

Remark: TODO: Change `\iftrue` in `commands.tex` `\def\rem` to `\iffalse` in final copy

Remark: In case we don't have enough content, just add more code examples, and maybe compare specific unit testing frameworks in some frameworks.

Remark: Talking about CI (continuous integration) platforms like Travis-CI is also possible

COMP2123 self-learning report

Unit testing

Chan Kwan Yin, Lee Chun Yin

December 23, 2018

Abstract

This report discusses the motivation, available techniques and difficulties of unit testing.

Contents

1	What is unit testing?	1
2	Motivation	1
2.1	Discover bugs early	1
2.2	As a method of specification	1
3	Unit testing methods	1
3.1	Testing for expected result	1
3.2	Increasing the test size	2
3.3	Generating test cases randomly	3
3.3.1	Random output, inverse-evaluated input	3
3.3.2	Property-based testing	3
3.3.3	Reimplement algorithm by brute-force	4
3.3.4	Testing for edge cases	5
4	Unit testing techniques	5
4.1	Code coverage	5
4.2	Test-driven development (TDD)	6
4.3	Behaviour-driven development (BDD)	7
4.4	Summary: TDD vs BDD	8
4.5	Dependency mocking	9
5	Difficulties	10
5.1	Coupling	10
5.2	Test case selection	11

1 What is unit testing?

Instead of just running the program and seeing if the output is correct, unit testing splits down a software into small components and checks if the desired behaviour of each component is correct.

2 Motivation

2.1 Discover bugs early

As a software grows in scale, there would be some "stale" components of the software unmodified for a long time. New components are built upon these older components, but if the older components are not stable enough, their bugs would propagate to the new components, creating confusion for new developers.

In addition, if a bug is only discovered a long time after writing the corresponding code, the developer is very likely to forget about what the code was intended to do.

Unit testing allows identification of bugs as soon as possible with little impact, overall reducing the development cost.

2.2 As a method of specification

Unit tests can also be used as a means of project requirement specification. It is common for programmers and users to misunderstand each other's requirements and waste a lot of time. Unit tests can provide a technical and strict definition of behaviour.

3 Unit testing methods

3.1 Testing for expected result

The intuitive way is to write a test that tests the output of each function.

If we have a `fooBar.cpp` with the following definition:

```
1 #include <string>
2
3 std::string fooBar() {
4     return "qux";
5 }
```

To ensure `fooBar()` always return "qux", we can write a test to test this behaviour:

```
1 #include "fooBar.cpp"
2 #include "assert.h"
3
4 void testFooBar() {
5     ASSERT_EQUAL(fooBar(), "qux");
6 }
```

The `ASSERT_EQUAL` macro function would compare the result of `fooBar()` with "qux" and trigger an error if they are not equal. This macro function can be implemented very easily:

```
1 #ifndef ASSERT_H
2 #define ASSERT_H
3
4 #include <iostream>
5
6 #define ASSERT_EQUAL(actual, expect) {\
7     auto actualValue = actual; \
8     auto expectValue = expect; \
9     bool eq = actualValue == expectValue; \
10    if (eq) { \
11        std::cout << std::string("Successful: ") + #actual + " = " + actualValue << std::endl; \
12    } else { \
13        throw std::string("Expected ") + #actual + " to return " + expectValue + ", got " +
        actualValue; \
14    }
15 }
```

```

14     } \
15 }
16
17 #endif

```

In this implementation, if the values are not equal, an exception string is thrown.

Different unit testing frameworks may have different error behaviour, and some are able to integrate with IDEs for advanced analysis.

Some other common assertions include:

- Null checks
- Arithmetic comparisons $< \leq > \geq$
- That an expected exception must be thrown

By running a series of similar tests every time before moving to another project subcomponent, bugs can be identified before it spreads to other components. This is particularly helpful when certain bugfixes might result in prototype changes, resulting in incompatibility with other components during bugfixes.

3.2 Increasing the test size

If a function accepts parameters, multiple calls with different values should be passed to the function.

Suppose we want to test a function that converts a number to scientific notation rounded to 3 significant figures:

```

1 #include <cmath>
2 #include "SciNotation.h"
3
4 SciNotation sciNot(double number) {
5     int exp = (int) floor(log10(number)) - 2;
6     int digits = (int) round(number * pow(10.0, -exp));
7     return SciNotation{digits, exp};
8 }

```

Testing the simple case of rounding 1235 to 124×10^1 is successful:

```

1 #include "sciNot1.cpp"
2 #include "assert.h"
3
4 void testSciNot() {
5     ASSERT_EQUAL(sciNot(1235), SciNotation(124, 1))
6 }

```

```

1 Executing test testSciNot
2 Successful: sciNot(1235) = 124e1

```

However, adding more test parameters, such as negative inputs, would turn out that the function does not always function as expected:

```

1 #include "sciNot1.cpp"
2 #include "assert.h"
3
4 void testSciNot() {
5     ASSERT_EQUAL(sciNot(1235), SciNotation(124, 1))
6     ASSERT_EQUAL(sciNot(-1234), SciNotation(-123, 1))
7 }

```

```

1 Executing test testSciNot
2 Successful: sciNot(1235) = 124e1
3 Failed test for testSciNot: Expected sciNot(-1234) to return -123e1, got 0e2147483646

```

From this, we can identify that a negative input results in a negative value, apparently because of the `log10` function call. So a simple workaround of fixing the input to a positive number can be made:

```

1 #include <cmath>
2 #include "SciNotation.h"
3
4 SciNotation sciNot(double number) {
5     int sign = number >= 0 ? 1 : -1;
6     number *= sign;
7     int exp = (int) floor(log10(number)) - 2;
8     int digits = (int) round(number * pow(10.0, -exp));
9     return SciNotation{digits * sign, exp};
10 }

```

```

1 Executing test testSciNot
2 Successful: sciNot(1235) = 124e1
3 Successful: sciNot(-1234) = -123e1

```

3.3 Generating test cases randomly

Since a larger test size is more likely to catch bugs, it is a good idea to generate a large sample of test cases.

Generating random test cases also avoids the case where a developer is writing specific test cases (e.g. only multiples of 7) to avoid unit tests from reflecting a known bug in the code.

3.3.1 Random output, inverse-evaluated input

But a contradictory condition arises: If the test cases are generated, how to test if the result is correct? It is not possible to calculate the value in the test generator, because that would involve reimplementing the tested function.

Instead, if an inverse of the tested function can be written (or to generate any of the possible inverse values if the function is not an injection), the random generator can be used to generate random results instead.

Nevertheless, the inverse function is often harder to write and likely to have bugs than the function to test for, and is usually unintuitive to be used as a way of software specification. A better approach is to use property-based testing.

3.3.2 Property-based testing

Instead of comparing if the output is equal to an expected result, property-based testing compares the characteristics of the output to determine if it is reasonable.

In C++, the `rapidcheck` library provides a property-based testing framework that generates a random sequence of data by type for assertion.

For example, in the 3-significant-figure example above, the requirements are:

1. Given an input n , the output $\hat{n} = d \times 10^x$ is returned.
2. The result significand d must have exactly 3 significant figures, i.e. $d \in [100, 999]$
3. The error must be correct to the 3rd significant figure: $\hat{n} - n \in [-0.5 \times 10^x, 0.5 \times 10^x]$

The 1st requirement is defined through the function prototype, and the 2nd and 3rd requirements can be specified by the following test:

```

1 #include <cmath>
2 #include <vector>
3 #include <rapidcheck.h>
4 #include "sciNot5.cpp"
5
6 using namespace std;
7 using namespace rc;
8
9 void testSciNot() {
10     check("requirement 2: 3 significant figures", [] (const vector<double> &vec) {
11         for (double input : vec) {

```

```

12         SciNotation sci = sciNot(input);
13         int digits = sci.getDigits();
14         RC_ASSERT(100 <= digits && digits < 999);
15     }
16 });
17 check("requirement 3: error within bounds", [] (const vector<double> &vec) {
18     for(double input : vec) {
19         SciNotation sci = sciNot(input);
20         double error = sci.toDouble() - input;
21         error /= pow(10, sci.getExp());
22         RC_ASSERT(-0.5 <= error && error < 0.5);
23     }
24 });
25 }

```

```

1 Executing test testSciNot
2 Using configuration: seed=96209734982424326
3
4 - requirement 2: 3 significant figures
5 Falsifiable after 4 tests and 2 shrinks
6
7 std::tuple<std::vector<double>>>:
8 ([0])
9
10 sciNot5.test.cpp:14:
11 RC_ASSERT(100 <= digits && digits < 999)
12
13 Expands to:
14 100 <= 0 && true
15
16 - requirement 3: error within bounds
17 OK, passed 100 tests
18 Some of your RapidCheck properties had failures. To reproduce these, run with:
19 RC_PARAMS="reproduce=BQicLFXdpJXZtVmb0BiM6AyMgMXan5Waml2Yh5GdgYWanVnclNnBX_Yh+4
    cVBYw1PW0POXVAGc9jF6jzVFgBX_Yh+4cVBgAAEMAAAgACEA"

```

The data generated by rapidcheck (or any other property-based testing framework) should be further filtered in the test to eliminate incorrect values, or a custom input generator can be written to generate inputs appropriate for the tested unit. The rapidcheck framework provides a **Gen** API to customize input data generation.

3.3.3 Reimplement algorithm by brute-force

Even though the test only expresses the software specification, it may still be difficult to implement. Suppose the following function is specified:

- Parameter `int n`: the number of nodes
- Parameter `vector<pair<int, int> > edges`: each pair is two numbers indicating two node IDs of a directed edge.
- Node IDs are integers in the range $[0, n)$.
- If a Hamiltonian circuit is present, return any possible `vector<int> *` such that the node IDs in the vector represent a Hamiltonian circuit for the graph.
- If a Hamiltonian circuit is not present, return `NULL`.

Assume the function is implemented in a certain polynomial-time algorithm (which is possible given more specifications on the graph, but it is omitted as Hamiltonian circuit algorithm is not the focus of this project).

While it is simple to test for the correctness for returned vectors, it is very difficult to test if returning `NULL` is correct. The specification "Hamiltonian circuit is not present" is very well-defined, but it is still not possible to validate. In this case, a brute-force approach can be used to test for the correctness.

Although the function was implemented in an efficient polynomial-time algorithm, tests do not strictly require a high performance. It is therefore appropriate to test the new algorithm with a potentially less error-prone and more reliable algorithm using brute-force, i.e. NP time. This is because unit tests do not have to be capable of tackling with large data sets; light-weightedness is one of the core concepts of unit testing. (Stress testing is another topic not related to unit testing)

3.3.4 Testing for edge cases

Some common special cases should be included in the test data set explicitly, since they are common sources of error:

- Empty strings, multi-line strings, strings with Unicode, strings with emoji (The emoji issue actually caused [major bugs in iOS](#) for several times)
- Small numbers: 0, 1, 2, -1, -2, 0.5, -0.5
- Large numbers: $2^{31} - 1$, $2^{32} - 1$, $2^{63} - 1$, $2^{64} - 1$
- Special floats: $1/0$ ($+\infty$), $-1/0$ ($-\infty$), $0/0$ (NaN)

Remark: E.g. test for empty strings, `Float.INFINITY`, 0, etc.

4 Unit testing techniques

4.1 Code coverage

Test coverage is a criterion to assess the representativeness of the unit tests of a project by counting the number of lines executed in the test.

In GNU C++, code coverage is supported by the `gcov` utility. First, the `-fprofile-arcs -ftest-coverage` flag should be provided both when compiling source files of interest into intermediate object files and when compiling the object files into final executable files.

```
1 g++ -c -o fooBar.o -fprofile-arcs -ftest-coverage fooBar.cpp
2 g++ -o fooBar fooBar.o qux.o -fprofile-arcs -ftest-coverage
```

When the executable `fooBar` is executed, the file `fooBar.gcd` is also generated, containing the statistics during the execution. Coverage reports can hence be generated:

```
1 gcov fooBar
```

Applying code coverage into `sciNot5` from the previous section, 70 `gcov` reports are generated:

```
1 Any.hpp.gcov
2 ApplyTuple.h.gcov
3 Arbitrary.hpp.gcov
4 Assertions.hpp.gcov
5 Capture.h.gcov
6 Check.hpp.gcov
7 Common.hpp.gcov
8 Container.hpp.gcov
9 Create.hpp.gcov
10 ExecRaw.hpp.gcov
```

And way more. This is because every header or source file included in the compilation of `sciNot5.a` is flagged for coverage reporting. The file `sciNot5.cpp.gcov` might be of more interest:

```
1      -: 0:Source:sciNot5.cpp
2      -: 0:Programs:6
3      -: 1:#include <cmath>
4      -: 2:#include "SciNotation.h"
5      -: 3:
6      2568: 4:SciNotation sciNot(double number) {
```



```

7      2568:    5:      int sign = number >= 0 ? 1 : -1;
8      2568:    6:      number *= sign;
9      2568:    7:      int exp = (int) floor(log10(number)) - 2;
10     2568:    8:      int digits = (int) round(number * pow(10.0, -exp));
11     2568:    9:      return SciNotation{digits * sign, exp};
12     -:   10:}
13     -:   11:

```

All lines in `sciNot(double)` got called 2629 times, and the file has been tested with full coverage. While this is obvious, coverage may be more useful when the file has multiple conditions, loops or even error conditions:

```

1  #include <string>
2  #include "fibo.h"
3
4  int fibo(int times) {
5      int a = 1, b = 1;
6      if(times <= 0) throw std::string("input out of bounds"); // error value
7      if(times <= 2) return 1; // init condition
8      int i = 3;
9      while(true) {
10         int tmp = a + b;
11         a = b;
12         b = tmp;
13         if(i == times) break;
14         i++;
15     }
16     return b;
17 }

```

```

1  #include "assert.h"
2  #include "assert_int.h"
3  #include "fibo.h"
4
5  void testFibo() {
6      ASSERT_EQUAL(fibo(2), 1)
7      ASSERT_EQUAL(fibo(3), 2)
8  }

```

```

1      -:   0:Source:fibo.cpp
2      -:   0:Programs:6
3      -:   1:#include <string>
4      -:   2:#include "fibo.h"
5      -:   3:
6      4:   4:int fibo(int times) {
7      4:   5:       int a = 1, b = 1;
8      4*:  6:       if(times <= 0) throw std::string("input out of bounds"); // error value
9      4:   7:       if(times <= 2) return 1; // init condition
10     2:   8:       int i = 3;
11     -:   9:       while(true) {
12     2:  10:           int tmp = a + b;
13     2:  11:           a = b;
14     2:  12:           b = tmp;
15     2:  13:           if(i == times) break;
16     #####: 14:           i++;
17     #####: 15:       }
18     2:  16:       return b;
19     -:  17:}

```

Apparently, the error condition is never run (8 runs unchanged), 4 runs

4.2 Test-driven development (TDD)

As a consequence of unit testing, development flow becomes more fluent if each development subtask is based on certain test cases. By using TDD, it ensures that every part of the program can be tested properly, as nearly all code is written based on pre-written test cases. It also prevents writing duplicated code in large collaborative projects, and TDD can avoid waste of time for implementing already done features/requirements.

The workflow of TDD is as follows:

1. Create a new unit test for the software for a requirement that has not yet been implemented.
2. Run the unit test. If the unit test passes, this means the new test overlaps with previously written tests, and that a new test should be written to cover the new requirement
3. Write code to pass (only) the newly added test, while ensuring previous tests also pass
4. Run the tests
5. Repeat the above process for every requirement

Suppose we want to write a calculator with a function `int add(string numbers)` that has the following requirements:

1. The method can take in any amount of space-separated integers as an input, and return their sum.
2. If the string is empty, the method should return 0.
3. If non-space-and-digit characters are passed within the string, the method throw an "Invalid character" exception.

We begin the development process by writing the test for the first requirement as follows:

Remark: C++
test

As the `sum` method is currently empty, the test fails as expected. So we can proceed to write the code which solves this test:

```
1 #include <string>
2 #include <sstream>
3 #include "Calculator.h"
4 using namespace std;
5
6 int add(string numbers){
7     stringstream ss (numbers); int current, ans = 0;
8     while (ss >> current)
9         ans += current;
10    return ans;
11 }
```

To verify our implementation, we run the test we have written in test 1:

Remark: Run test, success

And we can continue with implementing the next requirement.

Exercise: Write the test and implementation for the second test.

4.3 Behaviour-driven development (BDD)

While TDD provides a reliable method to test individual function executions, it cannot precisely and concisely define functions that depend on global/object states.

The main difference between TDD and BDD is that while TDD ensures that every sub-component of the software is correct, BDD provides a more "high-level" view of what the entire software should provide to the user. It is often regarded as a "bridge" between product designers and developers, as through BDD, product designers can effectively communicate their needs to developers who specialize in implementing those requirements.

BDD defines the behaviour of a function using some specific terms or keywords. Each of these is written in the format of a "user story", and follow special "syntax" via the usage of certain keywords. The usage of these keywords are important, as a user story can be automatically phrased into "testable" code with a framework (such as Cucumber, on next section), given that it is formatted correctly with those "grammar". Here we present some of the keywords, along with the syntax general to most BDD frameworks:

The following is an example BDD of Cucumber:

```

1 Feature: Is it Friday yet?
2   Everybody wants to know when it's Friday
3
4   Scenario: Sunday isn't Friday
5     Given today is Sunday
6     When I ask whether it's Friday yet
7     Then I should be told "Nope"

```

Feature Describes a feature that the software should provide on the high-level, with some description

Scenario Describes a situation (similar to example) where that the software should be able to handle properly, defines the behavior of the program for a certain (more specific) situation.

Given The initial state of the program when this scenario takes place

When Similar to natural English, "When" implies that the following statement should be done "when" this statement is completed

Then Describes the expected output of the program after the preceeding "when" statement is fulfilled.

Why is such grammar important? As a software developer, we may often find it difficult to translate requirements given by product designers/users into code. By having the product designers write their requirements in such a fashion which is both "human readable" and "machine prasable", it is possible to translate those requirements in to unit tests **automatically**. For example, in Cucumber, the following can be generated automatically from the above BDD, which can be further manipulated to contain actual machine-runnable tests (in fact a BDD is just a combination of multiple unit tests, but written in a high-level manner):

```

1 @Given("^today is Sunday$")
2 public void today_is_Sunday() {
3     // Write code here that turns the phrase above into concrete actions
4     throw new PendingException();
5 }
6
7 @When("^I ask whether it's Friday yet$")
8 public void i_ask_whether_it_s_Friday_yet() {
9     // Write code here that turns the phrase above into concrete actions
10    throw new PendingException();
11 }
12
13 @Then("^I should be told \"([^\"]*)\"$")
14 public void i_should_be_told(String arg1) {
15     // Write code here that turns the phrase above into concrete actions
16     throw new PendingException();
17 }

```

If you are interested to learn more, head over to <http://docs.cucumber.io/guides/10-minute-tutorial/> to get started with the Cucumber framework!

4.4 Summary: TDD vs BDD

As you can see, the description of the program behavior in BDD is more "high level" compared to TDD. The person writing BDDs usually writes them from a user's perspective, where they might not have full knowledge on the internals of the software. On the other hand, TDDs focus more on the internal working on the software, and whether each part functions correctly.

Often in BDD, when a test fails, it is difficult to locate the root cause of the error, as the BDD usually involves the whole piece of software instead of a single modular component and it is difficult to locate a bug in the whole program as opposed to a single function. (in Cucumber, the framework subclassifies a task into multiple components to help to identify the exact source of error, but the subclassification is not as small as those in TDD). Furthermore, as BDD involves running the whole software to conduct behavior tests, it

takes much longer to conduct a single test (a BDD test can take up to a minute, while a TDD test may only take a few seconds).

The existence of BDDs is to keep developers on track on building products that directly address the user's requirements, instead of building functions that satisfy a low-level requirement. This means, that in BDD the software development is driven by how the user will use the software, and the program is built accordingly, while for TDD, it involves the concept where the correctness of a software implies that all sub-components are correctly written, meaning that code is built according to the feature of a certain component, but not the whole picture.

This means if a BDD test fails, it implies that there is a use case that the software will have a bug in regular use, and should not be released. On the other hand, if BDD passes and a TDD fails, this implies that although some components may contain bugs, in regular usage, they will not affect the running of the whole software, and the software can continue with its release, leaving the bug fix for later.

4.5 Dependency mocking

Often (if not always), sub-components in a software will have other components as dependencies. Take our calculator class as an example. Suppose we want to write a function `int eval(int a, int b, char operation)`

```
1 #include "Calculator.h"
2
3 int eval(int a, int b, char operation){
4     if (operation == '+') return a + b;
5     if (operation == '-') return a - b;
6     if (operation == '*') return a * b;
7     if (operation == '/') return a / b;
8 }
```

The implementation is done, and we have unit tests for this code. Now lets say that there is a new requirement, that if the result of the calculation exceeds the value of `int` ($2^{31} - 1$), it should throw an `OverflowException`. A trivial idea for implementation would be to implement the internal calculation of the operation via another data type that supports larger range (e.g. `long long`, `BigInteger`), and compare this value with $2^{31} - 1$. In other words, we can write the following instead:

```
1 #include "Calculator.h"
2 #include "BigInteger.h"
3
4 int eval(int a, int b, char operation){
5     BigInteger bigA(a), BigInteger bigB(b);
6     BigInteger result(0);
7
8     if (operation == '+') result = bigA + bigB;
9     if (operation == '-') result = bigA - bigB;
10    if (operation == '*') result = bigA * bigB;
11    if (operation == '/') result = bigA / bigB;
12
13    if (result > BigInteger((1<<31)-1))
14        throw std::overflow_error("too big");
15    else
16        return result.value();
17 }
```

The `Calculator` now depends on `BigInteger`. Assume that we have not implemented or tested our implementation of `BigInteger` yet, how are we going to test the component of `Calculator`? One of the ways is to substitute `BigInteger` with another "fake" `BigInteger`, which for the purpose of explanation, we assume the implementation is correct.

```
1 #ifndef BIGINTEGER_H
2 #define BIGINTEGER_H
3 #include <string>
4 using namespace std;
5
6 class BigInteger {
```

```

7   long long val;
8   public:
9       BigInteger(int a){
10          val = a;
11      }
12      BigInteger operator+(BigInteger rhs){
13          return val + rhs.val;
14      }
15      BigInteger operator-(BigInteger rhs){
16          return val - rhs.val;
17      }
18      BigInteger operator*(BigInteger rhs){
19          return val * rhs.val;
20      }
21      BigInteger operator/(BigInteger rhs){
22          return val / rhs.val;
23      }
24      bool operator>(BigInteger rhs){
25          // implementation omitted
26      }
27      int value(){
28          return (int) val;
29      }
30 };
31
32 #endif

```

Now our `Calculator` class can compile successfully, and we can test the implementation of `Calculator` without caring about the correctness of `BigInteger`.

This is exactly what "Dependency mocking" is. If inside the target we want to perform unit tests on, has a function call which either is untested, or we do not want to call (due to lack of resources, take long time to run), we can essentially substitute that function with another one which produces the "expected result" of the original function, assuming that the implementation of the original function was correct. Mocking is useful when it is impractical to call a function many time during the testing stage of development, such as those which consume a lot of time like database calls.

However, although mocking effectively isolates the target component from other components, making unit testing relevant to the implementation of only the target, mocking has several disadvantages for software development. By essentially "neglecting" the implementation of other components, it makes the target component prone to unexpected behavior of the target component. For instance, if there were any updates to the dependencies, such as a new exception being thrown, the current unit tests would not be aware of such new behavior, and the tests may not be accurate. For example, the unit test may run successfully with mocks, but with the actual "unmocked" implementation, it has a runtime error, and the developer may be unaware to this. This is particularly important in large projects with many dependencies, perhaps with dependencies not being written by the in-house developers, and it is difficult to update mocks to fully reflect the current behavior of the dependencies.

If you want to know more about mocking and its advantages/disadvantages, check out <https://enterprisecraftsmanship.com/2016/06/09/styles-of-unit-testing/>.

5 Difficulties

5.1 Coupling

Coupling is the cyclic dependency between units to be tested.

Remark: Explain what is coupling, why it is bad, how hard it is to prevent, and common practices e.g. abstraction, interfaces, mocking to prevent coupling

Remark: Consider cyclic dependency checks in golang, which is an excellent example of why it is hard

5.2 Test case selection

As the number of variables to a feature increases, the number of test cases might increase exponentially.