

Probabilistic Symbolic Execution

FIRST1 LAST1*, Institution1, Country1

FIRST2 LAST2†, Institution2a, Country2a and Institution2b, Country2b

TODO: Write Abstract

CCS Concepts: • **Software and its engineering** → **General programming languages**; • **Social and professional topics** → *History of programming languages*.

Additional Key Words and Phrases: keyword1, keyword2, keyword3

ACM Reference Format:

First1 Last1 and First2 Last2. 2022. Probabilistic Symbolic Execution. *Proc. ACM Program. Lang.* 1, POPL, Article 1 (January 2022), 12 pages.

1 INTRODUCTION

One of the goals of symbolic execution is to identify whether a “bad” state (e.g. an assert failure) is reached, signifying a bug in the program. In probabilistic programs, we are often more interested in either whether a “bad” state is reached too often, or, conversely, whether the “good” states are reached often enough.

1.1 Contributions

2 MOTIVATING EXAMPLE

The Monty Hall problem [?] is a classic probability puzzle based on the American television show, *Let’s Make a Deal*, which showcases how subtle probabilistic reasoning can be. The problem itself is simple:

You are a contestant on a gameshow and behind one of three doors there is a car and behind the others, goats. You pick a door and the host, who knows what is behind each of the doors, opens a different door, which has a goat. The host then offers you the choice to switch to the remaining door. Should you?

While it may seem unintuitive, regardless of the contestant’s original door choice, the contestant who always switches doors will win the car $\frac{2}{3}$ of the time, as opposed to a $\frac{1}{3}$ chance.

We can represent this problem as a probabilistic program, as shown in Fig. 1, where choice $\in [1, 3]$ and is the door which is originally chosen by the contestant and door_switch is true if the contestant wants to switch doors when asked, and false otherwise. Supposedly if door_switch = true, regardless of the value of choice, monty_hall should return true (i.e. the car is won) $\frac{2}{3}$ of the time. The problem thus becomes, given a probabilistic program with discrete sampling statements, and a program property, how do we verify that the program satisfies the property?

One typical solution would be to run a simulation on random inputs and compute the probability of winning if door_switch = true, and the probability of winning if door_switch = false. While

*with author1 note

†with author2 note

Authors’ addresses: First1 Last1, Department1, Institution1, Street1 Address1, City1, State1, Post-Code1, Country1, first1.last1@inst1.edu; First2 Last2, Department2a, Institution2a, Street2a Address2a, City2a, State2a, Post-Code2a, Country2a, first2.last2@inst2a.com, Department2b and Institution2b, Street3b Address2b, City2b, State2b, Post-Code2b, Country2b, first2.last2@inst2b.org.

Fig. 1. C code for the Monty Hall Problem

for this program this would work as the input space is tiny, the resulting output distributions would be absent any formal guarantees and only be an approximation.

3 PROBABILISTIC SYMBOLIC EXECUTION ALGORITHM

In this section, we present our technique for augmenting traditional symbolic execution to support probabilistic programs with discrete sampling instructions. We begin with a short review of a traditional symbolic execution algorithm and then discuss how we calculate exact path probabilities.

3.1 Background

Subhajit & Sumit: If you could add a high-level description of symb. exec., maybe an algorithm, that would be great!

Symbolic execution is a program analysis technique where a program is run on *symbolic* inputs and all program operations are replaced with those which manipulate these symbolic variables. During execution, program state is encoded symbolically in two parts: a path condition which is a conjunctive formula, φ , which records the branch conditions which are true for that particular path, and a mapping from program variables to symbolic expressions containing constants and *symbolic variables*, σ . When execution reaches an assignment of the form $x = e$, where e is a constant or program variable, $\sigma[x] = \sigma[e]$. Similarly, if $e = e_1 \oplus e_2$, where \oplus is an arbitrary binary operation, $\sigma[x] = \sigma[e_1] \oplus \sigma[e_2]$. When execution reaches a branch guarded by the condition c , execution proceeds down both branches, one where $\varphi = \varphi \wedge \sigma[c]$, and the other where $\varphi = \varphi \wedge \neg\sigma[c]$.

3.2 Adding Probabilistic Sampling

We now consider the problem of performing symbolic execution on a simple imperative probabilistic programming language, **pWhile**:

$$S := \text{skip} \mid x \leftarrow e \mid x \overset{\$}{\leftarrow} d \mid S_1; S_2 \mid \text{if } e \text{ then } S_1 \text{ else } S_2 \mid \text{while } e \text{ do } S$$

Above, x is a program variable, e is an expression, and d is a *discrete* distribution expression. In order to support sampling instructions we make the following additions to traditional symbolic execution:

- *Probabilistic symbolic variables*. We distinguish two types of symbolic variables: *universal* symbolic variables (identical to those in traditional symbolic execution), and *probabilistic* symbolic variables. For each unique sampling instruction a new probabilistic symbolic variable is created to denote the result of sample.
- *Distribution map*. We add a new mapping from probabilistic symbolic variables to distribution expressions, P , which tracks the distribution from which a probabilistic symbolic variable was originally sampled from.
- *Path probability*. For each path, we adjoin a path probability expression, p , which is parameterized by universal symbolic variables.

Algorithm 1 PSE Assignment Algorithm

```

1: function PSEASSIGNMENT( $x, e, \varphi, \sigma, P$ )
2:    $e_{sym} \leftarrow \sigma[e]$ 
3:    $\sigma[x] = e_{sym}$ 
4:   return ( $\varphi, \sigma, P$ )
5: end function

```

Assignment. For assignment statements of the form $x \leftarrow e$, where x is a *program* variable and e is an expression, probabilistic symbolic execution proceeds identically to traditional symbolic execution, as detailed in Alg. 1.

Algorithm 2 PSE Sampling Algorithm

```

1: function PSESAMPLE( $x, d, \varphi, \sigma, P$ )
2:    $\delta \leftarrow$  Generate a fresh probabilistic symbolic variable
3:    $\sigma[x] = \delta$ 
4:    $P[\delta] = d$ 
5:   return ( $\varphi, \sigma, P$ )
6: end function

```

Sampling. For sampling statements, $x \stackrel{\$}{\leftarrow} d$, where x is a *program* variable and d is a distribution expression, Alg. 2 is used. A fresh probabilistic symbolic variable, δ , is created, σ is updated to be $\sigma[x] = \delta$, and the original distribution d is recorded in P by setting $P[\delta] = d$.

Algorithm 3 PSE Branch Algorithm

```

1: function PSEBRANCH( $c, \varphi, \sigma, P$ )
2:    $c_{sym} \leftarrow \sigma[c]$ 
3:    $(\delta_1, \dots, \delta_n) \leftarrow \text{dom}(P)$ 
4:    $(d_1, \dots, d_n) \leftarrow (P[\delta_1], \dots, P[\delta_n])$ 
5:   
$$p_c \leftarrow \frac{\sum_{(v_1, \dots, v_n) \in \text{dom}(d_1) \times \dots \times \text{dom}(d_n)} [(c_{sym} \wedge \varphi)\{\delta_1 \mapsto v_1, \dots, \delta_n \mapsto v_n\}]}{\sum_{(v_1, \dots, v_n) \in \text{dom}(d_1) \times \dots \times \text{dom}(d_n)} [\varphi\{\delta_1 \mapsto v_1, \dots, \delta_n \mapsto v_n\}]}$$

6:   return ( $(\varphi \wedge c_{sym}, \sigma, P, p_c), (\varphi \wedge \neg c_{sym}, \sigma, P, 1 - p_c)$ )
7: end function

```

Branches. Note that with the inclusion of probabilistic symbolic variables we can now either branch on universal or probabilistic symbolic variables (or both). For guards whose symbolic version does not contain any probabilistic symbolic variables, probabilistic symbolic execution works nearly identically to traditional symbolic execution, save one detail: the probability of taking the branch. If c_{sym} is a symbolic expression which represents the guard to a branch, and c_{sym} does not contain any probabilistic symbolic variables, then we define the probability of taking the “true” branch is $[c_{sym}]$, and the probability of taking the “false” branch is $[\neg c_{sym}]$, where $[Q] = 1$ if Q is true, and 0 otherwise (known as Iverson brackets).

For probabilistic branches, i.e. guards which branch on probabilistic symbolic variables, Alg. 3 is used instead of the traditional symbolic execution branch algorithm. Without loss of generality,

Algorithm 4 Getting Path Constraints : Symbolic Execution

```

1: function SYMBEX( $P_{prog} : Program$ )
2:    $\phi_{paths} \leftarrow [], Ex_{stack} \leftarrow [], \Delta \leftarrow [], \phi_{paths} \leftarrow \phi$  ▷ Initialization
3:    $I_0 \leftarrow \text{GETSTARTINSTRUCTION}(P_{prog})$ 
4:    $S_0 \leftarrow [I_0, \phi_{path}, \Delta]$  ▷ Empty Initial State
5:    $Ex_{stack}.\text{PUSH}(S_0)$  ▷ Start with  $S_{cur}$  in Execution Stack
6:   while  $Ex_{stack} \neq \phi$  do
7:      $S_{cur} \leftarrow Ex_{stack}.\text{POP}(), I_{cur} \leftarrow S_{cur}[1]$  ▷ Start State,  $[I_0, \phi_{path}, \Delta]$ 
8:     switch  $\text{INSTYPE}(I_{cur})$  do
9:       case  $v := e$  ▷ Assignment Instruction
10:         $I_{cur} \leftarrow \text{GETNEXTINSTRUCTION}(I_{cur})$ 
11:         $S_{cur} \leftarrow [I_{cur}, \phi_{path}, \Delta[v \rightarrow \text{EVAL}(e, \Delta)]]$ 
12:         $Ex_{stack}.\text{PUSH}(S_{cur})$ 
13:        case if ( $c_{sym}$ ) then  $P_1$  else  $P_2$  ▷ Branch Instruction
14:          if ( $\text{ISAT}(c_{sym}) \wedge \text{ISAT}(\neg c_{sym})$ ) then ▷ Both cases SAT
15:             $I_1 \leftarrow \text{GETSTARTINSTRUCTION}(P_1)$ 
16:             $I_2 \leftarrow \text{GETSTARTINSTRUCTION}(P_2)$ 
17:             $S_{true} \leftarrow [I_1, \phi_{path} \wedge c_{sym}, \Delta]$ 
18:             $S_{false} \leftarrow [I_2, \phi_{path} \wedge (\neg c_{sym}), \Delta]$ 
19:             $Ex_{stack}.\text{PUSH}(S_{false})$ 
20:             $Ex_{stack}.\text{PUSH}(S_{true})$  ▷ Start with True State
21:          else if  $\text{ISAT}(c_{sym})$  then ▷ True case SAT
22:             $I_{cur} \leftarrow \text{GETSTARTINSTRUCTION}(P_1)$ 
23:             $S_{cur} \leftarrow [I_{cur}, \phi_{path} \wedge c_{sym}, \Delta]$ 
24:             $Ex_{stack}.\text{PUSH}(S_{cur})$ 
25:          else if  $\text{ISAT}(\neg c_{sym})$  then ▷ False case SAT
26:             $I_{cur} \leftarrow \text{GETSTARTINSTRUCTION}(P_2)$ 
27:             $S_{cur} \leftarrow [I_2, \phi_{path} \wedge (\neg c_{sym}), \Delta]$ 
28:             $Ex_{stack}.\text{PUSH}(S_{cur})$ 
29:          else ▷ Unconditional Branch
30:             $I_{cur} \leftarrow \text{GETNEXTINSTRUCTION}(I_{cur})$ 
31:             $S_{cur} \leftarrow [I_{cur}, \phi_{path}, \Delta]$ 
32:             $Ex_{stack}.\text{PUSH}(S_{cur})$ 
33:          end if
34:          case HALT ▷ Terminate Instruction
35:             $\phi_{paths} \leftarrow \phi_{paths}.\text{APPEND}(\phi_{path})$ 
36:        end while
37:      return  $\phi_{paths}$ 
38:    end function

```

consider a branch of the form **if** c **then** S_1 **else** S_2 where c is a probabilistic branch. As before, we define $c_{sym} = \sigma[c]$. The core of Alg. 3 is computing the probability of taking each branch. We interpret probabilistic branches as a conditioning operation on the distributions which are mentioned in the guard of the branch. Under this interpretation, we aim to compute the conditional probability $\Pr[c_{sym} \mid \varphi] = \frac{\Pr[c_{sym} \wedge \varphi]}{\Pr[\varphi]}$, where φ is the current path condition formula.

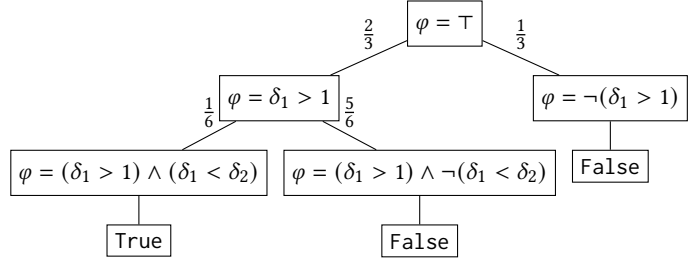
We use a form of model counting in order to compute p_c , the probability of c being true. Note that each probabilistic symbolic variable, δ , is mapped to exactly one distribution, d , and therefore, $\delta \in \text{dom}(d)$. So, assuming there are n probabilistic symbolic variables, $\delta_1, \dots, \delta_n$, and so n distributions, d_1, \dots, d_n , the set of all possible values $\delta_1, \dots, \delta_n$ be is $\mathcal{D} = \text{dom}(d_1) \times \dots \times \text{dom}(d_n)$. We then count the number of elements (or *assignments*) from \mathcal{D} which satisfy $c_{\text{sym}} \wedge \varphi$ and φ , and divide these two quantities as shown on line 5 of Alg. 3. Note that p_c is not necessarily a value, but rather a symbolic expression containing constants and universal symbolic variables. Additionally, we exploit the fact that the sum of the conditional probabilities of the branch outcomes is 1, which allows us to avoid computing the probability of taking the “false” branch directly.

```

1:  $x \xleftarrow{\$} \text{UniformInt}(1, 3)$ 
2:  $y \xleftarrow{\$} \text{UniformInt}(1, 3)$ 
3: if  $x > 1$  then
4:   if  $x < y$  then
5:     return True
6:   end if
7: else
8:   return False
9: end if

```

(a) Program



(b) Execution Tree

Fig. 2. An example program and its symbolic execution tree

Example. Consider the code in Fig. 2a and suppose you wish to calculate the probability of the program returning True. Following Alg. 2 for lines 1,2, we generate fresh probabilistic symbolic variables for x and y , δ_1 and δ_2 , respectively. We also store the distributions which δ_1 and δ_2 are samples from, namely the discrete uniform distribution $\mathcal{U}\{1, 3\}$. In our notation, we say that $\sigma = \{x \mapsto \delta_1, y \mapsto \delta_2\}$ and $P = \{\delta_1, \delta_2 \mapsto \mathcal{U}\{1, 3\}\}$. Now following Alg. 3 to process line 4 of Fig. 2a, note that $\mathcal{D} = \{1, 2, 3\} \times \{1, 2, 3\}$ and

$$\begin{aligned}
 p_c &= \Pr[\delta_1 < \delta_2 \mid \delta_1 > 1] = \frac{\Pr[(\delta_1 < \delta_2) \wedge (\delta_1 > 1)]}{\Pr[\delta_1 > 1]} \\
 &= \frac{\sum_{(v_1, v_2) \in \mathcal{D}} [(\delta_1 < \delta_2) \wedge (\delta_1 > 1) \{ \delta_1 \mapsto v_1, \delta_2 \mapsto v_2 \}]}{\sum_{(v_1, v_2) \in \mathcal{D}} [(\delta_1 > 1) \{ \delta_1 \mapsto v_1 \}]} \\
 &= \frac{1}{6}
 \end{aligned}$$

3.3 Query Generation

At this point, we have an execution tree where each branch is annotated with the probability of taking the true and false branch. From this tree we construct a set of paths which are represented by a path condition and set of probabilities. If there are n paths, and φ_i is the i^{th} path condition, and $p_{i,j}$ is the probability of the j^{th} branch in the i^{th} path, then we have $\{(\varphi_1, \{p_{1,1}, \dots, p_{1,k_1}\}), \dots, (\varphi_n, \{p_{n,1}, \dots, p_{n,k_n}\})\}$. To calculate the probability of an entire path, or equivalently, the probability of a path condition, φ ,

being true, we can simply multiple each of the branch probabilities. Recall that p_{i_j} is the *conditional* probability of the j^{th} branch condition being true along path i . Since φ_i is the conjunction of each of the k_i branch conditions, the probability of taking the entire path can be computed using the rule $\Pr[A \wedge B] = \Pr[A \mid B] \Pr[B]$ where A and B are any two events. In other words, since $\varphi_i = c_{i_1} \wedge \dots \wedge c_{i_{k_i}}$, then

$$\begin{aligned} \Pr[\varphi_i] &= \Pr[c_{i_1} \wedge \dots \wedge c_{i_{k_i}}] \\ &= \Pr[c_{i_1} (c_{i_2} \wedge \dots \wedge c_{i_{k_i-1}} \wedge c_{i_{k_i}})] \\ &= p_{i_1} \cdot \Pr[c_{i_2} \wedge \dots \wedge c_{i_{k_i}}] \\ &\vdots \\ &= \prod_{j=1}^{k_i} p_{i_j} \end{aligned}$$

So, if we let $p_i = \prod_{j=1}^{k_i} p_{i_j}$, then we can simplify our set of paths to $\Phi = \{(\varphi_1, p_1), \dots, (\varphi_n, p_n)\}$.

In our interpretation, the ultimate goal of probabilistic symbolic execution is to verify properties of probabilistic programs. In a probabilistic setting this often equates to either proving the upper bound of reaching some “bad” state, or proving the lower bound of reaching some “good” state. We achieve this through queries to an SMT solver, such as Z3 [?]. There are three components to our queries: 1) a universal quantification over the universal symbolic variables, 2) a filtering condition specifying what a “good” or “bad” state is, and 3) a desired upper/lower bound, potentially parameterized by universal symbolic variables.

- (1) If a program has m universal symbolic variables, $\alpha_1, \dots, \alpha_m$, we begin the query with a universal quantification, $\forall \alpha_1, \dots, \alpha_m$, in order to reason over any setting of the non-probabilistic program variables.
- (2) A *filtering condition*, ψ , is a predicate which determines whether a path is considered “bad” or “good”, depending on the property. Some example conditions are whether: a certain value is returned, a false positive (or negative) occurred, or a hash collision occurred. Out of all the paths, Φ , we keep only those which satisfy ψ , Φ' . We then sum over all of the paths in Φ' , giving us the probability of ψ occurring in program S :

$$\sum_{(\varphi, p) \in \Phi'} [\varphi] \cdot p.$$

Note that the probability expression p is multiplied by $[\varphi]$, as during some settings of $\alpha_1, \dots, \alpha_m$, the path represented by φ might not be reachable, and so we should exclude that probability from the sum. The inclusion of the Iverson brackets achieves this desired behavior.

- (3) Let δ (???) be the lower/upper bound *expression* which we want to prove the S does not violate.

A general query then takes the form of

$$\forall \alpha_1, \dots, \alpha_m. \delta \sim \sum_{(\varphi, p) \in \Phi'} [\varphi] \cdot p$$

where \sim is a binary relation (e.g. $>$, $<$, \leq , \geq).

TODO: Finish up once the overview section is done

For example, in the Monty Hall problem as described in Section 2, we were only concerned with

those paths that resulted in the contestant winning the car. So, we would filter on $\psi := \text{win} = \text{true}$, restricting the set of paths to be (TODO: Add on more after writing section 2). The probability of this occurring is then ... as ...

3.4 Formalization

In this section we present the formalization of our method. First, we will present our notation and definitions, and then provide proofs of correctness and soundness of our technique.

3.4.1 Notation & Definitions. The goal of this section is to describe how $R = (\varphi, \sigma, P)$, the inputs to Alg. 3 is an abstraction of a *distribution of program memories* before a branch guarded by a program expression c . To begin, we will define some notation:

- Let Vars be the set of all program variables, ForallSymVars be the set of all universal symbolic variables, ProbSymVars be the set of all probabilistic symbolic variables, $\text{SymVars} = \text{ForallSymVars} \cup \text{ProbSymVars}$ be the combined set of all symbolic variables, and Vals be the set of all values.
- Let $a_f : \text{ForallSymVars} \rightarrow \text{Vals}$ be an assignment of universal symbolic variables to values and let ForallAssign be the set of all such assignments.
- Similarly, let $a_p : \text{ProbSymVars} \rightarrow \text{Vals}$ be an assignment of probabilistic symbolic variables to values and let ProbAssign be the set of all such assignments.
- Let $m : \text{Vars} \rightarrow \text{Vals}$ be a program memory which translates program variables into values, and let Mems be the set of all program memories.
- Let $de : \text{Mems} \rightarrow (\text{Vals} \rightarrow [0, 1])$ be a distribution expression parameterized by program memories, and let DistExprs be the set of all distribution expressions.
- Let $d : \text{ForallAssign} \times \text{Mems} \rightarrow [0, 1]$ be a distribution of program memories parameterized by assignments to universal symbolic variables and let MemDists be the set of all parameterized distributions of program memories.

Additionally, we will use emphatic brackets for two purposes:

- If $e \in \text{ProgExprs}$ is a *program* expression containing the program variables $x_1, \dots, x_n \in \text{Vars}$, and $m \in \text{Mems}$, then

$$\llbracket e \rrbracket m = \text{eval}(e[x_1 \mapsto m(x_1), \dots, x_n \mapsto m(x_n)])$$

- If $e \in \text{SymExprs}$ is a *symbolic* expression containing the symbolic variables $\alpha_1, \dots, \alpha_n \in \text{ForallSymVars}$ and $\delta_1, \dots, \delta_m \in \text{ProbSymVars}$, and $a_f \in \text{ForallAssign}$ and $a_p \in \text{ProbAssign}$, then

$$\llbracket e \rrbracket a_f a_p = \text{eval}(e[\alpha_1 \mapsto a_f(\alpha_1), \dots, \alpha_n \mapsto a_f(\alpha_n), \delta_1 \mapsto a_p(\delta_1), \dots, \delta_m \mapsto a_p(\delta_m)])$$

With this notation in hand, we can now define what it means for R to be an abstraction of a distribution of programs memories.

Definition 3.1. Let $R = (\varphi, \sigma, P)$ be the abstraction generated by the symbolic execution algorithm where $\varphi : \text{ForallAssign} \times \text{ProbAssign} \rightarrow \{0, 1\}$ denotes whether the path condition is true or false under the given assignments, $\sigma : \text{Vars} \rightarrow \text{SymExprs}$ is the mapping from program variables to symbolic expressions generated through symbolic execution, and $P : \text{ForallAssign} \rightarrow \text{ProbSymVars} \rightarrow (\text{Vals} \rightarrow [0, 1])$ is the mapping from probabilistic symbolic variables to the distribution it is sampled from parameterized by assignments of forall symbolic variables. Additionally, for every assignment of forall symbolic variables, $a_f \in \text{ForallAssign}$, $\text{domain}(P(a_f)) = \{\delta_1, \dots, \delta_k\}$. Let $\alpha_1, \dots, \alpha_l \in \text{ForallSymVars}$ be the forall symbolic variables which correspond to the l parameters to the program. For every assignment of probabilistic and forall symbolic variables,

$a_f \in \text{ForallAssign}$, $a_p \in \text{ProbAssign}$, let $v : \text{ForallAssign} \rightarrow (\text{ProbAssign} \rightarrow [0, 1])$ be a distribution of assignments of probabilistic symbolic variables parameterized by assignments of forall symbolic variables, defined as

$$v(a_f, a_p) \triangleq \prod_{i=1}^k \Pr_{v \sim P(a_f, \delta_i)} [v = a_p(\delta_i)].$$

We say that a distribution d satisfies our abstraction R if, for all assignments of forall symbolic variables, $a_f \in \text{ForallAssign}$, $\Pr_{a'_p \sim v(a_f)} [\varphi(a_f, a'_p) = 1] > 0$, and if $\mu : \text{ForallAssign} \rightarrow (\text{ProbAssign} \rightarrow [0, 1])$ is defined as

$$\mu(a_f, a_p) = \frac{\Pr_{a'_p \sim v(a_f)} [a'_p = a_p \wedge \varphi(a_f, a'_p) = 1]}{\Pr_{a'_p \sim v(a_f)} [\varphi(a_f, a'_p) = 1]}.$$

Additionally, define $\text{convertToMem} : (\text{Vars} \rightarrow \text{SymExprs}) \rightarrow \text{ForallAssign} \rightarrow \text{ProbAssign} \rightarrow \text{Mem}$ as

$$\text{convertToMem}(\sigma, a_f, a_p) \triangleq \lambda(x : \text{Vars}) . \llbracket \sigma(x) \rrbracket a_f a_p,$$

and let $\text{convertFromMem}(\sigma, a_f, m) = (\text{convertToMem}(\sigma, a_f))^{-1}(m)$. Then,

$$d(a_f, m) = \sum_{a_p \in \text{convertFromMem}(\sigma, a_f, m)} \mu(a_f, a_p).$$

We additionally define the semantics for the three main types of statements which concerns probabilistic symbolic execution: assignments, probabilistic samples, and branches.

Definition 3.2 (Assignment Semantics). Let $d \in \text{MemDists}$ be a distribution of program memories parameterized by assignments to forall symbolic variables. Let $x = e$ be an arbitrary assignment of the program variable $x \in \text{Vars}$ to the program expression $e \in \text{ProgExprs}$. Let $\text{assign}_{x=e} : \text{Mems} \rightarrow \text{Mems}$ is defined as

$$\text{assign}_{x=e}(m) = \lambda(y : \text{Vars}) \begin{cases} \llbracket e \rrbracket m & \text{if } x = y \\ m(y) & \text{otherwise} \end{cases}$$

and let $\text{unassign}_{x=e} = \text{assign}_{x=e}^{-1}$. Then we define $d_{x=e}$ to be the distribution over program memories parameterized by assignments to forall symbolic variables after executing the assignment statement $x = e$ to be

$$d_{x=e}(a_f, m) = \sum_{m' \in \text{unassign}_{x=e}(m)} d(a_f, m').$$

Definition 3.3 (Sampling Semantics). Let $d \in \text{MemDists}$ be a distribution of program memories parameterized by assignments to forall symbolic variables. Let $x \sim de$ be an arbitrary sampling instruction which assigns the program variable $x \in \text{Vars}$ to a random element from the distribution of values parameterized by a memory, represented as a distribution expression $de \in \text{DistExprs}$. Let $\text{desample} : \text{Vars} \times \text{Mems} \rightarrow \mathcal{P}(\text{Mems})$ be defined as

$$\text{desample}(x, m) = \{m' \in \text{Mems} \mid \forall (y \in \text{Vars}) . (y \neq x \wedge m'(y) = m(y))\}.$$

Then, we define $d_{x \sim de}$ to be the distribution over program memories parameterized by assignments to forall symbolic variables after executing the sampling statement $x \sim de$ to be

$$d_{x \sim de}(a_f, m) = \sum_{m' \in \text{desample}(m)} (\llbracket de \rrbracket a_f)(m(x)) \cdot d(a_f, m').$$

Definition 3.4 (Conditional Distribution of Program Memories). Let $a_f \in \text{ForallAssign}$ be an arbitrary assignment of forall symbolic variables, c be a guard of an if condition, $d \in \text{MemDists}$ be a distribution of program memories parameterized by assignments to forall symbolic variables, and $x_1, \dots, x_n \in \text{Vars}$ be all of the program variables in c . Then for all program memories $m \in \text{Mems}$, d conditioned on a guard c being true, represented as d_c is defined as

$$d_c(a_f, m) = \frac{\Pr_{m' \sim d(a_f)} [m' = m \wedge \llbracket c \rrbracket m' = \text{true}]}{\Pr_{m' \sim d(a_f)} [\llbracket c \rrbracket m' = \text{true}]}$$

Similarly, for all program memories $m \in \text{Mems}$, d conditioned on a guard c being false, represented as $d_{\neg c}$ is defined as

$$d_{\neg c}(a_f, m) = \frac{\Pr_{m' \sim d(a_f)} [m' = m \wedge \llbracket c \rrbracket m' = \text{false}]}{\Pr_{m' \sim d(a_f)} [\llbracket c \rrbracket m' = \text{false}]}$$

3.4.2 Proofs.

4 IMPLEMENTATION

Subhajit & Sumit: Write something about KLEE and the implementation here!

We use a robust dynamic symbolic execution engine, KLEE to generate the *path constraints* corresponding to each path that our tool explores. For each path in the program, KLEE stores a list of *constraints* that encodes the whole path. These constraints are over *program* variables that have been marked `make_symbolic()` by the user and are stored in KLEE as *metadata* in the *state* data structure during the *dynamic* symbolic execution of the program. Apart from *path constraints*, the *state* also contains a mapping of *symbolic* variables to expressions that are either *symbolic* or *concrete* (in the case of concrete execution of the program) and a list of *instructions* currently getting executed as a part of the *state*.

During *dynamic* execution, KLEE executes each *instruction* in the program and updates the *state* along with the *symbolic* variables mapping (Δ) depending upon *type* of the *instruction* it executes as described in Algorithm 4. KLEE assigns values to *symbolic* variables by employing an SMT SOLVER (Z3 as in our case) for solving the current set of *constraints* seen so far in the program path and also by concretizing some of the values that get evaluated as a result of *concrete* execution (Eg. return value from external function calls).

For the purpose of our implementation, we modify KLEE to support creation of *probabilistic symbolic* variables whose values can be sampled from a distribution and dump the whole set of *path* constraints that is stored in the *state* corresponding to the current execution of the program at each of (1) **assignment** [Algorithm 4, Line 9], (2) **branch** [Algorithm 4, Line 13], & (3) **assume** [Algorithm 4, Line ??] statements. These dumps are later refined and used by the tool for further processing.

Upon reaching the *branch* instruction at Line 13 in Algorithm 4, based on whether the *path constraints* are solvable upon adding the *branch* [Line 13] condition, KLEE forks the current *state* by making two identical copies of it and then appending to each of *new* states one additional constraint encoding the *true* [Line 17] and *false* [Line 18] side of the branch and adds it to the *execution stack* containing *states* that need to be explored next. The process continues until all the *states* in the *execution stack* are explored.

5 CASE STUDIES

In this section we will briefly explain each of the case studies that we use in our evaluation (Section 6). For each case study, we will explain, (1) what the algorithm does, (2) which variables are concretized, universally quantified, and probabilistic, and (3) the property we aim to verify using

our technique. Note that we frame the queries as an existential query and hope to get “UNSAT” in order to reason over all the possible values for the universal symbolic variables.

5.1 Freivalds’ Algorithm

Freivalds’ algorithm [?] is a randomized algorithm used to verify matrix multiplication in $O(n^2)$ time. Given three $n \times n$ matrices A , B , and C , Freivalds’ algorithm checks whether $A \times B = C$ by generating a random $n \times 1$ vector containing 0s and 1s, \vec{r} and checks whether $A \times (B\vec{r}) - C\vec{r} = (0, \dots, 0)^T$. If so, the algorithm outputs “Yes”, and “No” otherwise. However, if $A \times B \neq C$, the probability that the algorithm returns “Yes” is at most $\frac{1}{2}$.

While the size of the matrices has to be concretized, the elements of the three matrices can be represented by universal symbolic variables and the elements of r as probabilistic symbolic variables. We want to verify the false positive error rate of $\frac{1}{2}$. To do this, we can ask Z3, for a fixed n , whether there exist any $n \times n$ matrices A , B , and C where $A \times B \neq C$ such that

$$\Pr[\text{freivalds}(A, B, C) = \text{Yes}] > \frac{1}{2}.$$

5.2 Randomized Response

Randomized response is a surveying technique which allows respondents to answer in a way that provides “plausible deniability” Before answering the query, a coin is flipped. If “tails”, then the respondent answers truthfully, if “heads”, a second coin is flipped and the respondent answers “Yes” if “heads” and “No” if tails. In fact, this method is $(\ln 3, 0)$ -differentially private.

The answer the respondent would give if they answered truthfully can be represented with a universal symbolic variable, and the results of the two coin flips as probabilistic symbolic variables. With this model, we can prove $(\ln 3, 0)$ differential privacy by asking our method if there exists a setting of the “truth” such that

$$\frac{\Pr[\text{Response} = \text{Yes} \mid \text{text} = \text{Yes}]}{\Pr[\text{Response} = \text{Yes} \mid \text{text} = \text{No}]} \neq 3$$

and

$$\frac{\Pr[\text{Response} = \text{No} \mid \text{text} = \text{No}]}{\Pr[\text{Response} = \text{No} \mid \text{text} = \text{Yes}]} \neq 3.$$

5.3 Reservoir Sampling

Algorithm 5 Reservoir Sampling

```

1: function RESERVOIRSAMPLING( $A[1..n], S[1..k]$ )
2:   for  $i = 1$  to  $k$  do
3:      $S[i] \leftarrow A[i]$ 
4:   end for
5:   for  $i = k + 1$  to  $n$  do
6:      $j \leftarrow \text{UniformInt}(1, i)$ 
7:     if  $j \leq k$  then
8:        $S[j] \leftarrow A[i]$ 
9:     end if
10:  end for
11:  return  $S$ 
12: end function

```

Reservoir sampling is an online, randomized algorithm to get a simple random sample of k elements from a population of n elements. It uses uniform integer samples to maintain a set of k elements drawn from the set of n elements. For the full algorithm, see Alg. 5.

The sizes of both the population, n , and the sample, k , need to be concretized. The elements from A , are universal symbolic variables, and each sample, j , are probabilistic symbolic variables. The property that we want to check is whether each sample has an equal probability of being returned by Alg. 5, namely $\frac{1}{\binom{n}{k}}$.

5.4 Schwartz-Zippel Lemma

The Schwartz-Zippel lemma is a probabilistic method of polynomial identity testing, that is, the problem of determining whether a given multivariate polynomial is identically equal to 0 or not. Given a non-zero polynomial of total degree d , P , over a field F , and r_1, \dots, r_n selected at random from a subset of F , say S , then the lemma states that

$$\Pr[P(r_1, \dots, r_n) = 0] \leq \frac{d}{|S|}.$$

The number of terms in P , n , has to be concretized, the degrees of each of the terms, and the size of the subset S all need to be concretized. The coefficients are represented as universal symbolic variables and r_1, \dots, r_n are probabilistic symbolic variables. The property that we want to check is whether there exists a non-zero polynomial, P , such that

$$\Pr[P(r_1, \dots, r_n) = 0] > \frac{d}{|S|}.$$

5.5 Bloom Filter

A Bloom filter is a space-efficient, probabilistic data structure used to rapidly determine whether an element is in a set. A Bloom filter is a bit-array of n bits and k associated hash functions, each of which maps elements in the set to places in the bit-array. To insert an element, x , into the filter, x is hashed using each of the k hash functions to get k array positions. All of the bits at these positions are set to 1. To check whether an element y is in the filter, y is again hashed by each of the hash functions to get k array positions. The bits at each of these positions are checked and the filter reports that y is in the filter if, and only if, each of the bits are set to 1. Note that false positives are possible due to hash collisions, but false negatives are not.

In order to bound the false positive error rate, most implementations of Bloom filters take in the expected number of elements to be inserted as well as the desired false positive error rate. From these two quantities, the optimal size of the bit-array, n , as well as the number of hash functions, k can be computed. If m is the expected number of elements and ε is the desired error rate, then

$$n = -\frac{m \ln \varepsilon}{(\ln 2)^2}$$

$$k = -\frac{\ln \varepsilon}{\ln 2}$$

We want to prove that for a given m and ε that the actual false positive rate does not exceed ε .

With our method, m and ε first need to be concretized. Then, we insert m elements x_i , where each x_i can be represented using a universal symbolic variable. We model each of the k uniform hash functions using the method described in Sec. (TODO: Insert bit about hash functions). We then want to check if there exist x_1, \dots, x_m such that the false positive rate exceeds ε .

5.6 Quicksort

Quicksort is a popular sorting algorithm which uses partitioning in order to achieve efficient sorting. One way to choose a pivot element is by way of a uniform random sample. Using this pivot method, the expected number of pivots required is $1.386n \log_2(n)$ where n is the length of the array. If we concretize the length of the array, but represent the elements as universal symbolic variables, we can compute the expected number of pivots required to sort the array.

6 EVALUATION

7 RELATED WORK

Subhajit & Sumit: If you could both start looking into related work (Mayhap, original PSE paper, Axprof, P4WN, PSI), that would be great!

8 CONCLUSION & FUTURE WORK

ACKNOWLEDGMENTS

A APPENDIX

Text of appendix ...