

---

# Hijacking Simulators with Universal Probabilistic Programming

---

Paper ID: 1010

## Abstract

Notes on how to develop a formalism to perform inference in population based simulators.

## 1 Introduction

One challenge in performing inference in such models is the combinatorial increase in *paths* that can be taken within the given probabilistic model. This is because each member of the population has a bounded, but large number of decisions to make. In addition to this, each member of the population can interact with other members of the population. To model a single population member quantum-based methods leveraging linear superposition and entanglement provide an ability to explore multiple paths simultaneously maintaining the whole program, *state*, in one iteration. Unfortunately, analytically writing the functional form of the state is, in this instance, impossible. The usability of current classical inference methods in this context is also problematic. However, by combining hijacking and amortizing parts of the simulator, in particular rejection sampling loops, we can not only correctly perform inference in such models, as the rejection sampling loops are now replaced with learnt functions, but we can also perform inference much more efficiently. This is because the rejection sampling loops are now replaced with learnt functions which enables us to correctly construct the density and as they are amortized, evaluation is cheap and only has to be performed once per call. This significantly reduces memory consumption and running time.

Without the amortization procedure, as updates are made sequentially MCMC methods are typically not-well suited this type of recursive estimation problem. As each time we get a new data point  $y_{t+1}$  we have to run new MCMC simulations for  $P(x_{0:t_1} | y_{0:t+1})$ , we also cannot directly reuse the samples  $\{x_{0:t}^i\}$ . However, as we can only run the simulator `forward()` we cannot return back to a previous state and re-run. So adapting existing inference techniques is the only viable option. MCMC methods can be effective with a well designed amortized inference procedure (**Note: This is not guaranteed**). Non-MCMC based methods such as importance sampling (IS) and sequential monte carlo (SMC), may also be sensible inference engines, although we shall discuss the limitations within the simulator context in the proceeding sections.

However, to implement this in practice we are required to do the following:

- Create a sensible training procedure
- Have 2 additional functions implemented within `pyprob`
  1. 1) That stores only the necessary raw samples from the desired rejection sampling addresses to be used for training.
  2. 2) A function at *test time* that replaces those given addresses with a learnt surrogate.
- Finer grained control over the actual latent variables of interest, to avoid unnecessary memory consumption and ensure that inference problems remain tractable.

There is no doubt that implementing some of these features will be more invasive to the underlying simulator code than previously just hijacking the raw-random number draws, but, it would enable

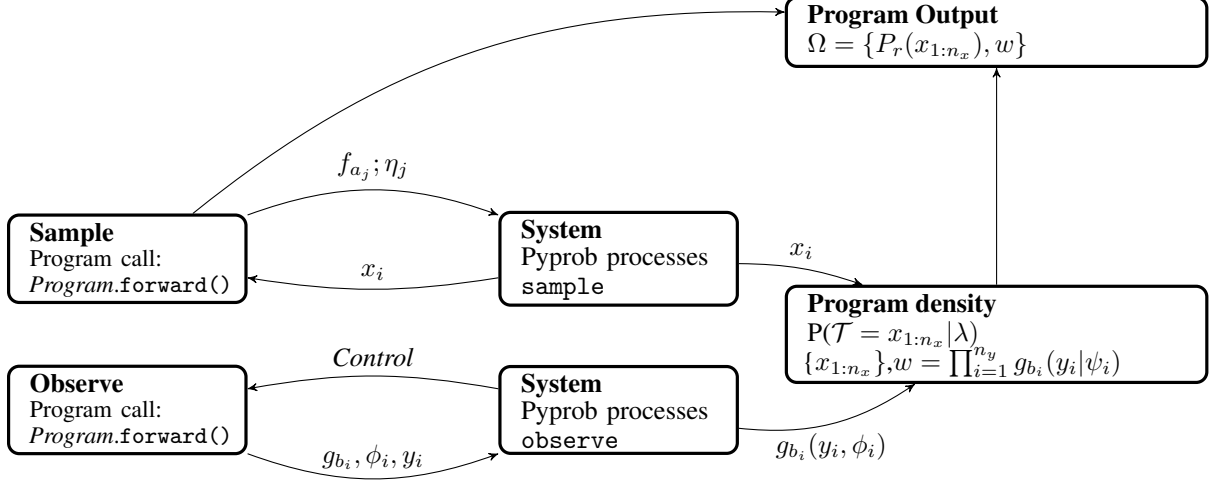


Figure 1: An overview of the hijacking system.  $x_i$  represent latent variables.  $f_{a_i}$  represent raw sample calls,  $g_{b_j}$  represent conditioning statements.  $\mathcal{T}$  corresponds to the program trace.  $w$  represents the collective weights.  $r$  represents run number.  $n_x$  represents the number of latent variables and  $n_y$  represents the number of observations, which will always be fixed, but is model dependent.  $\lambda$  represents all the other variables generated by running the simulator `.forward()`.

us to attack difficulties in the inference procedures that are currently not amenable to the current approach.

## 2 Limitations with different inference engines

### 2.1 Prior based sampling

In prior based sampling we take our existing program and look directly at the product of all the raw sample calls, denoted  $f_{a_i}(x_i|\eta)$  and construct a proposal that is dependent on those. In this sampling scheme our proposal is defined as:

$$q(x_{1:n_x}) = \begin{cases} \prod_{i=1}^{n_x} f_{a_i}(x_i|\eta) & \text{if } \mathcal{B}(x_{1:n_x}) = 1 \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

where  $\mathcal{B}(\cdot)$  is satisfied by the simulator and  $f_{a_i}(\cdot)$  represents a raw random sample,  $\eta$  is a subset of the variables in-scope at the point of sampling and  $n_x$  is the number of latent variables. This is true in all inference engines used in this hijacking framework. However, this makes use of no conditioning and is not sufficient for performing inference in complex, population-based simulations.

## 3 Non-prior Based Sampling

In non-prior based sampling we utilise conditioning too, which enables us to explore spaces more efficiently. In the current hijacking set-up, we can only do conditioning at the end of each `forward()` run of the simulator. This does inhibit certain inference engines and makes aspects of performing inference within this setting more complicated. Nonetheless, there are several inference schemes that we can exploit and utilise. We shall now give an Introduction to each of these schemes.

### 3.1 Importance Sampling

In the generalised probabilistic programming setting an importance sampling scheme over the program raw samples and observations can be defined as follows:

$$p(\mathcal{T} = x_{1:n_x} | \lambda) = \begin{cases} \prod_{i=1}^{n_x} f_{a_i}(x_i|\eta) \prod_{j=1}^{n_y} g_{b_j}(y_j|\phi) & \text{if } \mathcal{B}(x_{1:n_x}) = 1 \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

where the observe statements are defined by  $g_{b_j}(\cdot)$ ,  $\phi$  plays the same role as  $\eta$  and  $n_y$  represents the number of observations. The expectation over the whole program can be defined as:

$$\mathbb{E}_{p(\Omega|\lambda)}[h(\Omega)] = \mathbb{E}_{p(x_{1:n_x}|\lambda)}[h(P_r(x_{1:n_x}))] \quad (3)$$

**Note: It is a little weird how we define this, as the the program output  $\Omega$  should also include the product of the weights**

where  $r$  represents a run of the program,  $P_r(\cdot)$  is the entirety of the program output,  $h(\cdot)$  is just a deterministic function.

As the importance sampler updates the weights according to the following formula:

$$w(x_{1:n_x}) = \prod_{i=1}^{n_x} \frac{f_{a_i}(x_i|\eta_i)}{q_i(x_i|f_{a_i}, \eta_i)} \prod_{j=1}^{n_y} g_{b_j}(y_j|\phi_j) \quad (4)$$

When the number of observations is fixed, the combined weight is equal to  $w(x_{1:n_x}) = \prod_{j=1}^{n_y} w_j$ . To generate the weight up to each observe we can define  $K(j)$  to be a function that returns the least  $i$  before hitting an observe  $j$ . This enables us to write the  $j$ -th weight as follows:

$$w_j(x_1 : K(j)) = g_{b_j}(y_j|\psi_j) \prod_{K(j-1)+1}^{K(j)} v_i \quad (5)$$

where we define  $v_j = \frac{f_{a_j}(x_j|\eta_j)}{q_j(x_j|f_{a_j}, \eta_j)}$

Taking these definitions we can construct a target density over the latents that is proportional to the products of sample and observe statements **Note: Add update posterior here**. Thus, given the proposal  $q(x_{1:K(t)})$  we define the importance weight as:

$$w(x_{1:t}) = \prod_{j=1}^t w_j(x_{1:K(j)}) = \frac{p_t(x_{1:K(t)})}{q(x_{1:K(t)})} = \begin{cases} \prod_{j=1}^t g_{b_j}(y_j|w_i) \prod_{K(j-1)+1}^{K(j)} v_i & \text{if } B_t(x_1 : K(t)) = 1 \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

One problem with importance sampling is in its ability to scale and condition upon multiple observations. | Choosing a *good* prior remains challenging.

### 3.2 Sequential Monte Carlo

Given multiple observations, we can ake the inference more computationally efficient by employing sequential monte carlo (SMC), although as we can only condition at the end of one full run of the simulation, computational speed increases are minimal. SMC allows for conditioning on multiple observations, but the number of observations must be fixed. **Note: If we could condition on observations within the simulator, we could get dramatic speed increases over standard importance sampling.**

To implement SMC in the current set-up we round require a forking method and the

### 3.3 Random-walk Metropolis Hastings

Due to the current set-up performing Random-walk Metropolis Hastings (RMH) is not feasible as the rejection sampling blocks within the simulator mean that we cannot correctly perform RMH as we do not have a functional form for all parts of the simulator to construct a target density that we can evaluate. **Note: Ask tom what: one can sometimes achieve improvements by instead using the type of the distribution object associated with  $x_m$  (what does this first part mean?) to automatically construct a valid random walk proposal that allows for improved hill climbing behavior on the individual updates. .**

---

**Algorithm 1** SMC for Population Based Simulators - well any general program - not new
 

---

```

1: procedure SMC( $T, M, \mathbf{x}, \mathbf{w}, \mathbf{v}, \mathbf{g}$ )
2:   for  $t = 1 : T$  do
3:     for  $m = 1 : M$  do
4:        $w^m = 1$ 
5:       for  $k = K(t-1) + 1 : K(t)$  do
6:          $x_k^m \sim q_k(x_k^m | f_{a_k^m}, \eta_j^m)$ 
7:          $w_t^m \leftarrow w_t^m v_j^m(x_k^m)$ 
8:       end for
9:     end for
10:     $w_t^m \leftarrow w_t^m g_{b_j^m}(y_t^m | \psi_t^m)$ 
11:     $\{x_{1:t}^{1:m}\} \leftarrow \text{resample}(\{x_{1:t}^{1:m}\}, \{w_{1:t}^{1:m}\})$ 
12:  end for
13: end procedure

```

---

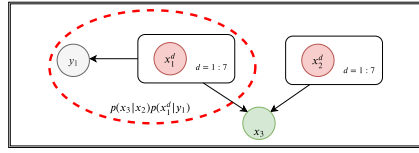


Figure 2: The DAG of the Parasite model.

## 4 An Example: Modelling Parasite Densities in Simulation

The parasite model by Smith et al. [11] is a simple stochastic model for determining the number of infections within a given population demographic. The models hypothesized in [11] were found, after extensive testing and development, to not characterise well the number of infections among different demographics, especially those under 15 years old. By using probabilistic programming we can immediately understand the expressiveness of the proposed model and can extract important insights for extending the proposed models, which increases one's ability to develop improved models to estimate parasite densities, given our observations of the seasonal entomological inoculation rate (EIR). In addition to this we can greatly reduce the computational costs, whilst improving model comprehension and expressivity.

**Common features of the models** The models used are all discrete time micro-simulations of malaria in humans, based on the published model used in [1], which we refer to as the base model. The model for infection of the vector as a function of parasite densities [2] (and hence for the effect of vaccination on onward transmission), the case-management model [3], the model of vaccination by pre-erythrocytic vaccines [4], the models for clinical outcomes and mortality [5,6], and the basic design of the scenarios used for predicting impacts [3,7] all correspond to the previous descriptions, which are accompanied by detailed rationales for the formulations used.

For all models, the number of infections introduced in individual  $i$  in five-day time step  $t$ , is distributed as:

$$h(i, t) \sim \text{Pos}(\lambda(i, t)) \quad (7)$$

where the expected number  $\lambda(i, t) = S(i, t)E_a(i, t)$  and  $S(i, t)$  is the susceptibility (proportion of inoculations that result in infections) and  $E_a(i, t)$  is the expected number of entomological inoculations, adjusted for age and individual factors.

**Base model** Although a statistical model is not formally defined within [11] we define the unknown *disease dynamics* parameters with a set of latent variables  $\{x_i\}^N$  and construct the model as follows:

In the base model  $E_a(i, t) = E_{max}(t) \frac{A(a(i, t))}{A_{max}}$  where  $E_{max}(t)$  refers to the usual measure of the EIR computed from human bait collections on adults  $A(\cdot)$  and, the availability to mosquitoes, is assumed to be proportional to average body surface area and to depend only on age  $a(i, t)$ . The function used to compute  $A(a(i, t))$  increases with age up to age 20 years where it reaches a value of  $A_{max}$ .

To capture observed relationships between infection and measured entomological exposure, the susceptibility is then:

$$S(i, t) = (S_\infty + \frac{1 - S_\infty}{1 + \frac{E_a(i, t)}{E^*}})(S_{imm} + \frac{1 - S_{imm}}{1 + (\frac{X_p(i, t)}{X_p^*})^{\gamma_p}}) \quad (8)$$

where  $S_{imm}, X_p^*, E^*, \gamma_p, S_\infty$  are latent variables and  $X_p(i, t) = \int_{-a(i, t)}^t E_a(i, \tau) d\tau$ . Fitting this model to data of [3] using their method learns the values  $S_\infty = 0.049$  and  $E^* = 0.032$  inoculations/person-night. The model for each individual infection  $j$  in host  $i$  comprises a time series of parasite densities. In a naïve host the expectation of the logarithm of the density of the infection at time point  $\tau$  in the course of the infection,  $\ln(y_0(i, j, \tau))$ , is taken from a statistical description of parasite densities in malaria therapy patients [6]. In the presence of previous exposure and co-infection, the expected log density for each concurrent infection is:

$$E(\ln(y(i, j, \tau))) = D_y(i, t)D_m(i, t)D_h(i, t)\ln(y_0(i, j, \tau)) + \ln(\frac{D_x}{M(i, t)} + 1 - D_x) \quad (9)$$

where  $M(i, t)$  is the total multiplicity of infection and

$$D_y(i, t) = \frac{1}{1 + \frac{X_y(i, j, t)}{X_y^*}} \quad (10)$$

where  $X_y(i, j, t) = \sum_{t-a}^t Y(i, t) - \sum_{t_0, j}^t y(i, j, \tau)$  (note that a continuous time approximation to this is given in the original publications [10, 6] and hence measures the cumulative parasite load, where

$$D_y(i, t) = \frac{1}{1 + \frac{X_h(i, t)}{X_h^*}} \quad (11)$$

where  $X_h(i, t) = \sum_{t-a}^t h(i, \tau) - 1$ , the number of inoculations since birth, excluding the one under consideration, which measures the inocula experienced by the host up to the time point under consideration.

$$D_m(i, t) = 1 - \alpha_m \exp(-\frac{0.693a(i, t)}{a_m^*}) \quad (12)$$

which measures the effect of maternal immunity.  $X_y^*, X_h^*, D_x, a_m^*$  and  $\alpha_m$  are essentially latent variables to be learnt given the field data.

**Mass action model of the force of infection** The base model assumes that susceptibility is a sigmoidal function of the exposure, which conflicts with the usual mass action assumption for infectious disease transmission. However, if entomological exposure is over-dispersed, infection-exposure relationships similar to those observed in the field can be reconciled with mass action principles, thus providing explanation for the observed non-linear relationship between infection and entomological exposure and hence a more justifiable representation of the infection process than that used in the base model [12]. We simulate this by including random variation in the availability of the human host to mosquitoes. This variation comprises an individual-specific component which is captured by defining for each simulated human a baseline availability,  $A_b(i)$ , which is sampled at birth using:

$$\log(A_b(i)) \sim \mathcal{N}(-\frac{\sigma_b^2}{2}, \sigma_b^2) \quad (13)$$

so that  $A_b(i)$  is log-normally distributed with arithmetic mean 1. The age adjusted EIR at time  $t$ , adjusted to the individual is then:

$$E_b(i, t) = E_{max}(t)A_b(t)\frac{A(a(i, t))}{A_{max}} \quad (14)$$

Additional log normal variation is introduced at each time step in the simulation to make the expected number of entomological inoculations a function both of the individual and of log-normal noise measured by the variance parameter,  $\sigma_n^2$ :

$$\ln(E_a(i, t)) \sim \text{Normal}(\ln(E_b(i, t)) - \frac{\sigma_n^2, \sigma_n^2}{}) \quad (15)$$

The susceptibility in these models is independent of EIR, i.e.:

$$S(i, t) = S_{\infty} \left( S_{imm} + \frac{1 - S_{imm}}{1 + \left( \frac{X_p(i, t)}{X_p^*} \right)^{\gamma_p}} \right) \quad (16)$$

Three different parameterisations, XML files, of this model were considered in [8] (E0063, R0065 and R0068). In each case, the overall variability in  $E_a(i, t)$  was constrained to the same total, estimated by fitting to the data of [3]. A probabilistic programming approach to determining these parameters would have provided a way to leverage state of the art inference algorithms that enable conditioning on arbitrary models. The parameters learnt from the original approach taken by the Malaria community generated values of  $\sigma_b^2 = 0.5, 0.3$  and  $0.05$  for each model respectively, corresponding to  $\sigma_n^2 = 0.187, 0.567$  and  $0.704$ . **TODO: To finish - Add prob prog example + plots**

**Models of decay of natural immunity To complete**

**TODO: Add some of the plots generated from this model and an accompanying analysis to them.**

## 5 Amortized Rejection Sampling

One common problem in traditional simulators, *programs*, outside of the standard probabilistic programming framework is that they rely heavily on rejection sampling to generate samples from the target density. Not only is rejection sampling computationally inefficient, [especially if criteria in which samples are generated has been poorly calibrated] (may leave out this statement as it may make a reviewer question why would amortization work any better.), but it creates complications for inference. As the rejection sampling loops provide no explicit form of the target density, we can not simply evaluate the density and leverage statistically and computationally efficient algorithms such as MCMC sampling. This inhibits our ability to transform an arbitrary stochastic simulator into a probabilistic programming framework [2]. To this end, we derive both a mathematical and engineering based framework to facilitate the amortization of rejection sampling loops in arbitrary simulators, which provides a path way to both leverage a larger, more efficient class of inference methods and convert arbitrary stochastic simulators into probabilistic programs systems.

## 6 Rejection Sampling

Rejection samplers are heavily used when designing simulators, due to their ease of use and ability to model any arbitrary distribution on  $\mathbb{R}^D$ , with a sufficiently well designed rejection sampling scheme [8, 7, 4] and have their origins in 1940 - 1950's, during the development of other prominent Monte carlo techniques at Los Alamos.

In the standard setting we let  $X$  be a random variable whose density is  $p(x)$  that we cannot easily sample from, but can evaluate. Using a random variable  $Z$ , whose density is  $g(z)$  that we can efficiently sample from, then we can generate a sample  $X$  by sampling  $Z$  and accept  $Z$  with probability  $p_{accept} = \frac{p(x)}{Mg(x)}$ , which we can do until acceptance.  $M$  is constant that has a finite bound,

$M \in (1, \infty)$ , on the likelihood ratio  $\frac{p(x)}{g(x)}$  over the support of  $X$ .  $M$  must also satisfy the following inequality  $p(x) \leq Mg(x) \forall x$ . This implies that the support of  $Z$  must also contained the support of  $X$ , that is  $g(x) > 0$  whenever  $p(x) > 0$ . The standard algorithm is defined as follows:

As useful rejection sampling is, it has many short comings. In particular, to simulate proposals for certain distributions, such as truncated normal distributions, regardless of the easy to sample from proposal chosen the computational running time of such rejection statements can be shown to have an expectation equal to infinity. Thus, the ability to learn and replace such rejection sampling blocks with learnt surrogate functions, that are expressive enough to reciprocate the behavior of the rejection sampling block buys the user two things. A) It reduces the computational complexity of the simulator. B) As we now have a direct functional form of the density, we can correctly leverage more computationally and sample efficient inference schemes.

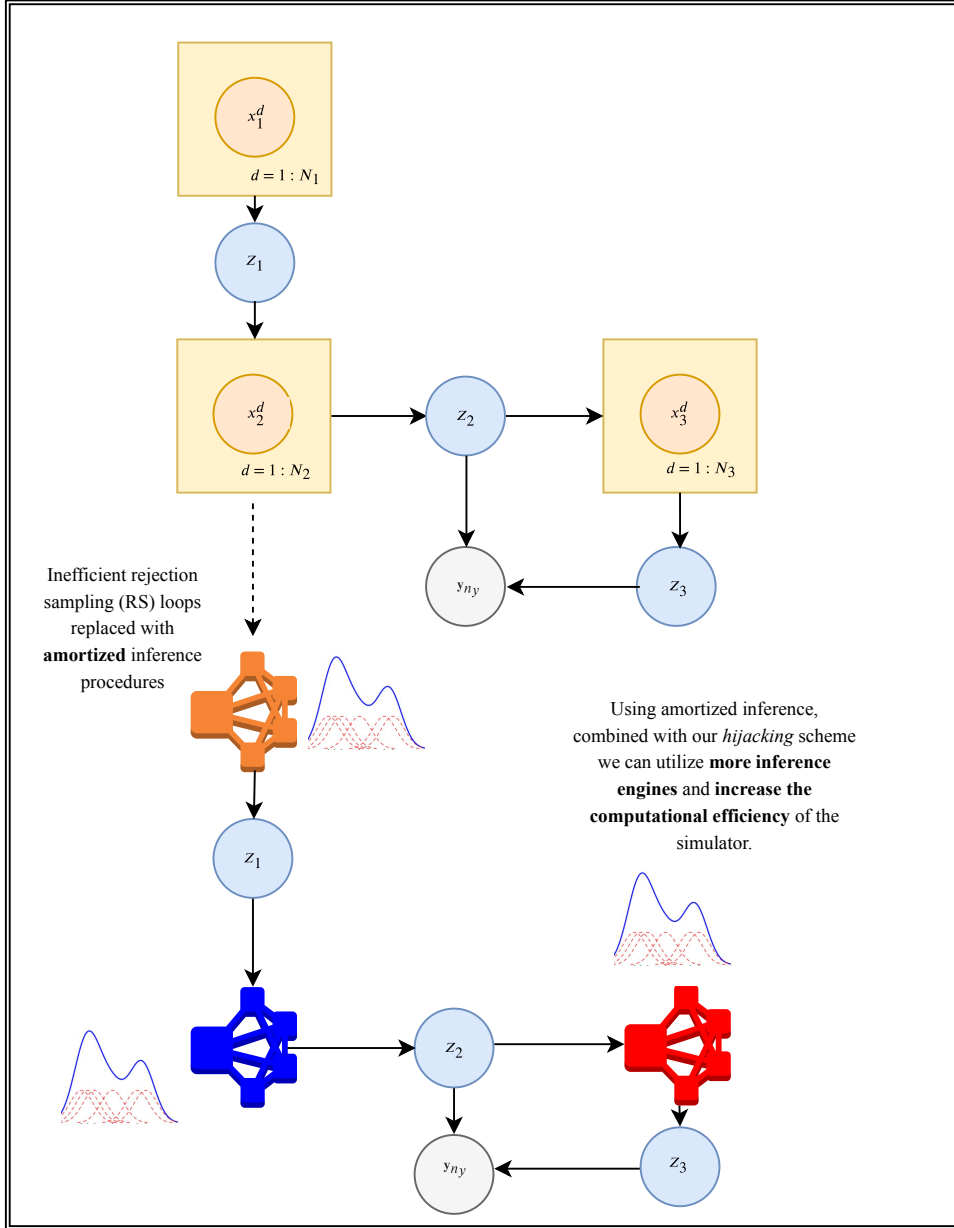


Figure 3: Using PyProb we can locate all regions of inefficient sampling and replace those points in the code with learnt functions that represent the distributions of the rejection sampling statements.  $Z_i$  represents the marginal and by learning surrogate functions for each of the rejection sampling loops. In doing so we can learn approximate raw sample distributions  $f(x_i|\eta)$  for each rejection sampling loop. This means that we can evaluate the function at each part in the simulator, enabling us to construct a target density that can be evaluated when using MCMC inference schemes such as RMH.

---

**Algorithm 2** Rejection Sampling

---

```

At iteration  $i$  s.t  $i \geq 1$ 
 $X_i \sim g(X_i)$ 
 $U_i \sim \mathcal{U}(0, 1)$ 
if  $U_i \leq M \frac{f(X_i)}{g(X_i)}$  then
    accept  $X_i \sim f$ 
else
     $i \leftarrow i + 1$ , repeat
end if

```

---

## 7 Details

We do this by Replacing the rejection sampling blocks with learnt functions from an amortized inference procedure Figure 3. This then enables us to replace rejection sampling loops with amortized functions that are quick to evaluate, provide a functional form to the density and significantly reduce memory over-heads, as compared to [2], due to not having to track the thousands of raw sample calls  $f_{a_i}$  generated in standard simulators.

Consider the following program, that is representative of a simple stochastic simulator and two types of rejection sampler:

Stochastic Simulator	Rejection Sampler 1	Rejection Sampler 2
<pre> def f():     z1 ~ N(0, 1)     z2 ← R1(z1, 1)     z3 ~ U(0, 2)     z4 ← R2(z2, z3, 1)      return g(z1:4) </pre>	<pre> def R1(z1):     x ~ N(z1, 1)     if x &gt; 0:         return x     else:         return R1(z1) </pre>	<pre> def R2(z2, z3):     y ← 0     while y &lt; z2:         y ~ N(z3, 1)     end     return y </pre>
$P(\mathcal{T} = x_{1:n_x}   \lambda)$	$q_1(z_2   z_1)$	$q_2(z_4   z_2, z_3)$

Here the rejection sampling functions take an additional argument 1, or 0, to specify whether or not we want to collect samples to train, 1, else replace with the learnt function, 0.  $P$  represents the program density and  $q_1$  and  $q_2$  represents the unknown densities that we wish to learn through an amortized approach.

However, in order to do this, we must correctly specify our variational objective.

## 8 An overview of Amortized Inference schemes

**Bradley: Todo: Tom has provided me with all the information to conduct this approach**

- Whilst the variational posterior we choose as our proposal has a functional form and can be sampled, the functional form of the rejection sampling loops is not explicitly known, which leads to problems when calculating the  $KL$ -divergence.
- There are two approaches that we can take to calculating some divergence metric.
- 1) Using something like a non-parametric divergence estimator that makes use of Renyi and  $L_2$  divergences (see [1] for details)
- 2) Use some form of kernel density to estimation to approximate the functional form of the rejection sampling block, which can be used in place. From which we learn a variational approximation characterized by a trained network that we can sample from quickly.
- An issue with approach 2) is that if I can approximate the density in functional form aren't there more efficient methods that we can use? (**Ask tom**)
- Or, if assuming that the rejection sampling block can be arbitrary dimension, are we first learning a flow to represent that density and then learn a variational approximation of something that we can easily sample from.



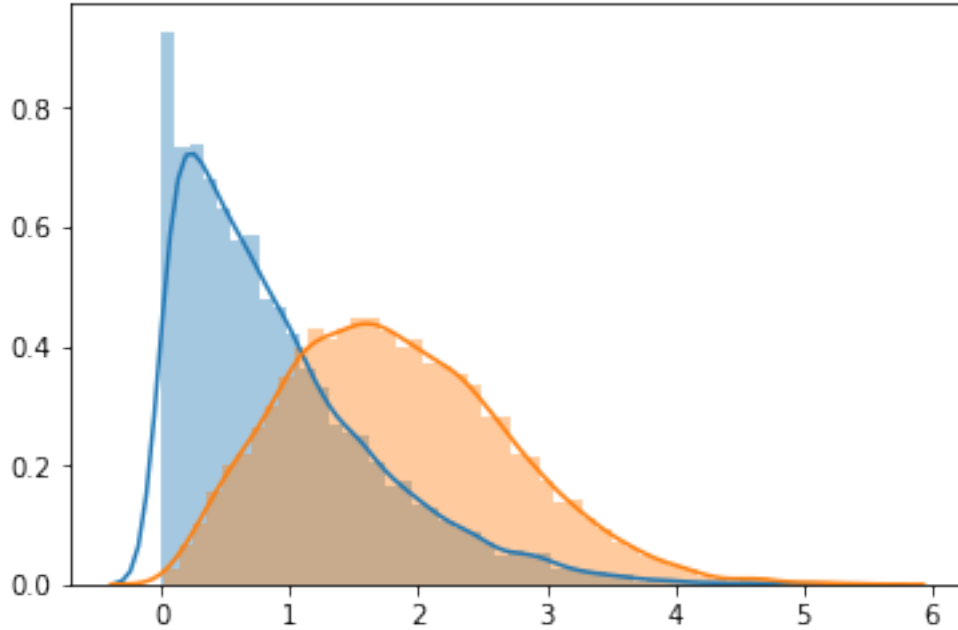


Figure 4: Simulated densities for rejection sampling block 1, *blue* and rejection sampling block 2, *orange*.

- But, this is very much the idea behind normalising flows, take a simple distribution that you can sample and then map it to the “complicated” density.
- Stein Variational inference may be quite nice to use here [5], but it does make explicit assumptions in requiring the densities to be differentiable.
- For an overview of all variational methods, including the use of different divergence metrics see here [14]

### 8.1 Using a simple Neural Network regressor

As a first approach, we use a regressor to learn a function to map between the input / output pairs that are outputted from the rejection sampling blocks with in the simulator.

**Architecture** For each rejection sampling block we learn a proposal  $q(z_1, \dots, z_i | \eta)$

### 8.2 Using variants of VAEs

### 8.3 Density estimation using Normalising flows

Notes from models that define an explicit, tractable density are highly effective, because they permit the use of an optimization algorithm directly on the log-likelihood of the training data. However, the family of models that provide a tractable density is limited, with different families having different disadvantage. The rejection sampling loops, if proposed correctly, provide a way of generating samples from some density, but do not provide a way of constructing the target density (**Bradley: Ask Tom, is this true?**). Within event and population-based simulators, the rejection sampling loops have been designed with care, to ensure that they are representative to the events that they are simulating. So we can, with confidence, assume that they are representative of some density. However, to utilise powerful inference engines, we must be able to construct the full density of the program.

One recent technique for density estimation, in the amortized inference setting, is normalising flows combined with variation inference [9]. Flow-based methods work by leveraging a simple property of

differential geometry applied to random variables. That being the change of variables. Let  $z \sim r(z)$  where  $r(z)$  represents the posterior defined by the rejection sampling loop, and define  $\pi$  to be a smooth and invertible function. When  $\pi$  acts  $z$  we transform the random variable from one domain, to another. Thus, if  $z' = \pi(z)$  then we can construct the density of the probability density of the new random variable  $z'$  under the flow initiated by  $\pi$ :

$$r(z') = r(\pi^{-1}(z')) \left| \det \frac{\partial \pi^{-1}}{\partial z'} \right| = r(\pi^{-1}(z')) \left| \det \frac{\partial \pi}{\partial z} \right|^{-1} \quad (17)$$

- Do we need to know the base distribution that the RS block is sampling from? i.e there maybe be multiple  $r(z)$  in each RS block.
- Is it right to say here that  $z'$  is the output of the RS block, we do not know  $\pi$  here, or is  $r(z)$  the raw sample  $f_{a_i}$  and  $\pi$  the distribution specified by the rejection sampling loop?
- We need to learn  $\pi(\cdot)$

## 9 Constructing the Objective

The framework we present in this section is generally applicable to stochastic simulator environment, where there exists rejection sampling loops.

We start by noting that the prior, defined in the context of the outer samples is defined as:

$$p(x) = \begin{cases} \prod_{j=1}^{n_x} f_{a_j}(x_j; \psi_j) & \text{if } \mathcal{B}(x) = 1 \\ 0 & \text{otherwise} \end{cases} \quad (18)$$

We start with the information projection we construct the variational objective as follows:

$$\begin{aligned} J(\eta) &= KL(p(x) || q(x; \eta)) \\ &= \int p(x) \ln \left( \frac{p(x)}{q(x; \eta)} \right) dx \\ &\quad \text{Let } a_j \in \{1, \dots, L\}, S_R \subset \{1, \dots, L\} \\ &= \mathbb{E}_{p(x)} \left[ \sum_{j=1, a_j \notin S_R}^{n_x} \ln(p(x)) \right] - \mathbb{E}_{p(x)} \left[ \sum_{j=1, a_j \in S_R}^{n_x} \ln(q_{a_j}(x_j; \eta)) \right] \end{aligned}$$

no dependence on  $\eta$

Thus, to define the objective for each subproblem  $r \in S_R$  as:

$$J_r(\eta_r) = -\mathbb{E}_{p(x)} \left[ \sum_{j=1}^{n_x} \mathbb{I}(a_j = r) \ln(q_r(x_j; \eta_r, \phi_j)) \right] \quad (19)$$

where during training we extract samples from each forward run and train each rejection sampler separately.

## 10 Additional TODOs

- Before implementing such tools into the simulator, spend time looking for simpler stochastic simulators that can be hijacked but have target densities that can be constructed analytically.
- If we are not utilizing the additional information stored in the trace, which is only used during the inference compilation training procedure, then we do not need it all, as for importance sampling we only require the weights after each `forward()` update.
- If we cannot find any pre-existing simulators that are amenable to this approach then, we need to create our own, with a rejection sampling loop, that can be exploited by the new machinery. This would enable us to debug, provide a simple example for the paper and understand the feasibility of the entire approach.

## References

- [1] Michele Basseville. Divergence measures for statistical data processing—an annotated bibliography. *Signal Processing*, 93(4):621–633, 2013.
- [2] Atilim Gunes Baydin, Lukas Heinrich, Wahid Bhimji, Bradley Gram-Hansen, Gilles Louppe, Lei Shao, Kyle Cranmer, Frank Wood, et al. Efficient probabilistic inference in the quest for physics beyond the standard model. *arXiv preprint arXiv:1807.07706*, 2018.
- [3] John C Beier, Charles N Oster, Fred K Onyango, James D Bales, James A Sherwood, Peter V Perkins, David K Chumo, Davy V Koech, Richard E Whitmire, Clifford R Roberts, et al. Plasmodium falciparum incidence relative to entomologic inoculation rates at a site proposed for testing malaria vaccines in western kenya. *The American journal of tropical medicine and hygiene*, 50(5):529–536, 1994.
- [4] George Casella, Christian P. Robert, and Martin T. Wells. *Generalized Accept-Reject sampling schemes*, volume Volume 45 of *Lecture Notes–Monograph Series*, pages 342–347. Institute of Mathematical Statistics, Beachwood, Ohio, USA, 2004.
- [5] Yihao Feng, Dilin Wang, and Qiang Liu. Learning to draw samples with amortized stein variational gradient descent. *arXiv preprint arXiv:1707.06626*, 2017.
- [6] Nicolas Maire, Thomas Smith, Amanda Ross, Seth Owusu-Agyei, Klaus Dietz, and Louis Molineaux. A model for natural immunity to asexual blood stages of plasmodium falciparum malaria in endemic areas. *The American journal of tropical medicine and hygiene*, 75(2\_suppl):19–31, 2006.
- [7] Luca Martino, Jesse Read, and David Luengo. Improved adaptive rejection metropolis sampling algorithms. *arXiv preprint arXiv:1205.5494*, 2012.
- [8] Renate Meyer, Bo Cai, and François Perron. Adaptive rejection metropolis sampling using lagrange interpolation polynomials of degree 2. *Computational Statistics and Data Analysis*, 52(7):3408 – 3423, 2008.
- [9] Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. *arXiv preprint arXiv:1505.05770*, 2015.
- [10] Thomas Smith, Gerry F Killeen, Nicolas Maire, Amanda Ross, Louis Molineaux, Fabrizio Tediosi, Guy Hutton, Jürg Utzinger, Klaus Dietz, and Marcel Tanner. Mathematical modeling of the impact of malaria vaccines on the clinical epidemiology and natural history of plasmodium falciparum malaria: Overview. *The American journal of tropical medicine and hygiene*, 75(2\_suppl):1–10, 2006.
- [11] Thomas Smith, Nicolas Maire, Klaus Dietz, Gerry F Killeen, Penelope Vounatsou, Louis Molineaux, and Marcel Tanner. Relationship between the entomologic inoculation rate and the force of infection for plasmodium falciparum malaria. *The American journal of tropical medicine and hygiene*, 75(2\_suppl):11–18, 2006.
- [12] Thomas A Smith. Estimation of heterogeneity in malaria transmission by stochastic modelling of apparent deviations from mass action kinetics. *Malaria journal*, 7(1):12, 2008.
- [13] Howard G Tucker. A generalization of the glivenko-cantelli theorem. *The Annals of Mathematical Statistics*, 30(3):828–830, 1959.
- [14] C. Zhang, J. Bütepage, H. Kjellström, and S. Mandt. Advances in variational inference. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(8):2008–2026, Aug 2019.

## A The unconditional acceptance probability

Given a collection of i.i.d random variables  $\{X_i\}_{i=1:n}$  with cumulative density function (CDF)  $F_X$ . Then by the Glivenko-Cantelli theorem [13] we can define the empirical distribution function for  $x \in \mathbb{R}$  as:

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{[X_i, \infty)}(x) \quad (20)$$

And for any fixed  $x \in \mathbb{R}$ , the Strong law of large numbers ensures that  $F_n(x) \rightarrow F_X(x)$  as  $n \rightarrow \infty$  as

$$\mathbb{E}[\mathbb{I}_{[X_i, \infty)}(x)] = P[\mathbb{I}_{[X_i, \infty)}(x) = 1] = P[X_i \leq x] = F_X(x) \quad (21)$$

where  $\mathbb{I}_A(x)$  is the indicator function for a set  $A$ .

$$\text{Prob}(U \leq \frac{p_X(z)}{Mg_Z(z)}) = \mathbb{E}(\cdot) \quad (22)$$

## B Law of total expectations

If  $X$  is a random variable whose expected value is  $\mathbb{E}[X]$  and  $Z$  is a random variable on the same probability space then:

$$\mathbb{E}[X] = \mathbb{E}_{p(z)}[\mathbb{E}_{p(x|z)}[X|Z]] \quad (23)$$

That is, the expected value of  $X$  given  $Z$  is the same as the expectation of  $X$ .