
Supplementary Material for "Localizing Bugs in Program Executions with Graphical Models"

Laura Dietz

Max-Planck Institute for Computer Science
Saarbruecken, Germany
dietz@mpi-inf.mpg.de

Valentin Dallmeier

Saarland University
Saarbruecken, Germany
dallmeier@cs.uni-saarland.de

Andreas Zeller

Saarland University
Saarbruecken, Germany
zeller@cs.uni-saarland.de

Tobias Scheffer

Potsdam University
Potsdam, Germany
scheffer@cs.uni-potsdam.de

Derivation of the Predictive Distribution

Given a collection \mathcal{G}_m of previously seen execution graphs for method m and a new execution $G_m = (V_m, E_m, L_m)$, Bayesian inference determines the likelihood $p((u, v) \in E_m | V_m, \mathcal{G}_m, \alpha_\psi, \beta_\psi)$ of each of the edges (u, v) , thus indicating unlikely transitions in the new execution of m represented by execution graph G_m . Since we employ independent models for all methods m , inference can be carried out for each method separately.

In order to infer the probability of an edge, Equation 1 integrates over the model space.

$$p((u, v) \in E_m | V_m, \mathcal{G}_m, \alpha_\psi, \beta_\psi) = \int p((u, v) \in E_m | V_m, \Psi) p(\Psi | \mathcal{G}_m, \alpha_\psi, \beta_\psi) d\Psi \quad (1)$$

According to the Bernoulli graph model, the likelihood of the existence of an edge (u, v) given the parameter vector Ψ is a Bernoulli distribution. The distribution is conditioned on existence of the start vertex u and yields zero probability if the labels do not overlap appropriately:

$$p((u, v) \in E_m | V_m, \Psi) = \begin{cases} \psi_{m, s_1 \dots s_n} & \text{if } u \in V_m, L_m(u) = s_1 \dots s_{n-1}, \text{ and } L_m(v) = s_2 \dots s_n \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

Corresponding to Figure 3a), the random variable $b_{G, u, s} = \text{true}$ iff an edge $(u, v) \in G$ exists in the graph such that $L(u) = s_1 \dots s_{n-1}$ and $L(v) = s_2 \dots s_{n-1} s$. The likelihood of a graph is proportional to the likelihood that all edges in E_G are generated ($b_{G, u, s} = \text{true}$) and all edges in the complementary set \bar{E}_G are not (Equation 3). If the start vertex u is not contained, random variables b are false with probability 1, thus we can omit them from the product over \mathcal{G} in Equation 4. This product with shared parameter ψ yields a Binomial distribution. Since the product of Binomial and Beta distributions yields a reparametrized Beta distribution, we arrive at Equation 5 using counts of successful and failed trials. Shorthand $\#_{(u, s_n)}^{\mathcal{G}}$ abbreviates the number of graphs $G \in \mathcal{G}$ with $b_{G, u, s_n} = \text{true}$ which is the case if an edge (u, v) between vertices labeled $L(u) = s_1 \dots s_{n-1}$ and $L(v) = s_2 \dots s_n$ exists; and $\#_u^{\mathcal{G}}$ refers to the number of graphs $G \in \mathcal{G}$ that have a vertex u labeled $s_1 \dots s_{n-1}$, in which case a draw from ψ is issued.

$$p(\Psi|\mathcal{G}_m, \alpha_\psi, \beta_\psi) \propto p(\Psi|\alpha_\psi, \beta_\psi) \prod_{G \in \mathcal{G}_m} \underbrace{p(V_G)}_{\text{const}} p(E_G|V_G, \Psi) (1 - p(\bar{E}_G|V_G, \Psi)) \quad (3)$$

$$\propto \prod_{s_1 \dots s_{n-1}} \prod_{s \in S} \left(p_\beta(\psi_{m, s_1 \dots s_{n-1} s} | \alpha_\psi, \beta_\psi) \prod_{\substack{G \in \mathcal{G}_m | u \in V_G \\ L(u) = s_1 \dots s_{n-1}}} p(b_{G, u, s} | \psi_{m, s_1 \dots s_{n-1} s}) \right) \quad (4)$$

$$\propto \prod_{s_1 \dots s_n \in (S_m)^n} p_\beta \left(\psi_{m, s_1 \dots s_{n-1} s_n} | \#_{(u, s_n)}^{\mathcal{G}} + \alpha_\psi, \#_u^{\mathcal{G}} - \#_{(u, s_n)}^{\mathcal{G}} + \beta_\psi \right) \quad (5)$$

The predictive distribution in Equation 1 using a Beta-distributed posterior has an analytic solution:

$$p((u, v) \in E_m | V_m, \mathcal{G}_m, \alpha_\psi, \beta_\psi) = \frac{\#_{(u, s_n)}^{\mathcal{G}} + \alpha_\psi}{\#_u^{\mathcal{G}} + \alpha_\psi + \beta_\psi}. \quad (6)$$

By definition, an execution graph G for an execution contains a vertex if its label is a substring of the execution's trace t . Likewise, an edge is contained if an aggregation of the vertex labels is a substring of t . It follows that $\#_u^{\mathcal{G}} = \#\{t \in T | s_1 \dots s_{n-1} \in t\}$ and $\#_{(u, s_n)}^{\mathcal{G}} = \#\{t \in T | s_1 \dots s_n \in t\}$. Equation 6 can be reformulated as in Equation 7 to predict the probability of seeing the code position $\tilde{s} = s_n$ after a fragment of preceding statements $\tilde{f} = s_1 \dots s_{n-1}$ using the trace representation of an execution. Thus, it is not necessary to represent execution graphs G explicitly.

$$p(\tilde{s} | \tilde{f}, T, \alpha_\psi, \beta_\psi) = \frac{\#\{t \in T | \tilde{f}\tilde{s} \in t\} + \alpha_\psi}{\#\{t \in T | \tilde{f} \in t\} + \alpha_\psi + \beta_\psi} \quad (7)$$