



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Semantics for probabilistic programming: higher-order functions, continuous distributions, and soft constraints

Citation for published version:

Staton, S, Yang, H, Heunen, C, Kammar, O & Wood, F 2016, Semantics for probabilistic programming: higher-order functions, continuous distributions, and soft constraints. in LICS '16 Proceedings of the 31st Annual ACM/IEEE Symposium on Logic in Computer Science. ACM, New York, USA, pp. 525-534, 31st Annual ACM/IEEE Symposium on Logic in Computer Science, New York City, United States, 5/07/16. DOI: 10.1145/2933575.2935313

Digital Object Identifier (DOI):

[10.1145/2933575.2935313](https://doi.org/10.1145/2933575.2935313)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

LICS '16 Proceedings of the 31st Annual ACM/IEEE Symposium on Logic in Computer Science

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



Semantics for probabilistic programming: higher-order functions, continuous distributions, and soft constraints

Sam Staton Hongseok Yang
Frank Wood
University of Oxford

Chris Heunen
University of Edinburgh

Ohad Kammar
University of Cambridge

Abstract

We study the semantic foundation of expressive probabilistic programming languages, that support higher-order functions, continuous distributions, and soft constraints (such as Anglican, Church, and Venture). We define a metalanguage (an idealised version of Anglican) for probabilistic computation with the above features, develop both operational and denotational semantics, and prove soundness, adequacy, and termination. They involve measure theory, stochastic labelled transition systems, and functor categories, but admit intuitive computational readings, one of which views sampled random variables as dynamically allocated read-only variables. We apply our semantics to validate nontrivial equations underlying the correctness of certain compiler optimisations and inference algorithms such as sequential Monte Carlo simulation. The language enables defining probability distributions on higher-order functions, and we study their properties.

1. Introduction

Probabilistic programming is the idea to use programs to specify probabilistic models; probabilistic programming languages blend programming constructs with probabilistic primitives. This helps scientists express complicated models succinctly. Moreover, such languages come with generic inference algorithms, relieving the programmers of the nontrivial task of (algorithmically) answering queries about their probabilistic models. This is useful in *e.g.* machine learning.

Several higher-order probabilistic programming languages have recently attracted a substantial user base. Some languages (such as Infer.net [21], PyMC [26], and Stan [33]) are less expressive but provide powerful inference algorithms, while others (such as Anglican [34], Church [12], and Venture [20]) have less efficient inference algorithms but more expressive power. We consider the more expressive languages, that support higher-order functions, continuous distributions, and soft constraints. More precisely, we consider a programming language (§3) with higher-order functions (§6) as well as the following probabilistic primitives.

Sampling The command `sample(t)` draws a sample from a distribution described by *t*, which may range over the real numbers.

Soft constraints The command `score(t)` puts a score *t* (a positive real number) on the current execution trace. This is typically used to record that some particular datum was observed as being drawn from a particular distribution; the score describes how surprising the observation is.

Normalization The command `norm(u)` runs a simulation algorithm over the program fragment *u*. This takes the scores into account and returns a new, normalized probability distribution. The argument to sample might be a primitive distribution, or

a distribution defined by normalizing another program. This is called a *nested query*, by analogy with database programming.

Here is a simple example of a program. We write $gauss(\mu, \sigma)$ for the Gaussian probability distribution whose density function is $density_gauss(a, (\mu, \sigma)) = \frac{1}{\sigma\sqrt{2\pi}} \exp(-\frac{(a-\mu)^2}{2\sigma^2})$.

```

1  norm(
2    let x = sample(gauss(0.0, 3.0)) in
3    score(density_gauss(5.0, (x, 1.0)));
4    return(x < 4.5))

```

(1)

Line 2 samples *x* from a prior Gaussian distribution. The soft constraint on Line 3 expresses the likelihood of the observed data, 5.0, coming from a Gaussian given the prior *x*. Line 4 says that what we are actually interested in is a boolean random variable over the sample space. Line 1 calculates a posterior distribution for the return value, using the prior and the likelihood. In this example we can precisely calculate that the posterior distribution on {true, false} has $p(\text{true}) = 0.5$.

Languages like this currently lack formal exact semantics. The aim of this paper is to provide just such a foundation as a basis for formal reasoning. Most expressive probabilistic programming languages are explained in terms of their Monte Carlo simulation algorithms. The simplest such algorithm, using importance and rejection sampling, is the *de facto* semantics against which other algorithms are ‘proved approximately correct’. Such ‘semantics’ are hard to handle and extend.

We provide two styles of semantics, operational and denotational. For first-order probabilistic programs, the denotational semantics is straightforward: types are interpreted as measurable spaces, and terms are interpreted as measurable functions (§4). Operational semantics is more complicated. For discrete distributions, an operational semantics might be a probabilistic transition system, but for continuous distributions, it must be a stochastic relation (labelled Markov process). We resolve this by equipping the set of configurations with the structure of a measurable space (§5).

The advantage to the operational semantics is that it is easily extended to higher-order programs (§7). Denotational semantics for higher-order programs poses a problem, because measurable spaces do not support the usual β/η theory of functions: they do not form a Cartesian closed category (indeed, $\mathbb{R}^{\mathbb{R}}$ does not exist as a measurable space [3]). Earlier semantics deal with this by either excluding higher-order functions or considering only discrete distributions. We resolve this by moving from the category of measurable spaces, where standard probability theory takes place, to a functor category based on it (§8). The former embeds in the latter, so we can still interpret first-order concepts. But the functor category does have well-behaved function spaces, so we can also interpret higher-order concepts. Moreover, by lifting the monad of probability distribu-

tions [11] to the functor category, we can also interpret continuous distributions. Finally, we can interpret observations by considering probability distributions with continuous density, irrespective of the categorical machinery (§9).

The denotational semantics is sound and adequate with respect to the operational semantics (§5.3,8.3), which means one can use the denotational model to directly check program equations whilst respecting computational issues. For example:

- we demonstrate a key program equation for sequential Monte Carlo simulation (§4.1);
- we show that every term of first-order type is equal to one without λ -abstractions or application, and hence is interpreted as a measurable function (Proposition 8.3).

2. Preliminaries

We recall basic definitions and facts of measure theory.

Definition 2.1. A σ -algebra on a set X is a family Σ of subsets of X , called *measurable (sub)sets*, which contains X and is closed under complements and countable unions. A *measurable space* is a set with a σ -algebra.

A *measure* on a measurable space (X, Σ) is a function $p: \Sigma \rightarrow [0, \infty]$ such that $p(\emptyset) = 0$ and $p(\bigcup_{i \in \mathbb{N}} U_i) = \sum_{i \in \mathbb{N}} p(U_i)$ for each sequence U_1, U_2, \dots of disjoint measurable sets. A *probability measure* or *probability distribution* is a measure p with $p(X) = 1$.

In the paper we use a few important constructions for measurable spaces. The first is to make a set X into a measurable space by taking the full powerset of X as Σ , yielding a *discrete* measurable space. When X is countable, a probability distribution on (X, Σ) is determined by its values on singleton sets, that is, by specifying a function $p: X \rightarrow [0, 1]$ such that $\sum_{x \in X} p(x) = 1$.

The second construction is to combine a collection of measurable spaces $(X_i, \Sigma_i)_{i \in I}$ by *sum* or *product*. The underlying sets are the disjoint union $\sum_{i \in I} X_i$ and product $\prod_{i \in I} X_i$ of sets. The measurable sets in the sum are $\sum_{i \in I} U_i$ for $U_i \in \Sigma_i$. The σ -algebra of the product is the smallest one containing all the subsets $\prod_{i \in I} U_i$ where $U_i \in \Sigma_i$ equals X_i but for a single index i .

For the third, the real numbers form a measurable space $(\mathbb{R}, \Sigma_{\mathbb{R}})$ under the smallest σ -algebra that contains the open intervals; the measurable sets are called *Borel* sets. Restricting to any measurable subset gives a new measurable space, such as the space $\mathbb{R}_{\geq 0}$ of nonnegative reals and the unit interval $[0, 1]$.

The fourth construction is to make the set $P(X)$ of all probability measures on a measurable space (X, Σ_X) into a measurable space, by letting $\Sigma_{P(X)}$ be the smallest σ -algebra containing the sets $\{p \in P(X) \mid p(U) \in V\}$ for all $U \in \Sigma_X$ and $V \in \Sigma_{[0,1]}$.

Definition 2.2. Let $(X, \Sigma_X), (Y, \Sigma_Y)$ be measurable spaces. A function $f: X \rightarrow Y$ is *measurable* if $f^{-1}(U) \in \Sigma_X$ for $U \in \Sigma_Y$.

We can *push forward* a measure along a measurable function: if $p: \Sigma_X \rightarrow [0, 1]$ is a probability measure on (X, Σ_X) and $f: X \rightarrow Y$ is a measurable function, then $q(U) = p(f^{-1}(U))$ is a probability measure on (Y, Σ_Y) .

Definition 2.3. A *stochastic relation* between measurable spaces (X, Σ_X) and (Y, Σ_Y) is a function $r: X \times \Sigma_Y \rightarrow [0, 1]$ such that $r(x, -): \Sigma_Y \rightarrow [0, 1]$ is a probability distribution for all $x \in X$, and $r(-, V): X \rightarrow [0, 1]$ is measurable for all $V \in \Sigma_Y$.

Giving a stochastic relation from (X, Σ_X) to (Y, Σ_Y) is equivalent to giving a measurable function $(X, \Sigma_X) \rightarrow (P(Y), \Sigma_{P(Y)})$. Stochastic relations $r: X \times \Sigma_Y \rightarrow [0, 1]$ and $s: Y \times \Sigma_Z \rightarrow [0, 1]$ compose associatively to $(s \circ r): X \times \Sigma_Z \rightarrow [0, 1]$ via the formula

$$(s \circ r)(x, W) = \int_Y s(y, W) r(x, dy).$$

Finally, for a predicate φ , we use the indicator expression $[\varphi]$ to denote 1 if φ holds, and 0 otherwise.

3. A first-order language

This section presents a first-order language for expressing Bayesian probabilistic models. The language forms a first-order core of a higher-order extension in Section 6, and provides a simpler setting to illustrate key ideas. The language includes infinitary type and term constructors, constant terms for all measurable functions between measurable spaces, and constructs for specifying Bayesian probabilistic models, namely, operations for sampling distributions, scoring samples, and normalizing distributions based on scores. This highly permissive and slightly unusual syntax is not meant to be a useful programming language itself. Rather, its purpose is to serve as a semantic metalanguage to which a practical programming language compiles, and to provide a mathematical setting for studying high-level constructs for probabilistic computation.

Types The language has types

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

where I ranges over countable sets. A type \mathbb{A} stands for a measurable space $[\![\mathbb{A}]\!]$. For example, \mathbb{R} denotes the measurable space of reals, $\mathbb{P}(\mathbb{A})$ is the space of probability measures on \mathbb{A} , and 1 is the (discrete) measurable space on the singleton set. The other type constructors correspond to products and sums of measurable spaces. Notice that countable sums are allowed, enabling us to express usual ground types in programming languages via standard encoding. For instance, the type for booleans is $1 + 1$, and that for natural numbers $\sum_{i \in \mathbb{N}} 1$.

Terms We distinguish typing judgements: $\Gamma \vdash_d t: \mathbb{A}$ for deterministic terms, and $\Gamma \vdash_p t: \mathbb{A}$ for probabilistic terms (see also e.g. [19, 25, 29]). In both, \mathbb{A} is a type, and Γ is a list of variable/type pairs. Variables stand for deterministic terms.

Intuitively, a probabilistic term $\Gamma \vdash_p t: \mathbb{A}$ may have two kinds of effects: during evaluation, t may sample from a probability distribution, and it may score the current execution trace (typically according to likelihood of data). Evaluating a deterministic term $\Gamma \vdash_d t: \mathbb{A}$ does not have any effects.

Sums and products The language includes variables, and standard constructors and destructors for sum and product types.

$$\frac{}{\Gamma, x: \mathbb{A}, \Gamma' \vdash_d x: \mathbb{A}} \quad \frac{\Gamma \vdash_d t: \mathbb{A}_i}{\Gamma \vdash_d (i, t): \sum_{i \in I} \mathbb{A}_i}$$

$$\frac{\Gamma \vdash_d t: \sum_{i \in I} \mathbb{A}_i \quad (\Gamma, x: \mathbb{A}_i \vdash_d u_i: \mathbb{B})_{i \in I}}{\Gamma \vdash_d \text{case } t \text{ of } \{(i, x) \Rightarrow u_i\}_{i \in I}: \mathbb{B}} \quad (z \in \{d, p\})$$

$$\frac{}{\Gamma \vdash_d *: 1} \quad \frac{\Gamma \vdash_d t_j: \mathbb{A}_j \text{ for all } j \in \{0, 1\}}{\Gamma \vdash_d (t_0, t_1): \mathbb{A}_0 \times \mathbb{A}_1} \quad \frac{\Gamma \vdash_d t: \mathbb{A}_0 \times \mathbb{A}_1}{\Gamma \vdash_d \pi_j(t): \mathbb{A}_j}$$

In the rules for sums, I may be infinite. In the last rule, j is 0 or 1. We use some standard syntactic sugar, such as *false* and *true* for the injections in the type $\text{bool} = 1 + 1$, and *if* for case in that instance.

Sequencing We include the standard constructs (e.g. [19, 22]).

$$\frac{\Gamma \vdash_d t: \mathbb{A}}{\Gamma \vdash_p \text{return}(t): \mathbb{A}} \quad \frac{\Gamma \vdash_p t: \mathbb{A} \quad \Gamma, x: \mathbb{A} \vdash_p u: \mathbb{B}}{\Gamma \vdash_p \text{let } x = t \text{ in } u: \mathbb{B}}$$

Where $\mathbb{A} = 1$, we write $(t; u)$ for $\text{let } x = t \text{ in } u$.

Language-specific constructs The language has constant terms for all measurable functions.

$$\frac{\Gamma \vdash_d t: \mathbb{A}}{\Gamma \vdash_d f(t): \mathbb{B}} \quad (f: [\![\mathbb{A}]\!] \rightarrow [\![\mathbb{B}]\!] \text{ measurable}) \quad (2)$$

In particular, all the usual distributions are in the language, including the Dirac distribution $\text{dirac}(x)$ concentrated on outcome x , the Gaussian distribution $\text{gauss}(\mu, \sigma)$ with mean μ and standard deviation σ , the Bernoulli distribution $\text{bern}(p)$ with success probability p , the exponential distribution $\text{exp}(r)$ with rate r , and the Beta distribution $\text{beta}(\alpha, \beta)$ with parameters α, β .¹ For example, from the measurable functions $42.0: 1 \rightarrow \mathbb{R}$, $e^{(-)}: \mathbb{R} \rightarrow \mathbb{R}$, $\text{gauss}: \mathbb{R} \times \mathbb{R} \rightarrow P(\mathbb{R})$ and $<: \mathbb{R} \times \mathbb{R} \rightarrow 1 + 1$ we can derive:

$$\frac{\Gamma \vdash 42.0: \mathbb{R}}{\Gamma \vdash 42.0: \mathbb{R}} \quad \frac{\Gamma \vdash t: \mathbb{R}}{\Gamma \vdash e^t: \mathbb{R}} \\ \frac{\Gamma \vdash \mu: \mathbb{R} \quad \Gamma \vdash \sigma: \mathbb{R}}{\Gamma \vdash \text{gauss}(\mu, \sigma): P(\mathbb{R})} \quad \frac{\Gamma \vdash t: \mathbb{R} \quad \Gamma \vdash u: \mathbb{R}}{\Gamma \vdash t < u: \text{bool}}$$

The following terms form the core of our language.

$$\frac{\Gamma \vdash t: P(\mathbb{A})}{\Gamma \vdash \text{sample}(t): \mathbb{A}} \quad \frac{\Gamma \vdash t: \mathbb{R}}{\Gamma \vdash \text{score}(t): 1}$$

The first term samples a value from a distribution t . The second multiplies the score of the current trace with t . Since both of these terms express effects, they are typed under \vdash instead of \vdash .

The reader may think of the score of the current execution trace as a state, but it cannot be read, nor changed arbitrarily: it can only be multiplied by another score. The argument t in $\text{score}(t)$ is usually the density of a probability distribution at an observed data point. For instance, recall the example (1) from the Introduction:

```
norm(let x = sample(gauss(0.0, 3.0)) in
  score(density_gauss(5.0, (x, 1.0)); return(x < 4.5)))
```

An execution trace is scored higher when x is closer to the datum 5.

When the argument t in $\text{score}(t)$ is 0, this is called a *hard constraint*, meaning ‘reject this trace’; otherwise it is called a soft constraint. In the discrete setting, hard constraints are more-or-less sufficient, but even then, soft constraints tend to be more efficient.

Normalization Two representative tasks of Bayesian inference are to calculate the so-called *posterior distribution* and *model evidence* from a prior distribution and likelihood. Programs built from sample and score can be thought of as setting up a prior and a likelihood. Consider the following program:

```
let x = sample(bern(0.25)) in
let y = (if x then return 5.0 else return 2.0) in
score(density_exp(0.0, y)); return(x)
```

Here the prior y comes from the Bernoulli distribution, and the likelihood concerns datum 0.0 coming from an exponential distribution with rate y . Recall that $\text{density_exp}(0.0, y) = y$. So there are two execution traces, returning either true, with probability 0.25 and score 5.0, or false, with probability 0.75 and score 2.0.

The product of the prior and likelihood gives an unnormalized posterior distribution on the return value: ($\text{true} \mapsto 0.25 \cdot 5 = 1.25$, $\text{false} \mapsto 0.75 \cdot 2 = 1.5$). The normalizing constant is the average score: $(0.25 \cdot 5 + 0.75 \cdot 2) = 2.75$, so the posterior is ($\text{true} \mapsto \frac{1.25}{2.75} \approx 0.45$, $\text{false} \mapsto \frac{1.5}{2.75} \approx 0.55$). The average score is called the model evidence. It is a measure of how well the model encoded by the program matches the observation.

Note that the sample $x = \text{true}$ matches the datum better, so the probability of true goes up from 0.25 to 0.45 in the posterior.

In our language we have a term $\text{norm}(t)$ that will usually convert a probabilistic term t into a deterministic value, which is its posterior distribution together with the model evidence.

$$\frac{\Gamma \vdash t: \mathbb{A}}{\Gamma \vdash \text{norm}(t): (\mathbb{R} \times P(\mathbb{A})) + 1 + 1}$$

¹ Usually, $\text{gauss}(0.0, 0.0)$ is undefined, but we just let $\text{gauss}(0.0, 0.0) = \text{gauss}(0.0, 1.0)$, and so on, to avoid worrying about this sort of error.

If the model evidence is 0 or ∞ , the conversion fails, and this is tracked by the ‘+1+1’.

4. Denotational semantics

This section discusses the natural denotational semantics of the first-order language. The basic idea can be traced back a long way (e.g. [17]) but our treatment of score and norm appear to be novel. As described, types \mathbb{A} are interpreted as measurable spaces $[\![\mathbb{A}]\!]$. A context $\Gamma = (x_1: \mathbb{A}_1, \dots, x_n: \mathbb{A}_n)$ is interpreted as the measurable space $[\![\Gamma]\!] \stackrel{\text{def}}{=} \prod_{i=1}^n [\![\mathbb{A}_i]\!]$ of its valuations.

- Deterministic terms $\Gamma \vdash t: \mathbb{A}$ are interpreted as measurable functions $[\![t]\!]: [\![\Gamma]\!] \rightarrow [\![\mathbb{A}]\!]$, providing a result for each valuation of the context.
- Probabilistic terms $\Gamma \vdash t: \mathbb{A}$ are interpreted as measurable functions $[\![t]\!]: [\![\Gamma]\!] \rightarrow P(\mathbb{R}_{\geq 0} \times [\![\mathbb{A}]\!])$, providing a probability measure on (score, result) pairs for each valuation of the context.

Informally, if $(\Omega, p: \Omega \rightarrow [0, 1])$ is the prior sample space of the program (specified by sample), and $l: \Omega \rightarrow \mathbb{R}_{\geq 0}$ is the likelihood (specified by score), and $r: \Omega \rightarrow [\![\mathbb{A}]\!]$ is the return value, then a measure in $P(\mathbb{R}_{\geq 0} \times [\![\mathbb{A}]\!])$ is found by pushing forward p along (l, r) .

Sums and products The interpretation of deterministic terms follows the usual pattern of the internal language of a distributive category (e.g. [27]). For instance, $[\![\Gamma, x: \mathbb{A}, \Gamma' \vdash x: \mathbb{A}]\!](\gamma, a, \gamma') \stackrel{\text{def}}{=} a$, and $[\![\Gamma \vdash f(t): \mathbb{A}]\!](\gamma) \stackrel{\text{def}}{=} f([\![t]\!](\gamma))$ for measurable $f: [\![\mathbb{A}]\!] \rightarrow [\![\mathbb{B}]\!]$. This interpretation is actually the same as the usual set-theoretic semantics of the calculus, as one can show by induction that the induced functions $[\![\Gamma]\!] \rightarrow [\![\mathbb{A}]\!]$ are measurable.

Sequencing For probabilistic terms, we proceed as follows.

$$[\![\text{return}(t)]\!](\gamma)(U) \stackrel{\text{def}}{=} [(1, [\![t]\!](\gamma)) \in U],$$

which denotes a Dirac distribution, and $[\![\text{let } x = t \text{ in } u]\!](\gamma)(V)$ is

$$\int_{\mathbb{R}_{\geq 0} \times [\![\mathbb{A}]\!]} \left([\![u]\!](\gamma, x) (\{(s, b) \mid (r \cdot s, b) \in V\}) \right) \left([\![t]\!](\gamma)(d(r, x)) \right).$$

As we will explain shortly, these interpretations come from treating $P(\mathbb{R}_{\geq 0} \times (-))$ as a commutative monad, which essentially means the following program equations hold.

$$[\![\text{let } x = \text{return}(x) \text{ in } u]\!] = [\![u]\!] \quad [\![\text{let } x = t \text{ in return}(x)]\!] = [\![t]\!] \\ [\![\text{let } y = (\text{let } x = t \text{ in } u) \text{ in } v]\!] = [\![\text{let } x = t \text{ in let } y = u \text{ in } v]\!] \\ [\![\text{let } x = t \text{ in let } y = u \text{ in } (x, y)]\!] = [\![\text{let } y = u \text{ in let } x = t \text{ in } (x, y)]\!]$$

The last two equations assume the standard conditions on free variables of u, v and t . The last equation justifies a useful program optimisation technique [34, §5.5].

Language-specific constructs We use the monad:

$$[\![\text{sample}(t)]\!](\gamma)(U) \stackrel{\text{def}}{=} [\![t]\!](\gamma)(\{a \mid (1, a) \in U\}) \\ [\![\text{score}(t)]\!](\gamma)(U) \stackrel{\text{def}}{=} ((\max([\![t]\!](\gamma), 0), *) \in U)$$

Here are some program equations to illustrate the semantics so far.

$$[\![\text{score}(7.0); \text{score}(6.1)]\!] = [\![\text{score}(42.7)]\!] \\ [\![\text{let } x = \text{sample}(\text{gauss}(0.0, 1.0)) \text{ in return}(x > 0.0)]\!] = [\![\text{sample}(\text{bern}(0.5))]\!] \\ [\![\text{let } x = \text{sample}(\text{gauss}(0.0, 1.0)) \text{ in return}(x > x)]\!] = [\![\text{return}(\text{false})]\!]$$

Normalization We interpret $\text{norm}(t)$ as follows. Each probability measure p on $(\mathbb{R}_{\geq 0} \times X)$ induces an unnormalized posterior measure \bar{p} on X : let $\bar{p}(U) = \int_{\mathbb{R}_{\geq 0} \times U} p(d(r, x))$. As long as the average score $\bar{p}(X)$ is not 0 or ∞ , we can normalize \bar{p} to build

a posterior probability measure on X . This construction forms a natural transformation $\iota_X : P(\mathbb{R}_{\geq 0} \times X) \rightarrow (\mathbb{R} \times P(X)) + 1 + 1$

$$\iota_X(p) \stackrel{\text{def}}{=} \begin{cases} (1, *) & \text{if } \bar{p}(X) = 0 \\ (2, *) & \text{if } \bar{p}(X) = \infty \\ (0, (\bar{p}(X), \lambda U. \frac{\bar{p}(U)}{\bar{p}(X)})) & \text{otherwise} \end{cases} \quad (3)$$

Let $\llbracket \text{norm}(t) \rrbracket(\gamma) \stackrel{\text{def}}{=} \iota(\llbracket t \rrbracket(\gamma))$. Here are some examples:

$$\begin{aligned} & \llbracket \text{norm}(\text{let } x = \text{sample}(\text{gauss}(0.0, 3.0)) \text{ in} \\ & \quad \text{score}(\text{density_gauss}(5.0, (x, 1.0)); \text{return}(x < 4.5)) \rrbracket \\ & = (0, (xxx, \text{bern}(0.5))) \\ & \llbracket \text{norm}(\text{let } x = \text{sample}(\text{bern}(0.25)) \text{ in} \\ & \quad (\text{if } x \text{ then score}(5.0) \text{ else score}(2.0)); \text{return}(x)) \rrbracket \\ & = (0, (2.75, \text{bern}(\frac{5}{11}))) \\ & \llbracket \text{norm}(\text{let } x = \text{sample}(\text{exp}(1.0)) \text{ in score}(e^x)) \rrbracket = (2, *) \\ & \llbracket \text{norm}(\text{let } x = \text{sample}(\text{beta}(1, 3)) \rrbracket = \llbracket \text{norm}(\text{score}(\frac{1}{1+3}); \\ & \quad \text{in score}(x); \text{return}(x)) \rrbracket = \llbracket \text{norm}(\text{sample}(\text{beta}(2, 3))) \rrbracket \end{aligned}$$

The fourth equation shows how infinite model evidence errors can arise when working with infinite distributions. In the last equation, the parameter x of $\text{score}(x)$ represents the probability of true under $\text{bern}(x)$. The equation expresses the so called conjugate-prior relationship between Beta and Bernoulli distributions, which has been used to optimise probabilistic programs [35].

Monads The interpretation of `let` and `return` given above arises from the fact that $P(\mathbb{R}_{\geq 0} \times (-))$ is a commutative monad on the category of measurable spaces and measurable functions (see also [6, §2.3.1], [? , §6]). Recall that a commutative monad (T, η, μ, σ) in general comprises an endofunctor T together with natural transformations $\eta_X : X \rightarrow T(X)$, $\mu_X : T(T(X)) \rightarrow T(X)$, $\sigma_{X,Y} : X \times T(Y) \rightarrow T(X \times Y)$ satisfying some laws [16]. Using this structure we interpret `return` and `let` following Moggi [22]:

$$\begin{aligned} & \llbracket \Gamma \vdash \text{return}(t) : \mathbb{A} \rrbracket(\gamma) \stackrel{\text{def}}{=} \eta_{\llbracket \mathbb{A} \rrbracket}(\llbracket t \rrbracket(\gamma)) \\ & \llbracket \Gamma \vdash \text{let } x = t \text{ in } u : \mathbb{B} \rrbracket(\gamma) \stackrel{\text{def}}{=} \mu_{\llbracket \mathbb{B} \rrbracket}(T(\llbracket u \rrbracket)(\sigma_{\llbracket \Gamma \rrbracket, \llbracket \mathbb{A} \rrbracket}(\gamma, \llbracket t \rrbracket(\gamma)))) \end{aligned}$$

Concretely, we make $P(\mathbb{R}_{\geq 0} \times (-))$ into a monad by combining the standard commutative monad structure [11] on P and the commutative monoid $(\mathbb{R}_{\geq 0}, \cdot, 1)$ with the monoid monad transformer.

(We record a subtle point about ι (3). It is *not* a monad morphism, and we cannot form a commutative monad of *all* measures because the Fubini property does not hold in general.)

4.1 Sequential Monte Carlo simulation

The program equations above justify some simple program transformations. For a more sophisticated one, consider *sequential Monte Carlo simulation* [7]. A key idea of its application to probabilistic programs is: ‘Whenever there is a score, it is good to renormalize and resample’. This increases efficiency by avoiding too many program executions with low scores [24, Algorithm 1].

The denotational semantics justifies the soundness of this transformation. For a term with top-level score, *i.e.* a term of the form $(\text{let } x = t \text{ in } (\text{score}(u); v))$ where u and v may have x free:

$$\begin{aligned} & \llbracket \text{norm}(\text{let } x = t \text{ in } (\text{score}(u); v)) \rrbracket \\ & = \llbracket \text{norm}(\text{case } (\text{norm}(\text{let } x = t \text{ in } \text{score}(u); \text{return}(x))) \text{ of} \quad 1 \\ & \quad (0, (e, d)) \Rightarrow \text{score}(e); \text{let } x = \text{sample}(d) \text{ in } v \quad 2 \\ & \quad | (1, *) \Rightarrow \text{score}(0); \text{return}(w) \quad 3 \\ & \quad | (2, *) \Rightarrow \text{let } x = t \text{ in } (\text{score}(u); v)) \rrbracket \quad 4 \end{aligned}$$

Let us explain the right hand side. Line 1 renormalizes the program after the score, and in non-exceptional execution returns the model evidence e and a new normalized distribution d . Line 2 immediately records the evidence e as a score, and then resamples d , using

the resampled value in the continuation v . Line 3 propagates the error of 0: w is a deterministic term of the right type whose choice does not matter. Finally, line 4 detects an infinite evidence error, and undoes the transformation. This error does not arise in most applications of sequential Monte Carlo simulation.

5. Operational semantics

In this section we develop an operational semantics for the first-order language. There are several reasons to consider this, even though the denotational semantics is arguably straightforward. First, extension to higher-order functions is easier in operational semantics than in denotational semantics. Second, operational semantics conveys computational intuitions that are obscured in the denotational semantics. We expect these computational intuitions to play an important role in studying approximate techniques for performing posterior inference, such as sequential Monte Carlo, in the future.

Sampling from probability distributions complicates operational semantics. Sampling from a discrete distribution can immediately affect control flow. For example, in the term

`let $x = \text{sample}(\text{bern}(0.5))$ in if x then return(1.1) else return(8.0)`

the conditional depends on the result of sampling the Bernoulli distribution. The result is 1.1 with probability 0.5 (*cf.* [5, §2.3]).

Sampling a distribution on \mathbb{R} introduces another complication. Informally, there is a transition

$$\text{sample}(\text{gauss}(0.0, 1.0)) \longrightarrow \text{return}(r)$$

for every real r , but any single transition has zero probability. We can assign non-zero probabilities to sets of transitions; informally:

$$\Pr(\text{sample}(\text{gauss}(0.0, 1.0)) \longrightarrow \{\text{return}(r) \mid r \leq 0\}) = 0.5.$$

To make this precise we need a σ -algebra on the set of terms, which can be done using *configurations* rather than individual terms. A configuration is a *closure*: a pair $\langle t, \gamma \rangle$ of a term t with free variables and an environment γ giving values for those variables as elements of a measurable space. (See also [13, §3], [? , §3].)

Sampling a distribution p on $\mathbb{R} + \mathbb{R}$ exhibits both complications:

$$\begin{aligned} \text{let } x = \text{sample}(p) \text{ in case } x \text{ of } (0, r) & \Rightarrow \text{return}(r + 1.0) \\ | (1, r) & \Rightarrow \text{return}(r - 1.0) \end{aligned} \quad (4)$$

The control flow in the *case* distinction depends on which summand is sampled, but there is potentially a continuous distribution over the return values. We handle this by instantiating the choice of summand in the syntax, but keeping the value of the summand in the environment, so that expression (4) can make a step to the closure

$$\begin{aligned} & \langle \text{let } x = \text{return}(0, y) \text{ in} \quad , y \mapsto 42.0 \rangle. \\ & \text{case } x \text{ of } (0, r) \Rightarrow \text{return}(r + 1.0) \\ & | (1, r) \Rightarrow \text{return}(r - 1.0) \end{aligned}$$

A type is *indecomposable* if it has the form \mathbb{R} or $P(\mathbb{A})$, and a context Γ is *canonical* if it only involves indecomposable types.

Configurations Let $z \in \{d, p\}$. A z -*configuration* of type \mathbb{A} is a triple $\langle \Gamma, t, \gamma \rangle$ comprising a canonical context Γ , a derivable term $\Gamma \vdash t : \mathbb{A}$, and an element γ of the measurable space $\llbracket \Gamma \rrbracket$. We identify contexts that merely rename variables, such as

$$\begin{aligned} & \langle (\langle u : \mathbb{R}, y : P(\mathbb{R}) \rangle), f(x, y), (x \mapsto 42.0, y \mapsto \text{gauss}(0.0, 1.0)) \rangle \\ & \approx \langle (\langle u : \mathbb{R}, v : P(\mathbb{R}) \rangle), f(u, v), (u \mapsto 42.0, v \mapsto \text{gauss}(0.0, 1.0)) \rangle. \end{aligned}$$

We call d -configurations *deterministic configurations*, and p -configurations *probabilistic configurations*; they differ only in typing. We will abbreviate configurations to $\langle t, \gamma \rangle$ when Γ is obvious.

Values v in a canonical context Γ are well-typed deterministic terms of the form

$$v, w ::= x_i \mid * \mid (v, w) \mid (i, v) \quad (5)$$

where x_i is a variable in Γ . Similarly, a probabilistic term t in context Γ is called *probabilistic value* or *p-value* if $t \equiv \text{return}(v_0)$ for some value v_0 . Remember from Section 4 that the denotational semantics of values is simple and straightforward.

Write $\text{Con}_d(\mathbb{A})$ and $\text{Con}_p(\mathbb{A})$ for the sets of deterministic and probabilistic configurations of type \mathbb{A} , and make them into measurable spaces by declaring $U \subseteq \text{Con}_z(\mathbb{A})$ to be measurable if the set $\{\gamma \in \llbracket \Gamma \rrbracket \mid \langle t, \gamma \rangle \in U\}$ is measurable for all judgements $\Gamma \vdash t : \mathbb{A}$.

$$\text{Con}_z(\mathbb{A}) = \sum_{\substack{(\Gamma, t) \\ \Gamma \text{ canonical,} \\ \Gamma \vdash t : \mathbb{A}}} \llbracket \Gamma \rrbracket \quad (6)$$

Further partition $\text{Con}_z(\mathbb{A})$ into $\text{ConV}_z(\mathbb{A})$ and $\text{ConN}_z(\mathbb{A})$ based on whether a term in a configuration is a value or not:

$$\begin{aligned} \text{ConV}_d(\mathbb{A}) &= \{\langle \Gamma, t, \gamma \rangle \in \text{Con}_d(\mathbb{A}) \mid t \text{ is a value}\} \\ \text{ConN}_d(\mathbb{A}) &= \{\langle \Gamma, t, \gamma \rangle \in \text{Con}_d(\mathbb{A}) \mid t \text{ is not a value}\} \\ \text{ConV}_p(\mathbb{A}) &= \{\langle \Gamma, t, \gamma \rangle \in \text{Con}_p(\mathbb{A}) \mid t \text{ is a p-value}\} \\ \text{ConN}_p(\mathbb{A}) &= \{\langle \Gamma, t, \gamma \rangle \in \text{Con}_p(\mathbb{A}) \mid t \text{ is not a p-value}\} \end{aligned}$$

Particularly well-behaved values are the *ordered values* $\Gamma \vdash v : \mathbb{A}$, where each variable appears exactly once, and in the same order as in Γ .

Lemma 5.1. *Consider a canonical context Γ , a type \mathbb{A} , an ordered value $\Gamma \vdash v : \mathbb{A}$, and the induced measurable function*

$$\llbracket v \rrbracket : \llbracket \Gamma \rrbracket \rightarrow \llbracket \mathbb{A} \rrbracket.$$

The collection of all such functions for given \mathbb{A} is countable, and forms a coproduct diagram.

Proof. By induction on the structure of types. The key fact is that every type is a sum of products of indecomposable ones, because the category of measurable spaces is distributive, *i.e.* the canonical map $\sum_{i \in I} (\mathbb{A} \times \mathbb{B}_i) \rightarrow \mathbb{A} \times \sum_{i \in I} \mathbb{B}_i$ is an isomorphism. \square

For example, $\mathbb{A} = (\mathbb{R} \times \text{bool}) + (\mathbb{R} \times \mathbb{R})$ has 3 ordered values, first $(x : \mathbb{R} \vdash (0, (x, \text{true})) : \mathbb{A})$, second $(x : \mathbb{R} \vdash (0, (x, \text{false})) : \mathbb{A})$, and third $(x : \mathbb{R}, y : \mathbb{R} \vdash (1, (x, y)) : \mathbb{A})$, inducing a canonical measurable isomorphism $\mathbb{R} + \mathbb{R} + \mathbb{R} \times \mathbb{R} \cong \llbracket \mathbb{A} \rrbracket$.

Evaluation contexts We distinguish three kinds of *evaluation contexts*: $\mathcal{C}[-]$ is a context for a deterministic term with a hole for deterministic terms; $\mathcal{D}[-]$ and $\mathcal{E}[-]$ are contexts for probabilistic terms, the former with a hole for probabilistic terms, the latter with a hole for deterministic terms.

$$\begin{aligned} \mathcal{C}[-] ::= & (-) \mid \pi_j \mathcal{C}[-] \mid (\mathcal{C}[-], t) \mid (v, \mathcal{C}[-]) \mid (i, \mathcal{C}[-]) \\ & \mid \text{case } \mathcal{C}[-] \text{ of } \{(i, x) \Rightarrow t_i\}_{i \in I} \mid f(\mathcal{C}[-]) \end{aligned}$$

$$\mathcal{D}[-] ::= (-) \mid \text{let } x = \mathcal{D}[-] \text{ in } t \quad (7)$$

$$\begin{aligned} \mathcal{E}[-] ::= & \mathcal{D}[\text{return}[-]] \mid \mathcal{D}[\text{sample}[-]] \mid \mathcal{D}[\text{score}[-]] \\ & \mid \text{case } \mathcal{D}[-] \text{ of } \{(i, x) \Rightarrow t_i\}_{i \in I} \end{aligned}$$

where t, t_i are general terms and v is a value.

5.1 Reduction

Using the tools developed so far, we will define a measurable function for describing the reduction of d-configurations, and a stochastic relation for describing that of p-configurations:

$$\begin{aligned} \longrightarrow & : \text{ConN}_d(\mathbb{A}) \rightarrow \text{Con}_d(\mathbb{A}), \\ \longrightarrow & : \text{ConN}_p(\mathbb{A}) \times \Sigma_{\mathbb{R}_{\geq 0} \times \text{Con}_p(\mathbb{A})} \rightarrow [0, 1], \end{aligned}$$

parameterised by a family of measurable functions

$$\nu_{\mathbb{A}} : \text{Con}_p(\mathbb{A}) \rightarrow (\mathbb{R}_{\geq 0} \times P(\llbracket \mathbb{A} \rrbracket)) + 1 + 1 \quad (8)$$

indexed by types \mathbb{A} .

Reduction of deterministic terms Define a type-indexed family of relations $\longrightarrow \subseteq \text{ConN}_d(\mathbb{A}) \times \text{Con}_d(\mathbb{A})$ as the least one that is closed under the following rules.

$$\begin{aligned} \langle \Gamma, \pi_j(v_0, v_1), \gamma \rangle &\longrightarrow \langle \Gamma, v_j, \gamma \rangle \\ \langle \Gamma, \text{case } (i', v) \text{ of } \{(i, x) \Rightarrow t_i\}_{i \in I}, \gamma \rangle &\longrightarrow \langle \Gamma, t_{i'}[v/x], \gamma \rangle \\ \langle \Gamma, f(w), \gamma \rangle &\longrightarrow \langle (\Gamma, \Gamma'), v, (\gamma, \gamma') \rangle \\ &\quad (w \text{ a value } \wedge \Gamma' \vdash v : \mathbb{A} \text{ an ordered value } \wedge f(\llbracket w \rrbracket(\gamma)) = \llbracket v \rrbracket(\gamma')) \\ \langle \Gamma, \text{norm}(t), \gamma \rangle &\longrightarrow \langle (\Gamma, x : \mathbb{R}, y : P(\mathbb{B})), (0, (x, y)), \gamma[x \mapsto r, y \mapsto p] \rangle \\ &\quad (\mathbb{A} = P(\mathbb{B}) \wedge \nu_{\mathbb{B}}(\langle \Gamma, t, \gamma \rangle) = (0, (r, p)) \wedge x, y \notin \text{dom}(\Gamma)) \\ \langle \Gamma, \text{norm}(t), \gamma \rangle &\longrightarrow \langle \Gamma, (i, *), \gamma \rangle \\ &\quad (\mathbb{A} = P(\mathbb{B}) \wedge \nu_{\mathbb{B}}(\langle \Gamma, t, \gamma \rangle) = (i, *), i \in \{1, 2\}) \\ \frac{\langle \Gamma, t, \gamma \rangle \longrightarrow \langle \Gamma', t', \gamma' \rangle}{\langle \Gamma, \mathcal{C}[t], \gamma \rangle \longrightarrow \langle \Gamma', \mathcal{C}[t'], \gamma' \rangle} &\quad (\mathcal{C}[-] \text{ is not } (-)) \end{aligned}$$

The rule for $f(w)$ keeps the original context Γ and the closure γ because they might be used in the continuation, even though they are not used in v . The rules obey the following invariant.

Lemma 5.2. *If $\langle \Gamma, t, \gamma \rangle \longrightarrow \langle \Gamma', t', \gamma' \rangle$, then $\Gamma' = (\Gamma, \Gamma'')$ and $\gamma' = (\gamma, \gamma'')$ for some Γ'' and $\gamma'' \in \llbracket \Gamma'' \rrbracket$.*

Proof. By induction on the structure of derivations. \square

This lemma allows us to confirm that our specification of a relation $\longrightarrow \subseteq \text{ConN}_d(\mathbb{A}) \times \text{Con}_d(\mathbb{A})$ is well-formed ('type preservation').

Proposition 5.3. *The induced relation is a measurable function.*

Proof. There are three things to show: that the relation is entire ('progress'); that the relation is single-valued ('determinacy'); and that the induced function is measurable. All three are shown by induction on the structure of terms. The case of application of measurable functions crucially uses Lemma 5.1. \square

Reduction of probabilistic terms Next, we define the stochastic relation \longrightarrow for probabilistic terms, combining two standard approaches: for indecomposable types, which are uncountable, use labelled Markov processes, *i.e.* give a distribution on the measurable set of resulting configurations; for decomposable types (sums, products *etc.*), probabilistic branching is discrete and so a transition system labelled by probabilities suffices.

Proposition 5.4. *Let $(X_i)_{i \in I}$ be an indexed family of measurable spaces. Suppose we are given:*

- a function $q : I \rightarrow [0, 1]$ that is only nonzero on a countable subset $I_0 \subseteq I$, and such that $\sum_{i \in I_0} q(i) = 1$;
- a probability measure q_i on X_i for each $i \in I_0$.

This determines a probability measure p on $\sum_{i \in I} X_i$ by

$$p(U) = \sum_{i \in I} q(i) q_i(\{a \mid (i, a) \in U\})$$

for U a measurable subset of $\sum_{i \in I} X_i$,

We will use three entities to define the desired stochastic relation $\longrightarrow : \text{ConN}_p(\mathbb{A}) \times \Sigma_{\mathbb{R}_{\geq 0} \times \text{Con}_p(\mathbb{A})} \rightarrow [0, 1]$.

$$\begin{aligned}
& \Pr(\langle \Gamma, \mathcal{E}[t], \gamma \rangle \longrightarrow (\Gamma', \mathcal{E}[t'])) \stackrel{\text{def}}{=} [\langle \Gamma, t, \gamma \rangle \longrightarrow (\Gamma', t', \gamma')] \\
& \Pr(\langle \Gamma, \mathcal{D}[t], \gamma \rangle \longrightarrow (\Gamma', \mathcal{D}[t'])) \stackrel{\text{def}}{=} \Pr(\langle \Gamma, t, \gamma \rangle \longrightarrow (\Gamma', t')) \\
& \Pr(\langle \Gamma, \text{let } x = \text{return}(v) \text{ in } t, \gamma \rangle \longrightarrow (\Gamma, t[v/x])) \stackrel{\text{def}}{=} 1 \\
& \Pr(\langle \Gamma, \text{case } (j, v) \text{ of } \{(i, x) \Rightarrow t_i\}_{i \in I}, \gamma \rangle \longrightarrow (\Gamma, t_j[v/x])) \stackrel{\text{def}}{=} 1 \\
& \Pr(\langle \Gamma, \text{score}(v), \gamma \rangle \longrightarrow (\Gamma, \text{return}(*))) \stackrel{\text{def}}{=} 1 \\
& \Pr(\langle \Gamma, \text{sample}(v), \gamma \rangle \longrightarrow ((\Gamma, \Gamma'), \text{return}(v'))) \\
& \stackrel{\text{def}}{=} \begin{cases} \llbracket v \rrbracket(\gamma)(\{\llbracket v' \rrbracket(\gamma') \mid \gamma' \in \llbracket \Gamma' \rrbracket\}) & \text{if } \Gamma' \Vdash v' : \mathbb{A} \text{ ordered} \\ 0 & \text{otherwise} \end{cases} \\
& \Pr(\langle \Gamma, \mathcal{E}[t], \gamma \rangle \longrightarrow_{(\Gamma', \mathcal{E}[t'])} U) \\
& \stackrel{\text{def}}{=} [\langle \Gamma, t, \gamma \rangle \longrightarrow \langle \Gamma', t', \gamma' \rangle \wedge \gamma' \in U] \\
& \Pr(\langle \Gamma, \mathcal{D}[t], \gamma \rangle \longrightarrow_{(\Gamma', \mathcal{D}[t'])} U) \stackrel{\text{def}}{=} \Pr(\langle \Gamma, t, \gamma \rangle \longrightarrow_{(\Gamma', t')} U) \\
& \Pr(\langle \Gamma, \text{let } x = \text{return}(v) \text{ in } t, \gamma \rangle \longrightarrow_{(\Gamma, t[v/x])} U) \stackrel{\text{def}}{=} [\gamma \in U] \\
& \Pr(\langle \Gamma, \text{case } (j, v) \text{ of } \{(i, x) \Rightarrow t_i\}_{i \in I}, \gamma \rangle \longrightarrow_{\Gamma, t_j[v/x]} U) \stackrel{\text{def}}{=} [\gamma \in U] \\
& \Pr(\langle \Gamma, \text{score}(v), \gamma \rangle \longrightarrow_{(\Gamma, \text{return}(*))} U) \stackrel{\text{def}}{=} [\gamma \in U] \\
& \Pr(\langle \Gamma, \text{sample}(v), \gamma \rangle \longrightarrow_{((\Gamma, \Gamma'), \text{return}(v'))} U) \\
& \stackrel{\text{def}}{=} \frac{\llbracket v \rrbracket(\gamma)(\{\llbracket v' \rrbracket(\gamma') \mid \gamma' \in \llbracket \Gamma' \rrbracket \wedge (\gamma, \gamma') \in U\})}{\llbracket v \rrbracket(\gamma)(\{\llbracket v' \rrbracket(\gamma') \mid \gamma' \in \llbracket \Gamma' \rrbracket\})} \\
& \text{Sc}(\langle \Gamma, \mathcal{E}[t], \gamma \rangle) \stackrel{\text{def}}{=} 1 \quad \text{Sc}(\langle \Gamma, \mathcal{D}[t], \gamma \rangle) \stackrel{\text{def}}{=} \text{Sc}(\langle \Gamma, t, \gamma \rangle) \\
& \text{Sc}(\langle \Gamma, \text{sample}(v), \gamma \rangle) \stackrel{\text{def}}{=} 1 \\
& \text{Sc}(\langle \Gamma, \text{score}(v), \gamma \rangle) \stackrel{\text{def}}{=} \max(\llbracket v \rrbracket(\gamma), 0) \\
& \text{Sc}(\langle \Gamma, \text{let } x = \text{return}(v) \text{ in } t, \gamma \rangle) \stackrel{\text{def}}{=} 1 \\
& \text{Sc}(\langle \Gamma, \text{case } (j, v) \text{ of } \{(i, x) \Rightarrow t_i\}_{i \in I}, \gamma \rangle) \stackrel{\text{def}}{=} 1
\end{aligned}$$

Figure 1. Entities used to define reduction of probabilistic terms

1. A countably supported probability distribution on the set $\{(\Gamma, t) \mid \Gamma \Vdash t : \mathbb{A}\}$ for each $C \in \text{ConN}_p(\mathbb{A})$. We write $\Pr(C \longrightarrow (\Gamma, t))$ for the probability of (Γ, t) .
2. A probability measure on the space $\llbracket \Gamma \rrbracket$ for each $C \in \text{ConN}_p(\mathbb{A})$ and (Γ, t) with $\Pr(C \longrightarrow (\Gamma, t)) \neq 0$. Write $\Pr(C \longrightarrow_{\Gamma, t} U)$ for the probability of a measurable subset $U \subseteq \llbracket \Gamma \rrbracket$.
3. A measurable function $\text{Sc} : \text{ConN}_p(\mathbb{A}) \rightarrow \mathbb{R}_{\geq 0}$, representing the score of the one-step transition relation. (For one-step transitions, the score is actually deterministic.)

These three entities are defined by induction on the structure of the syntax of \mathbb{A} -typed p-configurations in Figure 1. We combine them to define a stochastic relation as follows.

Proposition 5.5. *The map $\text{ConN}_p(\mathbb{A}) \times \Sigma_{\mathbb{R}_{\geq 0} \times \text{Con}_p(\mathbb{A})} \rightarrow [0, 1]$ that sends (C, U) to $\Pr(C \longrightarrow U)$, defined as*

$$\sum_{(\Gamma, t)} \Pr(C \longrightarrow (\Gamma, t)) \Pr(C \longrightarrow_{\Gamma, t} \{\gamma \mid (\text{Sc}(C), \langle \Gamma, t, \gamma \rangle) \in U\}),$$

is a stochastic relation.

Proof. For each p-configuration $C = \langle -, t, - \rangle$, use induction on t to see that the probability distribution $\Pr(C \longrightarrow (-))$ on pairs

(Γ', t') and the distribution $\Pr(C \longrightarrow_{(-)} (-))$ indexed by such pairs satisfy the conditions in Proposition 5.4. It follows that the partial evaluation $\Pr(C \longrightarrow (-))$ of the function in the statement is a probability measure, so it suffices to establish measurability of the other partial evaluation $\Pr((-) \longrightarrow U)$. Recall that $\text{ConN}_p(\mathbb{A})$ is defined in terms of the sum of measurable spaces, and that all p-configurations in each summand have the same term. Finally, use induction on the term shared by all p-configurations in the summand to see that the restriction of $\Pr((-) \longrightarrow U)$ to each summand is measurable. \square

5.2 Termination

To see that the reduction process terminates, we first define the transitive closure. This is subtle, as sampling can introduce countably infinite branching; although each branch will terminate, the expected number of steps to termination can be infinite.

We use the deterministic transition relation to define an evaluation relation $\Downarrow \subseteq \text{Con}_d(\mathbb{A}) \times \text{ConV}_d(\mathbb{A})$, by setting $C \Downarrow D$ if $\exists n. C \Downarrow^n D$, where

$$\frac{}{C \Downarrow^0 C} \quad (C \in \text{ConV}_d(\mathbb{A})) \quad \frac{D \Downarrow^n E \quad C \longrightarrow D}{C \Downarrow^{n+1} E}$$

To define evaluation for probabilistic configurations, we need *sub-stochastic relations*: functions $f : X \times \Sigma_Y \rightarrow [0, 1]$ that are measurable in X , satisfy $f(x, Y) \leq 1$ for every $x \in X$, and are countably additive in Y , i.e. $f(x, \bigcup_{i \in \mathbb{N}} U_i) = \sum_{i \in \mathbb{N}} f(x, U_i)$ for a sequence U_1, U_2, \dots of disjoint measurable sets. Thus a stochastic relation (as in Definition 2.3) is a sub-stochastic relation with $f(x, Y) = 1$. Define a sub-stochastic relation

$$\Pr(- \Downarrow -) : \text{Con}_p(\mathbb{A}) \times \Sigma_{\mathbb{R}_{\geq 0} \times \text{ConV}_p(\mathbb{A})} \rightarrow [0, 1]$$

by $\Pr(C \Downarrow U) \stackrel{\text{def}}{=} \sum_n \Pr(C \Downarrow^n U)$, where $\Pr(C \Downarrow^0 U)$ is given by $[(1, C) \in U]$, and $\Pr(C \Downarrow^{n+1} U)$ is

$$\int_{(r, D)} \Pr(D \Downarrow^n \{(s, E) \mid (r \cdot s, E) \in U\}) \Pr(C \longrightarrow d(r, D)).$$

Proposition 5.6 (Termination). *Evaluation of deterministic terms is a function: $\forall C. \exists D. C \Downarrow D$. Evaluation of probabilistic terms is a stochastic relation: $\forall C. \Pr(C \Downarrow (\mathbb{R}_{\geq 0} \times \text{ConV}_p(\mathbb{A}))) = 1$.*

Proof. By induction on the structure of terms. \square

Termination comes from the omission of recursion in our language. We do so for now, because our semantic model does not yet handle higher-order recursion, and probabilistic while-languages are already well-understood (e.g. [17]). (See also the discussion about domain theory in §8.)

5.3 Soundness

For soundness, extend the denotational semantics to configurations:

- define $s_d : \text{Con}_d(\mathbb{A}) \rightarrow \llbracket \mathbb{A} \rrbracket$ by $\langle \Gamma, t, \gamma \rangle \mapsto \llbracket t \rrbracket(\gamma)$;
- define $s_p : \text{Con}_p(\mathbb{A}) \times \Sigma_{\mathbb{R}_{\geq 0} \times \llbracket \mathbb{A} \rrbracket} \rightarrow [0, 1]$ similarly by $(\langle \Gamma, t, \gamma \rangle, U) \mapsto \llbracket t \rrbracket(\gamma)(U)$. We may also use this stochastic relation as a measurable function $s_p : \text{Con}_p(\mathbb{A}) \rightarrow P(\mathbb{R}_{\geq 0} \times \llbracket \mathbb{A} \rrbracket)$;
- define $s_{v_p} : \text{ConV}_p(\mathbb{A}) \rightarrow \llbracket \mathbb{A} \rrbracket$ by $\langle \Gamma, \text{return}(v), \gamma \rangle \mapsto \llbracket v \rrbracket(\gamma)$. Note that in this first-order language, s_{v_p} is a surjection which equates two value configurations iff they are related by weakening, contraction or exchange of variables.

Assumption 5.7. Throughout this section we assume that the normalization function ν on configurations (8) is perfect, i.e. it corresponds to ι , the semantic normalization function (3):

$$\nu_{\mathbb{A}}(\langle \Gamma, t, \gamma \rangle) = \iota_{\llbracket \mathbb{A} \rrbracket}(s_p(\langle \Gamma, t, \gamma \rangle)).$$

Lemma 5.8 (Context extension). *Let $z \in \{d, p\}$. Suppose that $\langle \Gamma, t, \gamma \rangle$ and $\langle \Gamma, \Gamma', t, (\gamma, \gamma') \rangle$ are configurations in $\text{Con}_z(\mathbb{A})$. Then $s_z(\langle \Gamma, t, \gamma \rangle) = s_z(\langle \Gamma, \Gamma', t, (\gamma, \gamma') \rangle)$.*

Proposition 5.9 (Soundness). *The following diagrams commute (in the category of measurable functions, and stochastic relations, respectively).*

$$\begin{array}{ccc}
 \text{ConN}_d(\mathbb{A}) & \xrightarrow{s_d} & \text{ConN}_p(\mathbb{A}) \xrightarrow{s_p} \mathbb{R}_{\geq 0} \times \llbracket \mathbb{A} \rrbracket \\
 \text{one-step reduction} \downarrow & \nearrow & \downarrow \text{one-step reduction} \\
 \text{Con}_d(\mathbb{A}) & \xrightarrow{s_d} & \mathbb{R}_{\geq 0} \times \text{Con}_p(\mathbb{A}) \xrightarrow{\text{id} \times s_p} \mathbb{R}_{\geq 0} \times \mathbb{R}_{\geq 0} \times \llbracket \mathbb{A} \rrbracket \\
 & & \uparrow \text{multiplication}
 \end{array}$$

Proof. By induction on the structure of syntax. The inductive steps with evaluation contexts use the extension Lemma 5.8, which applies by the invariant Lemma 5.2. \square

Adequacy The denotational semantics is adequate, in the sense:

$$\llbracket t \rrbracket(*) = P(\mathbb{R}_{\geq 0} \times s_{V_P})(\Pr(\langle \emptyset, t, * \rangle \Downarrow (-))) \quad \text{for all } \vdash_P t : \mathbb{A}.$$

That is, the denotation $\llbracket t \rrbracket(*)$ is nothing but pushing forward the probability measure $\Pr(\langle \emptyset, t, * \rangle \Downarrow (-))$ from the operational semantics along the function s_{V_P} . This adequacy condition holds because Proposition 5.9 ensures that

$$\left(\sum_{k \leq n} \Pr(\langle \emptyset, t, * \rangle \Downarrow^k \{(r, C) \mid (r, s_{V_P}(C)) \in U\}) \right) \leq \llbracket t \rrbracket(*) (U)$$

for all n and U , and Proposition 5.6 then guarantees that the left-hand side of this inequality converges to the right-hand side as n tends to infinity.

6. A higher-order language

This short section extends the first-order language with functions and thunks [18], allowing variables to stand for program fragments. In other words, ‘programs are first-class citizens’.

Types Extend the grammar for types with two new constructors.

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

Informally, $\mathbb{A} \Rightarrow \mathbb{B}$ contains deterministic functions, and $T(\mathbb{A})$ contains thunked (*i.e.* suspended) probabilistic programs. Then $\mathbb{A} \Rightarrow T(\mathbb{B})$ contains probabilistic functions. A type is *measurable* if it does not involve \Rightarrow or T , *i.e.* if it is in the grammar of Section 3.

Terms Extend the term language with the following rules. First, the usual abstraction and application of deterministic functions:

$$\frac{\Gamma, x : \mathbb{A} \vdash t : \mathbb{B}}{\Gamma \vdash \lambda x. t : \mathbb{A} \Rightarrow \mathbb{B}} \quad \frac{\Gamma \vdash t : \mathbb{A} \Rightarrow \mathbb{B} \quad \Gamma \vdash u : \mathbb{A}}{\Gamma \vdash t u : \mathbb{B}}$$

Second, we have syntax for thunking and forcing (*e.g.* [18, 22, 25]).

$$\frac{\Gamma \vdash_P t : \mathbb{A}}{\Gamma \vdash \text{thunk}(t) : T(\mathbb{A})} \quad \frac{\Gamma \vdash t : T(\mathbb{A})}{\Gamma \vdash_P \text{force}(t) : \mathbb{A}}$$

All the rules from Section 3 are also still in force. For rule (2) to still make sense, we only include constant terms for measurable functions $f : \llbracket \mathbb{A} \rrbracket \rightarrow \llbracket \mathbb{B} \rrbracket$ between measurable types \mathbb{A} and \mathbb{B} .

Examples One reason for higher types is to support code structuring. The separate function types and thunk types allow us to be flexible about calling conventions. For example, sampling can be reified as the ground term

$$\vdash \lambda x. \text{thunk}(\text{sample}(x)) : P(\mathbb{A}) \Rightarrow T(\mathbb{A}),$$

which takes a probability measure and returns a suspended program that will sample from it. On the other hand, to reify the normalization construction, we use a different calling convention.

$$\vdash \lambda x. \text{norm}(\text{force}(x)) : T(\mathbb{A}) \Rightarrow \mathbb{R} \times P(\mathbb{A}) + 1 + 1$$

This function takes a suspended probabilistic program and returns the result of normalizing it.

Example: higher-order expectation Higher types also allow us to consider probability distributions over programs. For an example, consider this term.

$$\begin{aligned}
 E_h &\stackrel{\text{def}}{=} \lambda(d, f) : T(\mathbb{A}) \times (\mathbb{A} \Rightarrow \mathbb{R}). \\
 &\quad \text{case } (\text{norm}(\text{let } a = \text{force}(d) \text{ in } \text{score}(f(a)))) \text{ of} \\
 &\quad \quad (0, (e, y)) \Rightarrow e \\
 &\quad \quad \mid (1, *) \Rightarrow 0.0 \mid (2, *) \Rightarrow 0.0
 \end{aligned}$$

It has type $T(\mathbb{A}) \times (\mathbb{A} \Rightarrow \mathbb{R}) \Rightarrow \mathbb{R}$. Intuitively, given a thunked probabilistic term t and a function f that is nonnegative, E_h treats t as a probability distribution on \mathbb{A} , and computes the expectation of f on this distribution. Notice that \mathbb{A} can be a higher type, so E_h generalises the usual notion of expectation, which has not been defined for higher types because the category of measurable spaces is not Cartesian closed.

7. Higher-order operational semantics

In this section we consider operational semantics for the higher-order extension of the language. In an operational intuition, $\text{force}(t)$ forces a suspended computation t to run. For example,

$$\vdash \text{thunk}(\text{sample}(\text{gauss}(0.0, 1.0))) : T(\mathbb{R})$$

is a suspended computation that, when forced, will sample the normal distribution.

Assumption 7.1. From the operational perspective it is unclear how to deal with sampling from a distribution over functions. For this reason, in this section, we only allow the type $P(\mathbb{A})$ when \mathbb{A} is a measurable type. We still allow probabilistic terms to have higher-order types, and we still allow $T(\mathbb{A})$ where \mathbb{A} is higher-order.

7.1 Reduction

We now extend the operational semantics from Section 5 with higher types. Values (5) are extended as follows.

$$v ::= \dots \mid \lambda x. t \mid \text{thunk}(t)$$

Evaluation contexts (7) are extended as follows.

$$C[-] ::= \dots \mid C[-] t \mid v C[-] \quad \mathcal{E}[-] ::= \dots \mid \mathcal{D}[\text{force}[-]]$$

There are two additional redexes: $(\lambda x. t) v$ and $\text{force}(\text{thunk}(t))$. The deterministic transition relation is extended with this β -rule:

$$\langle \Gamma, (\lambda x. t) v, \gamma \rangle \longrightarrow \langle \Gamma, t[v/x], \gamma \rangle.$$

Extend the probabilistic transition relation with the following rules.

$$\begin{aligned}
 \Pr(\langle \Gamma, \text{force}(\text{thunk}(t)), \gamma \rangle \longrightarrow \langle \Gamma, t \rangle) &= 1 \\
 \Pr(\langle \Gamma, \text{force}(\text{thunk}(t)), \gamma \rangle \longrightarrow_{(\Gamma, t)} U) &= [\gamma \in U] \\
 \text{Sc}(\langle \Gamma, \text{force}(\text{thunk}(t)), \gamma \rangle) &= 1
 \end{aligned}$$

7.2 Termination

The evaluation relations for deterministic and probabilistic configurations of the higher-order language are defined as in Subsection 5.2. The resulting rewriting system still terminates, even though configurations may now include higher-order terms.

Proposition 7.2 (Termination). *Evaluation of deterministic terms is a function: $\forall C. \exists D. C \Downarrow D$. Evaluation of probabilistic terms is a stochastic relation: $\forall C. \Pr(C \Downarrow (\mathbb{R}_{\geq 0} \times \text{ConV}_P(\mathbb{A}))) = 1$.*

Proof. We sketch an invariant of higher-order terms that implies the termination property, formulated as unary logical relations via sets

$$\begin{aligned} R(\Gamma \Vdash \mathbb{A}) &\subseteq \{t \mid \Gamma \Vdash t : \mathbb{A}\}, \\ R_v(\Gamma \Vdash \mathbb{A}) &\subseteq \{t \mid \Gamma \Vdash t : \mathbb{A} \wedge t \text{ a z-value}\}, \end{aligned}$$

for each canonical context Γ , type \mathbb{A} , and $z \in \{d, p\}$, defined by:

$$\begin{aligned} R(\Gamma \Vdash \mathbb{A}) &= \{t \mid \forall \gamma. \langle \Gamma, t, \gamma \rangle \Downarrow \langle \Gamma', t', \gamma' \rangle \wedge t' \in R_v(\Gamma' \Vdash \mathbb{A})\} \\ R(\Gamma \Vdash \mathbb{A}) &= \{t \mid \forall \gamma. \Pr(\langle \Gamma, t, \gamma \rangle \Downarrow \\ &\quad (\mathbb{R}_{\geq 0} \times \sum_{\Gamma'} R_v(\Gamma' \Vdash \mathbb{A}) \times \llbracket \Gamma' \rrbracket)) = 1\} \\ R_v(\Gamma \Vdash \mathbb{A}) &= \{\text{return}(v) \mid v \in R_v(\Gamma \Vdash \mathbb{A})\} \\ R_v(\Gamma \Vdash \mathbb{A}) &= \{x \mid (x : \mathbb{A}) \in \Gamma\} \quad \text{for } \mathbb{A} \text{ indecomposable} \\ R_v(\Gamma \Vdash 1) &= \{*\} \\ R_v(\Gamma \Vdash \mathbb{A}_1 \times \mathbb{A}_2) &= \{(v_1, v_2) \mid \forall j. v_j \in R_v(\Gamma \Vdash \mathbb{A}_j)\} \\ R_v(\Gamma \Vdash \sum \mathbb{A}_i) &= \{(i, v) \mid v \in R_v(\Gamma \Vdash \mathbb{A}_i)\} \\ R_v(\Gamma \Vdash T(\mathbb{A})) &= \{\text{thunk}(t) \mid t \in R(\Gamma \Vdash \mathbb{A})\} \\ R_v(\Gamma \Vdash \mathbb{A} \Rightarrow \mathbb{B}) &= \{\lambda x. t \mid \forall \Gamma' \supseteq \Gamma, u \in R_v(\Gamma' \Vdash \mathbb{A}). \\ &\quad t[u/x] \in R(\Gamma' \Vdash \mathbb{B})\} \end{aligned}$$

Induction on the structure of a term $\Gamma, x_1 : \mathbb{A}_1, \dots, x_n : \mathbb{A}_n \Vdash t : \mathbb{B}$ for $z \in \{d, p\}$ now proves that $v_i \in R_v(\Gamma \Vdash \mathbb{A}_i)$ for $i = 1, \dots, n$ implies $t[\vec{v}/\vec{x}] \in R(\Gamma \Vdash \mathbb{B})$. \square

8. Higher-order denotational semantics

This section gives denotational semantics for the higher-order language, without using Assumption 7.1. We are to interpret the new constructs $T(\mathbb{A})$, thunk , and force . We will interpret probabilistic judgements as Kleisli morphisms $\llbracket \Gamma \rrbracket \rightarrow T(\llbracket \mathbb{A} \rrbracket)$ for a certain monad T , and set $\llbracket T(\mathbb{A}) \rrbracket \stackrel{\text{def}}{=} T(\llbracket \mathbb{A} \rrbracket)$, so thunk and force embody the correspondence of maps $\llbracket \Gamma \rrbracket \rightarrow T(\llbracket \mathbb{A} \rrbracket)$ and $\llbracket \Gamma \rrbracket \rightarrow \llbracket T(\mathbb{A}) \rrbracket$.

On which category can the monad T live? Interpreting λ -abstraction and application needs a natural ‘currying’ bijection between morphisms $\llbracket \Gamma \rrbracket \times \mathbb{R} \rightarrow \mathbb{R}$ and morphisms $\llbracket \Gamma \rrbracket \rightarrow \llbracket \mathbb{R} \Rightarrow \mathbb{R} \rrbracket$. But measurable functions cannot do this: it is known that no measurable space $\llbracket \mathbb{R} \Rightarrow \mathbb{R} \rrbracket$ can support such a bijection [3].

We resolve the problem of function spaces by embedding the category of measurable spaces in a larger one, where currying is possible, and that still has the structure to interpret the first order language as before. As the larger category we will take a category of functors $\mathbf{Meas}^{\text{op}} \rightarrow \mathbf{Set}$ from a category \mathbf{Meas} of measurable spaces and measurable functions to the category \mathbf{Set} of sets and functions. This idea arises from two traditions. First, we can think of a variable of type \mathbb{R} as a read-only memory cell, as in the operational semantics, and functor categories have long been used to model local memory (e.g. [23]). Second, the standard construction for building a Cartesian closed category out of a distributive one is based on functor categories (e.g. [28]).

Other models of higher-order programs Semantics of higher-order languages with *discrete* probability are understood well. For terminating programs, there are set-theoretic models based on a distributions monad, and for full recursion one can use probabilistic powerdomains [14] or coherence spaces [9]; our higher-order denotational semantics could be enriched in domains to support recursion. It is also plausible one could model continuous distributions in domain theory, since it supports computable real analysis (e.g. [8]); this could be interesting because computability is subtle for probabilistic programming (e.g. [1]). Nonetheless, we contend it is often helpful to abstract away computability issues when studying probabilistic programming languages, to have access to standard theorems of probability theory to justify program transformations.

8.1 Semantic model

Fix a category \mathbf{Meas} of measurable spaces and measurable functions that is essentially small but large enough for the purposes of Section 4. For example, \mathbf{Meas} could be the category of *standard Borel spaces* [4, 32]: one can show that $\llbracket \mathbb{A} \rrbracket$ is standard Borel by induction on \mathbb{A} , and the class of all standard Borel spaces is countable up to measurable isomorphism.

In Section 4 we interpreted first-order types \mathbb{A} as measurable spaces $\llbracket \mathbb{A} \rrbracket$. We will interpret higher-order types \mathbb{A} as functors $\llbracket \mathbb{A} \rrbracket : \mathbf{Meas}^{\text{op}} \rightarrow \mathbf{Set}$. Informally, when \mathbb{A} is a first-order type and Γ is a first-order context, we will have $\llbracket \mathbb{A} \rrbracket(\llbracket \Gamma \rrbracket) \cong \mathbf{Meas}(\llbracket \Gamma \rrbracket, \llbracket \mathbb{A} \rrbracket) \approx \{t \mid \Gamma \Vdash t : \mathbb{A}\}$. For a second order type $(\mathbb{A} \Rightarrow \mathbb{B})$, we will have

$$\llbracket (\mathbb{A} \Rightarrow \mathbb{B}) \rrbracket(\llbracket \Gamma \rrbracket) \cong \mathbf{Meas}(\llbracket \Gamma \rrbracket \times \llbracket \mathbb{A} \rrbracket, \llbracket \mathbb{B} \rrbracket) \approx \{t \mid \Gamma, x : \mathbb{A} \Vdash t : \mathbb{B}\}$$

so that β/η equality is built in. To put it another way, $\llbracket \mathbb{A} \rrbracket(\mathbb{R}^n)$ models terms of type \mathbb{A} having n read-only real-valued memory cells.

Lemma 8.1. *For a small category \mathbf{C} with countable sums, consider the category $\overline{\mathbf{C}}$ of countable-product-preserving functors $\mathbf{C}^{\text{op}} \rightarrow \mathbf{Set}$, and natural transformations between them.*

- $\overline{\mathbf{C}}$ has all colimits;
- $\overline{\mathbf{C}}$ is Cartesian closed if \mathbf{C} has products that distribute over sums;
- There is a full and faithful embedding $\mathbf{y} : \mathbf{C} \rightarrow \overline{\mathbf{C}}$, given by $\mathbf{y}(c) \stackrel{\text{def}}{=} \overline{\mathbf{C}}(-, c)$, which preserves limits and countable sums.

Proof. See e.g. [28, §7], or [15, Theorems 5.56 and 6.25]. The embedding \mathbf{y} is called the *Yoneda embedding*. \square

For a simple example, consider the category \mathbf{CSet} of countable sets and functions. It has countable sums and finite products, but is not Cartesian closed. Because every countable set is a countable sum of singletons, the category \mathbf{CSet} is equivalent to \mathbf{Set} .

Our semantics for the higher-order language will take place in the category $\overline{\mathbf{Meas}}$. Note that products in $\overline{\mathbf{Meas}}$ are pointwise, e.g. $(F \times G)(X) = F(X) \times G(X)$ for all $F, G \in \overline{\mathbf{Meas}}$ and all $X \in \mathbf{Meas}$, but sums are not pointwise, e.g. $(1 + 1) \in \overline{\mathbf{Meas}}$ is the functor that assigns a measurable space X to the set of its measurable subsets. This is essential for \mathbf{y} to preserve sums.

Distribution types We have to interpret distribution types $P(\mathbb{A})$ in our functor category $\overline{\mathbf{Meas}}$. How can we interpret a probability distribution on the type $\mathbb{R} \Rightarrow \mathbb{R}$? We can answer this pragmatically, without putting σ -algebra structure on the set of all functions. If $\llbracket \mathbb{R} \Rightarrow \mathbb{R} \rrbracket$ were a measurable space, a random variable valued in $\llbracket \mathbb{R} \Rightarrow \mathbb{R} \rrbracket$ would be given by a measurable space (X, Σ_X) , a probability distribution on it, and a measurable function $X \rightarrow \llbracket \mathbb{R} \Rightarrow \mathbb{R} \rrbracket$. Despite there being no such measurable space $\llbracket \mathbb{R} \Rightarrow \mathbb{R} \rrbracket$, we can speak of uncurried measurable functions $X \times \mathbb{R} \rightarrow \mathbb{R}$. Thus we might define a probability distribution on $\llbracket \mathbb{R} \Rightarrow \mathbb{R} \rrbracket$ to be a triple

$$((X, \Sigma_X), f : X \times \mathbb{R} \rightarrow \mathbb{R}, p : \Sigma_X \rightarrow [0, 1])$$

of a measurable space (X, Σ_X) of ‘codes’, a measurable function f where we think of $f(x, r)$ as ‘the function coded x evaluated at r ’, and a probability distribution p on the codes. These triples should be considered modulo renaming the codes. This is exactly the notion of probability distribution that arises in our functor category.

Lemma 8.2. *For a small category \mathbf{C} with countable sums:*

- any functor $F : \mathbf{C} \rightarrow \mathbf{C}$ extends to a functor $\overline{F} : \overline{\mathbf{C}} \rightarrow \overline{\mathbf{C}}$ satisfying $\overline{F} \circ \mathbf{y} \cong \mathbf{y} \circ F$, given by

$$\overline{F}(G)(b) = \left(\sum_a G(a) \times \mathbf{C}(b, F(a)) \right) / \sim$$

where the equivalence relation \sim is the least one satisfying $(a', x, Fg \circ f) \sim (a, Gg(x), f)$;

- similarly, any functor $F: \mathbf{C} \times \mathbf{C} \rightarrow \mathbf{C}$ in two arguments extends to a functor $\bar{F}: \bar{\mathbf{C}} \times \bar{\mathbf{C}} \rightarrow \bar{\mathbf{C}}$, with $\bar{F} \circ (\mathbf{y} \times \mathbf{y}) \cong \mathbf{y} \circ F$:

$$\bar{F}(G, H)(c) = \left(\sum_{a,b} G(a) \times H(b) \times \mathbf{C}(c, F(a, b)) \right) / \sim$$

- any natural transformation $\alpha: F \rightarrow G$ between functors $F, G: \mathbf{C} \rightarrow \mathbf{C}$ lifts to a natural transformation $\bar{\alpha}: \bar{F} \rightarrow \bar{G}$, and similarly for functors $\mathbf{C} \times \mathbf{C} \rightarrow \mathbf{C}$;
- and this is functorial, i.e. $\overline{G \circ F} \cong \bar{G} \circ \bar{F}$ and $\overline{\beta \circ \alpha} = \bar{\beta} \circ \bar{\alpha}$.

Proof. $\bar{F}(G)$ is the left Kan extension of G along F , see e.g. [15]. Direct calculation shows $\bar{F}(G)$ preserves products if G does. \square

Thus the commutative monads P and $T = P(\mathbb{R}_{\geq 0} \times (-))$ on **Meas** lift to commutative monads \bar{P} and $\bar{T} \cong \bar{P}(\mathbf{y}\mathbb{R}_{\geq 0} \times (-))$ on **Meas**. The latter monad captures the informal importance-sampling semantics advocated by the designers of Anglican [34].

8.2 Conservativity

We interpret the types of the higher order language as objects in **Meas** using its categorical structure.

$$\begin{aligned} (\sum_{i \in I} \mathbb{A}_i) &\stackrel{\text{def}}{=} \sum_{i \in I} \langle \mathbb{A}_i \rangle & \langle \mathbb{R} \rangle &\stackrel{\text{def}}{=} \mathbf{y}\mathbb{R} & \langle P(\mathbb{A}) \rangle &\stackrel{\text{def}}{=} \bar{P}(\langle \mathbb{A} \rangle) \\ \langle \mathbb{A} \times \mathbb{B} \rangle &\stackrel{\text{def}}{=} \langle \mathbb{A} \rangle \times \langle \mathbb{B} \rangle & \langle 1 \rangle &\stackrel{\text{def}}{=} 1 & \langle T(\mathbb{A}) \rangle &\stackrel{\text{def}}{=} \bar{T}(\langle \mathbb{A} \rangle) \\ \langle \mathbb{A} \Rightarrow \mathbb{B} \rangle &\stackrel{\text{def}}{=} \langle \mathbb{A} \rangle \Rightarrow \langle \mathbb{B} \rangle \end{aligned}$$

We extend this interpretation to contexts.

$$\langle x_1 : \mathbb{A}_1, \dots, x_n : \mathbb{A}_n \rangle \stackrel{\text{def}}{=} \prod_{i=1}^n \langle \mathbb{A}_i \rangle$$

Deterministic terms $\Gamma \vdash_d t : \mathbb{A}$ are interpreted as natural transformations $\langle \Gamma \rangle \rightarrow \langle \mathbb{A} \rangle$ in **Meas**, and probabilistic terms $\Gamma \vdash_p t : \mathbb{A}$ as natural transformations $\langle \Gamma \rangle \rightarrow \bar{T}(\langle \mathbb{A} \rangle)$, by induction on the structure of terms as in Section 4. Application and λ -abstraction are interpreted as usual in Cartesian closed categories [27]. Thunk and force are trivial from the perspective of the denotational semantics, because $\langle T\mathbb{A} \rangle = T\langle \mathbb{A} \rangle$. To interpret $\text{norm}(t)$, use Lemma 8.2 to extend the normalization functions $T(X) \rightarrow \mathbb{R}_{\geq 0} \times P(X) + 1 + 1$ between measurable spaces (3) to natural transformations $\bar{T}(F) \rightarrow \mathbf{y}(\mathbb{R}_{\geq 0}) \times \bar{P}(F) + 1 + 1$.

Note that all measurable types \mathbb{A} have a natural isomorphism $\langle \mathbb{A} \rangle \cong \mathbf{y}[\mathbb{A}]$. This interpretation conserves the first-order semantics of Section 4:

Proposition 8.3. For $z \in \{d, p\}$, and first-order Γ and \mathbb{A} :

- for first-order $\Gamma \vdash_z t, u : \mathbb{A}$, $\llbracket t \rrbracket = \llbracket u \rrbracket$ if and only if $\langle t \rangle = \langle u \rangle$;
- every term $\Gamma \vdash_z t : \mathbb{A}$ has $\langle t \rangle = \langle u \rangle$ for a first-order $\Gamma \vdash_z u : \mathbb{A}$.

Proof. We treat $z = d$; the other case is similar. By induction on the structure of terms, $\langle t \rangle = \mathbf{y}[\llbracket t \rrbracket] : \langle \Gamma \rangle \rightarrow \langle \mathbb{A} \rangle$. The first point follows from faithfulness of \mathbf{y} ; the second from fullness and (2). \square

One interesting corollary is that the interpretation of a term of first-order type is always a measurable function, even if the term involves thunking, λ -abstraction and application. This corollary gives a partial answer to a question by Park *et al.* on the measurability of all λ -definable ground-type terms in probabilistic programs [25] (partial because our language does not include recursion).

8.3 Soundness

The same recipe as in Section 5.3 will show that the higher-order denotational semantics is sound and adequate with respect to the higher-order operational semantics. This needs Assumption 7.1.

A subtle point is that configuration spaces (6) involve uncountable sums: the set of terms of a given type is uncountable, but \mathbf{y} only preserves countable sums. This is not really a problem because only countably many terms are reachable from a given program.

Definition 8.4. For a type \mathbb{A} , the binary *reachability* relation \rightsquigarrow_d^* on $\{(\Gamma, t) \mid \Gamma \vdash_d t : \mathbb{A} \wedge \Gamma \text{ canonical}\}$ is the least reflexive and transitive relation with $(\Gamma, t) \rightsquigarrow_d^* (\Gamma', u)$ if $(\Gamma, t, \gamma) \rightarrow (\Gamma', u, \gamma')$ for $\gamma \in \llbracket \Gamma \rrbracket$, $\gamma' \in \llbracket \Gamma' \rrbracket$. Similarly, \rightsquigarrow_p^* is the least reflexive and transitive relation on $\{(\Gamma, t) \mid \Gamma \vdash_p t : \mathbb{A} \wedge \Gamma \text{ canonical}\}$ with $(\Gamma, t) \rightsquigarrow_p^* (\Gamma', u)$ if $\text{Pr}(\langle \Gamma, t, \gamma \rangle \rightarrow \langle \Gamma', u \rangle) \neq 0$ for $\gamma \in \llbracket \Gamma \rrbracket$.

Proposition 8.5. Let $z \in \{d, p\}$. For any closed term $t : \mathbb{A}$, the set of reachable terms $\{(\Gamma, u) \mid (\emptyset, t) \rightsquigarrow_z^* (\Gamma, u)\}$ is countable.

Proof. One-step reachability is countable by induction on terms. Since all programs terminate by Proposition 7.2, the reachable terms form a countably branching well-founded tree. \square

We may thus restrict to the configurations built from a countable set \mathcal{U} of terms that is closed under subterms and reachability. Extend the denotational semantics in **Meas** to configurations by defining $s_d : \mathbf{y}(\text{Con}_d(\mathbb{A})) \rightarrow \langle \mathbb{A} \rangle$, $s_p : \mathbf{y}(\text{Con}_p(\mathbb{A})) \rightarrow \bar{T}(\langle \mathbb{A} \rangle)$, and $s_{V_p} : \mathbf{y}(\text{Con}_{V_p}(\mathbb{A})) \rightarrow \bar{T}(\langle \mathbb{A} \rangle)$; use the isomorphisms

$$\mathbf{y}(\text{Con}_z(\mathbb{A})) \cong \sum_{(\Gamma \vdash_z t : \mathbb{A}) \in \mathcal{U}} \langle \Gamma \rangle$$

to define s_d, s_p, s_{V_p} by copairing the interpretation morphisms $\langle \Gamma \vdash_d t : \mathbb{A} \rangle : \langle \Gamma \rangle \rightarrow \langle \mathbb{A} \rangle$ and $\langle \Gamma \vdash_p u : \mathbb{A} \rangle : \langle \Gamma \rangle \rightarrow \bar{T}(\langle \mathbb{A} \rangle)$.

Proposition 8.6 (Soundness). *The following diagrams commute.*

$$\begin{array}{ccc} \mathbf{y}(\text{Con}_{N_d}(\mathbb{A})) & \xrightarrow{s_d} & \mathbf{y}(\text{Con}_{N_p}(\mathbb{A})) \\ \mathbf{y}(\text{reduction}) \downarrow & \searrow s_d & \downarrow \mathbf{y}(\text{reduction}) \\ \mathbf{y}(\text{Con}_d(\mathbb{A})) & \xrightarrow{s_d} & \mathbf{y}(\text{Con}_p(\mathbb{A})) \end{array} \quad \begin{array}{ccc} & & \xrightarrow{s_p} \\ & & \bar{T}(\langle \mathbb{A} \rangle) \\ & \nearrow \bar{\mu} \cdot \bar{T}(s_p) & \end{array}$$

Adequacy It follows that the higher denotational semantics remains adequate, in the sense that for all probabilistic terms $t : \mathbb{A}$,

$$\langle t \rangle_1(*) = (\bar{T}(s_{V_p}))_1(\text{Pr}(\langle \emptyset, t, * \rangle \Downarrow (-))).$$

Adequacy is usually only stated for first-order types. At first-order types \mathbb{A} the function s_{V_p} does very little, since global elements of $\langle \mathbb{A} \rangle$ correspond bijectively with value configurations modulo weakening, contraction and exchange in the context. At higher types, the corollary still holds, but s_{V_p} is not so trivial because we do not reduce under thunk or λ .

9. Continuous densities

In most examples, the argument t to $\text{score}(t)$ is a density function for a probability distribution. When scores are based on density functions, this makes the relationship between score and sample tighter than we have expressed so far. Our language easily extends to accommodate such distributions. We just add a collection *density* types to the syntax.

$$\mathbb{D} ::= \mathbb{R} \mid \text{bool} \mid \mathbb{N} \mid 1 \mid \mathbb{D} \times \mathbb{D} \quad \mathbb{A} ::= \dots \mid \mathbb{D}(\mathbb{D})$$

The \mathbb{D} in this grammar denotes a measurable space that: (i) carries a separable metrisable topology that generates the σ -algebra; and that (ii) comes with a chosen σ -finite measure $\mu_{\mathbb{D}}$. An example is \mathbb{R} with its usual Euclidean topology and the Lebesgue measure (which maps each interval to its size).

The type $\mathbb{D}(\mathbb{D})$ denotes a measurable space $\llbracket \mathbb{D}(\mathbb{D}) \rrbracket$ of continuous functions $f : \llbracket \mathbb{D} \rrbracket \rightarrow \mathbb{R}_{\geq 0}$ with $\int_{\mathbb{D}} f d\mu_{\mathbb{D}} = 1$. The σ -algebra of $\llbracket \mathbb{D}(\mathbb{D}) \rrbracket$ is the least one making $\{f \mid f(x) \leq r\}$ measurable for all $(x, r) \in \llbracket \mathbb{D} \rrbracket \times \mathbb{R}_{\geq 0}$.

A density type $D(\mathbb{D})$ comes with two measurable functions

$$\begin{aligned} ev : [D(\mathbb{D})] \times [\mathbb{D}] &\rightarrow \mathbb{R}_{\geq 0} & ev(f, x) &= f(x) \\ dist : [D(\mathbb{D})] &\rightarrow [P(\mathbb{D})] & dist(f)(U) &= \int_U f \, d\mu_{\mathbb{D}} \end{aligned}$$

Note that measurability of ev relies on continuity of the densities [3]. The usual way of imposing a soft constraint based on a likelihood can be encoded in our first-order language as $\text{score}(ev(f, x))$, where the datum x is observed with a probability distribution with density f . Thus the categorical machinery used to interpret higher-order functions is not needed for such soft constraints.

There is a limit to this use of continuity: probability measures produced by $\text{norm}(t)$ need not have continuous density. For example, $\text{norm}(\text{return}(42.0))$ produces a discontinuous Dirac measure.

Probability densities are often used by *importance samplers*. The importance sampler generates samples of $f \in [D(\mathbb{R})]$ by first sampling from a proposal distribution g where sampling is easy, and then normalizing those samples x from g according to their importance weight $f(x)/g(x)$. This works provided the support of g is \mathbb{R} , e.g. $g = \text{density_gauss}(x, (0.0, 1.0))$.

$$\begin{aligned} &[\text{norm}(\text{sample}(dist(f)))] = \\ &[\text{norm}(\text{let } x = \text{sample}(dist(g)) \text{ in } \text{score}(\frac{ev(f, x)}{ev(g, x)}); \text{return}(x))] \end{aligned}$$

Density types can be incorporated into the higher order language straightforwardly. The only subtlety is that denotational semantics now needs the base category to contain $[D(\mathbb{D})]$.

10. Conclusion and future work

We have defined a metalanguage for higher-order probabilistic programs with continuous distributions and soft constraints, and presented operational and denotational semantics, together with useful program equations justified by the semantics. One interesting next step is to use these tools to study other old or new language features and concepts (such as recursion, function memoisation [30], measure-zero conditioning [5], disintegration [2, 31], and exchangeability [10, 20, 34]) that have been experimented with in the context of probabilistic programming. Another future direction is to formulate and prove the correctness of inference algorithms, especially those based on Monte Carlo simulation, following [13].

Acknowledgments

We thank T. Avery, I. Garnier, T. Le, K. Sturtz, and A. Westerbaan. This work was supported by the EPSRC, a Royal Society University Fellowship, An Institute for Information & communications Technology Promotion (IITP) grant funded by the Korea government (MSIP, No. R0190-15-2011), DARPA PPAML, and the ERC grant ‘causality and symmetry — the next-generation semantics’.

References

- [1] N. L. Ackerman, C. E. Freer, and D. M. Roy. Noncomputable conditional distributions. In *LiCS*, 2011.
- [2] N. L. Ackerman, C. E. Freer, and D. M. Roy. On computability and disintegration, 2015. URL <http://arxiv.org/abs/1509.02992>.
- [3] R. J. Aumann. Borel structures for function spaces. *Illinois Journal of Mathematics*, 5:614–630, 1961.
- [4] S. K. Berberian. *Borel spaces*. World Scientific, 1988.
- [5] J. Borgström, A. D. Gordon, M. Greenberg, J. Margetson, and J. van Gael. Measure transformer semantics for Bayesian machine learning. *LMCS*, 9(3):11, 2013.
- [6] E.-E. Doberkat. *Stochastic Relations: Foundations for Markov Transition Systems*. Chapman & Hall, 2007.
- [7] A. Doucet, N. de Freitas, and N. Gordon. *Sequential Monte Carlo Methods in Practice*. Springer, 2001.
- [8] A. Edalat and M. H. Escardó. Integration in Real PCF. *Inf. Comput.*, 160:128–166, 2000.
- [9] T. Ehrhard, C. Tasson, and M. Pagani. Probabilistic coherence spaces are fully abstract for probabilistic PCF. In *POPL*, 2014.
- [10] C. E. Freer and D. M. Roy. Computable de Finetti measures. *Ann. Pure Appl. Logic*, 163(5):530–546, 2012.
- [11] M. Giry. A categorical approach to probability theory. *Categorical Aspects of Topology and Analysis*, 915:68–85, 1982.
- [12] N. Goodman, V. Mansinghka, D. M. Roy, K. Bonawitz, and J. B. Tenenbaum. Church: a language for generative models. In *UAI*, 2008.
- [13] C. Hur, A. V. Nori, S. K. Rajamani, and S. Samuel. A provably correct sampler for probabilistic programs. In *FSTTCS*, 2015.
- [14] C. Jones and G. D. Plotkin. A probabilistic powerdomain of evaluations. In *LiCS*, 1989.
- [15] G. M. Kelly. *Basic concepts of enriched category theory*. CUP, 1982.
- [16] A. Kock. Monads on symmetric monoidal closed categories. *Archiv der Mathematik*, XXI:1–10, 1970.
- [17] D. Kozen. Semantics of probabilistic programs. *Journal of Computer and System Sciences*, 22:328–350, 1981.
- [18] P. B. Levy. Call-by-push-value: A subsuming paradigm. In *TLCA’02*.
- [19] P. B. Levy, J. Power, and H. Thielecke. Modelling environments in call-by-value programming languages. *Inf. Comput.*, 185(2), 2003.
- [20] V. K. Mansinghka, D. Selsam, and Y. N. Perov. Venture: a higher-order probabilistic programming platform with programmable inference. 2014. URL <http://arxiv.org/abs/1404.0099>.
- [21] T. Minka, J. Winn, J. Guiver, and D. Knowles. Infer.NET 2.4, Microsoft Research Cambridge, 2010.
- [22] E. Moggi. Notions of computation and monads. *Inf. Comput.*, 93(1): 55–92, 1991.
- [23] F. J. Oles. Type algebras, functor categories and block structure. In *Algebraic methods in semantics*. CUP, 1984.
- [24] B. Paige and F. Wood. A compilation target for probabilistic programming languages. In *ICML*, 2014.
- [25] S. Park, F. Pfenning, and S. Thrun. A probabilistic language based on sampling functions. *ACM TOPLAS*, 31(1):171–182, 2008.
- [26] A. Patil, D. Huard, and C. J. Fonnesbeck. PyMC: Bayesian stochastic modelling in Python. *Journal of Statistical Software*, 35, 2010.
- [27] A. M. Pitts. Categorical logic. In *Handbook of Logic in Computer Science*, volume 5. OUP, 2000.
- [28] J. Power. Generic models for computational effects. *TCS*, 364(2): 254–269, 2006.
- [29] N. Ramsey and A. Pfeffer. Stochastic lambda calculus and monads of probability distributions. In *POPL*, 2002.
- [30] D. M. Roy, V. Mansinghka, N. Goodman, and J. Tenenbaum. A stochastic programming perspective on nonparametric Bayes. In *ICML Workshop on Nonparametric Bayesian*, 2008.
- [31] C.-C. Shan and N. Ramsey. Symbolic Bayesian inference by lazy partial evaluation, 2016.
- [32] S. M. Srivastava. *A course on Borel sets*. Springer, 1998.
- [33] Stan Development Team. Stan: A C++ library for probability and sampling, version 2.5.0, 2014. URL <http://mc-stan.org/>.
- [34] F. Wood, J. W. van de Meent, and V. Mansinghka. A new approach to probabilistic programming inference. In *AISTATS*, 2014.
- [35] H. Yang. Program transformation for probabilistic programs, 2015. Presentation at DALI.