

## Simple R

R is a language for statistical analysis, machine learning, and data visualization. Despite its stable and rapid growth over the years, we still don't have an instrument for rigorous inspection and analysis of R programs.

Since the behavior of R can be fairly ambiguous, it's difficult to identify and fix vulnerabilities of a given program.

Simple R is a subset of R. It supports the core R operations, but excludes programs that use unsafe and hard-to-analyze operations such as super-assignments.

## Approach

The symbolic execution engine executes a given Simple R program symbolically, finding and reporting errors of a specified kind. It also produces a concrete test case for a discovered error.

$$\text{IFTRUE} \frac{H(m_v) = (\text{true}, A_v)}{\langle\langle \text{C}[\text{if } m_v \text{ then } e_t \text{ else } e_f], m_\Gamma \rangle\rangle : S; H; C \rangle \Rightarrow \langle\langle \text{C}[e_t], m_\Gamma \rangle\rangle : S; H; C \rangle}$$

$$\text{IFFALSE} \frac{H(m_v) = (\text{false}, A_v)}{\langle\langle \text{C}[\text{if } m_v \text{ then } e_t \text{ else } e_f], m_\Gamma \rangle\rangle : S; H; C \rangle \Rightarrow \langle\langle \text{C}[e_f], m_\Gamma \rangle\rangle : S; H; C \rangle}$$

$$\text{IFSYMBOLIC} \frac{H(m_v) = (X_{\text{bool}}, A_v) \quad C' = C \cup \{X_{\text{bool}} = \text{true}\} \quad C'' = C \cup \{X_{\text{bool}} = \text{false}\}}{\langle\langle \text{C}[\text{if } m_v \text{ then } e_t \text{ else } e_f], m_\Gamma \rangle\rangle : S; H; C \rangle \Rightarrow \langle\langle \text{C}[e_t], m_\Gamma \rangle\rangle : S; H; C' \rangle, \langle\langle \text{C}[e_f], m_\Gamma \rangle\rangle : S; H; C'' \rangle}$$

$$\text{DIVSYMBOLIC} \frac{m_v \text{ fresh} \quad Z_{\text{int}} \text{ fresh} \quad H(m_{v_1}) = (v_1, A_{v_1}) \quad H(m_{v_2}) = (v_2, A_{v_2})}{C' = C \cup \{Z_{\text{int}} = v_1/v_2\} \cup \{v_2 \neq 0\} \quad H' = H[m_v \mapsto (Z_{\text{int}}, A)] \quad C'' = C \cup \{v_2 = 0\} \quad H'' = H[m_v \mapsto (\text{DivError}, 0)]} \langle\langle m_{v_1}/m_{v_2}, m_\Gamma \rangle\rangle; H; C \rangle \mapsto \langle\langle m_v; H'; C' \rangle\rangle, \langle\langle m_v; H''; C'' \rangle\rangle}$$

## The Symbolic Execution Engine

```
skew <- function(x)
{
  sum2 <- (sum(x-mean(x)))^2
  sum3 <- sum((x-mean(x))^3)
  skew <- (sqrt(length(x)*sum3)/(sum2^(1.5)))
  return(skew)
}
```

$$\frac{\sqrt{n} \sum_{i=1}^n (x_i - \bar{x})^3}{(\sum_{i=1}^n (x_i - \bar{x})^2)^{3/2}}$$

```
x <- scan("data.txt")
if (!(skew(x) == 0)) cat ("skewed")
```

```
skew <- function(x)
{
  sum2 <- (sum(x-mean(x)))^2
  sum3 <- sum((x-mean(x))^3)
  skew <- (sqrt(length(x)*sum3)/(sum2^(1.5)))
  return(skew)
}
```

```
x <- SimpleR.int()
if (!(skew(x) == 0)) cat ("skewed")
```

Error: Missing value where TRUE/FALSE needed

Example 1: x = c(-1, 0, 1)

Example 2: x = c(-4, -2, 0, 1, 5)

Example 3: x = c(4, 3, 5, 7, 8, 12, 10)

Unit Test 1:  
Division by Zero

Unit Test 2:  
Negative Indices

Unit Test 3:  
Error in if

## Challenges

We used SMT-LIB to represent symbolic vectors and path constraints. Symbolic vectors are placeholders for usual R vectors, which are the basic objects of R.

Finding a certain kind of error in a unit of R code boils down to solving Satisfiability Modulo Theories formulae using a theorem prover.

```
- (forall ((a (Array s1 s2)) (i s1) (e s2))
  (= (select (store a i e) i) e))

- (forall ((a (Array s1 s2)) (i s1) (j s1) (e s2))
  (=> (distinct i j)
    (= (select (store a i e) j) (select a j))))

- (forall ((a (Array s1 s2)) (b (Array s1 s2)))
  (=> (forall ((i s1)) (= (select a i) (select b i)))
    (= a b)))
```

Boolean, Integers and Real theories are directly supported in SMT-LIB v.2, while symbolic vectors were implemented using the theory of arrays. Thus, an X Y Array in SMT-LIB is a map from X to Y, where X,Y can be Int, Real or Boolean.

## Future Work

Make a user-friendly tool  
Relational Symbolic Execution

## Contact

[lzharmukhametova@college.harvard.edu](mailto:lzharmukhametova@college.harvard.edu)