$.

randomized algorithms in software engineering, using software testing examples. This choice was motivated by the fact that half the publications in search-based software engineering are on software testing [47].

Random testing, when used to find a test case for a specific testing target (e.g., a test case that triggers a failure or covers a particular branch/path) follows a geometric distribution. When there is more than one testing target, e.g., full structural coverage, it follows a coupon’s collector problem distribution [11]. Given θ the probability of sampling a test case that covers the desired testing target, then the expectation of random testing is $\mu = 1/\theta$ and its variance is $\delta^2 = (1 - \theta)/\theta^2$ [31]. Figure 2 plots the density function of a geometric distribution with $\theta = 0.01$ and a normal distribution with same μ and δ^2 . In this context, the density function represents the probability that, for a given number of sampled test cases l , we cover the target after sampling exactly l test cases. For random testing, the most likely outcome is $l = 1$, whereas for a normal distribution it is $l = \mu$. Notice that the geometric distribution is discrete (i.e., it is defined only on integer values), whereas a normal distribution is continuous. Furthermore, the density function of the normal distribution is always positive for any value, whereas for the geometric distribution it is equal to 0 for negative values, where in this context the values are the number of sampled test cases. Therefore, a testing technique can *never* follow a normal distribution in a strict way, although it might be a reasonable approximation.

As it is easily visible from Figure 2, the geometric distribution has a very strong departure from normality! Comparisons of novel techniques versus random testing (and this is the practice when search algorithms are evaluated [47]) using t -tests are hence very arguable. In general, in contrast to many physical and behavioral phenomena, in terms of their effectiveness, the probability distributions for search algorithms may strongly depart from normality. A common example is when the search landscape of the addressed problem has trap-like regions [79].

The Central Limit theorem states that the *sum* of n random variables converges to a normal distribution [31] as n increases. For example, consider the result of throwing a dice. There are only six possible outcomes, each one with probability $1/6$. If we consider the *sum* of two dice (i.e., $n = 2$), we have 11 possible outcomes, from value 2 to 12. Figure 3 shows that with $n = 2$, in the case of dice, we already obtain a distribution that resembles the normal one, even though with $n = 1$ it is very far from normality. In our context, these random variables are the results of the n runs of the analyzed algorithm. This theorem has three assumptions: the n variables should be independent and their mean μ and variance δ^2 should exist (i.e., they should be different from infinity). When using randomized algorithms, having n independent runs is usually trivial to achieve (we just need to use different seeds for the pseudo-random generators). But the existence of the mean and variance requires more scrutiny. As shown before, those values μ and δ^2 exist for random testing. A well known “paradox” in statistics in which mean and variance do not exist is the Petersburg Game [31]. Similarly, the existence of mean and variance in search algorithms is not always guaranteed, as discussed next.

To put this discussion on a more solid ground, let us briefly describe the Petersburg Game. Assume a player tosses an unbiased coin until a head is obtained. The player first gives an amount of money to the opponent which needs to be negotiated, and then she receives from the opponent an amount of money (Kroner) equal to $k = 2^t$, where t is the number of times the coin was tossed. For example, if the player obtains two tails and

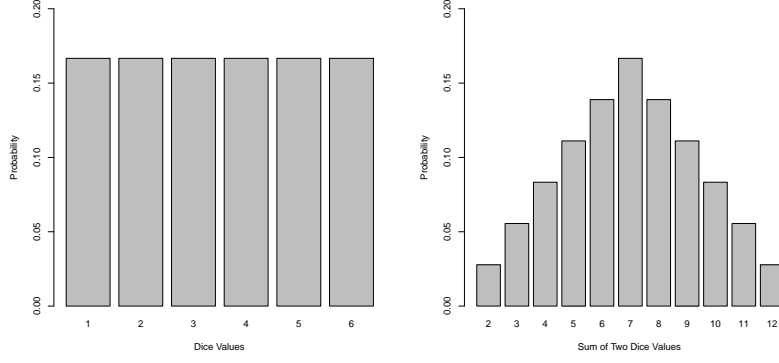


Figure 3: Density functions of the outputs of one dice and the sum of two dice.

then a head, then she would receive from the opponent $k = 2^3 = 8$ Kroner. *On average*, how many Kroner k will she receive from the opponent in a single match? The probability of having $k = 2^x$ is equivalent to get first $x - 1$ tails and then one head, so $p(2^x) = 2^{-(x-1)} \times 2^{-1} = 2^{-x}$. Therefore, the average reward is $\mu = E[k] = \sum_k k p(k) = \sum_t 2^t p(2^t) = \sum_t 2^t \times 2^{-t} = \sum_t 1 = \infty$. Unless the player gives an *infinite* amount of money to the opponent before starting tossing the coin, then the game would not be fair *on average* for the opponent! This is a classical example of a random variable where it is not intuitive to see that it has no finite mean value. For example, obtaining $t > 10$ is very unlikely, and if one tries to repeat the game n times, the average value for k would be quite low and would be a very wrong estimate of the actual, theoretical average (infinity).

Putting the issue illustrated by the Petersburg Game principle in our context, if the performance of a randomized algorithm is bounded within a predefined range, then the mean and variance always exist. For example, if an algorithm is run for a predefined amount of time to achieve structural test coverage, and there are z structural targets, then the performance of the algorithm would be measured with a value between 0 and z . Therefore, we would have $\mu \leq z$ and $\delta^2 \leq z^2$, thus making the use of a t -test valid.

The problems arise if no bound is given on how the performance is measured. A randomized algorithm could be run until it finds an optimal solution to the addressed problem. For example, random testing could be run until the first failure is triggered (assuming an automated oracle is provided). In this case, the performance of the algorithm would be measured in the number of test cases that are sampled before triggering the failure and there would be no upper limit for a run. If we run a search algorithm on the same problem n times, and we have n variables X_i representing the number of test cases sampled in each run before triggering the first failure, we would estimate the mean with $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i$, and hence conclude that the mean exists. As the Petersburg Game shows, this can be wrong, because $\hat{\mu}$ is only an *estimation* of μ , which might not exist.

For most search algorithms convergence in finite time is proven under some conditions (e.g., [87]), and hence mean and variance exist. But in software engineering, when new problems are addressed, standard search algorithms with standard search operators may not be usable. For example, when testing for object-oriented software using search algorithms (e.g., [97]), complex non-standard search operators are required. Without formal proofs, it is not safe to speak about the existence of the mean in those cases.

However, the non-existence of the mean is usually not a problem from a practical standpoint. In practice, usually there are upper limits to the amount of computational resources a randomized algorithm can use. For example, a search algorithm can be prematurely stopped when reaching a time limit. Random testing could be stopped after 100,000 sampled test cases (for example) if it has found no failure so far. But in these cases, we are actually dealing with *censored* data [60] (in particular, right-censorship) and this requires proper care in terms of statistical testing and the interpretation of results, as discussed in Section 6.

Even under proper conditions for using a parametric test, one aspect that is often ignored is that t -test and U-test are two different approaches to analyze two different properties. Let us use a random testing example in which we count the number of test cases run before triggering a failure. Considering a failure rate θ , the mean value of test cases sampled by random testing is hence $\mu = 1/\theta$. Let us assume that a novel testing technique \mathcal{A} yields a normal distribution of the required number of test cases to trigger a failure. If we further consider the same variance as random testing and a mean that is 85% of that of random testing, which one is better? Random

testing with mean μ or \mathcal{A} with mean 0.85μ ? Assuming a large number of runs (e.g., n is equal to one million), a t -test would state that \mathcal{A} is better, whereas a Mann-Whitney U-test would state exactly the opposite. How come? This is not an error but the two tests are measuring different things: The t -test measures the difference in mean values whereas the Mann-Whitney U-test deals with their stochastic ranking, i.e., whether observations in one data sample are more likely to be larger than observations in the other sample. Notice that this latter concept is technically different from detecting difference in *median* values (which can be stated only if the two distributions have same shape). In a normal distribution, the median value is equal to the mean, whereas in a geometric distribution the median is roughly 70% of the mean [31]. On one hand, half of the data points for random testing would be lower than 0.7μ . On the other hand, with \mathcal{A} we have half of the data points above 0.85μ , and a significant proportion between 0.7μ and 0.85μ . This explains the apparent contradiction in results: though the average is higher for random testing, its median is lower than that of \mathcal{A} .

From a practical point of view, which statistical test should be used? Based on the discussions in this section, and in line with [64], we suggest to use Mann-Whitney U-test (to assess difference in stochastic order) rather than the t -test and Welch test (to assess difference in mean values). However, the full motivation will become clearer once we discuss censored data, effect size, and the choice of n in the next sections.

At any rate, there is an important aspect that needs to be considered: data can be “transformed” before given as input to a statistical test. As discussed in [88], a Welch test can be used instead of a U-test if the data are replaced by their rank. For example, consider the data set $\{24, 2, 274\}$ discussed in the introduction regarding random testing. Those values could be transformed into their ranks $\{2, 1, 3\}$ before given as input to a statistical test. What would be the motivation of doing so? The U-test might be negatively affected if the two compared distributions have “significantly” different variance, and in such case a Welch test on ranked data might be better (in the sense that it would have lower probability of Type I and II errors). However, the Welch test would still be negatively affected by violations of normality assumption (ranked data might not be normal). Ruxton [88] reports on some cases in which a Welch test on ranked data is better than a U-test, but the results of those *empirical* analyses might not generalize to the case of randomized algorithms applied to software engineering problems.

For simplicity and because it has widespread applications, we recommend to use a U-test rather than a Welch test on ranked data. There might be cases in which this latter test could be preferable, but it might be difficult, for a non-expert in statistics, to clearly identify those cases. Nevertheless, it is important to clarify that a Welch test on ranked data does not assess any more whether there is a statistical difference among the mean values of the two compared distributions. It assesses differences in mean values of the ranks and therefore determine whether there is any difference in stochastic ordering in the two distributions. For example, assume the two data sets $X = \{1, 2, 3, 4, 5, 6, 49\}$ and $Y = \{7, 8, 9, 10, 11, 12, 13\}$. If it were not for the “outlier” 49 in X , then all the values in Y would be greater than the values in X . Both data sets have mean value 10. A Welch test on the raw values would result in p-value equal to 1, which is not surprising considering that the two data sets have same mean. However, if we do a rank transformation, then the outlier 49 would be replaced by the value 14 (all the other values in X and Y remain the same). In this case, the resulting p-value of the Welch test would be 0.02, which suggests a strong difference in the stochastic ordering (i.e., ranks) between the two distributions.

In the discussion above, we have assumed that both \mathcal{A} and \mathcal{B} are randomized. If one of them is deterministic (e.g., \mathcal{B}), it is still important to use statistical testing. Consistent with the above recommendation, the non-parametric *One-Sample Wilcoxon* test should be used. Given $m_{\mathcal{B}}$ the performance measure of the deterministic algorithm, a one-sample Wilcoxon test would verify whether the performance of \mathcal{A} is symmetric about $m_{\mathcal{B}}$, i.e., whether by using \mathcal{A} one is as likely to obtain a value lower than $m_{\mathcal{B}}$ as otherwise.

6 Censored Data

Assume that the result of an experiment is dichotomous: either we find a solution to solve the software engineering problem at hand (*success*), or we do not (*failure*). For example, in software testing, if our goal is to cover a particular target (e.g., a specific branch), we can run a randomized algorithm with a time limit L . We will stop the algorithm as soon as we find a solution, otherwise we stop it after time L . The choice of L depends on the available computational resources. Another example is bug fixing [100] where we find a patch within time L , or we do not.

These types of experiments are dealing with *right-censored* data, and their properties are equivalent to survival/failure time analysis [60, 34]. Let X be the random variable representing the time a randomized algorithm takes to solve a software engineering problem, and let us consider n experiments in which we collect X_i values. We are dealing with right-censorship since, assuming a time limit L , we will not have observations X_i for the cases $X > L$. There are several ways to deal with this problem [60] and we will limit our discussion to simple solutions.

One interesting special case is when we cannot say for sure whether we have achieved our target, e.g., generation of test cases that achieve code branch coverage. Putting aside trivial cases, there are usually infeasible targets (e.g., unreachable code) and their number is unknown. As a result, such experiments are not dichotomous because we cannot know whether we have covered all feasible targets. Even when using a time limit L , in these cases we are not tackling censored data. However, if in the experiments the comparisons are made reusing artifacts from published studies in the literature, and if we want to know whether or not, within a given time, we can obtain better coverage than these reported studies, then such experiments can be considered dichotomous despite infeasible targets.

Let us consider the case in which we need to compare two randomized algorithms \mathcal{A} and \mathcal{B} on a software engineering problem with dichotomous outcome. Let X be the random variable representing the time \mathcal{A} takes to find a valid solution, and let Y be the same type of variable for \mathcal{B} . Let us assume that we run \mathcal{A} n times collecting observations X_i , and we do the same for \mathcal{B} . Using a time limit L , to evaluate which of the two algorithms is better, we can consider their *success rate* $\gamma = k/n$, i.e., the proportion of number of times k out of the n runs in which they find a valid solution. To evaluate whether there is statistical difference between the success rates of \mathcal{A} and \mathcal{B} , a test for differences in proportions is then appropriate, such as the Fisher exact test [60].

The Fisher exact test is a parametric test, which assumes that the analyzed data follows a binomial distribution. In contrast to other parametric tests (e.g., the t -test), its assumptions are always valid: the experiments are independent, then the success rate of a series of randomized experiments would always follow a binomial distribution, where γ represents the estimated probability of success. Furthermore, for values of n until roughly 100, the test is “exact”, because all the assumptions of the Fisher test are met. This means that the resulting p-values are precise, and not estimates based on how close the data is from satisfying the conditions of a test (e.g., normality and equal variance in a t -test). However, for larger values of n , the computational cost of the test would start to be too prohibitive, and approximations are then used to calculate the p-values.

Assume that out of $n = 100$ runs the success rate of \mathcal{A} is $\gamma_{\mathcal{A}} = 1\%$, whereas for \mathcal{B} we have $\gamma_{\mathcal{B}} = 5\%$. A Fisher exact test has a resulting p-value equal to 0.21, which might be considered high, i.e., there is a 21% probability that the success rates of the two algorithms are actually equal. In such cases, one can run more experiments (i.e., increase n) to obtain higher statistical power (i.e., decrease the p-value). Alternatively, if there is no statistically or practically significant difference between the success rates of \mathcal{A} and \mathcal{B} , a practical question is then to determine which technique uses *less* time. This is particularly relevant if the success rates of both techniques are high. There can be different ways to analyze such cases, such as considering artificial censorships at different times before L . For example, we can consider censorship at $L/2$, i.e., the success rate with half the time, and determine which technique still fares better and at an acceptable level. Note that such analysis does not require to run any further experiments as success rates can be computed at $L/2$ from existing runs. Another alternative to compare execution times is to apply a Mann-Whitney U-test, recommended above, using only the times of successful runs, which have X_i and Y_i values lower or equal to L .

A more complex situation is when one algorithm shows a significantly higher success rate, but takes more time to produce valid solutions than the other. This is a typical situation, that is not so uncommon, where a choice needs to be made. For example, on one hand, a *local search* [69] might be very fast in generating appropriate testing data if it starts from the right area of the search landscape. But, at the same time, it could yield a low success rate if most of the search landscape has gradient toward local optima, and if the number of these local optima is low. (Notice that this is just an example: it is not in the scope of the paper to give lengthy explanations of why that would be a problem for local search, see [8] for further details on this topic.) On the other hand, a population-based search algorithm, such as Genetic Algorithms, could avoid the problem of local optima, which in turn would result in higher γ than a local search. However, because an entire population is evolved at the same time, depending on the selection pressure of the algorithm (e.g., the value of the tournament size in tournament selection) and the population size, a Genetic Algorithm might take much longer than a local

search in its successful runs.

7 Effect Size

When comparing a randomized algorithm \mathcal{A} against another \mathcal{B} , given a large enough number of runs n , it is most of the time possible to obtain statistically significant results with a t -test or U-test. Indeed, two different algorithms are extremely unlikely to have exactly the same probability distribution. In other words, with large enough n we can obtain statistically difference even if that difference is so small as to be of no practical value.

Though it is important to assess whether an algorithm fares statistically better than another, it is in addition crucial to assess the magnitude of the improvement. To analyze such a property, *effect size* measures are needed [45, 52, 77]. In their systematic review of empirical analyses in software engineering, Kampenes *et al.* [52] found out that standardized effect sizes were reported in only 29% of the cases. In our review, we found only one [83], which uses the Vargha and Delaney’s \hat{A}_{12} statistics (which will be described later in this section).

Effect sizes can be divided in two groups: standardized and unstandardized. Unstandardized effect sizes are dependent from the unit of measurement used in the experiments. Let us consider the difference in mean between two algorithms $\Delta = \mu^{\mathcal{A}} - \mu^{\mathcal{B}}$. This value Δ has a measurement unit, that of \mathcal{A} and \mathcal{B} . For example, in software testing, μ can be the expected number of test executions to find the first failure. On one testing artifact we might have $\Delta_1 = \mu^{\mathcal{A}} - \mu^{\mathcal{B}} = 100 - 1 = 99$, whereas on another testing artifact we might have $\Delta_2 = \mu^{\mathcal{A}} - \mu^{\mathcal{B}} = 100,000 - 200,000 = -100,000$. Deciding based on Δ_1 and Δ_2 which algorithm is better is difficult to determine since the two scales of measurement are different. Δ_1 is very low compared to Δ_2 , but in that case \mathcal{A} is 100 times worse than \mathcal{B} , whereas it is only twice as fast in the case Δ_2 . Empirical analyses of randomized algorithms, if they are to be reliable and generalizable, require the use of large numbers of artifacts (e.g., programs). The complexity of these artifacts is likely to widely vary, such as the number of test cases required to fulfill a coverage criterion on various programs. The use of standardized effect sizes, that are independent from the evaluation criteria measurement unit, is therefore necessary to be able to compare results across artifacts and experiments.

In this section we first describe which is the most known standardized effect size measure and why it should *not* be used. We then describe two other standardized effect sizes, and how to apply them in practice. The most known effect size is the so called d family which, in the general form, it is $d = (\mu^{\mathcal{A}} - \mu^{\mathcal{B}})/\sigma$. In other words, the difference in mean is scaled over the standard deviation (several corrections exists to this formula, but for more details please see [45]). Though we obtain a measure that has no measurement unit, the problem is that it assumes normality of the data, and strong departures can make it meaningless [45]. For example, in a normal distribution, roughly 64% of the points lie within $\mu \pm \sigma$ [31], i.e., they are at most σ away from the mean μ . But for distributions with high skewness (as in the geometric distribution and as it is often the case for search algorithms), the results of scaling the mean difference by the standard deviation “would not be valid”, because “standard deviations can be very sensitive to a distribution’s shape” [45]. In this case, a non-parametric effect size should be preferred. Existing guidelines in [52, 77] briefly discuss the use of non-parametric effect sizes.

The Vargha and Delaney’s \hat{A}_{12} statistic is a non-parametric effect size measure [99, 45]. Its use has been advocated in [64], and one example of its use in software engineering in which randomized algorithms are involved can be found in [83]. In our context, given a performance measure M , \hat{A}_{12} measures the probability that running algorithm \mathcal{A} yields higher M values than running another algorithm \mathcal{B} . If the two algorithms are equivalent, then $\hat{A}_{12} = 0.5$. This effect size is easier to interpret compared to the d family. For example, $\hat{A}_{12} = 0.7$ entails we would obtain higher results 70% of the time with \mathcal{A} . Though this type of non-parametric effect size is not common in statistical tools, it can be very easily computed [64, 45]. The following formula is reported in [99]:

$$\hat{A}_{12} = (R_1/m - (m+1)/2)/n \quad (1)$$

where R_1 is the rank sum of the first data group we are comparing. For example, assume the data $X = \{42, 11, 7\}$ and $Y = \{1, 20, 5\}$. The data set X would have ranks $\{6, 4, 3\}$, whose sum is 13. The rank sum is a basic component in the Mann-Whitney U-test, and most statistical tools provide it. In Equation 1, m is the number of observations in the first data sample, whereas n is the number of observations in the second data sample. In most experiments, we would run two randomized algorithms the same number of times: $m = n$.

When dealing with dichotomous results (as discussed in Section 6), several types of effect size measures [45] can be considered. The *odds ratio* is the most used and “is a measure of how many times greater the odds are that a member of a certain population will fall into a certain category than the odds are that a member of another population will fall into that category” [45]. Given a the number of times algorithm \mathcal{A} finds an optimal solution, and b for algorithm \mathcal{B} , the odds ratio is calculated as $\psi = \frac{a+\rho}{n+\rho-a} / \frac{b+\rho}{n+\rho-b}$, where ρ is any arbitrary positive constant (e.g., $\rho = 0.5$) used to avoid problems with zero occurrences [45]. There is no difference between the two algorithms when $\psi = 1$. The cases in which $\psi > 1$ imply that algorithm \mathcal{A} has higher chances of success.

Both \hat{A}_{12} and ψ are standardized effect size measures. But because their calculation is based on a finite number of observations (e.g., n for each algorithm, so $2n$ when we compare two algorithms), they are only estimates of the real \hat{A}_{12}^* and ψ^* . If n is low, these estimations might be very inaccurate. One way to deal with this problem is to calculate *confidence intervals* (CI) for them [45]. A $(1 - \alpha)$ CI is a set of values for which there is $(1 - \alpha)$ probability that the value of the effect size lies in that range. For example, if we have $\hat{A}_{12} = 0.54$ and a $(1 - \alpha)$ CI with range $[0.49, 0.59]$, then with probability $(1 - \alpha)$ the real value \hat{A}_{12}^* lies in $[0.49, 0.59]$ (where $\hat{A}_{12} = 0.54$ is its most likely estimation). Such effect size confidence intervals can facilitate decision making as they enable the comparison of the costs of alternative algorithms while accounting for uncertainty in their estimates. To see how confidence intervals can be calculated, please see [45] and [99].

Notice that a confidence interval can replace a test of statistical difference (e.g., t -test and U-test). If the null hypothesis H_0 lies within the confidence interval, then there is insufficient evidence to claim there is a statistically significant difference. In the previous example, because 0.5 is inside the $(1 - \alpha)$ CI $[0.49, 0.59]$, then there is no statistical difference at the selected significance level α . For a dichotomous result, H_0 would be $\psi = 1$.

8 Number of Runs

How many runs do we need when we analyze and compare randomized algorithms? As many as necessary to show with high confidence that the obtained results are statistically significant and to obtain a small enough confidence interval for effect size estimates. In many fields of science (e.g., medicine and behavioral science), a common rule of thumb is to use at least $n = 30$ observations. In the many fields where experiments are very expensive and time consuming, it is in general not feasible to work with high values for n . Several new statistical tests have been proposed and discussed to cope with the problem of lack of power and violation of assumptions (e.g., normality of data) when smaller numbers of observations are available [102].

Empirical studies of randomized algorithms do not involve human subjects and the number of *runs* (i.e., n) is only limited by computational resources. When there is access to clusters of computers as this is the case for many research institutes and universities, and when there is no need for expensive, specialized hardware (e.g., in hardware-in-the-loop testing), then large numbers of runs can be carried out to properly analyze the behavior of randomized algorithms. Many software engineering problems are furthermore not highly computationally expensive, as for example code coverage at the unit testing level, and can therefore involve very large numbers of executions. There are however exceptions, such as the system testing of embedded systems (e.g., [10]) where each test case can be very expensive to run.

Whenever possible, in most cases, it is therefore recommended to use a very high number of runs. For most problems in software engineering, thousands of runs should not be a problem and would solve most of the problems related to the power and accuracy of statistical tests. For example, as illustrated in [71, 27] in Table 2, even when 100 runs are used the U-test might be not powerful enough to confirm a statistical difference at a 0.05 significance level, even when the data seems to suggest such a difference.

Most discussions in the literature about statistical tests focus on situations with small numbers of observations (e.g., as in [88]). However, with thousands of runs, one would detect statistically significant differences on practically any experiment (Section 4). It is hence essential to complement such analyses with a study of the effect size as discussed in Section 7. Even when having large numbers of runs is not necessary for a set α level (e.g., 0.05) if differences are large enough to show p-values less than α , additional runs would help tighten the confidence intervals for effect size estimates and would be of practical value.

In Section 4, we suggested to use U-test instead of t -test. For very large samples, such as $n = 1,000$, there would be no practical difference between them regarding power and accuracy. However, the choice of a non-

parametric test would be driven by its effect size measure. In Section 7 we argued that effect size measures based on the mean (i.e., the d family) were not appropriate for randomized algorithms in software engineering due to violations in distribution assumptions. It would then be inconsistent to investigate the statistical difference of mean values with a t -test if we cannot use a reliable measure for its effect size. In other words, it is advisable to use size measures that are consistent with the differences being tested by the selected statistical test.

9 Multiple Tests

In most situations, we need to compare several alternative algorithms. Furthermore, if we are comparing different algorithm settings (e.g., population size in a Genetic Algorithm), then each setting technically defines a different algorithm. This often leads to a large number of statistical comparisons. It is possible to use statistical tests that deal with multiple techniques (treatments, experiments) at the same time (e.g., Factorial ANOVA), and effect size has been defined for those cases [45]. There are several types of statistical tests regarding multiple comparisons, and the choice depends on which research question one is addressing. In this paper we only deal with the two most common research questions in our context, since several books are dedicated to this topic, and an exhaustive analysis would not be possible in this paper:

- Does the choice of a particular parameter affect the performance of a randomized algorithm?
- Among a set of randomized algorithms, which one is the best in solving the addressed problem?

Assume a parameter that can assume several different values $j \in J$, and that we have carried out a series of experiments for a set of parameter values $\{j_1, j_2, \dots, j_k\} \subseteq J$. For example, in a Genetic Algorithm, we might want to study whether applying different cross-over rates has any effect on the effectiveness of the algorithm. One could consider the values $\{0, 0.25, 0.5, 0.75, 1\}$, and have $n = 1,000$ independent experiments for each of these five rate values. If we are only interested to evaluate whether the choice of this rate has any effect on the effectiveness of a Genetic Algorithm, then an *omnibus* test such as ANOVA can be employed. The null hypothesis is that the choice of the parameter value has no effect on the mean effectiveness of the algorithm. However, ANOVA suffers of the same problems as the t -test, i.e., assumption about normality of the data and equal variance. A non-parametric equivalent is the so called Kruskal-Wallis test.

Assume that a Kruskal-Wallis test states that the choice of that crossover rate has a statistically significant effect (i.e., the resulting p-value is low, so we can reject the null hypothesis). A relevant question might then be which crossover rate should be used (i.e., which one gives the best performance?). An omnibus test is not able to answer such a research question. This situation is exactly equivalent to the case of identifying the best algorithm among $K = 5$ algorithms/variants. In this case, we would like to compare the performance of each algorithm against all other alternatives individually. Given a set of algorithms, we would not be interested in simply determining whether all of them have the same mean values. Rather, given K algorithms, we want to perform $Z = K(K - 1)/2$ pairwise tests and measure effect size in each case.

However, using several statistical tests inflates the probability of Type I error. If we have only one comparison, the probability of Type I error is equal to the obtained p-value. If we have many comparisons, even when all the p-values are low, there is usually a high probability that at least in one of the comparisons the null hypothesis is true as all these probabilities somehow add up. In other words, if in all the comparisons the p-values are lower than α , then we would normally reject all the null hypotheses. But the probability that at least one null hypothesis is true could be as high as $1 - (1 - \alpha)^Z$ for Z comparisons, which converges to 1 as Z increases.

One way to address this problem is to use the so called *Bonferroni adjustment* [82, 76]. Instead of applying each test assuming a significance level α , we would use an adjusted level α/Z . For example, if we want at most a 0.05 probability of Type I error and we have two comparisons, we would need to use two statistical tests with $\alpha = 0.025$, and then check whether both differences are significant (i.e., if both p-values are lower than 0.025). However, the Bonferroni adjustment has been repeatedly criticized in the literature [82, 76], and we largely agree with those critiques. For example, let us assume that for both those tests we obtain p-values equal to 0.04. If a Bonferroni adjustment is used, then both tests will not be statistically significant at $\alpha = 0.05$ level. A researcher could be tempted to publish the results of only one of them and claiming statistical significance because $0.04 < 0.05$. Such a practice can therefore hinder scientific progress by reducing

the number of published results [82, 76]. This would be particularly true in our application context in which many randomized algorithms can be compared to address the same software engineering problem: it would be very tempting to leave out the results of some of the poorly performing algorithms. Notice that there are other adjustment techniques that are equivalent to Bonferroni but that are less conservative [37]. However, the statistical significance of a single comparison would still depend on the number of performed and reported comparisons. Though we do not recommend the Bonferroni adjustment, it is important to always report the obtained p-values, not just whether a difference is significant or not at an arbitrarily chosen α level. If for some reasons the readers want to evaluate the results using a Bonferroni adjustment or any of its (less conservative) variants, then it is possible to do so. For a full list of other problems related to the Bonferroni adjustment, the reader is referred to [82, 76].

Instead of pairwise tests using Bonferroni-like corrections, another (less popular) approach is to use the so called *post-hoc* methods, such as the Tukey’s range test. This test is applied on each of the Z pairs, and it is very similar to a t -test. Similar to the Bonferroni method, it employs a p-value correction to handle possible inflation of probability of Type I error.

Alpha level adjustments can be very important when assessing the validity of behavioral/physical phenomena with high confidence. For example, the leading international journal *Nature* has the following *requirement*² for published research papers regarding multiple tests:

- Multiple comparisons: When making multiple statistical comparisons on a single data set, authors should explain how they adjusted the alpha level to avoid an inflated Type I error rate, or they should select statistical tests appropriate for multiple groups (such as ANOVA rather than a series of t-tests).

However, in Section 4 we stated that in software engineering in general, and for randomized algorithms in particular, we mostly deal with decision-making problems. For example, if we must test software, then we must choose one alternative among K different techniques. In this case, even if the p-values are higher than α , we need to test the software anyhow and we must make a choice. In this context, Bonferroni-like adjustments make even less sense. Just choosing one alternative at random because there is no statistically significant difference is not optimal as it ignores available information.

Assume that we have analyzed the performance of K algorithms using pairwise tests and effect sizes. How to visualize the results of such analyses to grasp the relations among their performance? There can be different ways, and here we just describe a simple but practical one, that for example was used by Fraser and Arcuri in [32]. In that work [32], the effects of six parameters of a search algorithm were investigated in the context of automated unit testing of object-oriented software. Five parameters are binary (**Bo**, **Xo**, **Ra**, **Pa** and **Be**) and one ternary (**W**), for a total of $2^5 \times 3 = 96$ configurations. Each configuration was compared against all the other 95 (i.e., a total of 96×95 comparisons, which can be divided by two due to the symmetric property of the comparisons). Pairwise comparisons were made using a U-test, where the α level was arbitrarily set to 0.05. Initially, a score zero is assigned to each configuration. For each comparison in which a configuration is statistically better, its score is increased by one, whereas it is reduced by one in case it is statistically worse. Therefore, in the end each configuration has a score between -95 and 95. The higher the score, the better the configuration is. After this first phase, these scores are ranked such that the highest score has the best rank, where better ranks have lower values. In case of ties, the ranks are averaged. For example, if we have five configurations with scores $\{10, 0, 0, 20, -30\}$, then their ranks will be $\{2, 3.5, 3.5, 1, 5\}$. In [32], this procedure was repeated for each artifact in the case study (i.e., for all the 100 branches used in that empirical study), and the average of these ranks over all artifacts were calculated for each configuration, for a total of $100 \times 96 \times 95/2 = 456,000$ statistical comparisons. After collecting all of these data, a table was made in which the configurations were ordered based on their average rank from top (best) to bottom (worst). The same table from [32] is reported in Table 4. From this table, not only it is clear which are the best configurations, but it also possible to visualize some trends in the data (e.g., configurations with **Ra** are always better and **Xo** does not seem particularly useful).

²<http://www.nature.com/nature/authors/gta/index.html#a5.9>, accessed February 2011.

Table 4: Results of empirical analysis performed in [32]. The table shows the performance of the the 96 configurations, ordered from top (best performance) to bottom (worst performance). Symbols are used to indicate whether a particular boolean parameter is activated.

Bo	Xo	Ra	Pa	Be	20	W 50	80	Av. Rank	Av. Success Rate
△		⊕	▽	⊞		W		31.475	0.464
△		⊕	▽	⊞		W		31.840	0.456
△		⊕		⊞		W		32.595	0.482
		⊕	▽	⊞		W		32.670	0.456
		⊕	▽			W		34.725	0.447
△		⊕				W		35.415	0.448
		⊕		⊞		W		36.070	0.442
△		⊕		⊞	W			37.335	0.423
△	⊗	⊕	▽	⊞		W		37.430	0.430
△		⊕		⊞			W	37.605	0.459
△	⊗	⊕		⊞	W			37.615	0.418
	⊗	⊕		⊞		W		38.080	0.422
△	⊗	⊕	▽	⊞		W		39.325	0.419
	⊗	⊕		⊞		W		39.455	0.423
	⊗	⊕	▽			W		39.580	0.413
△		⊕			W			39.790	0.431
		⊕		⊞	W			39.815	0.431
	⊗	⊕			W			40.050	0.414
△		⊕	▽		W			40.140	0.420
△	⊗	⊕	▽			W		40.330	0.425
△		⊕	▽	⊞	W			40.670	0.413
△		⊕	▽	⊞			W	40.700	0.432
△	⊗	⊕		⊞	W			40.835	0.405
		⊕		⊞			W	40.940	0.438
△		⊕	▽			W		41.200	0.455
△	⊗	⊕			W			41.350	0.410
		⊕	▽	⊞		W		41.695	0.423
		⊕	▽	⊞	W			41.890	0.405
		⊕	▽		W			41.925	0.413
	⊗	⊕	▽		W			42.150	0.399
	⊗	⊕	▽	⊞			W	42.195	0.401
	⊗	⊕	▽	⊞	W			42.470	0.388
△	⊗	⊕	▽		W			42.500	0.395
	⊗	⊕		⊞		W		42.800	0.422
		⊕			W			43.075	0.407
	⊗	⊕				W		43.095	0.421
△	⊗	⊕			W			43.255	0.420
△	⊗	⊕	▽	⊞	W			43.635	0.377
		⊕				W		45.160	0.398
	⊗	⊕	▽				W	45.205	0.393
		⊕	▽			W		45.285	0.412
△	⊗	⊕	▽			W		45.450	0.392
△		⊕				W		45.850	0.418
		⊕				W		46.460	0.401
△	⊗	⊕				W		46.625	0.388
△	⊗	⊕		⊞		W		46.700	0.409
△	⊗	⊕	▽	⊞		W		47.760	0.379
△	⊗	⊕				W		47.850	0.384
		⊕	▽	⊞		W		48.985	0.342
			▽	⊞		W		49.585	0.329
			▽	⊞		W		49.705	0.334
△			▽	⊞	W			49.995	0.369
△	⊗		▽	⊞		W		50.290	0.313
△			▽		W			50.740	0.356
△	⊗		▽			W		51.295	0.313
△			▽			W		51.350	0.340
△				⊞		W		51.570	0.327
△			▽	⊞			W	52.215	0.326
				⊞	W			52.800	0.330
			▽	⊞		W		53.260	0.330
	⊗		▽	⊞		W		53.610	0.309
△			▽	⊞		W		53.845	0.321
	⊗		▽	⊞	W			54.040	0.310
	⊗		▽			W		54.475	0.312
			▽	⊞	W			54.835	0.296
			▽		W			55.080	0.306
				⊞		W		55.290	0.317
	⊗		▽			W		55.390	0.313
	⊗		▽	⊞		W		55.605	0.304
△					W			55.635	0.305
			▽			W		55.695	0.324
△	⊗		▽		W			56.065	0.310
△						W		56.160	0.309
	⊗			⊞		W		56.200	0.304
△	⊗		▽	⊞			W	56.255	0.301
	⊗		▽		W			56.295	0.312
△	⊗		▽	⊞	W			56.655	0.312
△	⊗		▽			W		56.835	0.291
△	⊗					W		57.095	0.279
△	⊗			⊞		W		57.135	0.291
△				⊞			W	57.180	0.319
				⊞			W	57.390	0.306
						W		58.955	0.285
△	⊗			⊞		W		59.085	0.297
				⊞	W			59.190	0.297
△	⊗			⊞	W			59.270	0.285
	⊗					W		59.595	0.279
△						W		59.995	0.300
	⊗			⊞	W			60.145	0.281
	⊗					W		60.150	0.289
△	⊗				W			60.675	0.278
	⊗			⊞		W		60.705	0.289
△	⊗					W		60.975	0.292
						W		61.655	0.267
	⊗				W			65.220	0.238
					W			71.765	0.190

10 Choice of Artifacts

When assessing randomized algorithms, the choice of artifacts to which these algorithms are applied (e.g., source code or executable programs) is of paramount importance as it usually has a strong bearing on the evaluation results. When analyzing empirical analyses in the software engineering literature evaluating randomized algorithms, many of the studies are carried out on artificial and small artifacts. Empirical analyses on real industrial systems are rare, thus raising questions about the credibility of results and the usefulness of the proposed algorithms. However, achieving realism by using representative industrial systems is particularly challenging. We usually cannot precisely characterize the population of artifacts we are targeting in our studies. Even if we could, we usually do not have access to large collections of industrial artifacts that are readily available to be sampled. And even if that were the case, studies are necessarily limited in terms of resources and time, and the number of artifacts studied is typically much more restricted than one would like. As a result, studies about randomized algorithms in software engineering typically present threats to external validity, making it difficult to generalize the results to other systems than the ones under study. In this paper, because the focus is on how to apply statistical tests, we do not emphasize the details of how one should choose artifacts from a general standpoint. We rather concentrate on how this choice affects the statistical tests procedures and the number of runs required.

The first question one faces is whether the selected artifacts are *representative* of the type of problem that is being addressed. For example, assume one wants to evaluate a new tool for automatically generating unit tests for object-oriented software (e.g., Pex [96], Randoop [81] or EvoSuite [33]). Which (types of) classes should be selected for experimenting? Following common practice in many empirical studies (e.g., [5, 85, 13]), is only using “container classes” acceptable? Well, it all depends on what is the target set of classes for the evaluation. If the proposed testing techniques are aimed *only* at container classes (e.g., [13]), then this would likely be acceptable. On the other hand, if the goal is to propose a *general* tool for generating unit tests, then using only container classes would lead to *serious* threats to external validity. But then the question is which classes should ideally be used? Again, we do not have well defined populations of classes that we can target and sample. But one possible simple heuristic is to try to maximize the diversity in terms of the type of classes, their size and complexity, and various other properties that are deemed relevant given the objective of the randomized algorithm, e.g., number of tasks accessing a lock when investigating deadlocks or data races [93]. For example we could rely on a sample containing a mix of container classes, numerical applications and others coming from common benchmarks?

As a practical alternative, one could use open source repositories such as SourceForge³, and randomly select a subset of projects for experimenting among the 260,000 that are currently hosted. If one wants to evaluate the applicability of a general tool for unit testing, this would be much better than using only container classes or arbitrarily choosing some programs in a non-systematic way (as it is often the case in the literature). However, even if one randomly samples projects from SourceForge, the empirical analyses would likely have some sort of bias. For example, open source projects in general may not be representative of programs developed in industry. Embedded systems and financial applications, for example, are unlikely to be well represented among these open source projects.

Regarding randomized algorithms (in particular search and optimization algorithms), there are specific and rigorous theoretical reasons for which the choice of artifacts is extremely important. The *No Free Lunch* theorem states that, on average across all possible problems (artifacts in our case), all search algorithms have the same performance [104]. If one does not clearly define which is the *space* of artifacts being targeted, then any comparison among randomized algorithms is doomed to be arbitrary. For example, let us consider again the example of unit testing of object-oriented software. Assume that a case study involves 10 classes, and algorithm \mathcal{A} is statistically better on seven of them, whereas algorithm \mathcal{B} is statistically better on the other three. One could naively claim that algorithm \mathcal{A} is *on average* better than \mathcal{B} . But maybe, those seven classes for which \mathcal{A} is better are all container classes, whereas the other three classes are related to string manipulations (e.g., [4]). If one had chosen for the case study more classes of this latter type, then the conclusions could be different (i.e., \mathcal{B} would be considered *on average* better than \mathcal{A}). Though the problem of choosing *appropriate* artifacts is intrinsically difficult, it is important for researchers to define their target artifacts as well as possible and carefully attempt to provide plausible reasons for differences in results across artifacts, such as classes, based

³<http://sourceforge.net/>, accessed February 2011.

on a thorough analysis of their characteristics.

If for the addressed research question the considered artifacts can be considered representative of the target, it is meaningful to then use statistical tests for evaluating whether algorithm \mathcal{A} is significantly better than \mathcal{B} on all selected artifact instances. However, as we see below, which type of test is used of of the highest importance. Using again the same example, assume six classes have been selected for investigating the unit testing of object-oriented software. Each algorithm is run on each of these six classes n times (e.g., $n = 30$), and average values out of these runs are collected for each class. This makes up a total of $2 \times 6 \times 30 = 360$ runs. Assume that the algorithms are evaluated based on how many test cases they generate before reaching full coverage. For the first algorithm, we obtain the following average values $X = \{10k, 20k, 30k, 40k, 50k, 60k\}$, whereas for the second algorithm we obtain $Y = \{12k, 21k, 34k, 41k, 53k, 68k\}$. The average values are ordered by problem instance where $k = 1000$, i.e., in X , out of $n = 30$ runs on the first artifact the average number of test cases run equals 10,000. Further assume that the problem instances are ordered by difficulty (i.e., solving the first problem is much easier than solving the fifth, because on average it requires to generate/run less test cases). If one wants to evaluate whether there is any statistical difference between X and Y , an *unpaired test*, such as Mann-Whitney U-test, would yield a p-value equal to 0.699 (e.g., by using the R [84] command “`wilcox.test(X,Y)`”), thus suggesting the difference is not statistically significant. However, this would be technically incorrect since different artifacts present different levels of difficulty, and considering all data together at the same time would blur the relative performance of the compared algorithms. In other words, a run of an inefficient algorithm on an *easy* problem would likely result in a better value than a run of a more efficient algorithm that is run instead on a *difficult* problem. If the case study involves artifacts of different levels of difficulty (as it is usually the case, either by design or due to random sampling) then it might be challenging to detect any statistical difference with an unpaired test.

Alternatively, *paired tests* such as Wilcoxon T test can be used (e.g., “`wilcox.test(X,Y, paired=TRUE)`” in R [84]). In a paired T test, what is evaluated is whether the differences $Z_i = Y_i - X_i$ are centered around 0, i.e., the null hypothesis is $Z = 0$. In that example, we have $Z = \{2k, 1k, 4k, 1k, 3k, 8k\}$, i.e, on average the second algorithm is always better than the first. A Wilcoxon T test here gives p-value=0.035, which suggests a statistically significant difference among the performance of the two algorithms, a result in sharp contrast with the unpaired test results above. This highlights why it is extremely important to use paired tests when comparing randomized algorithms on a set of selected artifacts. In the above example, the second algorithm is better in six out of six cases, which is a clear case. But typically results are not that consistent, and several of the compared algorithms may perform best on different artifacts. For example, if we assume a case study involving 100 artifacts, if an algorithm fares better on 51 of these, then the difference among the two would not be statistically significant when using a paired test. Using the example where an algorithm \mathcal{A} is better than another \mathcal{B} on some artifacts and worse on other artifacts, a T test evaluates whether one algorithm is statistically better on a higher number of artifacts.

The above discussion on the use of appropriate statistical tests is incomplete as it considers the evaluation of a randomized algorithm as ternary, i.e. it is either better, equivalent or worse than another one. Consider the following example: algorithm \mathcal{A} is better on 60% of the case study, but only by a very limited amount. On the other hand, on the other 40% of the case study, it is much worse than algorithm \mathcal{B} . In this case, blindly applying a Wilcoxon T test would lead to the conclusion that \mathcal{A} is preferable, whereas a practitioner might prefer to use \mathcal{B} . Another option could be to collect standardized effect sizes for each problem instance, and then average them over all of problems instances. This would provide additional information, but it would not solve the problem of fully describing the relative performance of two randomized algorithms. Consider a case with five artifacts and the following \hat{A}_{12} measures $\{0.6, 0.6, 0.6, 0.6, 0.1\}$. One algorithm is better than the other on four artifacts ($\hat{A}_{12} = 0.6$), but worse on the last one ($\hat{A}_{12} = 0.1$). If we average those values on the entire case study, we would obtain $\hat{A}_{12} = 0.5$, thus suggesting there is no difference among the two algorithms! This example illustrates the fact that aggregate statistics on a set of artifacts are useful to summarize the comparisons of two (or more) algorithms, but that particular care needs to be taken to handle cases where sharp differences can be observed among artifacts. In general, one should report the performance of the algorithms on each problem instance separately and attempt, as discussed above, to explain differences. One useful way to show the relative performance of randomized algorithms on a set of artifacts is to use box-plots of the effect sizes, especially when dealing with many artifacts

Ideally, when realistic artifacts for a certain type of problems are difficult to find, one would like to be

able to generate large numbers of them automatically in a realistic fashion. However, this requires that the artifacts have a clear and predictable structure, that there exist heuristics to generate correct and meaningful instances of such artifacts. If this is possible, one strong advantage is that one can control and vary interesting properties of the artifacts (e.g., class size, number of test cases) to enable interesting sensitivity analyses and assess the performance of randomized algorithms as a function of these properties. For example, in our work with Hemmati [49], we analyzed different test suite reduction techniques for model-based testing of large systems. Obtaining real models from industry is difficult, and UML models of real systems are not common in open source repositories. Although our case study was based on two real industrial systems (e.g., one provided by Cisco Systems), to cope with possible threats to external validity, we also used a large set of artificially generated test suites following some specific rules and a randomized construction algorithm. For example, we wanted to vary the number of test cases in the test suites and the fault detection rate, in order to assess their impact on the effectiveness of the resulting selection technique. We wanted to do so while retaining as much as possible the realism of the test suites in the case studies. Such studies may be considered a type of simulation and may not generate fully realistic artifacts. But they may provide useful insights into the impact of some artifact properties on the effectiveness of a randomized algorithm.

For some types of software engineering problems, a large number of artifacts can be selected or generated (e.g., randomly selecting classes to investigate the unit testing of open source software). When evaluating randomized algorithms in this context one has to make the following decision: Assume a budget for experiments $b = n \times z$ for each algorithm, where n represents the times a randomized algorithm is run on each artifact, and z is the number of these artifacts. If we consider b to be fixed (e.g., depending on how long it takes to run b experiments), then a practical and important question is how to choose n and z ? Two extreme cases would be $(n = 1, z = b)$ and $(n = b, z = 1)$, but they would clearly lead to problems in terms of statistical testing and external validity, respectively. We have to strike a balance between two objectives: we want to analyze as many artifacts as possible to improve external validity and wish at the same time to retain enough runs (i.e., n) to check whether there is a statistically significant difference on any single artifact when applying and comparing two randomized algorithms. This would, for obvious reasons, not be possible if $n = 1$. Though in Section 8 we suggested as a rule of thumb to use $n = 1,000$ when possible, in certain circumstances this may not be an option. If one has the possibility to analyze a large number z of artifacts but has practical constraints regarding the number of experiments to be run (e.g., having experiments running on a PC for a couple of years would not be very practical), then it may be more appropriate to execute less runs, perhaps as low as $n = 30$ or even $n = 10$. But going lower than such values would make the use of standard statistical tests very difficult and very likely, depending on the actual effect size, bring statistical power to unacceptable low levels.

As we discussed in Section 3, there are cases in the literature (e.g., [78, 101]) in which a random instance generator is used, but then the algorithms are run only once (i.e., $n = 1$) on each artifact. For all the reasons discussed in this section, we do not consider those empirical studies as appropriate.

11 Practical Guidelines

Based on the above discussions, we propose a set of practical guidelines for the use of statistical tests in experiments comparing randomized algorithms. Though we expect exceptions, given the current state of practice (Section 3 and [3, 52]), we believe that it is important to provide practical guidance that will be valid in most cases and enable higher quality studies to be reported. We recommend that practitioners follow these guidelines and justify any necessary deviation.

There are many statistical tools that are available. In the following we will provide examples based on *R* [84], because it is a powerful tool that is freely available and supported by many statisticians. But any other professional tool would provide similar capabilities.

Practical guidelines are summarized as follows. Notice that often, for reasons of space, it is not possible to report all the data of the statistical tests. Based on the circumstances, authors need to make careful choices on what to report.

- When randomized algorithms are analyzed, clearly specify the number of runs and employed statistical tests. For example, they can be summarized in a threats to validity section, in which how randomness has been taken into account should be discussed and justified.

- On each artifact in the case study, run each randomized algorithm at least $n = 1,000$ times. If this is not possible, explain the reasons and report the total amount of time it took to run the entire case study. If for example 30 runs were performed and the total execution time was just one hour, then it is rather difficult to justify why a higher number of runs was not used to gain statistical power, lower p-values, and narrow the confidence interval of effect size estimates (Section 8).
- When a large number of artifacts can be used in the case study (e.g., for unit testing of open source software) but there are constraints in terms of execution time, then it is advisable to execute less runs per artifact (though at least $n = 10$) and use more artifacts (rather than having $n = 1,000$ but only few artifacts, see Section 10). The objective is to strike a balance between generalization and statistical power.
- For detecting statistical differences, use the two-tailed non-parametric Mann-Whitney U-test for interval-scale results and the Fisher exact test for dichotomous results (i.e., in the cases of censored data as discussed in Section 6). For the former case, in *R* you can use the function “w=wilcox.test(X,Y)” where *X* and *Y* are the data sets with the observations of the two compared randomized algorithms. If you are comparing a randomized algorithm against a deterministic one, use “w=wilcox.test(X,mu=D)”, where *D* is the resulting performance measure for the deterministic algorithm. When we have number of successes *a* for the first algorithm and *b* for the second, you can use “f=fisher.test(m)”, where *m* is a matrix derived in this way: “m=matrix(c(a,n-a,b,n-b),2,2)”. A constant $\rho = 0.5$ could be added to each cell of the matrix to address zero occurrence cases.
- Report all the obtained p-values, whether they are smaller than α or not, and not just whether differences are significant. The motivation is for the reader to choose the level of risk that is suitable in her application context. When reporting all p-values is not possible, one could report the proportion of significant test results: “*x* out of *y* tests were significant at α level . . .”.
- Always report standardized effect size measures. For dichotomous results, the odds ratio ψ (and its confidence interval) is automatically calculated with “f=fisher.test(m)”. For interval-scale results and the \hat{A}_{12} effect size, the rank sum R_1 used in Equation 1 can be calculated with “R1=sum(rank(c(X,Y))[seq_along(X)])”. It is also strongly advised to report effect size confidence intervals (but the support for \hat{A}_{12} is unfortunately limited). This is much easier to use than p-values for decision making as potential benefits can be compared to costs while accounting for uncertainty.
- To help the meta-analyses of published results across studies, report means and standard deviations (in case readers for some reasons want to calculate effect sizes in the *d* family). For dichotomous experiments, always report the values *a* and *b* (so that other types of effect sizes can be computed [45]).
- If space permits, provide full statistics for the collected data, as for example mean, median, variance, min/max values, skewness, median and absolute deviation. Box-plots are also useful to visualize them.
- When analyzing more than two randomized algorithms, use pairwise comparisons including pairwise statistical tests and effect size measures.
- Always state the employed statistical tool (there can be subtle differences on how the tests are computed).

12 Threats to Validity

The systematic review in Section 3 is based on only three sources, from which only 49 out of 223 papers were selected. A larger review might lead to different results, although we can safely argue that TSE and ICSE are representative of research trends in software engineering. Furthermore, that review is only used as a motivation for providing practical guidelines, and its results are in line with other larger systematic reviews [3, 52]. Last, papers sometimes lack precision and interpretation errors are always possible.

As already discussed in Section 11, our practical guidelines may not be applicable to all contexts. Therefore, in every specific context, one should always carefully assess them. For some specific cases, other statistical procedures could be preferable, especially when only few runs are possible.

13 Conclusion

In this paper we report on a systematic review to evaluate how the results of randomized algorithms in software engineering are analyzed. This type of algorithms (e.g., Genetic Algorithms) are widely used to address many software engineering problems, such as test case selection. Similar to previous systematic reviews on related topics [3, 52], we conclude that the use of rigorous statistical methodologies are somehow lacking when investigating randomized algorithms in software engineering.

To cope with this problem, we provide, discuss, and justify a set of *practical* guidelines targeting researchers in software engineering. In contrast to other guidelines in the literature for other scientific fields (e.g., [77] and [53]), the guidelines in this paper are tailored to the specific properties of randomized algorithms when applied to software engineering problems. The use of these guidelines is important in order to develop a reliable body of empirical results over time, by enabling comparisons across studies so as to converge towards generalizable results of practical importance. Otherwise, as in many other aspects of software engineering, unreliable results will prevent effective technology transfer and will inevitably limit the impact of research on practice.

Acknowledgments

We would like to thanks Lydie du Bousquet and Zohaib Iqbal for useful comments on an early draft of this paper. The work described in this paper was supported by the Norwegian Research Council. This paper was produced as part of the ITEA-2 project called VERDE.

References

- [1] R. Abraham and M. Erwig. Mutation Operators for Spreadsheets. *IEEE Transactions on Software Engineering (TSE)*, 35(1), 2009.
- [2] J. Aguilar-Ruiz, I. Ramos, J. C. Riquelme, and M. Toro. An evolutionary approach to estimating software development projects. *Information and Software Technology*, 43:875–882, 2001.
- [3] S. Ali, L. Briand, H. Hemmati, and R. Panesar-Walawege. A systematic review of the application and empirical investigation of search-based test-case generation. *IEEE Transactions on Software Engineering (TSE)*, 36(6):742–762, 2010.
- [4] M. Alshraideh and L. Bottaci. Search-based software test data generation for string data using program-specific search operators. *Software Testing, Verification and Reliability (STVR)*, 16(3):175–203, 2006.
- [5] J. H. Andrews, T. Menzies, and F. C. Li. Genetic algorithms for randomized unit testing. *IEEE Transactions on Software Engineering (TSE)*, 37(1), 2011.
- [6] J. Antunes, N. Neves, M. Correia, P. Verissimo, and R. Neves. Vulnerability discovery with attack injection. *IEEE Transactions on Software Engineering (TSE)*, 36(3):357–370, 2010.
- [7] A. Arcuri. Full theoretical runtime analysis of alternating variable method on the triangle classification problem. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 113–121, 2009.
- [8] A. Arcuri. Theoretical analysis of local search in software testing. In *Symposium on Stochastic Algorithms, Foundations and Applications (SAGA)*, pages 156–168, 2009.
- [9] A. Arcuri and L. Briand. A practical guide for using statistical tests to assess randomized algorithms in software engineering. In *ACM/IEEE International Conference on Software Engineering (ICSE)*, 2011.
- [10] A. Arcuri, M. Z. Iqbal, and L. Briand. Black-box system testing of real-time embedded systems using random and search-based testing. In *IFIP International Conference on Testing Software and Systems (ICTSS)*, pages 95–110, 2010.

- [11] A. Arcuri, M. Z. Iqbal, and L. Briand. Formal analysis of the effectiveness and predictability of random testing. In *ACM International Symposium on Software Testing and Analysis (ISSTA)*, pages 219–229, 2010.
- [12] A. Arcuri and X. Yao. A novel co-evolutionary approach to automatic software bug fixing. In *IEEE Congress on Evolutionary Computation (CEC)*, pages 162–168, 2008.
- [13] A. Arcuri and X. Yao. Search based software testing of object-oriented containers. *Information Sciences*, 178(15):3075–3095, 2008.
- [14] S. Artzi, A. Kiezun, J. Dolby, F. Tip, D. Dig, A. Paradkar, and M. D. Ernst. Finding bugs in web applications using dynamic test generation and explicit-state model checking. *IEEE Transactions on Software Engineering (TSE)*, 36(4):474–494, 2010.
- [15] F. Asadi, G. Antoniol, and Y. Gueheneuc. Concept Location with Genetic Algorithms: A Comparison of Four Distributed Architectures. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 153–162, 2010.
- [16] A. J. Bagnall, V. J. Rayward-Smith, and I. M. Whitley. The next release problem. *Information and Software Technology*, 43(14):883–890, 2001.
- [17] N. E. Beckman, A. V. Nori, S. K. Rajamani, R. J. Simmons, S. D. Tetali, and A. V. Thakur. Proofs from tests. *IEEE Transactions on Software Engineering (TSE)*, 36(4):495–508, 2010.
- [18] M. Bowman, L. C. Briand, and Y. Labiche. Solving the class responsibility assignment problem in object-oriented analysis with multi-objective genetic algorithms. *IEEE Transactions on Software Engineering (TSE)*, 36(6):817–837, 2010.
- [19] G. Canfora, M. D. Penta, R. Esposito, and M. L. Villani. An approach for qos-aware service composition based on genetic algorithms. In *Genetic and Evolutionary Computation Conference (GECCO)*, pages 1069–1075, 2005.
- [20] J. Cohen. Statistical power analysis for the behavioral sciences, 1988.
- [21] K. D. Cooper, P. J. Schielke, and D. Subramanian. Optimizing for reduced code space using genetic algorithms. In *Proceedings of the ACM SIGPLAN workshop on Languages, compilers, and tools for embedded systems*, pages 1–9, 1999.
- [22] M. Cowles and C. Davis. On the origins of the .05 level of statistical significance. *American Psychologist*, 37(5):553–558, 1982.
- [23] J. T. de Souza, C. L. Maia, F. G. de Freitas, and D. P. Coutinho. The Human Competitiveness of Search Based Software Engineering. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 143–152, 2010.
- [24] J. del Sagrado, I. M. del Aguila, and F. J. Orellana. Ant Colony Optimization for the Next Release Problem. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 67–76, 2010.
- [25] H. Do, S. Mirarab, L. Tahvildari, and G. Rothermel. The effects of time constraints on test case prioritization: A series of controlled experiments. *IEEE Transactions on Software Engineering (TSE)*, 36(5):593–617, 2010.
- [26] J. W. Duran and S. C. Ntafos. An evaluation of random testing. *IEEE Transactions on Software Engineering (TSE)*, 10(4):438–444, 1984.
- [27] J. Durillo, Y. Zhang, E. Alba, and A. Nebro. A Study of the Multi-objective Next Release Problem. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 49–58, 2009.

- [28] T. Dybå, V. Kampenes, and D. Sjøberg. A systematic review of statistical power in software engineering experiments. *Information and Software Technology (IST)*, 48(8):745–755, 2006.
- [29] P. Emberson and I. Bate. Stressing search with scenarios for flexible solutions to real-time task allocation problems. *IEEE Transactions on Software Engineering (TSE)*, 36(5):704–718, 2010.
- [30] M. Fay and M. Proschan. Wilcoxon-Mann-Whitney or t-test? On assumptions for hypothesis tests and multiple interpretations of decision rules. *Statistics Surveys*, 4:1–39, 2010.
- [31] W. Feller. *An Introduction to Probability Theory and Its Applications, Vol. 1*. Wiley, 3 edition, 1968.
- [32] G. Fraser and A. Arcuri. It is not the length that matters, it is how you control it. In *IEEE International Conference on Software Testing, Verification and Validation (ICST)*, 2011.
- [33] G. Fraser and A. Arcuri. Whole test suite generation. Technical report, Chair of Software Engineering, Saarland University, 2011.
- [34] G. Freitag, S. Lange, and A. Munk. Non-parametric assessment of non-inferiority with censored data. *Statistics in medicine*, 25(7):1201, 2006.
- [35] M. Gabel and Z. Su. Online inference and enforcement of temporal properties. In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 15–24, 2010.
- [36] V. Ganesh, T. Leek, and M. Rinard. Taint-based directed whitebox fuzzing. In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 474–484, 2009.
- [37] L. García. Escaping the Bonferroni iron claw in ecological studies. *Oikos*, 105(3):657–663, 2004.
- [38] V. Garousi. A genetic algorithm-based stress test requirements generator tool and its empirical evaluation. *IEEE Transactions on Software Engineering (TSE)*, 36(6):778–797, 2010.
- [39] B. Garvin, M. Cohen, and M. Dwyer. An improved meta-heuristic search for constrained interaction testing. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 13–22, 2009.
- [40] K. Ghani, J. Clark, and Y. Heslington. Widening the Goal Posts: Program Stretching to Aid Search Based Software Testing. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 122–131, 2009.
- [41] M. Gligoric, T. Gvero, V. Jagannath, S. Khurshid, V. Kuncak, and D. Marinov. Test generation through programming in udit. In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 225–234, 2010.
- [42] P. Godefroid, N. Klarlund, and K. Sen. Dart: directed automated random testing. In *ACM Conference on Programming language design and implementation (PLDI)*, pages 213–223, 2005.
- [43] S. Goodman. P values, hypothesis tests, and likelihood: implications for epidemiology of a neglected historical debate. *American Journal of Epidemiology*, 137(5):485–496, 1993.
- [44] S. Goodman. Toward evidence-based medical statistics. 1: The P value fallacy. *Annals of Internal Medicine*, 130(12):995–1004, 1999.
- [45] R. Grissom and J. Kim. *Effect sizes for research: A broad practical approach*. Lawrence Erlbaum, 2005.
- [46] Z. Gu, E. T. Barr, D. J. Hamilton, and Z. Su. Has the bug really been fixed? In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 55–64, 2010.
- [47] M. Harman, S. A. Mansouri, and Y. Zhang. Search based software engineering: A comprehensive analysis and review of trends techniques and applications. Technical Report TR-09-03, King’s College, 2009.

- [48] M. Harman and P. McMinn. A theoretical and empirical study of search based testing: Local, global and hybrid search. *IEEE Transactions on Software Engineering (TSE)*, 36(2):226–247, 2010.
- [49] H. Hemmati, A. Arcuri, and L. Briand. Empirical investigation of the effects of test suite properties on similarity-based test case selection. In *IEEE International Conference on Software Testing, Verification and Validation (ICST)*, 2011.
- [50] H. Hsu and A. Orso. MINTS: A general framework and tool for supporting test-suite minimization. In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 419–429, 2009.
- [51] S. Jha, S. Gulwani, S. A. Seshia, and A. Tiwari. Oracle-guided component-based program synthesis. In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 215–224, 2010.
- [52] V. Kampenes, T. Dybå, J. Hannay, and D. Sjøberg. A systematic review of effect size in software engineering experiments. *Information and Software Technology (IST)*, 49(11-12):1073–1086, 2007.
- [53] M. Katz. *Multivariable analysis: a practical guide for clinicians*. Cambridge Univ Pr, 2006.
- [54] K. Khan, R. Kunz, J. Kleijnen, and G. Antes. *Systematic reviews to support evidence-based medicine: how to review and apply findings of healthcare research*. RSM Press, 2004.
- [55] U. Khan and I. Bate. WCET analysis of modern processors using multi-criteria optimisation. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 103–112, 2009.
- [56] T. Khoshgoftaar, L. Yi, and N. Seliya. A multiobjective module-order model for software quality enhancement. *IEEE Transactions on Evolutionary Computation (TEC)*, 8(6):593–608, 2004.
- [57] A. Kieyzun, P. Guo, K. Jayaraman, and M. Ernst. Automatic creation of SQL injection and cross-site scripting attacks. In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 199–209, 2009.
- [58] D. Kim and S. Park. Dynamic Architectural Selection: A Genetic Algorithm Based Approach. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 59–68, 2009.
- [59] B. Kitchenham, O. Pearl Brereton, D. Budgen, M. Turner, J. Bailey, and S. Linkman. Systematic literature reviews in software engineering-A systematic literature review. *Information and Software Technology (IST)*, 51(1):7–15, 2009.
- [60] J. Klein and M. Moeschberger. *Survival analysis: techniques for censored and truncated data*. Springer Verlag, 2003.
- [61] S. Kpodjedo, F. Ricca, G. Antoniol, and P. Galinier. Evolution and Search Based Metrics to Improve Defects Prediction. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 23–32, 2009.
- [62] Z. Lai, S. Cheung, and W. Chan. Detecting atomic-set serializability violations in multithreaded programs through active randomized testing. In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 235–244, 2010.
- [63] K. Lakhota, M. Harman, and H. Gross. AUSTIN: A tool for Search Based Software Testing for the C Language and its Evaluation on Deployed Automotive Systems. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 101–110, 2010.
- [64] N. Leech and A. Onwuegbuzie. A Call for Greater Use of Nonparametric Statistics. Technical report, US Dept. Education, 2002.
- [65] F. Lindlar and A. Windisch. A Search-Based Approach to Functional Hardware-in-the-Loop Testing. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 111–119, 2010.

- [66] G. Lu, R. Bahsoon, and X. Yao. Applying Elementary Landscape Analysis to Search-Based Software Engineering. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 3–8, 2010.
- [67] A. Marchetto and P. Tonella. Search-based testing of Ajax web applications. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 3–12, 2009.
- [68] A. Masood, R. Bhatti, A. Ghafoor, and A. Mathur. Scalable and Effective Test Generation for Role-Based Access Control Systems. *IEEE Transactions on Software Engineering (TSE)*, pages 654–668, 2009.
- [69] P. McMinn. Search-based software test data generation: A survey. *Software Testing, Verification and Reliability*, 14(2):105–156, 2004.
- [70] P. McMinn. How Does Program Structure Impact the Effectiveness of the Crossover Operator in Evolutionary Testing? In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 9–18, 2010.
- [71] T. Menzies, S. Williams, B. Boehm, and J. Hihn. How to avoid drastic software process change (using stochastic stability). In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 540–550, 2009.
- [72] B. S. Mitchell and S. Mancoridis. On the automatic modularization of software systems using the bunch tool. *IEEE Transactions on Software Engineering (TSE)*, 32(3):193–208, 2006.
- [73] T. Mitchell. *Machine Learning*. McGraw Hill, 1997.
- [74] M. Motwani and P. Raghavan. *Randomized Algorithms*. Cambridge University Press, 1995.
- [75] P. A. Nainar and B. Liblit. Adaptive bug isolation. In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 255–264, 2010.
- [76] S. Nakagawa. A farewell to Bonferroni: the problems of low statistical power and publication bias. *Behavioral Ecology*, 15(6):1044–1045, 2004.
- [77] S. Nakagawa and I. Cuthill. Effect size, confidence interval and statistical significance: a practical guide for biologists. *Biological Reviews*, 82(4):591–605, 2007.
- [78] A. Ngo-The and G. Ruhe. Optimized Resource Allocation for Software Release Planning. *IEEE Transactions on Software Engineering (TSE)*, 35(1):109–123, 2009.
- [79] S. Nijssen and T. Back. An analysis of the behavior of simplified evolutionary algorithms on trap functions. *IEEE Transactions on Evolutionary Computation (TEC)*, 7(1):11–22, 2003.
- [80] A. Nori and S. K. Rajamani. An empirical study of optimizations in yogi. In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 355–364, 2010.
- [81] C. Pacheco, S. K. Lahiri, M. D. Ernst, and T. Ball. Feedback-directed random test generation. In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 75–84, 2007.
- [82] T. Perneger. What’s wrong with Bonferroni adjustments. *British Medical Journal*, 316:1236–1238, 1998.
- [83] S. Poulding and J. Clark. Efficient Software Verification: Statistical Testing Using Automated Search. *IEEE Transactions on Software Engineering (TSE)*, 36(6):763–777.
- [84] R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2008. ISBN 3-900051-07-0.

- [85] J. C. B. Ribeiro, M. A. Zenha-Rela, and F. F. de Vega. Test case evaluation and input domain reduction strategies for the evolutionary testing of object-oriented software. *Information and Software Technology*, 51(11):1534–1548, 2009.
- [86] J. A. Rice. *Mathematical Statistics and Data Analysis*. Duxbury Press, 2 edition, 1994.
- [87] G. Rudolph. Convergence analysis of canonical genetic algorithms. *IEEE transactions on Neural Networks*, 5(1):96–101, 1994.
- [88] G. Ruxton. The unequal variance t-test is an underused alternative to Student’s t-test and the Mann-Whitney U test. *Behavioral Ecology*, 17(4):688–690, 2006.
- [89] S. Sawilowsky and R. Blair. A more realistic look at the robustness and type II error properties of the t test to departures from population normality. *Psychological Bulletin*, 111(2):352–360, 1992.
- [90] C. A. Schaefer, V. Pankratius, and W. F. Tichy. Engineering parallel applications with tunable architectures. In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 405–414, 2010.
- [91] R. Sharma, M. Gligoric, A. Arcuri, G. Fraser, and D. Marinov. Testing container classes: Random or systematic? In *Fundamental Approaches to Software Engineering (FASE)*, 2011.
- [92] M. Shevertalov, J. Kothari, E. Stehle, and S. Mancoridis. On the Use of Discretized Source Code Metrics for Author Identification. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 69–78, 2009.
- [93] M. Shousha, L. Briand, and Y. Labiche. A uml/marte model analysis method for uncovering scenarios leading to starvation and deadlocks in concurrent systems. *IEEE Transactions on Software Engineering (TSE)*, 2010. 10.1109/TSE.2010.107.
- [94] C. L. Simons, I. C. Parmee, and R. Gwynllwy. Interactive, evolutionary search in upstream object-oriented class design. *IEEE Transactions on Software Engineering (TSE)*, 36(6):798–816, 2010.
- [95] T. Thum, D. Batory, and C. Kastner. Reasoning about edits to feature models. In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 254–264, 2009.
- [96] N. Tillmann and N. J. de Halleux. Pex — white box test generation for .NET. In *International Conference on Tests And Proofs (TAP)*, pages 134–253, 2008.
- [97] P. Tonella. Evolutionary testing of classes. In *ACM International Symposium on Software Testing and Analysis (ISSTA)*, pages 119–128, 2004.
- [98] P. Tonella, A. Susi, and F. Palma. Using Interactive GA for Requirements Prioritization. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 57–66, 2010.
- [99] A. Vargha and H. D. Delaney. A critique and improvement of the CL common language effect size statistics of McGraw and Wong. *Journal of Educational and Behavioral Statistics*, 25(2):101–132, 2000.
- [100] W. Weimer, T. Nguyen, C. L. Goues, and S. Forrest. Automatically finding patches using genetic programming. In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 364–374, 2009.
- [101] J. White, B. Dougherty, and D. Schmidt. Ascent: An algorithmic technique for designing hardware and software in tandem. *IEEE Transactions on Software Engineering (TSE)*, 36(6), 2010.
- [102] R. Wilcox. *Fundamentals of modern statistical methods: Substantially improving power and accuracy*. Springer Verlag, 2001.
- [103] C. Wohlin. *Experimentation in software engineering: an introduction*, volume 6. Springer Netherlands, 2000.

- [104] D. H. Wolpert and W. G. Macready. No free lunch theorems for optimization. *IEEE Transactions on Evolutionary Computation*, 1(1):67–82, 1997.
- [105] J. Xiao and W. Afzal. Search-based resource scheduling for bug fixing tasks. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 133–142, 2010.
- [106] Q. Yang and M. Li. A cut-off approach for bounded verification of parameterized systems. In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 345–354, 2010.
- [107] S. Yoo. A Novel Mask-Coding Representation for Set Cover Problems with Applications in Test Suite Minimisation. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 19–28, 2010.
- [108] X. Yuan and A. M. Memon. Generating event sequence-based test cases using gui runtime state feedback. *IEEE Transactions on Software Engineering (TSE)*, 36(1):81–95, 2010.
- [109] L. Zhang, S. Hou, J. Hu, T. Xie, and H. Mei. Is operator-based mutant selection superior to random mutant selection? In *ACM/IEEE International Conference on Software Engineering (ICSE)*, pages 435–444, 2010.
- [110] Y. Zhang and M. Harman. Search Based Optimization of Requirements Interaction Management. In *International Symposium on Search Based Software Engineering (SSBSE)*, pages 47–56, 2010.