

# Auxiliary variables in Probabilistic Programs

Ekansh Sharma   Daniel M. Roy

University of Toronto

January 9, 2018

## Motivation

In existing probabilistic programming languages (PPLs), programmers represent probability models by:

## Motivation

In existing probabilistic programming languages (PPLs), programmers represent probability models by:

- ▶ Giving an exact sampler for the prior distribution, using probabilistic primitive `sample`

## Motivation

In existing probabilistic programming languages (PPLs), programmers represent probability models by:

- ▶ Giving an exact sampler for the prior distribution, using probabilistic primitive `sample`
- ▶ Weighing the prior by data likelihood, using probabilistic primitive `observe` or `score`

## Motivation

In existing probabilistic programming languages (PPLs), programmers represent probability models by:

- ▶ Giving an exact sampler for the prior distribution, using probabilistic primitive `sample`
- ▶ Weighing the prior by data likelihood, using probabilistic primitive `observe` or `score`

We know, this approach is limiting!

## Motivation

In existing probabilistic programming languages (PPLs), programmers represent probability models by:

- ▶ Giving an exact sampler for the prior distribution, using probabilistic primitive `sample`
- ▶ Weighing the prior by data likelihood, using probabilistic primitive `observe` or `score`

We know, this approach is limiting!

- ▶ Certain computable distributions do **not** have an exact sampling representation
  - ▶ Aldous–Hoover representation for an exchangeable array is not necessarily computable (Ackerman et al., 2017)

## Motivation

In existing probabilistic programming languages (PPLs), programmers represent probability models by:

- ▶ Giving an exact sampler for the prior distribution, using probabilistic primitive `sample`
- ▶ Weighing the prior by data likelihood, using probabilistic primitive `observe` or `score`

We know, this approach is limiting!

- ▶ Certain computable distributions do **not** have an exact sampling representation
  - ▶ Aldous–Hoover representation for an exchangeable array is not necessarily computable (Ackerman et al., 2017)
- ▶ In many cases, it might be difficult to come up with an exact sampling representation.

## Motivation

In existing probabilistic programming languages (PPLs), programmers represent probability models by:

- ▶ Giving an exact sampler for the prior distribution, using probabilistic primitive `sample`
- ▶ Weighing the prior by data likelihood, using probabilistic primitive `observe` or `score`

We know, this approach is limiting!

- ▶ Certain computable distributions do **not** have an exact sampling representation
  - ▶ Aldous–Hoover representation for an exchangeable array is not necessarily computable (Ackerman et al., 2017)
- ▶ In many cases, it might be difficult to come up with an exact sampling representation.

## Contribution

Language construct—`slice-let`—that allows programmer to specify a distribution without giving an exact sampler



## Main Idea

Instead of specifying an exact sampler for a random variable  $x$ , introduce an auxiliary (random) variable  $u$ .

Programmer provides an exact representations for:

- ▶ conditional distribution  $u$  given  $x$
- ▶ conditional distribution  $x$  given  $u$

## Main Idea

Instead of specifying an exact sampler for a random variable  $x$ , introduce an auxiliary (random) variable  $u$ .

Programmer provides an exact representations for:

- ▶ conditional distribution  $u$  given  $x$
- ▶ conditional distribution  $x$  given  $u$

These two conditional distributions can be used to generate "exact–approximate" inference algorithms, in particular, algorithms called *slice samplers*.

## Main Idea

Instead of specifying an exact sampler for a random variable  $x$ , introduce an auxiliary (random) variable  $u$ .

Programmer provides an exact representations for:

- ▶ conditional distribution  $u$  given  $x$
- ▶ conditional distribution  $x$  given  $u$

These two conditional distributions can be used to generate "exact-approximate" inference algorithms, in particular, algorithms called *slice samplers*.

## Outline

- ▶ Syntax for `slice-let` in the first order language by Staton et al. (2016).
- ▶ Measure theoretic semantics for `slice-let`.
- ▶ Simple program transformation for an approximate inference algorithm.

# First Order Language (Staton et al., 2016)

## Syntax

*Types:*

$$\mathbb{A}, \mathbb{B} ::= \mathbb{R} \mid \mathbf{1} \mid P(\mathbb{A}) \mid \mathbb{A} \times \mathbb{B} \mid \sum_{i \in I} \mathbb{A}_i$$

*Deterministic Terms:*

$$a, b, c ::= x \mid * \mid (a, b) \mid (i, a) \mid \pi_j(a) \mid \text{case } a \text{ of } \{(i, x) \Rightarrow c_i\}_{i \in I} \\ \mid f(a) \mid \text{normalize}(t)$$

*Probabilistic Terms:*

$$t, u, v ::= \text{sample}(a) \mid \text{score}(a) \mid \text{case } a \text{ of } \{(i, x) \Rightarrow t_i\}_{i \in I} \\ \mid \text{let } x = t \text{ in } u \mid \text{return}(a)$$

*Program:*

$$\mathcal{P} ::= t$$

# First Order Language (Staton et al., 2016)

## Measure theoretic denotational semantics

1. Types are interpreted as measurable spaces  $\llbracket \mathbb{A} \rrbracket$
2. Context  $\Gamma = (x_1 : \mathbb{A}_1, \dots, x_n : \mathbb{A}_n)$  interpreted as measurable space,  $\llbracket \Gamma \rrbracket = \prod_{i=1}^n \llbracket \mathbb{A}_i \rrbracket$
3. Deterministic terms are interpreted as measurable functions  $\llbracket \Gamma \vdash_d a : \mathbb{A} \rrbracket : \llbracket \Gamma \rrbracket \rightarrow \llbracket \mathbb{A} \rrbracket$
4. Probabilistic terms are interpreted as sigma finite kernels  $\llbracket \Gamma \vdash_p t : \mathbb{A} \rrbracket : \llbracket \Gamma \rrbracket \times \Sigma_{\llbracket \mathbb{A} \rrbracket} \rightarrow [0, \infty]$ , where  $\Sigma_{\llbracket \mathbb{A} \rrbracket}$  denotes the  $\sigma$ -algebra on  $\llbracket \mathbb{A} \rrbracket$

## slice-let: Typing Rules

$$\frac{\begin{array}{l} \Gamma \vdash_p t_0 : \mathbb{A}_2 \\ \Gamma, x : \mathbb{A}_1 \vdash_p t_2 : \mathbb{A}_2 \end{array} \quad \begin{array}{l} \Gamma, u : \mathbb{A}_2 \vdash_p t_1 : \mathbb{A}_1 \\ \Gamma, x : \mathbb{A}_1, u : \mathbb{A}_2 \vdash_p t_3 : \mathbb{B} \end{array}}{\Gamma \vdash_p \text{slice-let } u = t_0 \text{ with } x = t_1 \text{ update } u = t_2 \text{ in } t_3 : \mathbb{B}}$$

- ▶  $\Gamma \vdash_p t_0 : \mathbb{A}_2$  specifies the initial distribution for auxiliary variable  $u$ , which may not be its marginal distribution
- ▶  $\Gamma, u : \mathbb{A}_2 \vdash_p t_1 : \mathbb{A}_1$  specifies the conditional distribution for random variable  $x$  given  $u$
- ▶  $\Gamma, x : \mathbb{A}_1 \vdash_p t_2 : \mathbb{A}_2$  specifies the conditional distribution for random variable  $u$  given  $x$

## slice-let: Denotational semantics

## slice-let: Denotational semantics

Semantics for let term is given by:

$$\llbracket \Gamma \vdash_{\rho} \text{let } x = t_1 \text{ in } t_2 : \mathbb{A} \rrbracket (\gamma)(A) = \int \llbracket t_2 \rrbracket (\gamma, x)(A) \llbracket t_1 \rrbracket (\gamma)(dx),$$

where  $A \in \Sigma_{[\mathbb{A}]}$ .



## slice-let: Denotational semantics

Semantics for `let` term is given by:

$$\llbracket \Gamma \vdash_p \text{let } x = t_1 \text{ in } t_2 : \mathbb{A} \rrbracket (\gamma)(A) = \int \llbracket t_2 \rrbracket (\gamma, x)(A) \llbracket t_1 \rrbracket (\gamma)(dx),$$

where  $A \in \Sigma_{[\mathbb{A}]}$ .

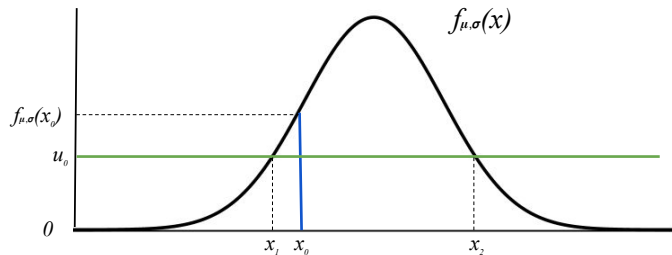
Semantics for `slice-let` term, similar to `let` term, is given by:

$$\begin{aligned} \llbracket \text{slice-let } u = t_0 \text{ with } x = t_1 \text{ update } u = t_2 \text{ in } t_3 \rrbracket (\gamma)(A) \\ = \int \llbracket t_3 \rrbracket (\gamma, x, u)(A) \mu(dx, du) \end{aligned}$$

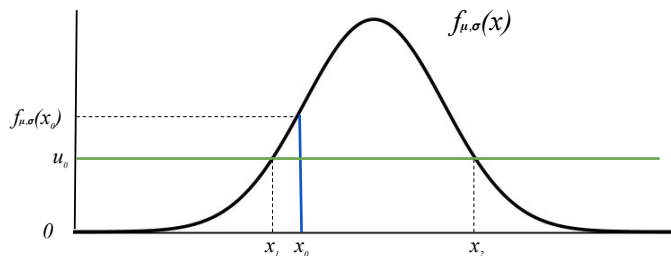
where  $\mu : \Sigma_{[\mathbb{A}]} \otimes \Sigma_{[\mathbb{R}_+]} \rightarrow [0, 1]$  is the unique measure such that, for all bounded measurable functions  $f$  and for all  $\gamma$

$$\begin{aligned} \iint f(x', u) \llbracket t_1 \rrbracket (\gamma, u)(dx') \mu(dx, du) &= \iint f(x, u') \llbracket t_2 \rrbracket (\gamma, x)(du') \mu(dx, du) \\ &= \int f(x, u) \mu(dx, du). \end{aligned}$$

## Example: Slice sampling Gaussian (Neal, 2003)



## Example: Slice sampling Gaussian (Neal, 2003)



For an auxiliary variable representation, programmer needs to specify the following:

- ▶  $u|x \sim \text{Uniform}[0, f_{\mu,\sigma}(x)]$ , where  $f_{\mu,\sigma}$  is the density of a Gaussian with mean  $\mu$  and standard deviation  $\sigma$
- ▶  $x|u \sim \text{Uniform}\{y \mid f_{\mu,\sigma}(y) > u\}$

## Example: Slice sampling Gaussian (Neal, 2003)

Program with `slice-let`

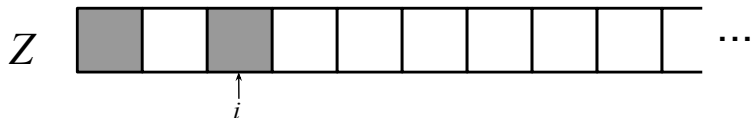
```
slice-let u = sample(Uniform[0,  $f_{\mu,\sigma}(0)$ ])  
  with x = sample(Uniform{y :  $u < f_{\mu,\sigma}(y)$ })  
  update u = sample(Uniform[0,  $f_{\mu,\sigma}(x)$ ])  
  in return(x)
```

## Example: Bernoulli Process with decaying weights

Let  $Z = (z_1, z_2, \dots)$  be an infinite dimensional  $\{0, 1\}$ -valued vector, such that

$$z_i \stackrel{ind}{\sim} \text{Bernoulli}(2^{-i})$$

Interested in representing  $Z$  and the last *active* index of  $Z$

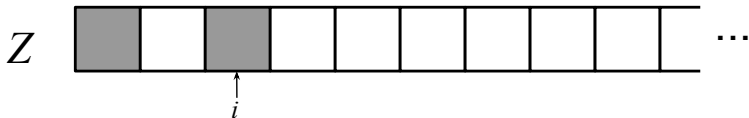


## Example: Bernoulli Process with decaying weights

Let  $Z = (z_1, z_2, \dots)$  be an infinite dimensional  $\{0, 1\}$ -valued vector, such that

$$z_i \stackrel{ind}{\sim} \text{Bernoulli}(2^{-i})$$

Interested in representing  $Z$  and the last *active* index of  $Z$



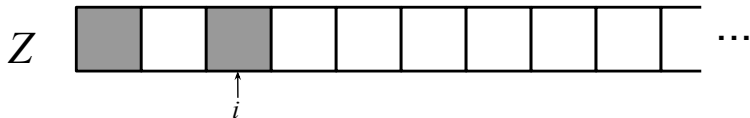
**Problem:** last-active-index:  $\{0, 1\}^{\mathbb{N}} \rightarrow \mathbb{Z}$  is not continuous.

## Example: Bernoulli Process with decaying weights

Let  $Z = (z_1, z_2, \dots)$  be an infinite dimensional  $\{0, 1\}$ -valued vector, such that

$$z_i \stackrel{ind}{\sim} \text{Bernoulli}(2^{-i})$$

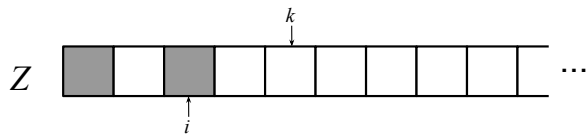
Interested in representing  $Z$  and the last *active* index of  $Z$



**Problem:** last-active-index:  $\{0, 1\}^{\mathbb{N}} \rightarrow \mathbb{Z}$  is not continuous.

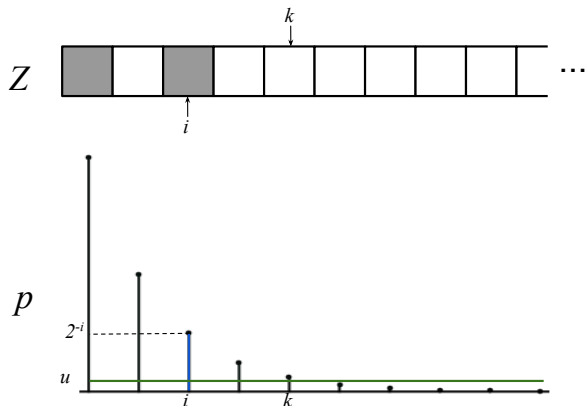
- ▶ Want a representation of  $Z$  and the last active index that exposes the conditional independence in  $Z$

## Example: Bernoulli Process with decaying weights





## Example: Bernoulli Process with decaying weights



Introduce an auxiliary variable  $u$  such that

- ▶  $u|Z, i \sim \text{Uniform}[0, 2^{-i}]$ , where  $i$  is the last active index
- ▶  $p(Z, i|u) \propto p(Z) \frac{1}{2^{-i}} \mathbb{I}(0 \leq u \leq 2^{-i})$

In effect, given  $u$ ,  $\forall j : 2^{-j} < u \implies z_j = 0$ . Thus, only need to sample  $Z$  up to an index  $k$ , given by  $u$ .

# Program Transformation

## Program Transformation for Inference via unfold

For the term  $\phi := \text{slice-let } u = t_0 \text{ with } x = t_1 \text{ update } u = t_2 \text{ in } t_3$ ,  
Add syntactic sugar `unfold`, defined inductively as follows:

$$\begin{aligned} \text{unfold}^1 \phi &:= \text{let } u = t_0 \text{ in let } x = t_1 \text{ in return}(x, u) \\ \text{unfold}^{n+1} \phi &:= \text{let } (x, u) = \text{unfold}^n v \text{ in} \\ &\quad \text{let } u = t_2 \text{ in} \\ &\quad \text{let } x = t_1 \text{ in return}(x, u) \end{aligned}$$

## Program Transformation for Inference via unfold

For the term  $\phi := \text{slice-let } u = t_0 \text{ with } x = t_1 \text{ update } u = t_2 \text{ in } t_3$ ,  
Add syntactic sugar `unfold`, defined inductively as follows:

$$\begin{aligned}\text{unfold}^1 \phi &:= \text{let } u = t_0 \text{ in let } x = t_1 \text{ in return}(x, u) \\ \text{unfold}^{n+1} \phi &:= \text{let } (x, u) = \text{unfold}^n \phi \text{ in} \\ &\quad \text{let } u = t_2 \text{ in} \\ &\quad \text{let } x = t_1 \text{ in return}(x, u)\end{aligned}$$

Sequence of terms  $\{\phi_n\}_{n \in \mathbb{N}}$  by applying the unfold transformation:

$$\phi_n := \text{let } (x, u) = \text{unfold}^n \phi \text{ in } t_3$$

## Assumption

For the term,

$\phi := \text{slice-let } u = t_0 \text{ with } x = t_1 \text{ update } u = t_2 \text{ in } t_3$ ,  
we assume that the terms  $t_0$ ,  $t_1$ , and  $t_2$  defined by the programmer satisfy following conditions:

1.  $\llbracket t_1 \rrbracket$  and  $\llbracket t_2 \rrbracket$  are *probability* kernels whose sequential composition define an ergodic Markov chain with a unique stationary distribution  $\mu$ .
2.  $\llbracket t_0 \rrbracket(\gamma)$  is a probability measure, such that  $\llbracket t_0 \rrbracket(\gamma) \ll \mu$ .

## Assumption

For the term,

$\phi := \text{slice-let } u = t_0 \text{ with } x = t_1 \text{ update } u = t_2 \text{ in } t_3$ ,  
we assume that the terms  $t_0$ ,  $t_1$ , and  $t_2$  defined by the programmer satisfy following conditions:

1.  $\llbracket t_1 \rrbracket$  and  $\llbracket t_2 \rrbracket$  are *probability* kernels whose sequential composition define an ergodic Markov chain with a unique stationary distribution  $\mu$ .
2.  $\llbracket t_0 \rrbracket(\gamma)$  is a probability measure, such that  $\llbracket t_0 \rrbracket(\gamma) \ll \mu$ .

Now, Markov chain theory tells us

$$\forall \gamma, \lim_{n \rightarrow \infty} \|\llbracket \text{unfold}^n \phi \rrbracket(\gamma) - \mu\|_{TV} = 0.$$

## Assumption

For the term,

$\phi := \text{slice-let } u = t_0 \text{ with } x = t_1 \text{ update } u = t_2 \text{ in } t_3$ ,  
we assume that the terms  $t_0$ ,  $t_1$ , and  $t_2$  defined by the programmer satisfy following conditions:

1.  $\llbracket t_1 \rrbracket$  and  $\llbracket t_2 \rrbracket$  are *probability* kernels whose sequential composition define an ergodic Markov chain with a unique stationary distribution  $\mu$ .
2.  $\llbracket t_0 \rrbracket(\gamma)$  is a probability measure, such that  $\llbracket t_0 \rrbracket(\gamma) \ll \mu$ .

Now, Markov chain theory tells us

$$\forall \gamma, \lim_{n \rightarrow \infty} \|\llbracket \text{unfold}^n \phi \rrbracket(\gamma) - \mu\|_{TV} = 0.$$

**Question:** Does the semantics of the program converge when we apply this transformation?

# Convergence Result

For the term

$\phi := \text{slice-let } u = t_0 \text{ with } x = t_1 \text{ update } u = t_2 \text{ in } t_3,$

## Theorem

Let  $\mathcal{P}$  be a program consisting of a **slice-let** term  $\phi$ . Let  $\{\phi_n\}$  be a sequence of terms obtained by applying the **unfold** transformation, such that replacing  $\phi$  with  $\phi_n$  in  $\mathcal{P}$  gives a program  $\mathcal{P}_n$ . If the following conditions hold:

1.  $\llbracket \Gamma \vdash_{\rho} \mathcal{P} : \mathbb{A} \rrbracket(\gamma)(\llbracket \mathbb{A} \rrbracket) < \infty$
  2. If  $\phi$  appears inside a term  $\text{normalize}(t')$ , then  $\llbracket \Gamma \vdash_{\rho} t' : \mathbb{A} \rrbracket(\gamma)(\llbracket \mathbb{A} \rrbracket) < \infty$
  3. If  $\phi$  appears inside  $f(a)$ , then  $f$  is a **continuous** function
- then,  $\llbracket \mathcal{P}_n \rrbracket(\gamma)$  converges strongly to  $\llbracket \mathcal{P} \rrbracket(\gamma)$ .



# Summary

- ▶ Probabilistic primitives, `sample`, `score`, and `normalize` are not enough for all probability models
- ▶ Introduced `slice-let` construct that enables programmer to give exact–approximate representations
- ▶ Discussed a program transformation for an inference algorithm.

## Appendix: Bernoulli Process with decaying weights

We are interested in representing the Bernoulli sequence,  $Z$ , and the last *active* index of  $Z$ .

Idea is that we define the term

$\Gamma, u : \mathbb{R}_+ \vdash_p$  truncated-sample :  $\{0, 1\}^{\mathbb{N}} \times \mathbb{N}$  such that required program becomes:

```
slice-let  $u = \text{sample}(\text{Uniform}[0, 0.5])$   
  with  $(x, i) = \text{return}(\text{normalize}(\text{truncated-sample}))$   
  update  $u = \text{sample}(\text{Uniform}[0, 2^{-i}])$   
  in  $\text{return}(x, i)$ 
```

## Appendix: Bernoulli Process with decaying weights

Define the terms  $\Gamma \vdash_d t_n : \{0, 1\}^n$  inductively as follows:

$$t_{n+1} := \text{let } x = \text{flip } (2^{-(n+1)}) \text{ in return}(t_n, x)$$
$$t_1 := \text{flip } (2^{-1})$$

Define the terms  $\Gamma, x : \{0, 1\}^n \vdash_d \text{last-active}_n : \mathbb{N}$  inductively as follows:

$$\text{last-active}_{n+1} := \text{if } (\pi_{n+1}(x) == 1) \text{ return}(n + 1)$$
$$\quad \text{else let } x = (\pi_1(x), \dots, \pi_n(x)) \text{ in last-active}_n$$
$$\text{last-active}_1 := \text{if } (\pi_1(x) == 1) \text{ return}(1) \text{ else return}(0)$$

Then, we can write the a term

$\Gamma, u : \mathbb{R}_+ \vdash_p \text{truncated-sample} : \{0, 1\}^{\mathbb{N}} \times \mathbb{N}$  as follows:

$$\text{truncated-sample} := \text{case } (\lfloor -\log_2(u) \rfloor, ()) \text{ of}$$
$$\quad (n, ()) \implies$$
$$\quad \text{let } x = t_n \text{ in}$$
$$\quad \text{score}(2^{-\text{last-active}_n});$$
$$\quad \text{return}((x, 0, 0, \dots), \text{last-active}_n)$$