Approximations in Probabilistic Programs

Ekansh Sharma ekansh@cs.toronto.edu University of Toronto Vector Institute

Abstract

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

We introduce a new language construct, stat, which converts the description of the Markov kernel of an ergodic Markov chain into a sample from its unique stationary distribution. Up to minor changes in how certain error conditions are handled, we show that language constructs for soft-conditioning and normalization can be compiled away from the extended language. We then explore the problem of approximately implementing the semantics of the language with potentially nested stat expressions, in a language without stat. For a single stat term, the natural unrolling yields provable guarantees in an asymptotic sense. In the general case, under uniform ergodicity assumptions, we are able to give quantitative error bounds and convergence results for the approximate implementation of the extended first-order language. We leave open the question of whether the same guarantees can be made assuming mere geometric ergodicity.

1 Introduction

28 Approximations are ubiquitous for any practical implemen-29 tation of a probabilistic programming language (PPL) for 30 Bayesian modelling; This is because computing the normal-31 ized posterior of a Bayesian model is intractable. Broadly 32 speaking, there are two types of implementations for com-33 puting the "approximate" posterior: 1) Languages like Stan, 34 Church, and Venture use versions of Markov chain Monte 35 Carlo (MCMC) algorithm to approximate the posterior; 2) 36 Languages like Tensorflow Probability and Pyro use varia-37 tional inference to approximate the posterior [Abadi et al. 38 2015; Bingham et al. 2019; Carpenter et al. 2017; Goodman 39 et al. 2012; Tolpin et al. 2016]. A reasonable question to ask is: 40 Can probabilistic programming systems quantify the error 41 induced by these approximations? Also, do we know how 42 the error scales under composition of multiple "approximate" 43 programs and nested queries?

44 The answer to both the questions is no. One reason is 45 that the semantics of probabilistic languages is not amenable 46 to approximations introduced by the compiler of these real 47 world languages. In this paper we bridge the gap between 48 ideal semantics of probabilistic programming languages and 49 approximations induced by compilers that use MCMC based 50 inference engines. We do this by proposing a new language 51 construct, stat, that takes as input an initial distribution and 52

54 2019.

55

Daniel M. Roy droy@utstat.toronto.edu University of Toronto Vector Institute

a Markov kernel, and outputs the unique stationary distribution corresponding to the Markov kernel, if there exists one. We show that having **stat** in the language is "essentially" equivalent to having constructs **norm** and **score** that compute the posterior distribution. Then in Section 5.1, we give the approximate compiler for the **stat** construct based on a unrolling scheme. We then identify some semantic constraints on the Markov kernel given as argument to the **stat** construct under which we can derive quantitative error bounds on a program.

2 Related Work

This work builds on top of the foundations laid by Staton [2017]; Staton et al. [2016] that gives semantics to the first order probabilistic languages with construct **norm** and **score**. Previously Borgström et al. [2016]; Hur et al. [2015]; Ścibior et al. [2017] have proves the asymptotic correctness of Markov Chain Monte Carlo based inference algorithms, but the do not quantify the error due to finite computation.

Rainforth [2018] gives quantitative bounds on the error due to nested Monte Carlo approximations in probabilistic programs; But this work assumes that we can produce exact samples when the queries are nested.

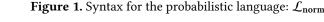
In the Markov chain literature, following papers study the convergence of Markov chain when the transition kernel is approximate [Medina-Aguayo et al. 2018; Roberts et al. 1998]. This is relevant for nested queries. Also, Medina-Aguayo et al. [2018] gives quantitative convergence bounds for Metropolis–Hastings algorithm when the acceptance probability can only be accessed in an approximate manner.

3 Language for MCMC inference

We first give an idealized first-order probabilistic language with the proposed construct **stat** that takes as input a transition kernel for a Markov chain on some state space and returns the stationary distribution associated with the Markov chain. The language we present is based on the first-order probabilistic language introduced and studied by Staton et al. [2016] and Staton [2017], which has constructs for sampling, soft constraints, and normalization. The key differences, which we highlight again below, are (i) a syntactic distinction between probabilistic terms with and without soft constraints, which affects also typing, and (ii) the introduction of the new construct, **stat**.

⁵³ PTML'19, December 14, 2019, Vancouver, BC, Canada

111	Tubou
112	Types:
113	$\mathbb{A}_0, \mathbb{A}_1 ::= \mathbb{R} \mid 1 \mid \mathbb{P}(\mathbb{A}) \mid \mathbb{A}_0 imes \mathbb{A}_1 \mid \sum \mathbb{A}_i$
114	$i \in \mathbb{N}$
115	Terms:
116	deterministic:
117	$a = a = m \left[$
118	$a_0, a_1 ::= x \mid * \mid (a_0, a_1) \mid (i, a) \mid \pi_j(a) \mid f(a)$
119	$ \text{ case } a \text{ of } \{(i, x) \Rightarrow a_i\}_{i \in I}$
120	purely probabilistic:
121	$t_0, t_1 ::= $ sample $(a) $ return $(a) $ let $x = t_0$ in t_1
122	
123	$ \text{ case } a \text{ of } \{(i, x) \Rightarrow t_i\}_{i \in I}$
124	$ $ stat $(t_0, \lambda x.t_1) $ norm (v)
125	probabilistic:
126	$v_0, v_1 ::= t \mid $ let $x = v_0$ in v_1
127	
128	case <i>a</i> of $\{(i, x) \Rightarrow v_i\}_{i \in I}$
129	score (<i>a</i>)
130	Program:
131	t is a program if t is a purely probabilistic
132	
133	with no free variables
134	



3.1 First order language with "stat"

We begin with types and syntax of the language, presented in Figure 1. For the remainder of this paper we call this language \mathcal{L}_{norm} . Along with standard types, this language has a type \mathbb{R} for real numbers and a type $P(\mathbb{A})$ as a type for the space of probability measures on A. The language has all the basic constructors, destructor, case statements, and sequencing. Along with standard programming language constructs, the language has probabilistic features including sample statements that takes in a probability measure as input and returns a sample from it, score statements that scales the prior program with likelihood, and norm term that takes an un-normalized measure and returns the normalized probability measure.

Here, we give a very brief review of the semantics of the language. For a detailed account, we refer the readers to Staton [2017]. Types in the language are interpreted as mea-surable spaces ($[A], \Sigma_{[A]}$). As in [Staton et al. 2016], each term in the language is either deterministic or probabilis-*tic*, satisfying typing judgments of the form $\Gamma \mid_{d} t \colon \mathbb{A}$ and $\Gamma \mid_{\overline{D}} t \colon \mathbb{A}$, respectively, given some environment/context $\Gamma = (x_1: \mathbb{A}_1, ..., x_n: \mathbb{A}_n)$. Letting $\llbracket \Gamma \rrbracket = \prod_{i=1}^n \llbracket \mathbb{A}_i \rrbracket$, a de-terministic term denotes a measurable function from the environment $\llbracket \Gamma \rrbracket$ to $\llbracket A \rrbracket$. As in [Staton 2017], a probabilistic term denotes an *S*-finite kernel from $\llbracket \Gamma \rrbracket$ to $\llbracket \mathbb{A} \rrbracket$. Different

from Staton [2017]; Staton et al. [2016], we distinguish a subset of probabilistic terms we call purely probabilistic, which satisfy an additional typing judgment $\Gamma \downarrow_{p1} t$: A. A purely probabilistic term denotes a probability kernel from $\llbracket \Gamma \rrbracket$ to [A]. Departing again from Staton [2017]; Staton et al. [2016], a program in our language is a purely probabilistic term with no free variables.

3.1.1 Sequencing and sampling terms

In addition to standard let statements and return statements for sequencing, the language has a construct for producing a random sample from a probability distribution.

As in [Staton 2017], the semantics of the let construct is defined in terms of integration as follows:

$$[[let x = t_1 in t_2]]_{Y,A} \stackrel{\text{def}}{=} \int_{[[A]]} [[t_2]]_{Y,x,A} [[t_1]]_{Y,dx}$$

Since both t_1 and t_2 are probabilistic terms, both are interpreted as S-finite kernels; The category of S-finite kernels is closed under composition, thus the term let $x = t_1$ in t_2 is also interpreted as an S-finite kernel.

The semantics of the return statement is given by the kernel $\llbracket \operatorname{return}(t) \rrbracket : \llbracket \Gamma \rrbracket \times \Sigma_{\llbracket \mathbb{A} \rrbracket} \to [0, 1]$

$$\llbracket \mathbf{return}(t) \rrbracket_{\gamma,A} \stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } \llbracket t \rrbracket_{\gamma} \in A \\ 0 & \text{otherwise} \end{cases}$$

Finally the sample statement takes in as argument a deterministic term of type $P(\mathbb{A})$ that is it takes in as argument a probability measure on the space [A]. Thus the semantics is given as:

$$\llbracket \mathbf{sample}(t) \rrbracket_{\gamma,A} = \llbracket t \rrbracket_{\gamma,A},$$

where $\llbracket t \rrbracket_{\gamma} \in \llbracket P(\mathbb{A}) \rrbracket$.

3.1.2 Soft constraints and normalization terms

We are studying a probabilistic language for Bayesian inference; We have terms in the language that is used to scale the prior by the likelihood of some observed data and a term that re-normalizes the scaled measure to return the posterior distribution over the return type. The constructs score and norm are the constructs that, respectively, scale the prior program, and normalize the program to return the posterior probability distribution on the output type, if there exists one. The semantics for the score construct are given by a *S*-finite kernel on the unit type, **1**, as follows:

$$\llbracket \mathbf{score}(a) \rrbracket_{\gamma,A} = \begin{cases} \left| \llbracket t \rrbracket_{\gamma} \right| & \text{if } A = \{()\} \\ 0 & \text{otherwise} \end{cases}$$

The main difference between this semantics and the denotational semantics of the language proposed in Staton [2017]; Staton et al. [2016] is in the semantics of norm. We interpret the semantics of **norm** terms as a probability kernel on the

sum space given as $\llbracket \mathbf{norm}(t) \rrbracket$: $\llbracket \Gamma \rrbracket \times \Sigma_{\mathbb{A}+1} \to [0, 1]$ defined as

$$\llbracket \mathbf{norm}(t) \rrbracket_{\gamma,A} = \begin{cases} \frac{\llbracket t \rrbracket_{\gamma, \lbrace u \rbrace (0, u) \in A \rbrace}}{\llbracket t \rrbracket_{\gamma, \llbracket A \rrbracket}} & \text{if } \llbracket t \rrbracket_{\gamma, \llbracket A \rrbracket} \in (0, \infty) \\ 0 & \text{else if } (1, ()) \notin A \\ 1 & \text{else if } (1, ()) \in A \end{cases}$$

where $A \in \llbracket A + 1 \rrbracket$. The key distinction is that we are not able to determine if the term $\llbracket t \rrbracket_{\gamma}$ is an infinite measure or a null measure.

3.1.3 Stationary terms

One of the main contributions of this paper is that we propose a new feature in the probabilistic language that takes as argument a Markov chain transition kernel on some state space and returns the stationary distribution associated with the kernel.

We do this by allowing the users to define a transition kernel on some measurable space using a standard lambda expression. Following is the syntax and typing rules for the stationary term:

$$\frac{\Gamma \mid_{p_1} t_0 \colon \mathbb{A} \quad \Gamma, x \colon \mathbb{A} \mid_{p_1} t_1 \colon \mathbb{A}}{\Gamma \mid_{p_1} \mathsf{stat}(t_0, \lambda x. t_1) \colon \mathbb{A} + 1}$$

To give the denotational semantics for the **stat**-term, we first introduce to the meta language to define the function $ST: \llbracket P(\mathbb{A}) \rrbracket \llbracket \mathbb{A} \rrbracket \to \llbracket P(\mathbb{A}) \rrbracket + 1$ as follows: For some $k: X \times \Sigma_X \to [0, 1]$,

$$ST(k) = \begin{cases} (0,\mu) & \text{if } \mu \text{ is unique and} \\ & \left\| \mu(\cdot) - \int_X k(x,\cdot)\mu(dx) \right\|_{\text{tv}} = 0 \\ (1,()) & \text{otherwise} \end{cases}$$

Now, we define the semantics to the stationary term as:

$$\llbracket \operatorname{stat}(t_0, \lambda x. t_1) \rrbracket_{\gamma, A} = \begin{cases} \mu_{\gamma}(A) \text{ if } \operatorname{ST}(\llbracket t_1 \rrbracket_{\gamma}) = (0, \mu_{\gamma}) \\ 0 & \text{ if } \operatorname{ST}(\llbracket t_1 \rrbracket_{\gamma}) = (1, ()) \text{ and} \\ (1, ()) \notin A \\ 1 & \text{ if } \operatorname{ST}(\llbracket t_1 \rrbracket_{\gamma}) = (1, ()) \text{ and} \\ (1, ()) \in A. \end{cases}$$

4 Removing soft constraints and normalization terms from \mathcal{L}_{norm}

So far we have modified the probabilistic language proposed by Staton et al. [2016] to include a new construct stat. Consider the language in Figure 2. We call this language \mathcal{L}_{stat} . The key difference between the \mathcal{L}_{stat} and \mathcal{L}_{norm} is that \mathcal{L}_{stat} does not have the language constructs for soft constraints and normalization. Following theorem says that \mathcal{L}_{stat} is equivalent to \mathcal{L}_{norm} .

Theorem 4.1 (Equivalence). For the languages \mathcal{L}_{norm} and \mathcal{L}_{stat} , there exists a function ϕ that maps \mathcal{L}_{norm} -phrases to

Types: $\mathbb{A}_{0}, \mathbb{A}_{1} :::= \mathbb{R} \mid \mathbf{1} \mid \mathbb{P}(\mathbb{A}) \mid \mathbb{A}_{0} \times \mathbb{A}_{1} \mid \sum_{i \in \mathbb{N}} \mathbb{A}_{i}$ Terms: deterministic: $a_{0}, a_{1} :::= x \mid * \mid (a_{0}, a_{1}) \mid (i, a) \mid \pi_{j}(a) \mid f(a)$ $\mid \mathbf{case} \ a \ \mathbf{of} \ \{(i, x) \Rightarrow a_{i}\}_{i \in I}$ $purely \ probabilistic:$ $t_{0}, t_{1} :::= \mathbf{sample}(a) \mid \mathbf{return}(a) \mid \mathbf{let} \ x = t_{0} \ \mathbf{in} \ t_{1}$ $\mid \mathbf{case} \ a \ \mathbf{of} \ \{(i, x) \Rightarrow t_{i}\}_{i \in I}$ $\mid \mathbf{stat}(t_{0}, \lambda x. t_{1})$

Program:

t is a program if *t* is a purely probabilistic with no free variables

Figure 2. Syntax for the probabilistic language: \mathcal{L}_{stat}

 \mathcal{L}_{stat} -phrases such that under the semantics described in Section 3 following statements hold:

- $\phi(t)$ is an $\mathcal{L}_{\text{stat}}$ -program for all $\mathcal{L}_{\text{norm}}$ programs t;
- $\varphi(\mathbb{F}(e_1, \ldots, e_a)) = \mathbb{F}(\varphi(e_1), \ldots, \varphi(e_a))$ for all features of $\mathcal{L}_{\text{stat}}$
- For all programs t, and for all γ , $\left\| \begin{bmatrix} t \end{bmatrix}_{\gamma} \begin{bmatrix} \phi(t) \end{bmatrix}_{\gamma} \right\|_{tv} = 0$

Proof idea. First, we identify that all **score** terms in a program in \mathcal{L}_{norm} is encapsulated in the **norm** statement. For every **norm** statement in a program, we can construct a Markov kernel that proposes a move with the prior probability distribution described by the program; Then accept/reject the move using the standard Metropolis–Hastings acceptance probability. This transition kernel has the stationary distribution that is semantically equivalent to the normalize term. Thus we can use the **stat** construct to compile away **norm** and **score**.

5 Approximate compilation of probabilistic programs

We saw in Section 3, to compute the the normalized measure one term **norm**(*t*) from an un-normalized probabilistic term *t*, we need to compute a normalization factor $\int_{[\![A]\!]} [\![t]\!]_{Y,dx}$. In the previous section we showed that the language \mathcal{L}_{norm} , that includes **norm** construct, is equivalent to the language \mathcal{L}_{stat} that doesn't have the **norm** construct. Unfortunately, \mathcal{L}_{stat} still includes the language construct **stat** and computing the stationary distribution for arbitrary Markov kernels is also computationally intractable. One advantage of using **stat** construct over **norm** is that there exists approximate compilation for the **stat** construct that, under some suitable assumptions are amenable to error quantification. In this

435

436

437

438

439

440

331 section we begin by giving a broad semantic constraints that 332 needs to be imposed on the Markov kernels programming 333 language under which we can give an approximate compila-334 tion scheme for a single stat construct which is asymptot-335 ically exact. We later impose stricter semantic restrictions that allows us to give quantitative bounds on error bounds 336 337 associated to the approximate compilation of program with multiple stat terms, including nested queries. 338

5.1 Approximate implementation for stat terms

Before giving an approximate compilation for a **stat** term, we first make an assumption on the inputs to the **stat** terms.

Assumption 1. For a term $\Gamma \mid_{p1} \operatorname{stat}(t_0, \lambda x.t_1) : \mathbb{A} + 1$ in the language, we make the following assumptions for all γ , $\llbracket t_1 \rrbracket_{\gamma}$ is an ergodic Markov kernel with stationary distribution $\pi_{\gamma} \in \mathcal{M}\llbracket \mathbb{A} \rrbracket$ and $\llbracket t_0 \rrbracket_{\gamma}$ is such that:

$$\lim_{N \to \infty} \left\| \int_{\llbracket \mathbb{A} \rrbracket} \llbracket t_1 \rrbracket_{\gamma, x}^N(\cdot) \llbracket t_0 \rrbracket_{\gamma, dx} - \pi_{\gamma}(\cdot) \right\|_{\text{tv}} = 0$$

Under Assumption 1, there exists a simple program transformation via iteration that approximates the **stat**-term. First, let's inductively define syntactic sugar **iterate** as follows:

iterate⁰ $(t_0, \lambda x.t_1) := t_0$

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

$$\operatorname{iterate}^{N}(t_0, \lambda x. t_1) := \operatorname{let} x = \operatorname{iterate}^{N-1}(t_0, \lambda x. t_1) \operatorname{in} t_1$$

The semantics of the approximate program transformation is given as:

$$\llbracket \text{iterate}^{N}(t_{0}, \lambda x.t_{1}) \rrbracket_{\gamma} = \int \llbracket t_{1} \rrbracket_{\gamma, x}^{N} \llbracket t_{0} \rrbracket_{\gamma}(dx)$$

Given these syntactic sugar, now we given the program transformation ϕ as follows:

 $\phi(\operatorname{stat}(t_0, \lambda x.t_1), n) \stackrel{\text{def}}{=} \operatorname{iterate}^n(t_0, \lambda x.t_1)$

Even though the Assumption 1 yields an asymptotically exact compilation scheme for the **stat** construct, a program with nested **stat** terms down not behave well in general. This is highlighted in the following problem.

Problem 1. Semantics of the stat construct is not continuous, i.e., for some term $\Gamma \mid_{p1} \operatorname{stat}(t_0, \lambda x. t_1) : \mathbb{A} + 1$ if we know that $\llbracket t_1 \rrbracket_{\gamma, x}$ is an ergodic kernel that has a unique stationary distribution, it is possible to construct an approximate implementation of the Markov transition kernel $\lambda x. t_1'$ such that

$$\exists \delta \in (0,1) \forall \gamma, x. \left\| \llbracket t_1 \rrbracket_{\gamma,x} - \llbracket t_1' \rrbracket_{\gamma,x} \right\| \leq \delta$$

but the Markov kernel $[[t'_1]]_{\gamma,x}$ does not have a stationary distribution. Such an example is given in proposition 1 of Roberts et al. [1998].

To side step the issue stated in Problem 1, we need to make further semantic restrictions on the Markov kernels passed in as input to the **stat** construct. In this paper we identify that if the transition kernel given to the **stat** construct is uniformly ergodic, then **stat** construct is continuous. **Definition 5.1.** A Markov chain with transition kernel P: $\llbracket \mathbb{A} \rrbracket \times \Sigma_{\llbracket \mathbb{A} \rrbracket} \to [0, 1]$ is uniformly ergodic with stationary distribution π if there exists $C \in [0, \infty)$, $\rho \in [0, 1)$ such that for all $x \in \llbracket \mathbb{A} \rrbracket$

$$\left\|P^{N}(x,\cdot) - \pi(\cdot)\right\|_{\mathrm{tv}} \le C\rho^{N}$$

5.2 Quantitative error bounds

We finally state the quantitative error bounds for an approximate probabilistic programs where each **stat** term is uniformly ergodic.

Theorem 5.2 (Quantitative error bound for probabilistic programs). Let *P* be a probabilistic program in the proposed language. Let $\{\operatorname{stat}(t_{0i}, \lambda x.t_{1i})\}_{i \in I}$ be the set of all stationary terms in the program $\emptyset \mid_{p_1} P : \mathbb{B}$ such that $\forall \gamma$, there exists constants $\{C_i\}$ and $\{\rho_i\}$ such that $[[t_{1i}]]_{\gamma}$ is uniformly ergodic with those constants. Let *P'* be a program where $\forall i \in I$, $\operatorname{stat}(t_{0i}, \lambda x.t_{1i})$ is replaced by $\phi(\operatorname{stat}(t_{0i}, \lambda x.t_{1i}), N_i)$ where $N_i \in \mathbb{N}$, then there exists constants $\{C'_i\}_{i \in I}$ such that

$$\left\| \llbracket P \rrbracket_{\gamma} - \llbracket P' \rrbracket_{\gamma} \right\|_{\mathrm{tv}} \leq \sum_{i \in I} C'_i \rho_i^{N_i}.$$

Proof. In appendix.

6 Summary and Discussion

Markov chain Monte Carlo algorithms are workhorses for approximate computation of the normalized posterior distribution. MCMC algorithms are popular because we they give us asymptotic convergence guarantees and are commonly used as "approximate" compilers for probabilistic programming languages. In this abstract we proposed a novel language construct stat that allows us give a formal description for such compilers. We then gave a simple compiler description that again at its core implements an MCMC algorithm. Typically quantifying the rate at which Markov chain with some given transition kernel converges is an open problem and the one we do not attempt to solve in this paper. We make a semantic assumption that the user using our language provides us with a description of the Markov kernel that converges uniformly to the corresponding target distribution. We show that under this uniform convergence property, we can derive rates at which the approximate compiler converges to the original program.

The assumption for uniform ergodicity is crucial for us to derive the quantitative bound. The main difficulty we found in relaxing the uniform ergodicity assumption is the fact that our language allows us to nest the **stat**. We leave as open problem if we can relax the uniform ergodicity assumption.

References

Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Van-houcke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden,

- Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2015.
- TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems.
- https://www.tensorflow.org/ Software available from tensorflow.org.

Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul A. Szerlip, Paul Hors-fall, and Noah D. Goodman. 2019. Pyro: Deep Universal Probabilis-tic Programming. J. Mach. Learn. Res. 20 (2019), 28:1-28:6. http:

//jmlr.org/papers/v20/18-403.html Johannes Borgström, Ugo Dal Lago, Andrew D. Gordon, and Marcin Szym-

czak. 2016. A Lambda-calculus Foundation for Universal Probabilistic

Programming. In Proceedings of the 21st ACM SIGPLAN International Conference on Functional Programming (ICFP 2016). ACM, New York, NY,

USA, 33-46. https://doi.org/10.1145/2951913.2951942 Bob Carpenter, Andrew Gelman, Matthew Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus Brubaker, Jiqiang Guo, Peter

- Li, and Allen Riddell. 2017. Stan: A Probabilistic Programming Language. Journal of Statistical Software, Articles 76, 1 (2017), 1-32. https: //doi.org/10.18637/jss.v076.i01
- Noah Goodman, Vikash Mansinghka, Daniel M Roy, Keith Bonawitz, and Joshua B Tenenbaum. 2012. Church: a language for generative models. arXiv preprint arXiv:1206.3255 (2012).

Chung-Kil Hur, Aditya V Nori, Sriram K Rajamani, and Selva Samuel. 2015.

- A provably correct sampler for probabilistic programs. In 35th IARCS Annual Conference on Foundations of Software Technology and Theoretical Computer Science (FSTTCS 2015). Schloss Dagstuhl-Leibniz-Zentrum fuer
- Informatik. Felipe Medina-Aguayo, Daniel Rudolf, and Nikolaus Schweizer. 2018. Per-
- turbation Bounds for Monte Carlo within Metropolis via Restricted Approximations. arXiv preprint arXiv:1809.09547 (2018).

Tom Rainforth. 2018. Nesting Probabilistic Programs. In UAI.

- Gareth O Roberts, Jeffrey S Rosenthal, and Peter O Schwartz. 1998. Con-vergence properties of perturbed Markov chains. Journal of applied probability 35, 1 (1998), 1-11.
- Adam Ścibior, Ohad Kammar, Matthijs Vákár, Sam Staton, Hongseok Yang, Yufei Cai, Klaus Ostermann, Sean K. Moss, Chris Heunen, and Zoubin Ghahramani. 2017. Denotational Validation of Higher-order Bayesian Inference. Proc. ACM Program. Lang. 2, POPL, Article 60 (Dec. 2017), 29 pages. https://doi.org/10.1145/3158148
- Sam Staton. 2017. Commutative semantics for probabilistic programming. In European Symposium on Programming. Springer, 855-879.
- Sam Staton, Frank Wood, Hongseok Yang, Chris Heunen, and Ohad Kam-mar. 2016. Semantics for probabilistic programming: higher-order func-tions, continuous distributions, and soft constraints. In 2016 31st Annual ACM/IEEE Symposium on Logic in Computer Science (LICS). IEEE, 1-10.
- David Tolpin, Jan Willem van de Meent, Hongseok Yang, and Frank Wood. 2016. Design and implementation of probabilistic programming language anglican. arXiv preprint arXiv:1608.05263 (2016).

A Proof of Theorem 5.2

Before we prove the Theorem 5.2, we prove the following proposition that characterizes the uniform continuity of let and case statements under uniform ergodicity assumption.

Proposition A.1. The following statements hold:

491
491
1. Let
$$\left[\Gamma \mid_{p_{1}} t_{0} : \mathbb{A}\right], \left[\Gamma \mid_{p_{1}} t_{0}' : \mathbb{A}\right], \left[\Gamma, x : \mathbb{A} \mid_{p_{1}} t_{1} : \mathbb{B}\right],$$

492
493
494
494
If for all $\gamma \in \Gamma, \|[t_{0}]]_{\gamma} - [[t_{0}]]_{\gamma}\|_{t_{v}} \leq \alpha$ and for all

PTML'19, December 14, 2019, Vancouver, BC, Canada

$$\begin{split} \llbracket x \rrbracket_{\gamma} \in \llbracket \mathbb{A} \rrbracket, \\ & \left\| \llbracket t_1' \rrbracket_{\gamma} (\llbracket x \rrbracket_{\gamma}) - \llbracket t_1 \rrbracket_{\gamma} (\llbracket x \rrbracket_{\gamma}) \right\|_{\text{tv}} \leq \beta \end{split}$$

then

 $\left\| \left[\left[\text{let } x = t_0 \text{ in } t_1 \right] \right]_Y - \left[\left[\text{let } x = t'_0 \text{ in } t'_1 \right] \right]_Y \right\|_{\text{ty}} \le \alpha + \beta$

2. Let
$$\left[\Gamma, x : \mathbb{A}_i \mid_{p_1} t_i : \mathbb{B}\right]$$
 and $\left[\Gamma, x : \mathbb{A}_i \mid_{p_1} t'_i : \mathbb{B}\right]$, such that

$$\forall i \in I, \forall x \in \llbracket \mathbb{A}_i \rrbracket, \left\| \llbracket t'_i \rrbracket_{\gamma, x} - \llbracket t_i \rrbracket_{\gamma, x} \right\|_{\mathrm{tv}} \leq \alpha_i$$

$$\begin{aligned} \left\| \begin{bmatrix} \text{case } a \text{ in } \left\{ (i, x) \Rightarrow t'_i \right\}_{i \in I} \end{bmatrix}_{\gamma} - \\ \left\| \begin{bmatrix} \text{case } a \text{ in } \left\{ (i, x) \Rightarrow t_i \right\}_{i \in I} \end{bmatrix}_{\gamma} \right\|_{\text{tv}} &\leq \sup \left\{ \alpha_i \right\}_{i \in I} \end{aligned}$$

Proof. 1. For the **let** construct:

$$\begin{split} \left\| \int \left\| t_{1} \right\|_{Y,x} \left\| t_{0} \right\|_{Y,dx} - \int \left\| t_{1}' \right\|_{Y,x} \left\| t_{0}' \right\|_{Y,dx} \right\|_{tv} \\ \leq \left\| \int \left\| t_{1} \right\|_{Y,x} \left\| t_{0} \right\|_{Y,dx} - \int \left\| t_{1}' \right\|_{Y,x} \left\| t_{0} \right\|_{Y,dx} \right\|_{tv} \\ + \left\| \int \left\| t_{1}' \right\|_{Y,x} \left\| t_{0} \right\|_{Y,dx} - \int \left\| t_{1}' \right\|_{Y,x} \left\| t_{0}' \right\|_{Y,dx} \right\|_{tv} \\ = \left\| \int \left(\left\| t_{1} \right\|_{Y,x} - \left\| t_{1}' \right\|_{Y,x} \right) \left\| t_{0} \right\|_{Y,dx} - \int \left\| t_{1}' \right\|_{Y} (x,A) \left\| t_{0}' \right\|_{Y,dx} \right\| \\ + \sup_{A} \left| \int \left\| t_{1}' \right\|_{Y} (x,A) \left\| t_{0} \right\|_{Y,dx} - \int \left\| t_{1}' \right\|_{Y} (x,A) \left\| t_{0}' \right\|_{Y,dx} \right| \\ \leq \int \sup_{A} \left\| \left\| t_{1} \right\|_{Y} (x') - \left\| t_{1}' \right\|_{Y} (x') \right\|_{tv} \left\| t_{0} \right\|_{Y,dx} \end{split}$$

$$+ \sup_{f \le 1} \left| \int f(x) [\![t_0]\!]_{\gamma, dx} - \int f(x) [\![t_0']\!]_{\gamma, dx} \right|$$

$$= \int \beta \llbracket t_0 \rrbracket_{\gamma, dx} + \alpha$$
$$= \alpha + \beta$$

2. For the case construct:

$$\begin{aligned} \left\| \begin{bmatrix} \operatorname{case} a \text{ in } \{(i, x) \Rightarrow t'_i\}_{i \in I} \end{bmatrix}_{Y} - \\ \left\| \operatorname{case} a \text{ in } \{(i, x) \Rightarrow t_i\}_{i \in I} \end{bmatrix}_{Y} \right\|_{\mathrm{tv}} \\ &= \left\| \begin{bmatrix} t'_i \end{bmatrix}_{v, Y} - \begin{bmatrix} t_i \end{bmatrix}_{v, Y} \right\|_{\mathrm{tv}} & \text{if } (i, v) = \llbracket a \rrbracket_{Y} \\ &\leq \alpha_i & \text{if } (i, v) = \llbracket a \rrbracket_{Y} \\ &\leq \sup_i \{\alpha_i\}_{i \in I} \end{aligned}$$

П

Remark A.2. A natural consequence of the sugar is that the following terms are syntactically equivalent.

iterate^N
$$(t_0, \lambda x.t_1) := iterate^{N-i}(iterate^i(t_0, \lambda x.t_1), \lambda x.t_1)$$

We characterize the error of an approximate iterate transformation.

Theorem A.3. Let $\Gamma, x \mid_{p_1} t_1$ and $\Gamma, x \mid_{p_1} t'_1$ be probabilistic terms. If $[t_1]_{\gamma}$ is uniformly ergodic with stationary distribution π and constants *C*, and ρ and

$$\left\| \llbracket t_1' \rrbracket_{\gamma} - \llbracket t_1 \rrbracket_{\gamma} \right\|_{\mathsf{tv}} \le \varepsilon,$$

then

$$\llbracket \text{iterate}^{N}(t_{0}, \lambda x.t_{1}') \rrbracket - \llbracket \text{stat}(t_{0}, \lambda x.t_{1}) \rrbracket \Big\|_{\text{tv}} \leq \frac{\varepsilon C}{1-\rho} + C\rho^{N}$$

We begin the proof of theorem Theorem A.3 by giving a simple contraction lemma that quantifies the error associated with the N step iteration transformation with a different initial distribution.

Lemma A.4 (Contraction). If $R(\llbracket t_1 \rrbracket_{\gamma}, ST(\llbracket t_1 \rrbracket_{\gamma}), C, \rho)$ then for all

$$\left\| \begin{bmatrix} [\text{iterate}^{N}(m_{1},\lambda x.t_{1})]]_{Y} - \\ [\text{iterate}^{N}(m_{2},\lambda x.t_{1})]]_{Y} \end{bmatrix}_{\text{tv}} \leq C\rho^{N} \left\| [m_{1}]]_{Y} - [m_{2}]]_{Y} \right\|.$$

Proof. The proof of this lemma follows directly from linearity of integration.

To prove Theorem A.3, we now give following theorem quantifies the distance between iteration transformation the transition kernels are close.

Lemma A.5. Let $\Gamma, x \mid_{p_1} t_1$ and $\Gamma, x \mid_{p_1} t'_1$ be probabilistic *terms.* If there exists $\pi \in \mathcal{M}(\mathbb{A})$, $C \in \mathbb{R}_+$, and $\rho \in [0, 1)$ such that $R(\llbracket t_1 \rrbracket_{\gamma}, \pi, C, \rho)$ and

$$\left\| \llbracket t_1' \rrbracket_{\gamma} - \llbracket t_1 \rrbracket_{\gamma} \right\|_{\mathrm{tv}} \le \varepsilon$$

then,

$$\left\| \llbracket \text{iterate}^{N}(t_{0}, \lambda x.t_{1}') \rrbracket - \llbracket \text{iterate}^{N}(t_{0}, \lambda x.t_{1}) \rrbracket \right\|_{\text{tv}} \leq \frac{\varepsilon C}{1 - \rho}$$

Proof. We show this by first noting that

$$\left\| \left[\left[\text{iterate}^{N}(t_{0}, \lambda x.t_{1}) \right] \right]_{\gamma} - \left[\left[\text{iterate}^{N}(t_{0}, \lambda x.t_{1}') \right] \right]_{\gamma} \right\|_{\text{tv}}$$

$$= \left\| \sum_{i=0}^{N-1} \llbracket \text{iterate}^{N-i}(\text{iterate}^{i}(t_{0},\lambda x.t_{1}'),\lambda x.t_{1}) \rrbracket_{Y} \\ - \llbracket \text{iterate}^{N-i-1}(\text{iterate}^{i+1}(t_{0},\lambda x.t_{1}'),\lambda x.t_{1}) \rrbracket_{Y} \right\|_{\text{tv}} \right\|_{\text{tv}}$$

$$= \left\| \sum_{i=0}^{N-1} \llbracket \text{iterate}^{N-i-1}(\text{let } x = \text{iterate}^{i}(t_{0},\lambda x.t_{1}') \text{ in } n,\lambda x.t_{1})) \rrbracket_{Y} \right\|_{\text{tv}}$$

$$= \left\| \sum_{i=0}^{N-1} \llbracket \text{iterate}^{N-i-1}(\text{let } x = \text{iterate}^{i+1}(t_{0},\lambda x.t_{1}') \text{ in } n,\lambda x.t_{1})) \rrbracket_{Y} \right\|_{\text{tv}}$$

$$= \left\| \sum_{i=0}^{N-1} C\rho^{N-i-1}(\lVert \llbracket \text{let } x = \text{iterate}^{i}(t_{0},\lambda x.t_{1}') \text{ in } n \rrbracket_{Y} \right\|_{\text{tv}}$$

$$= \left\| \sum_{i=0}^{N-1} C\rho^{i}\varepsilon \leq \frac{\varepsilon C}{1-\rho}.$$

Proof of Theorem A.3. Given the Lemma A.5 and the assumption in the theorem statement that $R(\llbracket t_1 \rrbracket_V, \pi, C, \rho)$ holds, the proof of theorem follows by triangle inequality.

Now we are in a position to prove the main theorem of this section.

Proof of Theorem 5.2. The proof is seen by induction on probabilistic terms.

- Base case:
 - Leaf node for the induction is the terms of the form $stat(t'_0, \lambda x.t'_1)$. By the assumption of uniform ergodicity there exists C', ρ' such that following holds

 $[[stat(t_0, \lambda x.t_1)]] - [[iterate^{N'}(t_0, \lambda x.t_1)]] \le C' \rho'^{N'}.$

Thus the base case holds.

• Inductive Step:

For the inductive step we show the hypothesis holds for for all constructor of probabilistic terms in our language.

- case terms: By the inductive hypothesis,

$$\| [t'_j] \|_{\gamma} - [t'_j] \|_{\gamma} \|_{\mathrm{tv}} \le \sum_{i \in I_i} C_i \rho^{N_i}.$$

Now, we show

$$\left\| \begin{bmatrix} \text{case } a \text{ in } \{(j, x) \Rightarrow t'_j\}_{j \in J} \end{bmatrix} \right\|_{\text{ty}} \le \sum_{i \in \bigcup_{i \in J} I_i} C'_i \rho_i^{N_i}$$

From Proposition A.1 and the inductive hypothesis, we know

$$\left\| \begin{bmatrix} \text{case } a \text{ in } \{(j, x) \Rightarrow t'_j\}_{j \in J} \end{bmatrix} - \\ - \begin{bmatrix} \text{case } a \text{ in } \{(j, x) \Rightarrow t_j\}_{j \in J} \end{bmatrix} \right\|_{\text{tv}} \leq \sup_{j \in J} \sum_{i \in I_j} C_i \rho^{N_i} \\ \leq \sum_{i \in J} \rho^{N_i}.$$

- let term: We need to show

$$\left\| \left[\left[\text{let } x = t_1 \text{ in } t_2 \right] \right] - \left[\left[\text{let } x = t'_1 \text{ in } t'_2 \right] \right] \right\|_{\text{tv}} \le \sum_{i \in I_1 \cup I_2} C_i \rho_i^{N_i}$$

This follows from inductive hypothesis and Proposition A.1.

 $i \in \bigcup_i I_i$

- stat term: From Theorem A.3 and inductive hypothesis it follows that

$$\left\| \begin{bmatrix} \operatorname{stat}(t_0, \lambda x. t_1) \end{bmatrix} - \\ \begin{bmatrix} \operatorname{iterate}^{N'}(t_0, \lambda x. t_1') \end{bmatrix} \right\|_{\operatorname{tv}} \leq C' \rho'^{N_i} + \sum_i \frac{C'C_i}{1 - \rho'} \rho_i^{N_i}.$$