

USING STATISTICAL DISTRIBUTIONS TO GENERATE RANDOM TEST DATA

by

JAMIE ZIMMERMAN

A THESIS

Presented to the Department of Computer and Information Science
and the Robert D. Clark Honors College
in partial fulfillment of the requirements for the degree of
Bachelor of Science

June 2018

An Abstract of the Thesis of

Jamie Zimmerman for the degree of Bachelor of Science
in the Department of Computer and Information Science to be taken June 2018.

Title: Using Statistical Distributions to Generate Random Test Data

Approved: _____

Dr. Michal Young

Many open-source software programs lack rigorous, wide-reaching test suites. Quality test cases provide assurance that the software functions properly and performs the job it is expected to, but they are often neglected primarily because the testing process is deeply dependent on human invention and manual writing. This has a substantial financial and temporal expense on the software development lifecycle. However, the data for some system tests can be generated from statistical distributions that exhibit certain trends or patterns when interpreted in the context of the program. Therefore, this approach can tease out specific behaviors in the program that warrant being checked. This project aims to automate the creation of those data points to streamline the process of writing large test suites.

Acknowledgements

It is with the deepest gratitude I salute Dr. Michal Young for his wisdom, guidance, forgiveness, and intellect. What an honor it has been to enjoy the time offered by a true polymath. Finally I thank him sincerely for a year of weekly meetings for which significant portions were devoted to chats about bicycles and coffee. I also thank Dr. Hank Childs for his seemingly inexhaustible generosity in taking on another student's senior project. I express to Dr. Rebecca Lindner appreciation for her constructive motivation to always improve my writing communication, her bright attention to detail, and her encouragement for me to remain a scholar in the Honors College, well before this project had materialized and when I thought I had no chance of completing it. I am forever thankful to my three advisors for encouraging me to strive for the highest.

Table of Contents

1	Introduction	1
2	Literature Survey	5
3	Proposed Argument	9
4	Background on Design Decisions	11
4.1	Parameters	11
4.2	Pairwise Testing	13
4.3	Context-Free Grammars	15
4.4	Law of Large Numbers	17
4.5	Cardioids	21
4.6	Preprocessing	22
5	Methods	23
6	Results	25
6.1	Celestial Body Simulation	25
6.2	Earthquake Analysis and Visualization Program	27
7	Concluding Thoughts	36
	Bibliography	40

List of Accompanying Materials

1. GenSequence: <https://github.com/TestCreator/GenSequence>
2. GenPairs: <https://github.com/TestCreator/GenPairs>

List of Terms

SUT - Software Under Test, the program to be tested

Production Rule - a substitution rule (which usually belongs to a grammar, a set of rules) that describes what the left-hand symbol could possibly be. For example, an A could be a B, denoted by $A \rightarrow B$. Whenever an A is seen, one can choose to replace it with a B.

Test Vector - A symbolic description of a test case, lacking exact data points but having an English language descriptor

Fault - a misbehavior in a software program, also known as a bug. The behavior can vary from a program crash to an unexpected, nonsensical output

CSV file format - a comma-separated value file that is viewable in Microsoft Excel

List of Figures

1.1	Basic Workflow of a software program	2
3.1	Basic Diagram of User's workflow with GenSequence	10
4.1	A CSV file of a test case input into the Team Builder program	12
4.2	Histogram of point values generated by normal distribution - sample size 100	19
4.3	Histogram of point values generated by normal distribution - sample size 10,000	20
4.4	Histogram of point values generated by normal distribution and re- duced to sample size of 300	20
6.1	Starting Locations of 70 celestial bodies	26
6.2	Ending Locations of 70 celestial bodies after 665 time steps	27
6.3	Plotting Magnitudes and Depths from the same test case	30
6.4	An Expected Earthquake Event – High Intensity, Deep Below Surface	32
6.5	Another Expected Earthquake Event – Low Intensity, Near Surface .	33
6.6	An Outlier Earthquake Event – High Intensity, Near Surface	34
6.7	K-Means Clustering of Earthquake events from case 13	35

List of Tables

4.1	Single Symbolic Test Vector	13
4.2	Single Concrete Test Case	14

Chapter 1

Introduction

Software testing is a crucial part of systems that run our lives. Industries that rely on predicting consumer habit trends or dispatching taxi cars are gravely affected when their software fails to provide accurate results. The gravity of this trust is even greater in safety-critical programs controlling gas-leak shut-off valves or anesthetic delivery machines. Consequently, before these pieces of software are deployed, they must be checked and tested rigorously. The breadth and completeness of testing is directly proportional to how much damage could occur if the software malfunctioned. The software must eventually reach production (used in the hands of the customer or in industrial practice “for real”) which means that it has to leave the testing stage eventually. When that time comes, testers must be assured that the software behaves properly. This only sometimes happens. A 2017 report found that 43% of software engineers did not have enough time to test their code before its release date [23]. Therefore, the tradeoff of cost and convenience is an important consideration in determining the rigor of testing. When the release date occurs, testers either have to ship their product without confidence it works or else invest extra time and energy into seeing it does work. Cost is never low, and convenience is never high, so software testing remains a lengthy and sometimes ad hoc practice in many industrial applications. Numerous papers agree that testing is a tremendously arduous part of

software development [25, 14, 26, 33, 13, 35, 10, 23, 18, 11]. In fact, Misailovic et. al report that 50% of the software development process is software testing. This is not surprising since testing the quality of software is itself riddled with difficulties on every front.

Firstly, testing the software means knowing what the software *should* do and what output it *should* generate. A key component of getting correct results is knowledge of what correctness looks like. Developers and testers must be careful not to confuse code that compiles and runs without errors with code that renders accurate results. To find accurate results, testers need an oracle, which is a way of determining if the actual result of a program is the expected result, typically a boolean pass or fail value. Figure 1.1 illustrates the divergence. Unsurprisingly, obtaining and checking this oracle is difficult. Given a particular input to a program, a tester must know what the expected result

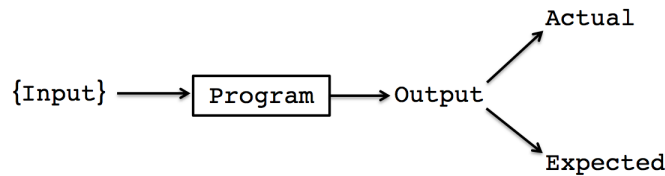


Figure 1.1: Basic Workflow of a software program

Then, she must know how to read or understand the actual result and be able to compare it to the expected.¹

In immensely complex scientific software or financial

programs swathed in accounting terms, the output of the program might be an unknown language to the tester. Even if the tester is a domain scientist, the output may be too complicated to read quickly and identify the trends that indicate she received a postive test result.

¹Consequently, many companies often build the role of tester into the role of developer, since the developer has the most intimate knowledge of the design and implementation of the software, and therefore has a better sense of what the output of the program exhibits. Even then, it may be impossible to describe the ideal expected output of the program.

Another major challenge in software testing is the need to know all the possible situations one may want to test for. In an automobile simulator, for example, a tester does not want to just test for braking and accelerating, she wants to test turn-signaling, beeping, and brake-lighting, and also combinations of those tests to make sure that the beeping sensor does not accidentally disable the brake pedal capability. Knowing and describing every single situation may be unknown to the tester; they may not even think that that is a situation they have to test for. But these situations must be tested. What is known as the “happy path”, the operational choice that most frequently represents the typical testing scenario, is certainly necessary for testing the programs reliability, but it does not find bugs. Outlier situations find many more bugs than the operational choice can. Lindvall et. al summarize the need for edge cases in their research in testing autonomous drones: “Just because the drone behaves in a safe manner for a set of test scenarios does not mean that it will always behave in a safe manner for other scenarios” [20]. Identifying all these possible outlier situations is **Problem One**.

With human’s prolific ability to create and store data quantifying the world around them, it is hard to imagine that a tester could not find data that they could use. Gotterbarn agrees that “Insufficient data is not a problem” [12]. For example, consider an ocean temperature monitoring software that has predictive power in future local hot spots or cold zones. Ocean temperature data points do exist, and the missing points can likely be interpolated quite easily. But this is a sample size of one. To be sure the temperature projection software is robust, the tester would want to study a variety of circumstances and possibilities situations that do not even exist. They might perhaps want to study temperature diffusion trends as a result of a hypothetical significant event 30 years from now. Data for that particular test case has to be created.

This is **Problem Two** – writing the actual test data itself. A symbolic description of the test must be turned into an actual concrete input to the software. A test vector checking that beeping does not disable braking must be turned into “beeping=5s&&braking=true” or whatever discrete format the SUT requires. This is not a step that should be done manually. For example, VisIt, a graphical visualization tool, is capable of handling several gigabyte files [17]. No tester wants to or even could write three gigabytes of data points modeling a real-world situation by hand just to satisfy one test case. Nearly every software testing paper agrees that writing data by hand seriously hinders the testing process [22, 25, 27, 28].

It is significantly easier to describe the trend of a certain test than it is to create the hundreds (if not more) data points that fit that description. Therefore, the goal of this project is to bridge that gap. It aims to aid the tester in designing and creating test suites so as to provide automation in the testing process. Automation, simplicity, and ease of use are key tenets of this project. Making a tester’s job easier is much more likely to help them do it well, and making this tool accessible and easy to use is a low-cost, high-convenience solution that not only codifies a standardized testing procedure but shortens the duration of the testing phase and ensures resilient software before its deployment date.

Chapter 2

Literature Survey

One key problem of software testing is knowing what to expect in the output of any given program. The essence of the oracle problem is best captured in the definition of non-testable programs, provided by Weyuker:

A program should be considered non-testable [if] (1) there does not exist an oracle; (2) it is theoretically possible, but practically too difficult to determine the correct output [7].

There are a variety of reasons why the oracle does not exist or cannot be exercised in a practical capacity. Weyuker describes that some software systems are like magic calculators – they are meant to inform us of the answer. She additionally mentions that some programs produce output that is impossible to read, either because of its volume or complexity. Here lies a core challenge of software testing. It is impossible to solve for every oracle or be able to describe every detail of a complex oracle - every system is entirely different. It goes without saying that this is even more impossible when there are time and resource constraints on the software development timeline. Consequently, a number of researchers have attempted to automate the creation of oracles [3]. For example, one research group used natural-language processing techniques to extract Java comments and design test cases that generate exceptions [11]. Another group identified and compared outputs of redundancies in software, using

them as pseudo-oracles [6]. However, the problem of generating complete, automatic oracles remains unsolved. While my project will not make any attempt to create oracles, it will provide information about the input. The test data will be constructed backwards in a sense. Instead of creating random data and then trying to discover its patterns, I will describe a trend and create data that fits that trend, in the hopes that knowing the input's shape will give meaning to the expected output.

Some researchers have done away almost entirely with creating end result oracles by developing techniques specifically used on non-testable programs [20, 30, 9, 8, 19]. Metamorphic testing exploits the relationships of different executions of different input data. For example, upon close inspection of a system's implementation, testers can see that outputs should relate directly to their inputs. They may not know much else about the system, or what its output means, but they know what differences they should see upon two different executions. Upon a certain execution they get $f(5) = 20$, and since they know the input-output relates directly, then they should predict that $f(6) > f(5)$. In and of itself, this is a partial oracle.

Other researchers have proposed methods that do not attempt to provide or calculate an oracle, but rather aid the tester, informed of the structure of the system, exactly what details she should watch and constrain about the oracle [28, 31, 15, 32]. The idea is that if an existing test suite can find artificial bugs in a system, it not only can find real bugs but also inform the tester what parts of the code must be correct in order for the system not to fail. This approach is called mutation testing. When a test suite finds artificial bugs, it determines which variables in a system are most effective in revealing faults, and then presents them to the tester to manually define. This makes the construction of an oracle easier.

The problem with metamorphic and mutation testing is that, while they skirt automatic creation of an oracle, they require just as much effort to develop as just creating an oracle by hand. Discovering metamorphic relationships requires a high

level of implementation knowledge. Mutation testing requires a substantial amount of setup work, and still turns to the tester for manually defined expected values. Even though these techniques are considered cutting-edge testing methods, they do not provide extra automation and therefore are not usable enough.

The need for automation is why random testing has developed such popularity. Random testing is not new, but it also is not very useful. Neither *purely* random testing nor purely random data generation is strong enough to adequately test all branches of the code, nor does it provide any useful information about the test case itself or the data it contains. Consequently, constrained randomness, or using randomness in conjunction with other techniques, has gained the benefits of randomness' automation with the structure of other established testing methodologies.

The biggest complaint of pure random testing is the distribution type typically implemented: uniform distribution, which is why Lu recommends sampling along sample-space partitions non-uniformly [21], preferring uneven numbers of “Random Node” selections and “Random Edge” selections. Another generalized technique is continuous feedback. Arcuri describes Adaptive Random Testing, which monitors the success of the result of a random test case, and then appropriately chooses the next input to be unlike the previous input to attempt to trigger a different result [2]. All input choices are random, but directed towards certain trends. Arcuri identified a primary problem with this technique: it often finds repeat bugs and does not extend to as many situations that pure random test cases otherwise could have found. However, recently, several groups similarly use continuous feedback loops to identify the next test cases but with greater control over random stepping so as to address Arcuri's concerns [34, 16, 29]. Finally, one group used random grammar derivations limited to a finite size to generate programs to test compilers [27]. Again this harnesses the automation of random while limiting its recklessness with additional constraints that guide the random selection.

Applying constraints on random data generation is particularly topical to areas of machine learning and database-driven applications. A team from Columbia University constructed a framework for generating test data for machine-learning models that produce ranking decisions on likely failure of a network of devices. Not only did their framework take in to account the number of examples and number of attributes of a dataset, it considered appearance of repeat values, missing entries, and partitioned categories of examples [25, 24]. However their tool only implemented uniform distribution and could only produce positive integers. This is a keystone feature that indicates the specialization of their framework, and makes extensibility to other applications of software testing nonexistent without refactoring. In regards to database-driven applications, there is a considerable lack of developing in testing. The typical approach of data generation is limited in its power. That data really only passes syntactic checks, like matching database schema descriptions, and general “not null, unique, primary, and foreign key constraints” [14], which indicates that the links between data tables makes sense. The data is database-compliant – it does not violate any database constraints – but it contains no semantic meaning. It does not exhibit any trends or patterns. The shortcomings of these applications influence the desired features of my project.

Chapter 3

Proposed Argument

I have designed an approach to these problems and have architected a straightforward solution that automates the steps described in the introduction. I hypothesize that Problem One is solved by the advent of combinatorial testing, which identifies most, if not all, of the test cases necessary for confidence in working software. My framework, Parmgen, has addressed the second problem, translating the English description of a test case into a concrete format that can be used as input into the program. I ask the user for a simple description of their input data, which requires less work than writing every data point by hand. Then I can use Python’s random statistical distributions to generate those points and write them to a local file convenient for the tester’s use. Figure 3.1 describes the user interaction with the architecture of tool. Red indicates user-written files. Blue indicates the output, the test data that can then be used for testing. This tool is written in the Python language, which is useful for prototyping the MVP (Minimum Viable Product). I call this entire pipeline GenSequence.

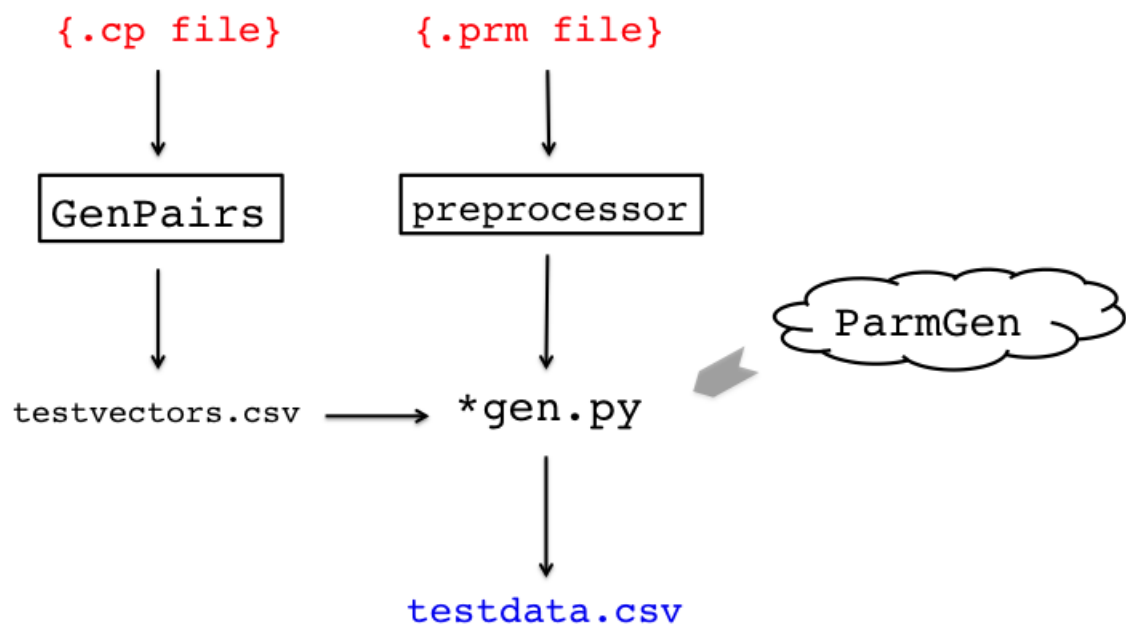


Figure 3.1: Basic Diagram of User's workflow with GenSequence

Chapter 4

Background on Design Decisions

4.1 Parameters

The “input” to a program is a set of parameters, named variables that represent specific groups of data points. How those parameters are generated, and what literal form they take, affects the outcome of the program. Large programs with larger inputs may benefit from having their input parameter data sets generated by statistical distributions.

I was inspired to create this tool by testing a project I helped build in my software engineering class, which took in the records of each student in a class of 40 students, describing things like their skills in certain technologies and their available free times during the week. This project made a fairly sophisticated decision about how to form optimal groups of students to work in teams, ensuring that each group has complementary skills and overlapping free times. During testing, it was important to consider bizarre scenarios, like what should happen when no student has any free time at all, or when half the class can only meet early in the week and the other half can only meet late in the week. If the tool can gracefully handle situations like this, it gives a tester confirmation that the software functions the way she expects it to. These “edge case” scenarios scarcely happen in real life. But it is important to test them because they might happen in real-life, and in a production environment the

software needs to be able to handle it. Moreover, correct behavior in a wide variety of situations increases confidence in both the likelihood of correct behavior in other untested situations and the infrequency of bugs.

These parameters can follow certain trends that can be mimicked by repeated trials of statistical distributions. The input to this program is merely an Excel csv file of a record table, such that every row is a student's entry and every column is all the entries for a certain parameter that really effects the outcome of the program. Figure 4.1 shows the CSV file of the literal input into the Team Building program. The software takes in the filename, opens the file, and parses the data. Each row is a student, and each column is a parameter, such as free time availability on Monday or self-rated HTML experience level.

</

4.2 Pairwise Testing

GenPairs is a pairwise testing tool [36] that creates symbolic test vectors that describe what each test case should generally look like, using English adjectives to describe what the input should be. Pairwise testing rests on the hypothesis that a majority of faults in a program are a result of combinations of parameters that work with one another to create or amplify a problem. Creating the minimum number of test cases that satisfactorily describe all combinations of possible parameters gives adequate code coverage. The success of finding all faults increase with greater numbers of combinations of parameters; that is to say that all combinations of four parameters will find more faults than all combinations of three. However, pairwise testing finds all combinations of two parameters and is sufficient for a diverse test suite while still reducing the time, and therefore cost, of creating or running potentially hundreds of other test cases. This is even more effective when the total number of parameters is also somewhat low. Combining only two parameters is a matter of convenience and time, since more test cases require more time to create and to execute.

While translating a symbolic test vector to concrete data may present a challenge, it is extremely valuable in reducing the uncertainty of knowing what sort of trend the input test data follows. Say, for example, a tester needs to generate a concrete point for this test vector (Table 4.1)

<i>Item Purchased</i>	<i>Price</i>	<i>Delivery Method</i>
Large Item	Expensive	Ultrafast

Table 4.1: Single Symbolic Test Vector

She knows that some test case will send a large item by ultrafast means. In our generation of concrete data (described in the next section) in this test vector, she might come across the test case in Table 4.2

Motorboat	18700.00	Jet Airplane
-----------	----------	--------------

Table 4.2: Single Concrete Test Case

The interaction of these two parameters, a very large object sent by ultrafast means, ought trigger a bug in the system, because her company does not have that capability. If this situation is not handled correctly in the software, the user of the software, the customer, will be misled about the quality of their shipping network. This example is very trivial, but it exemplifies the problem well, and shows how pairwise testing methods can bring this awareness to the tester and developer and help them improve their software. Now imagine a much longer and complicated test vector in a very advanced version of the company's ecommerce website; maybe now it has millions of possible item categories and hundreds of shipping methods and can also handle coupon codes and student discounts and special offerings: many more parameters to consider. It is helpful to know that the test vector was created from a description including Large Item - Expensive - Ultrafast. The symbolic test vector describes the input well enough for the tester to understand and use common sense to predict what the output ought to be. Then if she tests that test case and finds that the customer receives a message that sending a boat through lightning speed mail is impossible, she gains more confidence that her program generates the right output.

The pairwise testing tool is designed such that the tester describes all the possibilities for what each parameter could possibly be, and then the tool generates all the test vectors that describe what the parameter ends up becoming. When using this tool, testers should list every kind of statistical distribution that they would be interested in seeing, or what makes sense in the context of their program. This is so important because it makes GenSequence a more applicable tool for database-driven applications. If a tester would benefit from giving meaning to the data trend, they

can benefit from using the GenSequence framework. This is the semantic sense that is often missing in testing database programs.

Pairwise testing also satisfies the need to test “outlier” cases. The very methodology of pairwise testing ensures that a combination of classifications of parameters occurs at least once in the test suite. This is a test case that the tester might have otherwise forgotten or overlooked. It guarantees robust testing.

Finally, starting with the pattern first and then generating data can possibly give insight into the expected output of the SUT. Instead of generating a set of concrete points and trying to calculate its trend, working the opposite way makes understanding the input easier. Additionally, because there is a natural connection between input and output, understanding the input’s trend may indicate to the tester what to expect in the execution of their program on that test case, especially if their measure of correctness is performed by eyeballing the result.

4.3 Context-Free Grammars

We can use the ideas of context-free grammar derivations to generate concrete test data. We call this feature Makogram, a tool implemented in the Python templating tool Mako [1]. Say a test case has the following production rule:

$$\langle Case \rangle ::= \langle Item \rangle \langle Price \rangle \langle Delivery \ Method \rangle$$

What constitutes an item? Item might have a production rule

$$\langle Item \rangle ::= \text{‘Ship’} \mid \text{‘Book’} \mid \text{‘Loaf’}$$

The grammar tool used is constructed such that a non-terminal can be the result of a function, and that result can be a randomly generated data point. For example,

$$\langle Price \rangle ::= \text{fun } \{ \text{random_choice}(8,15) \} ()$$

Meaning that the data point for the Price parameter would be some number chosen at equal randomness from any number between 8 and 15.

A terminal symbol can be written to be the returned result of a much more sophisticated function – one perhaps that returns a Gaussian distribution of data points. For example, consider an earthquake modeling program, which takes a set of data points for its magnitudes parameter. It can have all sort of descriptions attached to it: Gaussian, uniform, triangular, etc. All of these descriptions will occur somewhere in a test vector, so each is mapped to the corresponding function. When the magnitudes parameter is set to be generated from a bell-curve distribution, the bell-curve random generator will generate that set of points.

Generating good test data for actual magnitudes of an earthquake must be based on the knowledge that a magnitude cannot be negative or greater than 10.0. This is the limit of the program. It cannot make logical guesses about the natural language of the parameter. But the tool is written to minimize the number of times the tester must identify this constraint – just once, in the beginning.

Moreover, the power of grammars makes creating the data very easy. GenSequence provides support for the *grammar* of the output csv files, which can be thought of as a template describing the “language” of a csv file. A derivation is a series of rule substitutions that eventually ends with a string. This “string” belongs to the language described by the grammar, and this string *is* the test data that is usable for testers because the production rules gave functional data points. Makogram was constructed not just to limit the number of substitutions but to reach exactly the desired number. So if a tester asks for 30 examples (30 rows of data), Makogram will hand back exactly 30 rows.

4.4 Law of Large Numbers

Given a symbolic test vector, concrete data points must now be generated for it by the random function associated with the description. However, asking any of Python's random number generator functions for exactly the number of points needed will provide a skewed distribution set because the sample size is just too small. The histogram of the data set doesn't actually capture the distribution that the probability prescribes. So GenSequence generates an exceptionally large data set, and strips it down to the desired size. This is an important step for multiple reasons. First, the SUT may not require an exceptionally large data set - it might ask for a specific number of points. Secondly, smaller sets are more convenient than larger ones, in terms of both disk storage space and the SUT's computational runtime. Finally, small is much easier to see with the naked eye. A tester often has to examine the program's outcome just by looking at it; a "small" execution of the program on a significantly fewer number of data points will be much easier to look at. Consequently it will be much easier to measure the program's correctness and be less susceptible to a human's error in interpretation. Ultimately we can guarantee that this smaller data set will have the statistical characteristics of a larger set because of the Law of Large Numbers.

The Law of Large Numbers states that the actual outcomes will approach the expected outcomes as the sample size increases to infinity . One canonical example is a series of coin flips, with an equal probability of flipping heads as tails. It is expected that exactly 50% of the samples will be heads and 50% tails. In a sample size of only one hundred flips, the percentage of heads to tails may only be 46.0% - 54.0%, but a sample size of ten thousand would come far closer to an even 50-50 split, possibly 50.37% - 49.63%.

This principle applies to all kinds of probability distribution types, not just coin flip probability. Students grades tend to figure towards a bell-curve distribution, but a graph of 50 students grades will look more misshapen than a graph of thousands of students grades. Moreover, the principle applies to many programming languages utilities that support randomness. The Python language contains a built-in module called `random` which provides a variety of generators for all flavors of probability distribution. Individual points are generated according to a particular distribution, and the generation of a large enough sample size is guaranteed to follow the expected trend.

How large is large? Bernoulli proved that the actual outcome of a simulation would approach the expected outcome as the sample size grows to infinity. So in general, the greater the sample size the more accurate the results.

If the use case arrives on the order of hundreds of data points, a default sampling size on the order of ten thousand randomly generated data points should suffice. Because the initially generated set of points is far too many, I have devised a scheme to reduce the sample set to the appropriate size while still encapsulating the distributions guaranteed by the Law of Large Numbers. For each parameter, I generate the too-large sample size, sort it in increasing order, and selectively choose every n th point to include in the reduced set of points, such that the size of the reduced set is the desired number - however many the test case needs.¹

I have yet to rigorously prove the reliability and strength of this selection scheme, but I hypothesize that it suffices for a first version of `GenSequence`, and does indeed capture the essence of the Law of Large Numbers. Pictorial representations show that

¹Given the statements of the previous section, some clarification is in order. Context-free grammars would imply that the random function would be called upon every derivation of a data point. But since the number of derivations is limited, that would mean that the number of function calls would be minimized, which would not produce the desired distribution. `GenSequence` declares the grammar, and then uses what I call a `Parm` object to create a large sample size and reduce it down into a small pool. When a derivation of the grammar is instantiated, it calls upon the pool, which yields singular points.

this paradigm is not perfect, but it certainly exhibits the desired statistical “shape” better than choosing exactly the desired number of points. Visibly a sample size of 300 systematically chosen from a large set follows the curve of 10,000 points much more closely than 100 initial points.

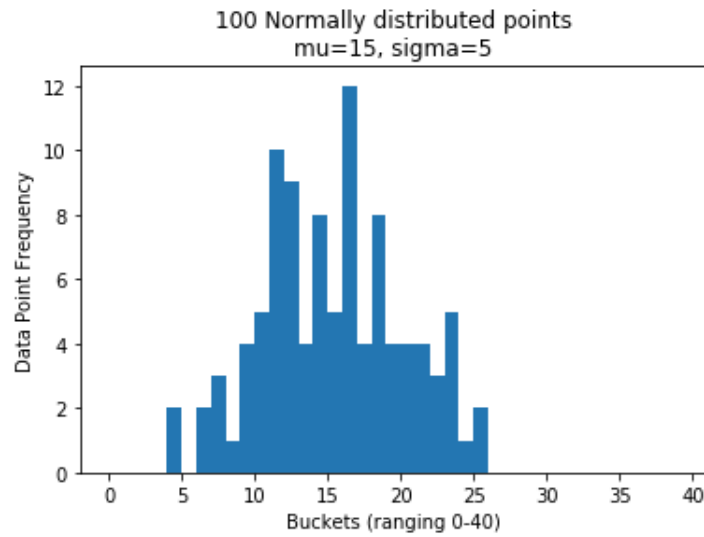


Figure 4.2: Histogram of point values generated by normal distribution - sample size 100

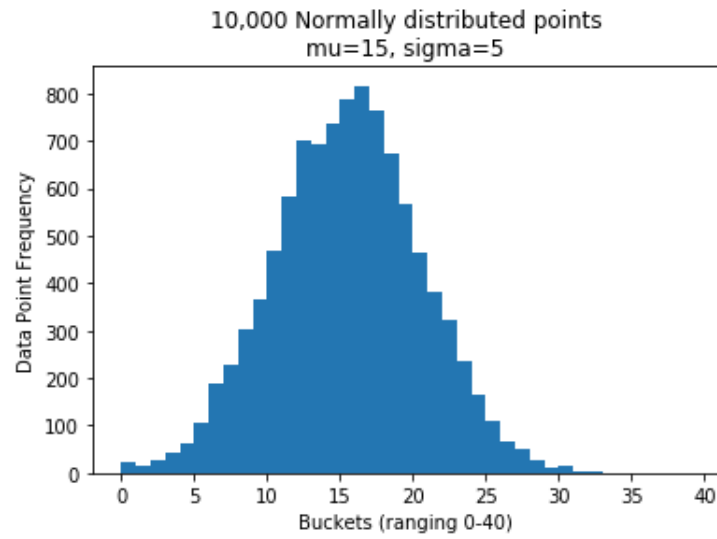


Figure 4.3: Histogram of point values generated by normal distribution - sample size 10,000

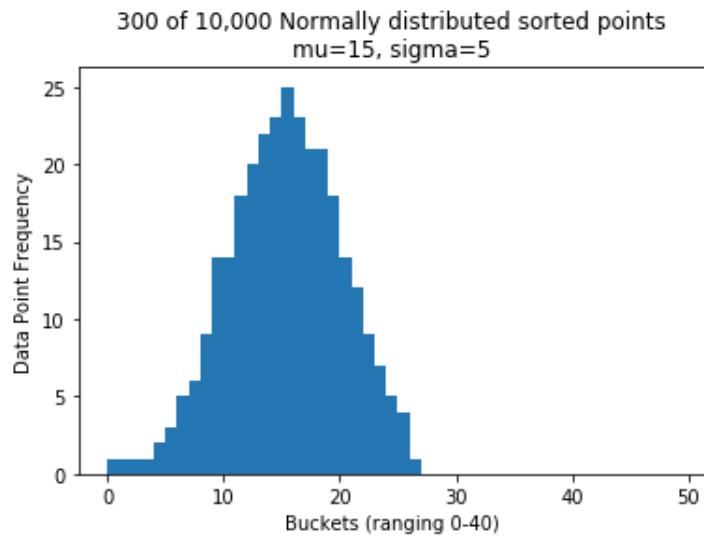


Figure 4.4: Histogram of point values generated by normal distribution and reduced to sample size of 300

4.5 Cardioids

In some testing situations, it may be useful to consider the relationship between two parameters, a relationship that mimics real life. For example, a physics-modeling simulation might have a loose property that larger objects move slower than smaller objects. These two parameters, size and velocity, are interdependent.

A Cardioid² object acts as a two-column parameter. The end user specifies the relationship between the two sub-parameters by providing descriptions of what data point range pairings are preferred (called favorites) and what data point range pairings should only occur infrequently (called outliers or non-favorites). The data generation step will look at these descriptions and generate the data point pairings such that 90% of the samples are favorites and 10% are outliers. For example, the tester will specify favorites for the cardioid relationship between size and velocity: large and slow, medium and medium, and small and fast. Non-favorites are outlying particles: small and slow, and large and fast.

The cardioid relationship spanning two parameters occurs only if the tester specifies it as a possibility in the test vector creation step. There are certainly situations where it does not make semantic sense to use cardioid relationships, and in that case, cardioids are not required for test data generation. On the other side, statistical distributions may not make semantic sense because they do not raise *every* test case one wants to check. Cardioids allow the user to include their own test data generation constraints if statistical randomness does not quite suffice, and allows for tests that maintain or exaggerate realistic sense. Again this is a tool that helps testers identify and execute edge cases.

²The name cardioid was chosen for its polarity. Its figure encompasses an area such that a majority of the space (about 90%) sits on one side of an axis and a small section (about 10%) sits on the other.

4.6 Preprocessing

An important arm of the automation pipeline is developing a little language [5], that is easy to write yet fully describes everything a tester wants out of this tool. Then a preprocessor can read this language and generate the code that, when executed, will create the data. Writing the data-generating code itself is complex enough that it limits automation, and building a machine that performs this step systematically significantly increases ease of use. I have fully developed the language into basic syntax and grammar rules, which the tester can then write in a .prm file.

I have also constructed a simple preprocessor to parse this information. It is very lightweight in that it compiles quickly, but does not have advanced error handling or give helpful compiler errors. The current iteration is implemented in Ply (Python Lex-Yacc) [4]. Ply works in two general steps. First it reads the input and identifies characters or words into tokens - essentially describing the type of each word or character. This is the lexing, or “tokenizing”, step. Next, the parsing step, parses the tokens altogether by identifying the grammar rules that combine them. Once the information has been identified as a derivation of a grammar rule, it is then available for interpretation. The preprocessor does not yet use Ply to read the entire .prm file, and mostly uses Python file-processing techniques, but the language is simple enough that I had little difficulty rendering the information into the data generator. For an MVP this tactic suffices.

Chapter 5

Methods

One key part of this project is deciding if it accomplished its goal and answered its research questions.

1. Is using random statistical distributions a valid approach to writing test data?
2. Do symbolic test vectors provide enough information about the input that it informs expectations of the output?
3. Does GenSequence make a tester's job easier or save her time?

GenSequence's intended use is for new projects, pieces of software that do not already have a monolithic code base. However, testers might use this tool at any stage during the software development cycle if they lack a substantial test suite. Perhaps their application is highly advanced, has been neatly designed, and has fixed bugs found through other means of testing or from end user reports. Therefore, that piece of software simply may be bug-free. That is unlikely, but it is possible. So measuring the usefulness of GenSequence on fault finding ability alone would depend on the state of the software it is testing, giving the indication that GenSequence is inadequate when really the software is just really well done.

It is important to determine how valid the random statistical distribution approach is. I will use a few case studies to determine if knowing how a parameterss

data was generated helps clarify the expected result. First I will test my testing tool against a planetary orbits simulation. Then I will observe my tool's data representation in an earthquake analysis program.

In response to both Question 2 & 3, the intention of this project was to streamline the testing process, and provided automated methods of data creation. The usefulness of this program is measured by the ease of test suite creation. I will consider these general metrics: the length of the script generating data, how much code a tester would have to write, and how quickly the data is generated (more data takes more time). If it is easier for a tester to write a GenSequence spec than to write a full test suite by hand, then GenSequence is a valuable tool, and it might encourage developers to write more thorough test suites.

Ultimately it is a tester's decision as to how useful this project is. More features added to the project only increase the ease of use, since a machine does more of the heavy lifting than a human brain and body otherwise would have. However, it is only useful when a test suite does not exist or in circumstances when it would be easier to generate a test suite by hand. There is further discussion later.

Chapter 6

Results

6.1 Celestial Body Simulation

Initial use of GenSequence on the planetary orbits program proved some promise in the capability of this tool. The tool generated 30 test cases. Test case 13 was programmatically named:

```
13-70-mass|right_slanted-position|right_-  
slanted-velocity|uniform-diameter|left_slanted.csv
```

13 is the test case number and 70 is the number of rows in the test case. Each row of the file represents an instance of a planetary body. GenSequence has encoded the way each parameter was generated in the filename. So `mass|right_slanted` means that planetary masses were generated with a right slanted distribution. Position, the location of each body, was generated with a right-slanted triangular probability. The velocities of the bodies were generated uniformly. The diameters, the visible size in the program's GUI, were generated by a left-slanted distribution. The starting point of this program's test data showed this visualization:

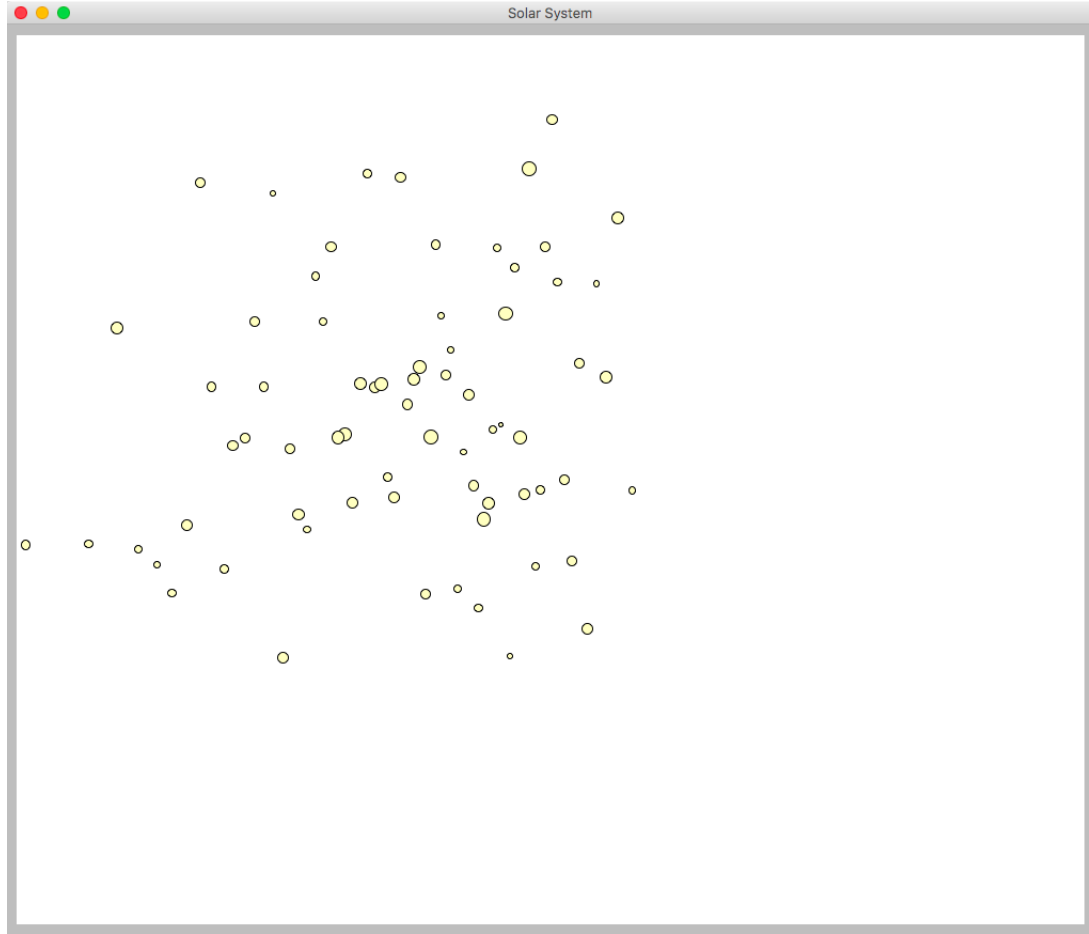


Figure 6.1: Starting Locations of 70 celestial bodies

Figure 6.1 accurately reflects the visible parameters position and diameter. The positions were generated by a right slant and, ostensibly so, the bodies are mostly located in one region. The diameters were generated by a left-slanted distribution, clearly slanted towards a smaller diameter, with most samples appearing small or medium size, but a notable number appearing fairly large (though from the naked eye one can agree it does require close inspection).

The program tracks and draws the motion of each planetary body through time, giving a useful summary of the program's execution by the end of the simulation (665 time steps). Running test case 13 resulted in Figure 6.2.

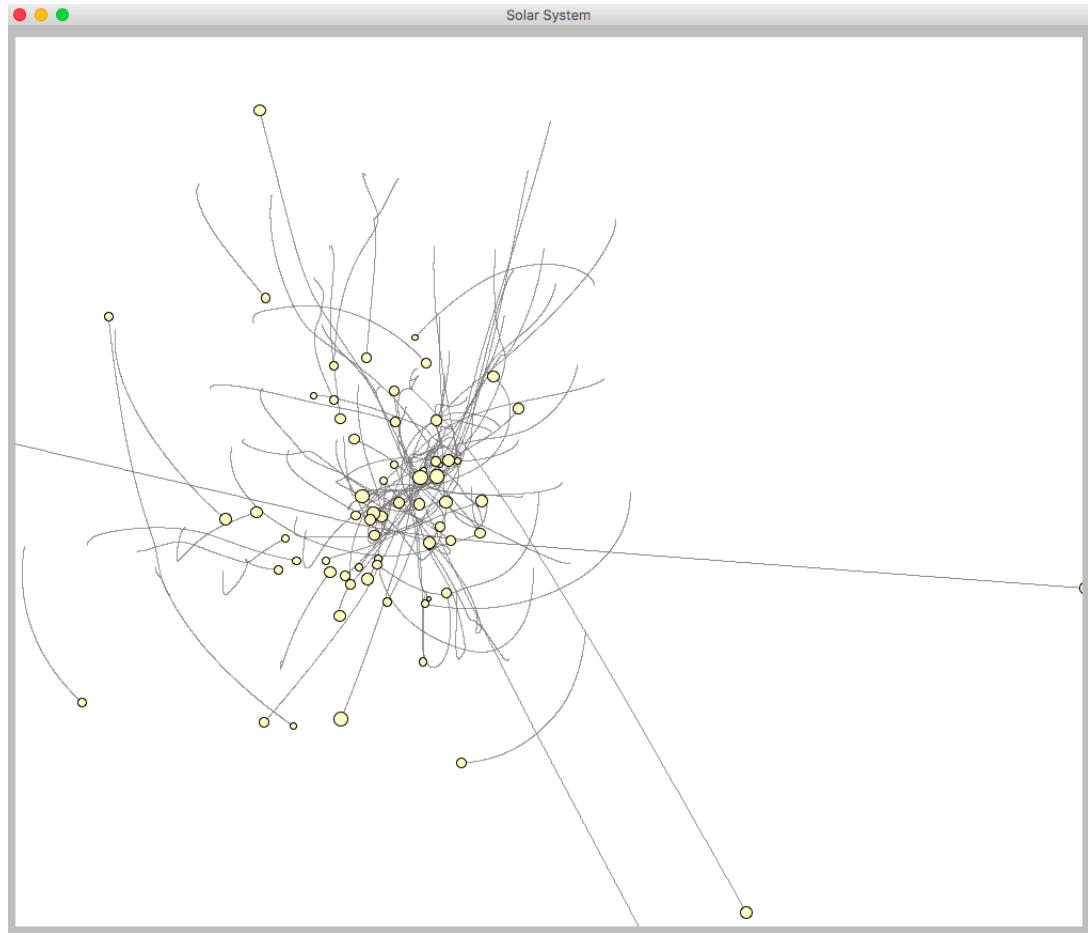


Figure 6.2: Ending Locations of 70 celestial bodies after 665 time steps

Spending a minute closely observing these paths registers some important observations about the natural behavior of gravity acting on planetary bodies. Some disobeyed their trajectory and were redirected by the cumulative mass in the center. The initial velocities manifested themselves as well. The paths are of all different lengths, which makes sense since they were generated by a uniform distribution.

6.2 Earthquake Analysis and Visualization Program

The earthquake analysis program performs basic statistical analysis of magnitudes, locations, and depths on a map, and also plots quake events on a map. Bigger

dots represent a greater data point value (This is important to note because dots representing magnitudes describe their intensity, not their destruction coverage). One test case had a file name (test vector) of:

```
6-70-magnitudes|cardioid-latitudes|left_-  
slanted-longitudes|right_slanted-depths|cardioid.csv
```

Observing the cardioid relationship between latitudes and longitudes is confirmation of GenSequence working. GenSequence was programmed with the following information:

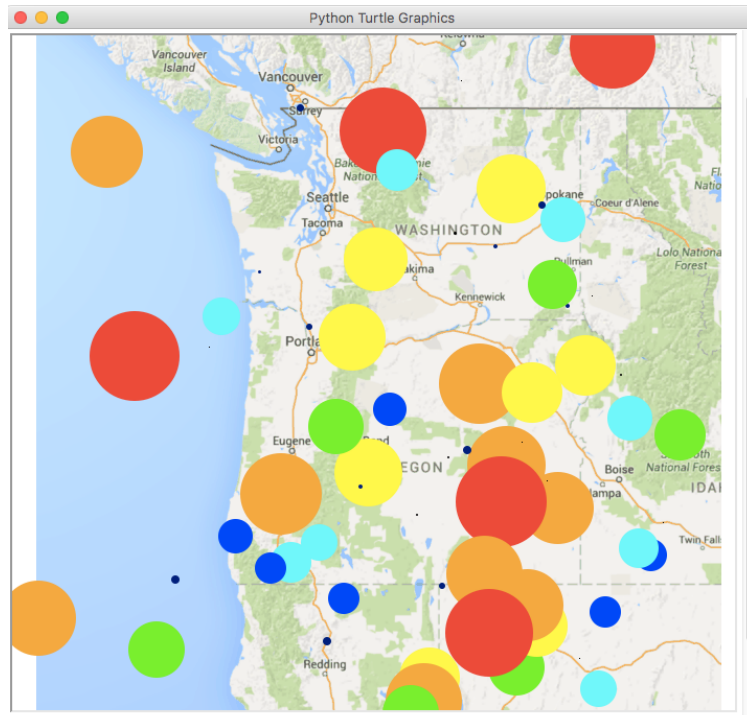
```
1 #magnitudes ranges  
2 Micro = Range(0.0, 2.0, exclusive_upper=True)  
3 Feelable = Range(4.5, 7.9, exclusive_lower=True)  
4 Great = Range(8.0, 9.5, exclusive_lower=True, exclusive_upper=True)  
5 #depths ranges  
6 Shallow = Range(0.0, 5.0, exclusive_lower=True, exclusive_upper=True)  
7 Mid = Range(5.0, 15.0)  
8 Deep = Range(15.0, 30.0, exclusive_lower=True)  
9 ...  
10 # specify the relationship between magnitudes and depths  
11 MagsDepths = Cardioid(Mags, Depths)  
12 MagsDepths.setFavorites([(Micro,Shallow), (Great,Deep), (Feelable,Mid)])  
13 MagsDepths.setNonFavorites([(Micro,Deep), (Great,Shallow),  
    (Feelable,Deep), (Feelable,Shallow)])
```

This is the literal implementation of the tool describing the most frequently occurring data pairs in this test case. Low intensity magnitudes should be near Earth's surface (values closer to 0) and very intense magnitudes should occur deep below the earth's surface (values closer to 30). Moreover, it is not very frequent that

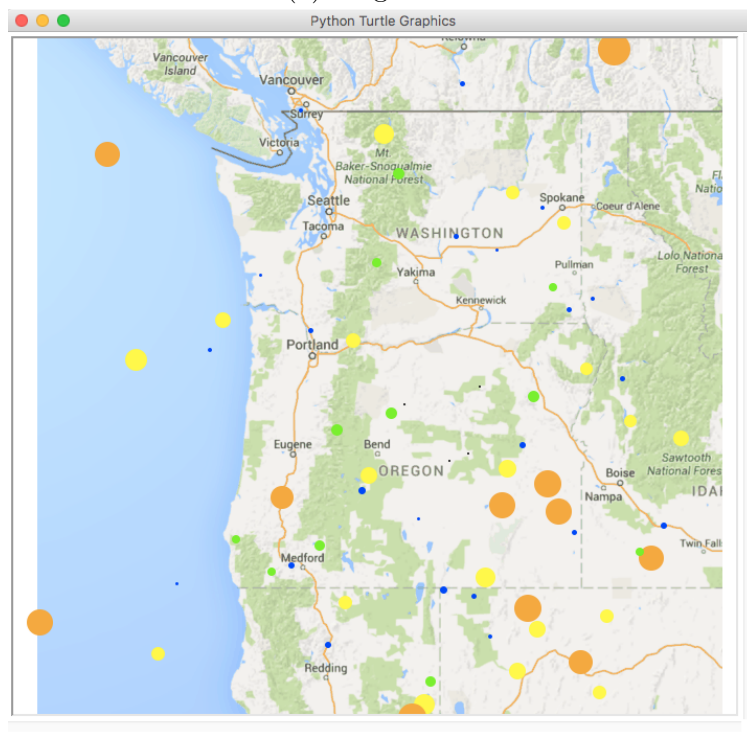
low-intensity earthquakes occur deeply, and high-intensity earthquakes occur near the surface.¹

The plots of magnitudes and depths (Figure 6.3) show this constraint propagating:

¹This assumption may be entirely false about the nature of earthquakes. This trend is one I created to add functionality to GenSequence, one that could easily be reversed by any seismologist user of the program.



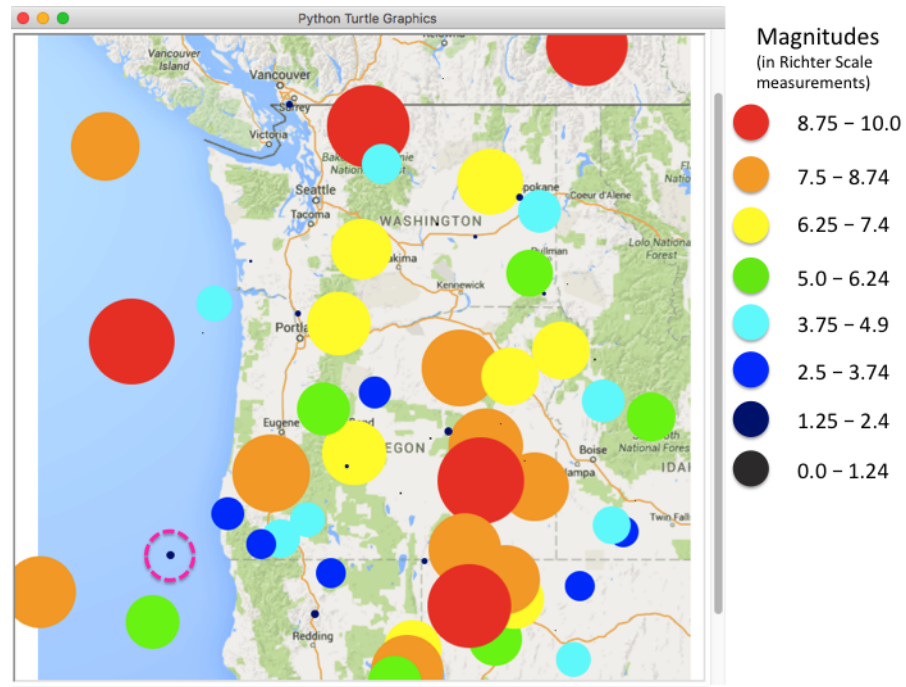
(a) Magnitudes



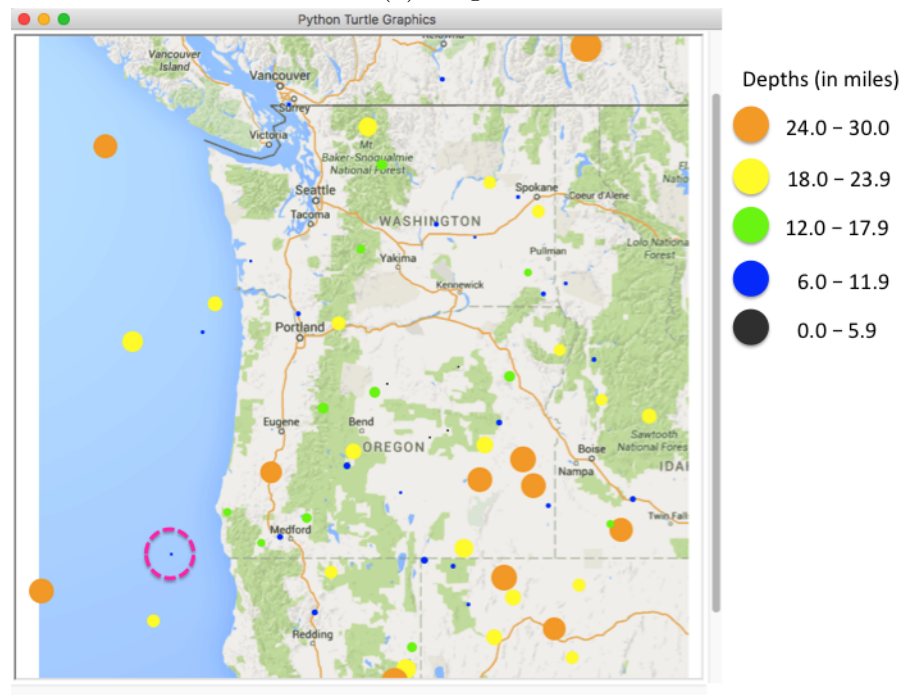
(b) Depths

Figure 6.3: Plotting Magnitudes and Depths from the same test case

The correlation reflects what is expected. The large red dot due west of Portland has a matching yellow depth dot, which means that a high-intensity quake was fairly deep (Figure 6.4). The small blue event due west of the Oregon-California border has a matching blue dot in the same place, indicating that a low-intensity quake was not that deep below the surface (Figure 6.5). One outlier is a sizable orange earthquake covering the Umatilla forest that has a very tiny blue depth dot (Figure 6.6).

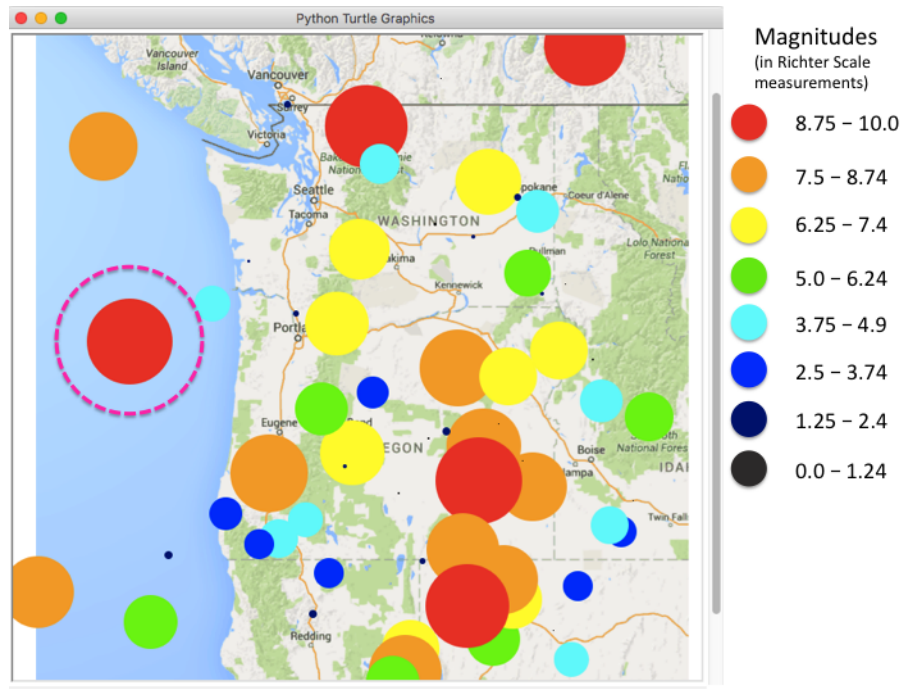


(a) Magnitude

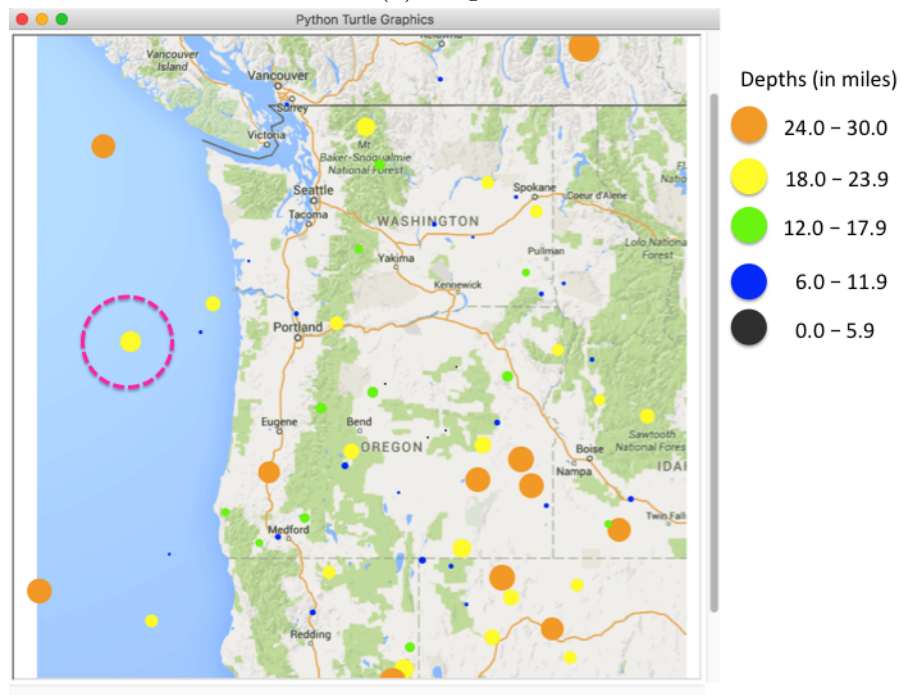


(b) Depth

Figure 6.4: An Expected Earthquake Event – High Intensity, Deep Below Surface

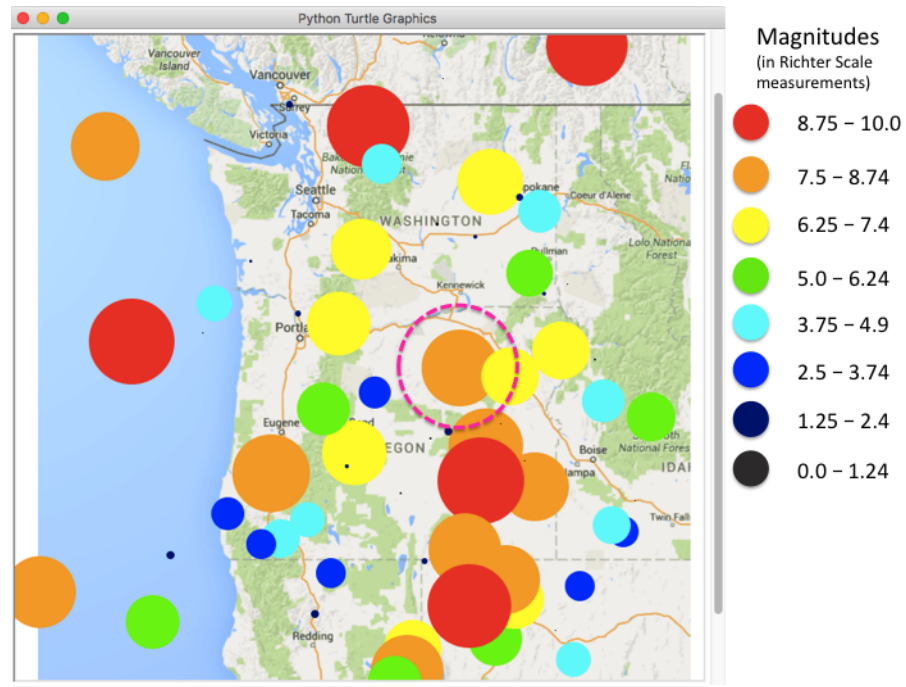


(a) Magnitude

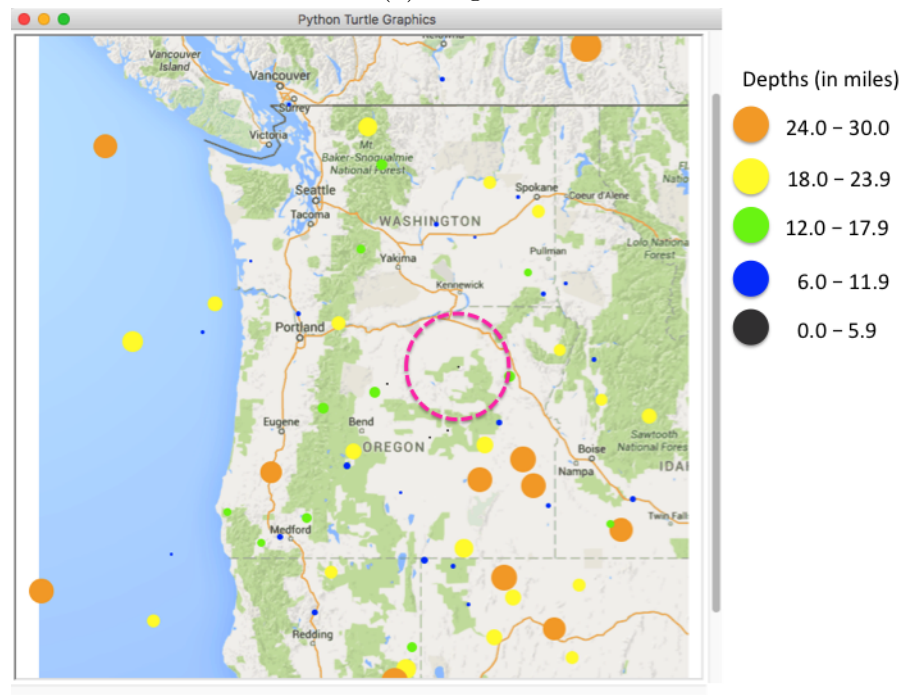


(b) Depth

Figure 6.5: Another Expected Earthquake Event – Low Intensity, Near Surface



(a) Magnitude



(b) Depth

Figure 6.6: An Outlier Earthquake Event – High Intensity, Near Surface

If we consider latitude-longitude oriented north-up, west-left, east-right, south-down, and given the left-slanted-ness of the latitudes and right-slanted-ness of longitudes, we should expect right-leaning longitudes and left-leaning latitudes. This test case should mark the locations of the earthquakes mostly drifting towards the lower-right corner. This is in fact what the program generates (Figure 6.7).

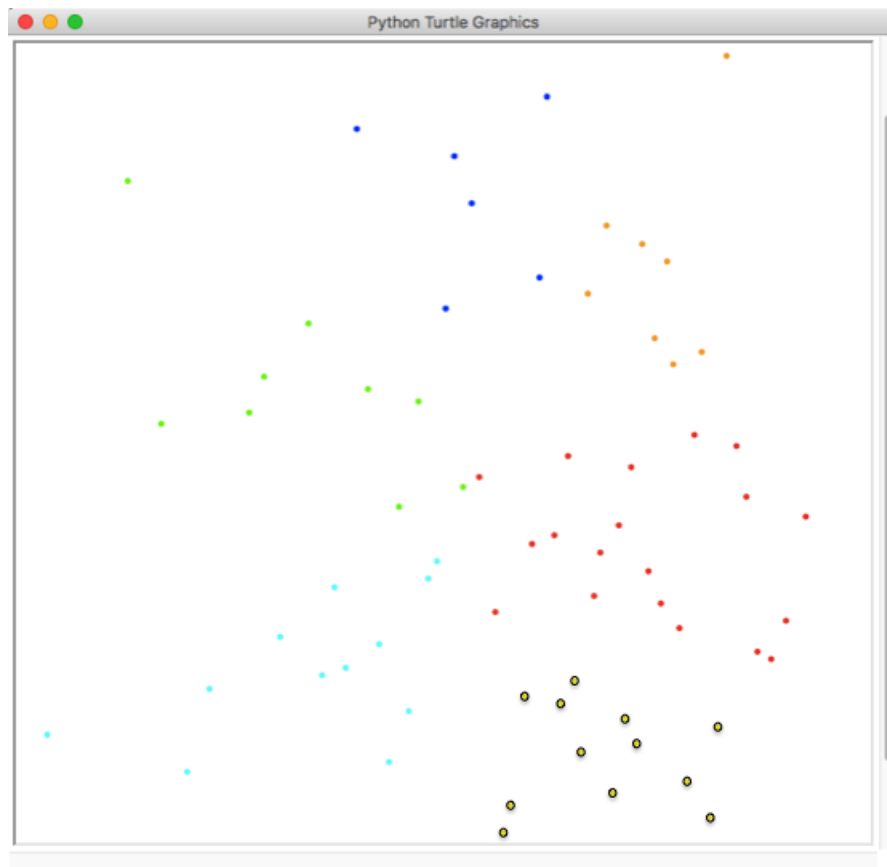


Figure 6.7: K-Means Clustering of Earthquake events from case 13

Chapter 7

Concluding Thoughts

I believe that GenSequence’s best application is for testing database-driven applications. These applications literally control the world. Consider financial applications alone: the record of every transaction, every payment, every credit report, and every bill can dictate a person’s entire life. The infrastructure surrounding that data must not expunge or fabricate any of it, and must maintain its integrity. Software that accesses, controls, and manipulates this huge amount of data carries substantial responsibility. Moreover, software that can read, interpret, and identify trends in a huge wealth of global data has amazing power in informing us of what happens in the world and how we can make better decisions. It is therefore of utmost importance to design that software well enough to trust its results.

GenSequence has addressed those demands. It functions simply but provides automation during the testing stage and generates reliable data that is what it says it is. If the context of the program warrants consideration of statistically unlikely but possible scenarios, GenSequence can make that happen. It primarily shows its power by its ease of use and nearly end-to-end automation. Moreover, the construction of data points does provide insight into the expected output. However, this insight is limited to eyeballing the result.

One question that remains unanswered is if knowledge of the input’s construction can produce any expectation of the output. With the smaller programs that GenSequence was tested against, it is difficult to answer that question since the programs tested against GenSequence are limited in breadth. Ideally I would test GenSequence against a program that has been largely untested at the system level and may still have bugs. In that situation, GenSequence’s symbolic test vectors might be extremely helpful in increasing confidence that the software works or else identify flaws in the program. Knowing the symbolic test vector might even help a developer debug their work because the pair of conflicting parameters is clearly displayed. Again, these claims are unproven, and would be the first step in improving GenSequence moving forward.

Given that the construction of GenSequence was a Minimum Viable Product, there is considerable room for improvement. What follows is a discussion of future work and extensibility.

Current Iteration

- Identify an open-source project in need of a test suite, and create a hypothetical test suite for it using GenSequence. The greatest indication that GenSequence is useful is cross-checking it against a program that someone else wrote.

Top of the Backlog

- Implement the entire preprocessor in Ply, develop an abstraction for a model of the Parns, and develop an interpreter that could execute those models. Currently the preprocessor performs most of its parsing via line-by-line file processing, which is an inelegant and clunky solution.

- Combine the specifications that enter GenPairs and the preprocessor into one; in other words, combine the roles of the .cp and .prm files. This would certainly reduce the work for the tester. In initial design of the prm language I had included distribution possibilities for each parm, but had removed it from prm since it was already specified in cp language that went into GenPairs. Ideally this specification could be parsed and found from the prm file and automatically sent into pairwise generation.
- Develop greater control over Double cardioids - relating multicolumns together. I had this idea when developing data for TeamBuilder. I wanted control over the relationship between general skills (columns Python and Java) and more niche, specialized skills (like Haskell and C pthreads library). This required generation of two cardioids, and then calculation to determine what pattern each data instance matched, and then re-sorting by implementing a hill climbing technique. This quickly fell out favor, as much more of the system needed implementation first.

Icebox

- Consider alternate formats for the output. Currently GenSequence will give the tester a csv file of raw data. However in database-driven applications, the test data must come in scripts of database insert statements, as the database state must be restored for every test case. Consequently it is worth implementing a utility feature that outputs a database creation script, say for example a shell script containing SQL insert statements.
- Rewrite documentation, and provide utilities that allow for easy initial use. A core component of automation is that it is easy to set up, if even for a proof-of-concept.

- Develop a greater understanding of how pseudo-random number generation works, and its influence on data generation according to distribution functions.
- Research and identify possible Machine Learning applications in which GenSequence could be used. I believe that GenSequence might have applications in some ML models, in which case knowing the demands of ML models might influence the way GenSequence is used.

Bibliography

- [1] Mako templates for python. <https://github.com/zszeek/mako>, 2017.
- [2] Andrea Arcuri and Lionel Briand. Adaptive random testing: An illusion of effectiveness? In *Proceedings of the 2011 International Symposium on Software Testing and Analysis*, ISSTA '11, pages 265–275, New York, NY, USA, 2011. ACM.
- [3] E. T. Barr, M. Harman, P. McMinn, M. Shahbaz, and S. Yoo. The oracle problem in software testing: A survey. *IEEE Transactions on Software Engineering*, 41(5):507–525, May 2015.
- [4] David Beazley. Python Lex-Yacc. <https://github.com/dabeaz/ply>, 2018.
- [5] Jon Bentley. *Programming Pearls*. ACM, New York, NY, USA, 1986.
- [6] Antonio Carzaniga, Alberto Goffi, Alessandra Gorla, Andrea Mattavelli, and Mauro Pezzè. Cross-checking oracles from intrinsic software redundancy. In *Proceedings of the 36th International Conference on Software Engineering*, ICSE 2014, pages 931–942, New York, NY, USA, 2014. ACM.
- [7] David Chays, Saikat Dan, Phyllis G. Frankl, Filippos I. Vokolos, and Elaine J. Weyuker. A framework for testing database applications. *SIGSOFT Softw. Eng. Notes*, 25(5):147–157, August 2000.
- [8] Junjie Chen, Yanwei Bai, Dan Hao, Lingming Zhang, Lu Zhang, Bing Xie, and Hong Mei. Supporting oracle construction via static analysis. In *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering*, ASE 2016, pages 178–189, New York, NY, USA, 2016. ACM.
- [9] Tsong Yueh Chen. Metamorphic testing: A simple method for alleviating the test oracle problem. In *Proceedings of the 10th International Workshop on Automation of Software Test*, AST '15, pages 53–54, Piscataway, NJ, USA, 2015. IEEE Press.

- [10] Vahid Garousi, Michael Felderer, Marco Kuhrmann, and Kadir Herkiloğlu. What industry wants from academia in software testing?: Hearing practitioners' opinions. In *Proceedings of the 21st International Conference on Evaluation and Assessment in Software Engineering*, EASE'17, pages 65–69, New York, NY, USA, 2017. ACM.
- [11] Alberto Goffi, Alessandra Gorla, Michael D. Ernst, and Mauro Pezzè. Automatic generation of oracles for exceptional behaviors. In *Proceedings of the 25th International Symposium on Software Testing and Analysis*, ISSTA 2016, pages 213–224, New York, NY, USA, 2016. ACM.
- [12] Don Gotterbarn. The creation of facts in the cloud: A fiction in the making. *SIGCAS Comput. Soc.*, 45(3):60–67, January 2016.
- [13] Priya Gupta and Prafullakumar Surve. Model based approach to assist test case creation, execution, and maintenance for test automation. In *Proceedings of the First International Workshop on End-to-End Test Script Engineering*, ETSE '11, pages 1–7, New York, NY, USA, 2011. ACM.
- [14] Klaus Haller. The test data challenge for database-driven applications. In *Proceedings of the Third International Workshop on Testing Database Systems*, DBTest '10, pages 6:1–6:6, New York, NY, USA, 2010. ACM.
- [15] Mark Harman, Yue Jia, and William B. Langdon. Strong higher order mutation-based test data generation. In *Proceedings of the 19th ACM SIGSOFT Symposium and the 13th European Conference on Foundations of Software Engineering*, ESEC/FSE '11, pages 212–222, New York, NY, USA, 2011. ACM.
- [16] John Hughes, Ulf Norell, Nicholas Smallbone, and Thomas Arts. Find more bugs with quickcheck! In *Proceedings of the 11th International Workshop on Automation of Software Test*, AST '16, pages 71–77, New York, NY, USA, 2016. ACM.
- [17] Lawrence Livermore National Laboratory. About VisIt.
- [18] William B. Langdon, Shin Yoo, and Mark Harman. Inferring automatic test oracles. In *Proceedings of the 10th International Workshop on Search-Based Software Testing*, SBST '17, pages 5–6, Piscataway, NJ, USA, 2017. IEEE Press.
- [19] Mikael Lindvall, Dharmalingam Ganesan, Ragnar Árdal, and Robert E. Wiegand. Metamorphic model-based testing applied on nasa dat: An experience report. In *Proceedings of the 37th International Conference on Software Engineering - Volume 2*, ICSE '15, pages 129–138, Piscataway, NJ, USA, 2015. IEEE Press.
- [20] Mikael Lindvall, Adam Porter, Gudjon Magnusson, and Christoph Schulze. Metamorphic model-based testing of autonomous systems. In *Proceedings of the*

- 2Nd International Workshop on Metamorphic Testing*, MET '17, pages 35–41, Piscataway, NJ, USA, 2017. IEEE Press.
- [21] Jianguo Lu, Hao Wang, and Dingding Li. Uniform random sampling not recommended. In *Companion Proceedings of the The Web Conference 2018*, WWW '18, pages 495–499, Republic and Canton of Geneva, Switzerland, 2018. International World Wide Web Conferences Steering Committee.
 - [22] Sasa Misailovic, Aleksandar Milicevic, Nemanja Petrovic, Sarfraz Khurshid, and Darko Marinov. Parallel test generation and execution with korat. In *Proceedings of the the 6th Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on The Foundations of Software Engineering*, ESEC-FSE '07, pages 135–144, New York, NY, USA, 2007. ACM.
 - [23] Philip A. Laplante Mohamad Kassab, Joanna F. DeFranco. Software testing: the state of the practice. *IEEE Software*, 34(5):46–52, 2017.
 - [24] Christian Murphy, Gail Kaiser, and Marta Arias. A framework for quality assurance of machine learning applications. Columbia University Computer Science Technical Reports, pages 1–10, New York, NY, USA, 2006. Department of Computer Science, Columbia University.
 - [25] Christian Murphy, Gail Kaiser, and Marta Arias. Parameterizing random test data according to equivalence classes. In *Proceedings of the 2Nd International Workshop on Random Testing: Co-located with the 22Nd IEEE/ACM International Conference on Automated Software Engineering (ASE 2007)*, RT '07, pages 38–41, New York, NY, USA, 2007. ACM.
 - [26] Kivanc Muslu, Yuriy Brun, and Alexandra Meliou. Preventing data errors with continuous testing. In *Proceedings of the 2015 International Symposium on Software Testing and Analysis*, ISSTA 2015, pages 373–384, New York, NY, USA, 2015. ACM.
 - [27] Michal H. Palka, Koen Claessen, Alejandro Russo, and John Hughes. Testing an optimising compiler by generating random lambda terms. In *Proceedings of the 6th International Workshop on Automation of Software Test*, AST '11, pages 91–97, New York, NY, USA, 2011. ACM.
 - [28] Matthew Patrick, Matthew D. Castle, Richard O. J. H. Stutt, and Christopher A. Gilligan. Automatic test image generation using procedural noise. In *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering*, ASE 2016, pages 654–659, New York, NY, USA, 2016. ACM.
 - [29] Korosh Koochekian Sabor and Stuart Thiel. Adaptive random testing by static partitioning. In *Proceedings of the 10th International Workshop on Automation of Software Test*, AST '15, pages 28–32, Piscataway, NJ, USA, 2015. IEEE Press.

- [30] Sergio Segura, Amador Durán, Javier Troya, and Antonio Ruiz Cortés. A template-based approach to describing metamorphic relations. In *Proceedings of the 2Nd International Workshop on Metamorphic Testing*, MET '17, pages 3–9, Piscataway, NJ, USA, 2017. IEEE Press.
- [31] Francisco Carlos M. Souza, Mike Papadakis, Yves Le Traon, and Márcio E. Delamaro. Strong mutation-based test data generation using hill climbing. In *Proceedings of the 9th International Workshop on Search-Based Software Testing*, SBST '16, pages 45–54, New York, NY, USA, 2016. ACM.
- [32] Matt Staats, Gregory Gay, and Mats P. E. Heimdahl. Automated oracle creation support, or: How i learned to stop worrying about fault propagation and love mutation testing. In *Proceedings of the 34th International Conference on Software Engineering*, ICSE '12, pages 870–880, Piscataway, NJ, USA, 2012. IEEE Press.
- [33] Rajeev Tiwari and Noopur Goel. Reuse: Reducing test effort. *SIGSOFT Softw. Eng. Notes*, 38(2):1–11, March 2013.
- [34] Kohsuke Yatoh, Kazunori Sakamoto, Fuyuki Ishikawa, and Shinichi Honiden. Feedback-controlled random test generation. In *Proceedings of the 2015 International Symposium on Software Testing and Analysis*, ISSTA 2015, pages 316–326, New York, NY, USA, 2015. ACM.
- [35] Andreas Zeller. Search-based testing and system testing: A marriage in heaven. In *Proceedings of the 10th International Workshop on Search-Based Software Testing*, SBST '17, pages 49–50, Piscataway, NJ, USA, 2017. IEEE Press.
- [36] Jamie Zimmerman and Michal Young. GenPairs. <https://github.com/TestCreator/GenPairs>, 2017.