# Delayed Sampling with Funsors

FRITZ OBERMEYER, Uber AI

ELI BINGHAM, Uber AI

Delayed sampling is an inference technique for automatic Rao-Blackwellization in sequential latent variable models. Funsors are a software abstraction generalizing Tensors and Distributions and supporting seminumerical computation including analytic integration. We demonstrate how to easily implement delayed sampling in a Funsor-based probabilistic programming language using effect handlers for `sample` statements and a `barrier` statement.

## 1 INTRODUCTION

Let us distinguish two types of inference strategies in probabilistic programming, call them *lazy* and *eager*. Let us say a strategy is lazy if it it first symbolically evaluates or compiles model code, then globally analyzes the code to create an inference update strategy. By contrast consider a strategy eager if it eagerly executes model code, drawing samples from each latent variable. For example the autograd-based inference algorithms in Pyro [Bingham et al. 2018] are eager in the sense that samples are eagerly created at each sample site.

However it is often advantageous to combine lazy and eager strategies, performing local exact computations within small parts of a probabilistic model, but drawing samples to communicate between those parts. Examples include Rao-Blackwellized SMC filters and their generalization as implemented in Birch [Murray et al. 2017], and reactive probabilistic programming [Baudart et al. 2019].

This work addresses the challenge of implementing boundedly-lazy inference in a Pyro-like language where samples are eagerly drawn and control flow may depend on those sample values. Our approach is to use Funsors, a software abstraction generalizing Tensors, Distributions, and lazy compute graphs. The core idea is to allow lazy sample statements during program execution, and to trigger sampling of lazy random variables only at user-specified `barrier` statements, typically before control flow. We implement our approach using two effect handlers [??].

## 2 DELAYED SAMPLING

Delayed sampling [Murray et al. 2017] is an inference technique for automatic Rao-Blackwellization in sequential latent variable models. Delayed sampling was introduced in the Birch probabilistic programming language [Murray and Schön 2018].

## 3 FUNSORS

Funsors [Obermeyer et al. 2019] are a software abstraction generalizing Tensors and Distributions and supporting seminumerical computation including analytic integration.

## 4 DELAYED SAMPLING WITH FUNSORS

Consider an embedded probabilistic programming language, extending a host language with two primitive statements and two effect handlers:

- The statement x = sample(name,dist) is a named stochastic statement, where x is a Funsor value, name is a unique identifier for the statement, and dist is a Funsor distribution.
- The statement x = barrier(x) eliminates any free variables from the recursive data structure x, which may contain Funsor values.

- The effect handler `condition({name:data})` conditions a model to observed data and affects only the single `sample` statement with matching `name`.
- The effect handler `log_joint()` records a representation of the cumulative log joint density of a model as a Funsor expression.

We use Python for the host language in this paper.

Consider a model of stochastic control system with piecewise control, attempting to keep a latent state z within the interval $[-10, 10]$

```
1  def model():
2      z = sample("z_init",Normal(0,1))  # latent state
3      k = 0                             # control
4      cost = 0                          # cumulative cost of controller
5      for t in range(1000):
6          if z > 10:                    # control flow depends on z
7              k -= 1
8          elif z < -10:
9              k += 1
10         else:
11             k = 0
12         cost += abs(k)
13         z = sample(f"z_{t}",Normal(z+k,1))
14         x = sample(f"x_{t}",Normal(z,1))
15     return cost
```

Now suppose we want to estimate the total controller cost given a sequence of observations x. One approach to inference is to apply Sequential Monte Carlo (SMC) filtering. To maintain a vectorized population of particles we can rewrite the model using a vectorized conditional (e.g. `where(cond,if_true,if_false)` as implemented in NumPy and PyTorch). Further we can support resampling of particle populations by adding a `barrier` statement; this is needed to communicate resampling decisions with the model's local state.

```
1  def model():
2      z = sample("z_init",Normal(0,1))    # latent state
3      k = 0 * z                           # control
4      cost = 0 * z                        # cumulative cost of controller
5      for t in range(1000):
6          z,k,cost = barrier((z,k,cost))
7          k = where(z > 10, k + 1, k)
8          k = where(z < 10, k + 1, k)
9          k = where(-10 <= z & z <= 10, 0 * k, k)
10         z = sample(f"z_{t}",Normal(z+k,1))
11         x = sample(f"x_{t}",Normal(z,1))
12     return cost
```

See appendix for details of the effect handlers `condition` and `log_joint` to implement SMC inference.

Notice that if there were no control k (or indeed if the control were linear), we could completely Rao-Blackwellize using a Kalman filter: inference via variable elimination would be exact. To implement linear-time exact inference, we require only:

- lazy versions of tensor operations such as *where*;

- a lazy interpretation for `sample` statements;
- a variable-eliminating interpretation of `barrier` statements; and
- a variable-eliminating computation of the final `log_joint` latent state.

See appendix for details of the effect handlers `condition` and `log_joint` to implement SMC inference. This is the approach taken by Pyro's discrete enumeration inference, which leverages broadcasting in the host tensor DSL to implement lazy sampling and lazy tensor ops.

To partially Rao-Blackwellize our SMC inference, we can combine the two approaches, emitting lazily sampled values from `sample` statements and eagerly sampling delayed samples at `barrier` statements. In contrast to the variable-elimination interpretation of `barrier`, this interpretation guarantees all local state is ground and hence can be inspected by conditionals. See appendix for details of effect handlers `condition` and `log_joint` to implement delayed sampling.

## 5 CONCLUSION

We demonstrated flexible inference algorithms in an embedded probabilistic programming language with a new `barrier` statement and support for lazy computations represented as Funsors.

## REFERENCES

Guillaume Baudart, Louis Mandel, Eric Atkinson, Benjamin Sherman, Marc Pouzet, and Michael Carbin. 2019. Reactive probabilistic programming. *arXiv preprint arXiv:1989.07563* (2019).

Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall, and Noah D. Goodman. 2018. Pyro: Deep Universal Probabilistic Programming. *arXiv preprint arXiv:1810.09538* (2018).

Lawrence M Murray, Daniel Lundén, Jan Kudlicka, David Broman, and Thomas B Schön. 2017. Delayed sampling and automatic rao-blackwellization of probabilistic programs. *arXiv preprint arXiv:1708.07787* (2017).

Lawrence M Murray and Thomas B Schön. 2018. Automated learning with a probabilistic programming language: Birch. *Annual Reviews in Control* (2018).

Fritz Obermeyer, Eli Bingham, Martin Jankowiak, Du Phan, and Jonathan P Chen. 2019. Functional tensors for probabilistic programming. *arXiv preprint arXiv:1910.10775* (2019).

## 6 EFFECT HANDLERS FOR INFERENCE VIA SEQUENTIAL MONTE CARLO

TODO

## 7 EFFECT HANDLERS FOR EXACT INFERENCE VIA VARIABLE ELIMINATION

TODO

## 8 EFFECT HANDLERS FOR INFERENCE VIA DELAYED SAMPLING

TODO