

When to Stop Testing: Statistical Testing & Reliability

CSCE 740 - Lecture 24 - 11/23/2015

The all-important question....

**When have we tested
enough?**

We Will Cover

- How do we know when we are done?
- Stopping Criteria
 - Requirements
 - Coverage
 - Budget
 - Plan
 - Mutation Analysis
 - Reliability

When All Black-Box Tests Pass?

- Write tests based on your requirements, then stop when they pass.
 - AKA: Stop when you can make an argument for verification.
- Are there problems with this?

When We Have Achieved Coverage?

- Set your sights on some structural coverage metric and test until that is achieved.
 - branch coverage, condition coverage, etc.
- Problems?

When We Run Out of Time?

- The “Budget Coverage Criterion”
 - The usual answer to when testing is done.
 - When we run out of time.
 - When the money dries up.
- Problems?

What if We Make a Plan?

- Plan a series of tests carefully, then test according to that plan.
 - Consider forms of black and white box testing.
 - Factor in the budget and the cost of test case creation in choosing how we test.
- When those tests are done, you are done.
- Problems?

Mutation Testing Revisited

Many testing techniques based on what we *think should happen*. Why not test based on *what we think could go wrong*?

Deliberately seed faults into a system, and see if you can distinguish the faulty system from the existing system.

Mutation Testing

1. Select *mutation operators* - code transformations that represent classes of faults that we are interested in.
2. Generate *mutants* by applying mutation operators to the program.
3. Execute tests against the program and mutants to *kill* mutants.

Mutation testing is used to judge adequacy of a test suite, based on the idea that mutations represent real faults.

- Relies on two assumptions:
 - **competent programmer hypothesis**
 - (the program is close to correct)
 - **coupling hypothesis**
 - (small mutations cascade into larger faults).

Mutation Coverage

Adequacy of the suite can be measured as:

$$\frac{(\# \text{ mutants killed})}{(\text{total mutants})}$$

A high “mutation score” can be used to argue that you wrote effective tests.

Mutation can subsume other coverage:

- Statement: apply statement deletion to all statements.
- Branch: Replace all predicates with constants T/F.

Mutation Testing?

- Deliberately seed faults into a system, and see if you can distinguish the faulty system from the existing system.
- Can provide evidence that real faults have been detected.
- Problems?

\$%\$*, I give up. When are we done?

- When we can argue that we're done.
- How do we do that? Provide evidence.
 - Make a plan.
 - Mutation testing can help assess adequacy.
 - One more method:
 - Examine the reliability of the system.

Analyzing Software Reliability

What is Reliability?

- Reliability is the probability of failure-free operation for a specified time in a specified environment for a given purpose.
- This means different things depending on the system and the users of that system.
- Informally, reliability is a measure of how well users think the system provides the services they require.

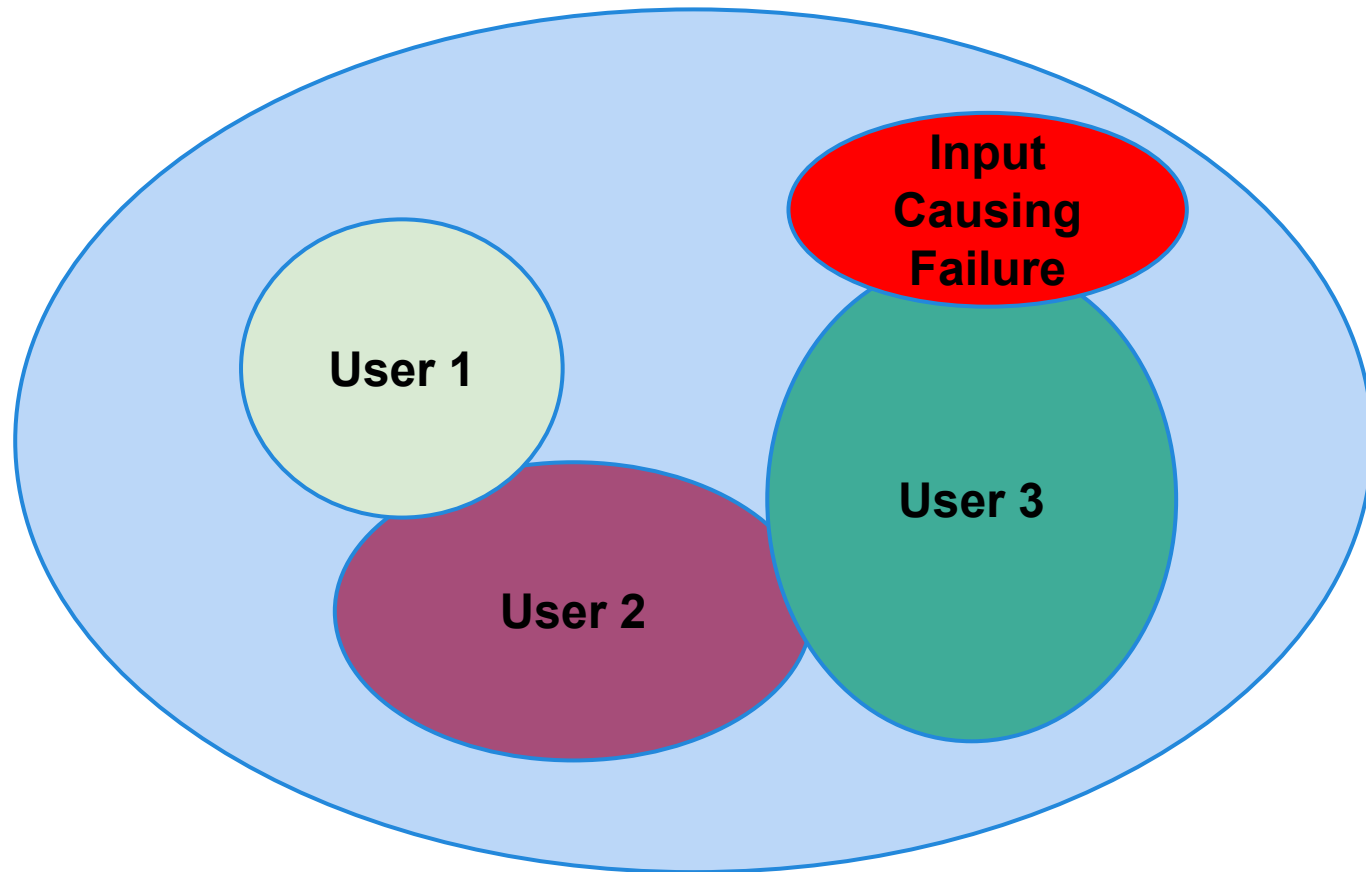
Reliability is Measurable

- Reliability can be defined and measured.
- Reliability requirements can be specified:
 - Non-functional requirements can define the number of failures that are acceptable during normal use of the system, or the time in which the system is allowed to be unavailable for use.
 - Functional requirements can define how the software avoids, detects, and tolerates faults to ensure they don't lead to failures.

Improving Reliability

- Reliability is improved when software faults that occur in the most frequently-used parts of the software are removed.
- Removing $X\%$ of the faults will not necessarily lead to an $X\%$ improvement in reliability.
 - In a study, removing 60% of the faults actually led to a 3% reliability improvement.
- Removing faults with serious consequences is the top priority.

Reliability Perception



Software Reliability

- Reliability cannot be defined objectively for all situations.
 - Reliability measurements quoted out of context are meaningless.
- Requires operational profile for its definition.
 - A profile of the expected pattern of software usage.
- Reliability must consider fault consequences.
 - Not all faults are equally serious.
 - System is perceived as unreliable if there are more serious faults.

How to Measure Reliability

- Measuring reliability is normal when building hardware, but hardware metrics often aren't suitable for software.
 - Based on component failures and the need to repair or replace a component once it has failed.
 - In hardware, the design is assumed to be correct.
- **Software failures are always design failures.**
 - Often, the system is available even though a failure has occurred.

Availability

- The availability of a system reflects its ability to deliver services when available (uptime/total time).
 - Takes repair and restart time into account.
 - Does not tend to take incorrect computations (partial failures) into account.
- Availability of 0.9999 means the system is available 99.99% of the time.
 - 0.9 = down for 144 minutes a day, 0.99 = down for 14.4 minutes, 0.999 = down for 84 seconds, 0.9999 = down for 8.4 seconds.

Probability of Failure on Demand (POFOD)

- The likelihood that a service request will result in a system failure (failures/requests over a period).
- POFOD = 0.001 means that 1 out of 1000 service requests result in a failure.
- Should be used as a reliability metric in situations where a failure on request can lead to a serious system failure.
 - Independent of the frequency of requests.
 - 1/1000 failure rate sounds risky, but if that becomes one failure per lifetime, it is pretty good.

Rate of Occurrence of Fault (ROCOF)

- Frequency of the occurrence of unexpected behavior.
 - Probable number of failures over a period of time or number of system executions.
- ROCOF of 0.02 means that 2 failures are likely per 100 time units.
- Most appropriate metric when requests are made on a regular basis (such as a shop).

Mean Time to Failure (MTTF)

- Measures the average length of time between observed failures.
- MTTF of 500 means that the time between failures is, on average, 500 time units (or requests).
- For systems with long user sessions, you want to require a long MTTF.

Data Needed for Measurements

To assess reliability, data must be captured from users' sessions with the system:

- Measure the number of failures per a given number of requests (used for POFOD).
- Measure the time or number of requests between failures, plus total elapsed time or request number (used for ROCOF and MTTF).
- Measure the time to restart after a failure (for availability).

Reliability Examples

- Provide software with 10000 requests.
 - Wrong result on 35 requests, crash on 5 requests.
 - What is the POFOD?
- $40 / 10000 = 0.0004$
- Run the software for 144 hours
 - (6 million requests). Software failed on 6 requests.
 - What is the ROCOF? The POFOD?
- $ROCOF = 6/144 = 1/24 = 0.04$
- $POFOD = 6/6000000 = (10^{-6})$

Reliability Examples

- You advertise a piece of software with a ROCOF of 0.001 failures per hour.
 - However, it takes 3 hours (on average) to get the system up again after a failure.
 - What is the availability per year?
- Failures per year:
 - approximately 8760 hours per year (24×365)
 - $0.001 \times 8760 = 8.76$ failures per year
- Availability
 - $8.76 \times 3 = 26.28$ hours of downtime per year.
 - Availability = $0.997 \left(\frac{8760 - 26.28}{8760} \right)$

Activity - Availability

- Your customers want an availability of at least 99% and a rate of fault occurrence of less than 2 failures per 8 hour work period.
- After testing your code for 6 full days, the product failed 27 times and it took an average of 32 minutes to restart after each failure.
 - What is the rate of fault occurrence?
 - What is the availability?
 - Is the product ready to ship?
 - If not, why not?

Activity Solution

- What is the rate of fault occurrence?
 - $27/144$ hours = 1.5/8 hour work day
- What is the availability?
 - Was down for 864 minutes out of of 144 hours.
 - ... or 864/8640 minutes.
 - ... or 10% of the time.
 - So, availability is 0.90.
- Is the product ready to ship? If not, why not?
 - ROCOF is fine, but availability is too low.
 - What should they do to improve?

Reliability Economics

- Raising reliability is expensive. It may be cheaper to accept unreliability and pay for failure costs.
- The balancing point depends on social and political factors and the system type.
 - A reputation for unreliable products may hurt more than the cost of improving reliability.
 - Cost of failure depends on risks of failure. For business systems, modest reliability may be fine.

Statistical Testing

- Rather than using tests to trigger faults, we can use tests to measure reliability.
- Test inputs should match the predicted usage profile of a user.
- By recording errors and other measurements, we can calculate ROCOF, POFOD, etc.
- An acceptable level of reliability should be specified and the software tested until that level is reached.

Operational Profiles

- Reflects how the software is used in practice.
- Consists of classes of input and the probability of their occurrence.
- Can be specified in advance if other systems exist that perform similar actions.
- For new systems, it is much harder to specify.
 - Conduct beta testing to gather initial usage data.
 - Remember that usage changes over time.

Statistical Testing Procedure

- Study existing systems and form an operational profile.
- Construct test input that reflects the profile.
- Apply inputs and count the frequency and type of failures that occur, along with the time between failures.
- After observing a statistically significant number of failures, compute the reliability.

Statistical Testing Challenges

- Operation profile uncertainty
 - A profile based on other systems may not be valid for your system.
- High cost of test input generation
 - Large volume of inputs needed. Can be expensive.
- Statistical uncertainty
 - Need to generate enough failures to estimate reliability. This is hard when the system is already reliable.
 - Hard to estimate confidence in operational profile.

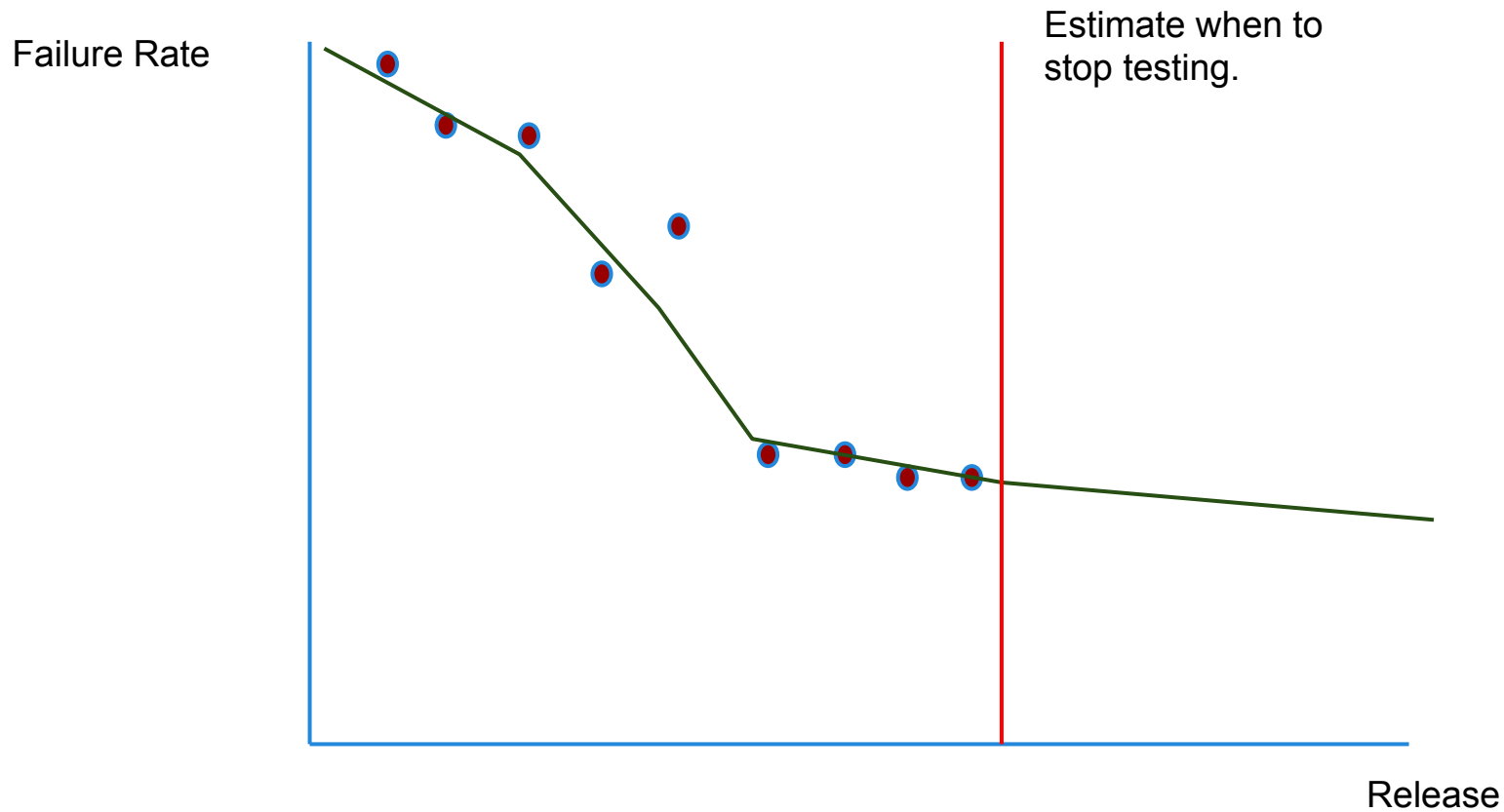
Getting the Most Out of Statistical Testing

- Statistical testing often reveals errors that do not emerge from other V&V activities.
 - As these emerge, fix the system and re-test.
 - As you gather more data, reliability growth can be modeled and used to plan testing.
- Mutation often used to plant “known faults” for reliability testing.
 - Can make an argument that you already rooted out most of the real faults.

Reliability Growth Modeling

- Can build mathematical model of the change in system reliability as changes are made to the code base.
- Used as a means of predicting additional reliability from additional testing and changes.
- Use statistical testing to measure reliability of each system version, and make a call on when to stop trying to raise reliability.

Reliability Prediction



Key Points

- Reliability is one of the most important software characteristics.
- We should aim to produce reliable software.
- Reliability depends on the pattern of usage of the software. Different users will interact differently.
 - Faulty software can be reliable for some users.

Key Points

- Reliability can be measured quantitatively.
 - ROCOF, POFOD, Availability, MTTF
- Statistical testing is used to estimate reliability without actual users.
- Reliability growth models may be used to predict when a required level of reliability may be achieved.

When Do We Stop Testing?

- Requirements-based tests may not exercise all code. Coverage criteria may not reflect requirements or produce good tests. “We ran out of money” just indicates poor planning.
- Come up with a plan that reflects your budget.
- You are done when you can present evidence that you have built a reliable system.

Next Time

- Software Evolution and Maintenance
 - Reading:
 - Sommerville, ch. 9
- Practice Final is up.
 - We will go over this during the last class.
- Homework 5 is up.
- Homework 4 due tomorrow.
 - Any questions?