

Conditional probability by lazy partial evaluation

Chung-chieh Shan

Indiana University
ccshan@indiana.edu

Abstract

This is the text of the abstract.

Categories and Subject Descriptors G.3 [Probability and Statistics]: distribution functions

Keywords keyword1, keyword2

1. Introduction

Probability distributions are a popular way to model and handle uncertainty. In particular, the typical *Bayesian* reasoner begins with a *prior* probability distribution on all possible worlds, *observes* the actual world, then *conditions* the prior distribution to obtain a *posterior* distribution on those possible worlds that match the observations. It is common and convenient to specify a distribution by composing a *generative story*, which is a procedure that picks a world randomly, usually by modeling the relevant aspects of how the world comes to be.

TODO our contributions

1.1 Background example

We illustrate the starting point of this paper with an example. Suppose we observed the players of a one-on-one game, perhaps to match them up to make future games more fun. For simplicity, suppose we saw just one game, in which the player Alice beat the player Bob. This outcome may be due to Alice’s skill level (whatever it means) being higher than Bob’s, or due to Alice being lucky in this particular game (whatever it means). Regardless, we can ask how Alice’s skill level compares to Bob’s, given that Alice beat Bob. To model the situation, we can write the following generative story:

```
do {a ← normal 10 3;           (1)
    b ← normal 10 3;
    l ← normal 0 2;
    let true = (l < a - b);
    return (a, b)}
```

The first line of this program means to pick a real number a from the normal distribution with mean 10 and standard deviation 3. The second line picks b from the same distribution. These numbers model the skill levels of Alice and Bob. The third line picks l , which models how much luck Bob has over Alice in this game. These

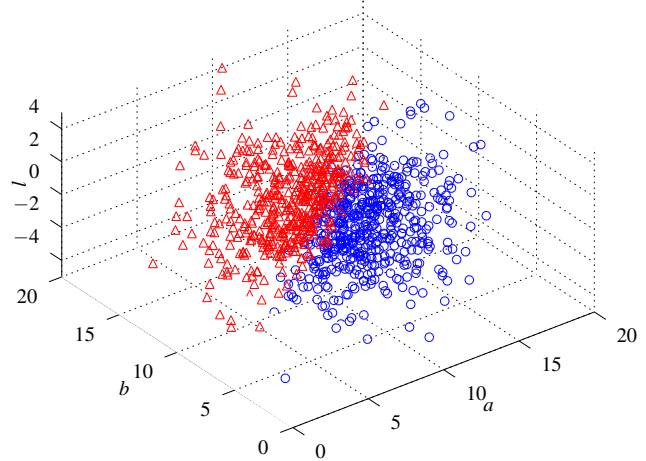


Figure 1. An approximate prior distribution: 1000 random samples generated by the first three lines of (1) (one line for each dimension). The red triangles are samples rejected by the fourth line of (1). The blue circles are samples accepted by that line.

three lines together define a *prior* distribution on \mathbb{R}^3 , in which each point (a, b, l) is a possible world. Figure 1 depicts this distribution approximately, as a cloud of points (both red triangles and blue circles). The fourth line performs a non-exhaustive pattern-match to restrict the distribution to the part where $l < a - b$. This step incorporates our observation that Alice beat Bob. The first four lines together define the *posterior* distribution on \mathbb{R}^3 , shown as blue circles in Figure 1. The last line projects this distribution from \mathbb{R}^3 to \mathbb{R}^2 .

A simple way to answer a question about the distribution is to interpret this program as a *sampler*—that is, a probabilistic algorithm that generates a random point. For instance, it is well-studied how to generate a pseudorandom number from a normal distribution. We run this sampler 1000 times, say (as shown in Figure 1). These resulting points approximate our uncertainty about Alice and Bob’s skill levels. For example, as shown in Figure 2, to estimate the probability that Alice is more skilled than Bob, we can compute the proportion of resulting points where $a > b$.

To make this estimation method more accurate, we can change the program (1) to one that denotes the same distribution but makes fewer random choices. For example, instead of picking l then testing $l < a - b$, the following equivalent program generate samples with uneven *importance weights*:

```
do {a ← normal 10 3;           (2)
    b ← normal 10 3;
    factor  $\frac{1 + \operatorname{erf}((a - b)/(2 \times \sqrt{2}))}{2}$ ;
    return (a, b)}
```

[Copyright notice will appear here once ‘preprint’ option is removed.]

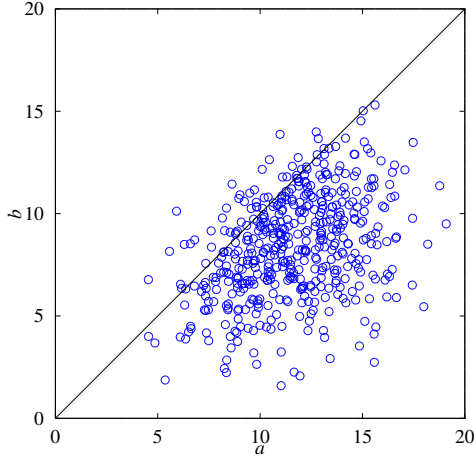


Figure 2. A boolean query on an approximate posterior distribution: fewer than 1000 random samples generated by (1). The samples to the right of the diagonal line are where $a > b$.

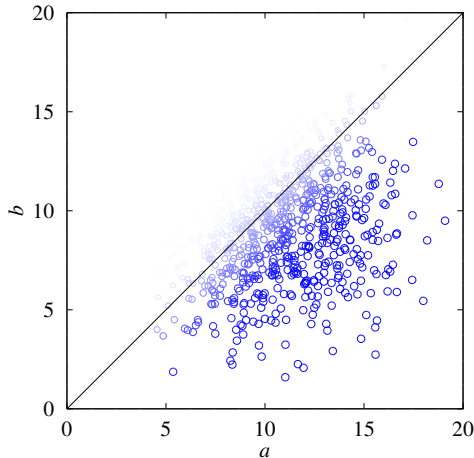


Figure 3. A boolean query on a better approximate posterior distribution: 1000 weighted random samples generated by (2). The size and darkness of each circle reflect its weight. As in Figure 2, the samples to the right of the diagonal line are where $a > b$.

Interpreted as a sampler, this program does not pick any concrete l , but rather attaches the weight $\frac{1+\text{erf}((a-b)/(2\times\sqrt{2}))}{2}$ (instead of the default weight 1) to each generated point (a, b) . This weight is the probability that, had we chosen a concrete l as in (1), we would have chosen an l that passes the test $l < a - b$. We can think of the weight as softening the rejection due to a failed pattern-match. As shown in Figure 3, to estimate the probability that Alice is more skilled than Bob, we should divide the total weight of points where $a > b$ by the total weight of all points. This estimate varies less from run to run than the estimate from (1) does, because this program makes fewer random choices while denoting the same distribution. That is good.

We wait until Section 3 to detail how expressions denote distributions. The gist of the present example is that we turned a prior distribution expressed as a program (namely the first three lines of (1)) into a posterior distribution expressed as a more efficient program (namely (2)), in two steps:

1. We expressed an observation as a non-exhaustive pattern-match (the fourth line of (1)).
2. We justified an optimization by denotation equality (between the third and fourth lines of (1) and the third line of (2)).

1.2 The need to generalize and mechanize conditioning

In general, just as it is useful to express any distribution as a sampler program, it is useful to express a posterior distribution as a sampler program such as (1) or (2). There are a variety of reasons:

1. A human can then read the code and understand the posterior as a distribution in its own right, then reason about it equationally.
2. A larger program can include the posterior. Posteriors are an essential part of many probabilistic inference algorithms (such as *Monte Carlo Markov chain* methods [13]) as well as probabilistic models (for reasoning about reasoning [17]).
3. A program transformation can generate the posterior from the prior. For example, the fourth line of (1) could have been added automatically.
4. A human or a machine such as a compiler pass can simplify and optimize the posterior. For example, we are building a probabilistic programming system Hakaru that can turn the third and fourth lines of (1) into the third line of (2).
5. A robot that receives a stream of observations as the world changes over time can update the posterior to reflect the evolution and uncertainty of its knowledge [10].

Before we can express the posterior distribution as a program, we need to specify what it is, so as to inform the design and justify the correctness of each of the applications listed above. In other words, we need to specify *conditioning*, the mathematical operation that turns the prior into the posterior. Unfortunately, the specification of conditioning in current probabilistic programming languages assumes that the condition has non-zero probability. For example, the probability of $l < a - b$ in (1) and Figure 1 is $1/2$. If this probability were zero, then we would get the empty distribution with no samples, which is unworthy of being called a posterior.

In practice, it is very common for the condition to have zero probability. TODO: Motivating example of zero-probability conditioning from July 2014 talk.

1.3 Our contributions

We define a probabilistic language whose types are measurable spaces that can be uncountable and whose terms can be interpreted simultaneously as samplers, measures, and functionals.

We generalize the specification of conditioning to when the condition has zero probability, by adapting to programming languages the notion of *disintegration* advertised by Chang and Pollard [5].

We implement this specification as a program transformation that generalizes both Bhat et al.’s density calculator [1, 2] and Fischer et al.’s sharing-preserving lazy partial evaluator [8].

2. Samplers, measures, and functionals

For its motivation, intuition, and correctness, our work relies on a three-way correspondence between *samplers*, *measures*, and *functionals*. Thus, we introduce these concepts and detail the correspondence in this section.

2.1 Samplers

A *sampler* is a program that uses randomness to produce a result. We can think of each run of the program as an experiment, and the result of the run is the *outcome* of the experiment. A typical sampler is composed of building blocks that are primitive samplers. For

example, one primitive sampler might be to choose a real number uniformly at random between 0 and 1. In our language, we can express this sampler by the primitive expression **random**. One way to flip a fair coin (that is, to choose between two outcomes with equal probability) is to perform **random** twice then see which result is bigger. Because our language is monadic, we can express this composed sampler by the expression

$$\mathbf{do} \{x \leftarrow \mathbf{random}; y \leftarrow \mathbf{random}; \mathbf{return} (x < y)\}. \quad (3)$$

Here **return** is the monad unit operation (`return` in Haskell) and $\mathbf{do} \{x \leftarrow \dots; \dots\}$ is the monad bind operation (`>>=` in Haskell). The type of **random** is $\mathbb{M}\mathbb{R}$ (where \mathbb{M} is the monad type constructor), because **random** is not a real number but rather a sampler that produces a real number. The type of (3) is $\mathbb{M}(\mathbb{1} + \mathbb{1})$, because our numeric comparison $<$ returns the sum of unit types $\mathbb{1} + \mathbb{1}$ (the value `inl ()` for **true** and the value `inr ()` for **false**).

To optimize and reason about programs, we want to consider many samplers equivalent that do not describe exactly the same procedure. For example, the following are two other perfectly correct ways to flip a fair coin in our language:

$$\mathbf{do} \{y \leftarrow \mathbf{random}; x \leftarrow \mathbf{random}; \mathbf{return} (x < y)\} \quad (4)$$

$$\mathbf{do} \{x \leftarrow \mathbf{random}; \mathbf{return} (x < 1/2)\} \quad (5)$$

If we are running these programs using a pseudo-random number generator, then we may well want to optimize the first two ways to the third way, despite (or due to) the third way using less randomness. Such an optimization is justified by the fact that the three ways denote the same *measure*, even though not the same sampler [15].

2.2 Measures

A *measure* is a mathematical function that maps sets to non-negative real numbers. In the slightly unfortunate standard terminology, the measure is said to *measure* sets and return their *measures*.

1. For example, **random** in our language denotes the *uniform probability measure on (0,1)*: given an interval (a,b) of real numbers, it returns the length of the intersection of the intervals $(0,1)$ and (a,b) . Thus the measure of the interval $(2/3,2)$ is $1/3$. This matches the fact that, a third of the time, choosing a real number uniformly at random between 0 and 1 produces a number between $2/3$ and 2.

A different measure may well give a different result from measuring the same set.

2. For example, **return** x denotes the *Dirac measure at x*: given the same set S to be measured, whether S is an interval (a,b) , it returns 1 if $x \in S$, and 0 otherwise. [TODO: Use x only for a variable and a,b only for atomic terms?]

The general idea is that each sampler corresponds to the measure that, given a set of outcomes, tells the *probability* that running the sampler will produce an outcome in the given set. The deterministic sampler **return** x always produces the same outcome x , so the probability that the outcome is in the measured set is either 1 or 0.

To avoid pathological cases, a measure is not required to measure every set of real numbers. Rather, we introduce the notion of a *measurable space*. A measurable space is a set A , equipped with a notion of what subsets of A can be measured. The complement of a measurable set and the union of a countable collection of measurable sets must be measurable. In particular, the empty set $\{\}$ and the full set A must both be measurable, because the empty set is the nullary union and the full set is the complement of the empty set.

Each type in our language is a measurable space. For example, the type $\mathbb{1}$ in our language is the singleton set $\{\}$, equipped with

the notion that the empty set $\{\}$ and the full set $\{\{\}\}$ are both measurable as required. Another base type \mathbb{R} in our language is the set of real numbers, equipped with the notion that every interval can be measured—so every set of real numbers made from intervals by complement and countable union must also be measurable.

It is easy to add other measurable spaces as base types. For example, we could add the measurable space of integers. Yet another useful measurable space is $[0, \infty]$, the set of non-negative real numbers augmented with positive infinity, again equipped with the notion that every interval can be measured. We don't need this type in our language, but it is an essential measurable space for semantics.

Mathematically speaking, then, a measurable space is a pair

$$\alpha = (\text{set}(\alpha), \text{measurable}(\alpha)), \quad (6)$$

where $\text{set}(\alpha)$ is a set and $\text{measurable}(\alpha)$ is a set of subsets of $\text{set}(\alpha)$ that is closed under complement and countable union. A measure μ on α is a function from $\text{measurable}(\alpha)$ to $[0, \infty]$, such that

$$\mu(S_1 \cup S_2 \cup \dots) = \mu(S_1) + \mu(S_2) + \dots \quad (7)$$

for every countable collection $S_1, S_2, \dots \in \text{measurable}(\alpha)$ of pairwise-disjoint measurable sets. (In particular, the empty collection forces $\mu(\{\})$ to be 0.) We sometimes write a measurable space to mean its underlying set. For example, we write $3 \in [0, \infty]$ to mean $3 \in \text{set}([0, \infty])$.

There is a category of measurable spaces [9]. Its morphisms are the *measurable functions*. A measurable function is a function such that the inverse image of every measurable set is measurable. Formally, given two measurable spaces α and β , a *measurable function* $f \in \alpha \rightarrow \beta$ is a function $f \in \text{set}(\alpha) \rightarrow \text{set}(\beta)$ such that $f^{-1}(T) \in \text{measurable}(\alpha)$ for every $T \in \text{measurable}(\beta)$. It is a standard exercise to prove that, if $S \in \text{measurable}(\alpha)$, then the indicator function $S_* \in \alpha \rightarrow [0, \infty]$ defined by

$$S_*(x) = \begin{cases} 1 & \text{if } x \in S \\ 0 & \text{if } x \notin S \end{cases} \quad (8)$$

is in fact measurable. Another useful exercise is to prove that, if $c_1, c_2, \dots \in \alpha \rightarrow [0, \infty]$ is a (countable) sequence of measurable functions, and the sequence is *increasing* in the sense that $c_i(x) \leq c_j(x)$ for all $x \in \alpha$ whenever $i \leq j$, then the pointwise limit function

$$\lambda x. \lim_{n \rightarrow \infty} c_n(x) \in \alpha \rightarrow [0, \infty] \quad (9)$$

is again measurable. [TODO: Don't use the letter c for both integers and metalanguage continuations.]

Three type constructors \times , $+$, and \mathbb{M} in our language build bigger measurable spaces out of smaller ones, as follows.

Given two measurable spaces α and β , we can take their product as well as disjoint union. The product $\alpha \times \beta$ is the Cartesian product of sets $\text{set}(\alpha) \times \text{set}(\beta)$, equipped with the notion that any Cartesian product $S \times T$ of two measurable sets $S \in \text{measurable}(\alpha)$ and $T \in \text{measurable}(\beta)$ can be measured (as is every set made from those Cartesian products by complement and countable union). Thus,

$$(x, y) \in \alpha \times \beta \quad \text{if} \quad x \in \alpha \quad \text{and} \quad y \in \beta. \quad (10)$$

The disjoint union $\alpha + \beta$ is the disjoint union of sets $\text{set}(\alpha) + \text{set}(\beta)$, equipped with the notion that any disjoint union $S + T$ of two measurable sets $S \in \text{measurable}(\alpha)$ and $T \in \text{measurable}(\beta)$ can be measured. As alluded to above, we write

$$\mathbf{inl} x \in \alpha + \beta \quad \text{if} \quad x \in \alpha, \quad (11)$$

$$\mathbf{inr} y \in \alpha + \beta \quad \text{if} \quad y \in \beta; \quad (12)$$

and we abbreviate $\mathbf{inl} ()$ as **true** and $\mathbf{inr} ()$ as **false**.

Given a measurable space α , it turns out we can turn the set of measures on it into a measurable space $\mathbb{M}\alpha$. To do so, we

equip the set with the notion that, for any two measurable sets $S \in \text{measurable}(\alpha)$ and $T \in \text{measurable}([0, \infty])$, the set of measures

$$\{\mu \mid \mu(S) \in T\} \quad (13)$$

is measurable. This construction \mathbb{M} is a monad on the category of measurable spaces [9]. A measurable function from α to $\mathbb{M}\beta$ (in other words, a morphism from α to β in the Kleisli category) is called a *kernel* from α to β .

2.3 Functionals

We have just described the measure built by the monad unit operation **return**, but not the monad bind operation. We need the monad bind operation not only to define the denotational semantics of our language, but also to define conditioning. To describe the monad bind operation, we turn to a certain class of *functionals* that are in one-to-one correspondence with measures.

A functional is just a higher-order function. Given a measure μ on a measurable space α , let's consider the functional μ^* that takes any function c from α and integrates it. As long as the integrand c is a *measurable* function from α to $[0, \infty]$ (that is, $c \in \alpha \rightarrow [0, \infty]$), we can define this integral $\mu^*(c) \in [0, \infty]$. This definition of integration is called *Lebesgue integration*. Its basic idea is to slice the integrand horizontally into sets measurable in α , then take the limit as the height of each slice approaches zero and the total height of all slices approaches infinity. [TODO: Slicing picture, from Wikipedia?]

Thus, to each measure $\mu \in \mathbb{M}\alpha$ corresponds an integration functional μ^* , which is a function from $\alpha \rightarrow [0, \infty]$ to $[0, \infty]$. This functional μ^* enjoys three important properties:

1. It extends μ , in that $\mu^*(S_*) = \mu(S)$ for all $S \in \text{measurable}(\alpha)$.
2. It is *linear*, in that

$$\mu^*(\lambda x. r \times c(x)) = r \times \mu^*(c) \quad (14)$$

$$\mu^*(\lambda x. c(x) + c'(x)) = \mu^*(c) + \mu^*(c') \quad (15)$$

for all $c, c' \in \alpha \rightarrow [0, \infty]$ and $r \in [0, \infty]$.

3. It satisfies *monotone convergence*, in that

$$\mu^*(\lambda x. \lim_{n \rightarrow \infty} c_n(x)) = \lim_{n \rightarrow \infty} \mu^*(c_n) \quad (16)$$

for all increasing $c_1, c_2, \dots \in \alpha \rightarrow [0, \infty]$.

Moreover, it turns out that μ^* is the *unique* functional with these properties. In other words, there is a one-to-one correspondence between measures μ and linear functionals satisfying monotone convergence μ^* . In fact, many useful measures are most easily defined by specifying their corresponding functionals—for example,

1. The uniform probability measure on $(0, 1)$ is the measure $\mu \in \mathbb{M}\mathbb{R}$ such that $\mu^*(c) = \int_0^1 c(x) dx$.
2. Given an element $x \in \alpha$, the Dirac measure at x is the measure $\mu \in \mathbb{M}\alpha$ such that $\mu^*(c) = c(x)$.
3. Given a measure $\mu \in \mathbb{M}\alpha$ and a kernel $\nu \in \alpha \rightarrow \mathbb{M}\beta$, the monad bind operation produces the measure $\xi \in \mathbb{M}\beta$ such that $\xi^*(c) = \mu^*(\lambda x. (\nu x)^*(c))$.

Therefore, de Finetti [7] and Pollard [16] advocate omitting the stars altogether. In other words, they advocate identifying each set $S \in \text{measurable}(\alpha)$ with its indicator function $S_* \in \alpha \rightarrow [0, \infty]$, and each measure $\mu \in \mathbb{M}\alpha$ with its integration functional $\mu^* \in (\alpha \rightarrow [0, \infty]) \rightarrow [0, \infty]$. Under this proposal, the monad operations are exactly those of the continuation monad.

Without going that far, we introduce two pieces of notation that make applying the correspondence more concise. First, just as we already abbreviate “the function f such that $f(x) = \dots$ ” to “ $\lambda x.$ ”, we abbreviate “the measure μ such that $\mu^*(c) = \dots$ ” to “ $\lambda^*c.$ ”. Second [TODO: drop this notation?], we write $\mu^*(c)$ as $\mu \star c$, so

that we can write $\mu^*(\lambda x. \dots)$ as $\mu \star \lambda x. \dots$ without parentheses. This way, we can write that

1. the uniform probability measure on $(0, 1)$ is $\lambda^*c. \int_0^1 c(x) dx$,
2. monad unit produces the measure $\lambda^*c. c(x)$, and
3. monad bind produces the measure $\lambda^*c. \mu \star \lambda x. \nu x \star c$.

2.4

For conditioning and density, it's useful to generalize from probability measures to all measures (actually, to sigma-finite measures/kernels), and from samplers to importance samplers.

Return to motivating examples.

State precisely the three sets, notions of equivalence, and correspondences preserving said notions of equivalence.

2.5 What is disintegration?

Given $\mu \in \mathbb{M}\alpha$ and $\xi \in \mathbb{M}(\alpha \times \beta)$, we say that a kernel $\nu \in \alpha \rightarrow \mathbb{M}\beta$ is a *disintegration* of ξ with respect to μ iff

$$\xi = \lambda^*c. \mu \star \lambda t. \nu t \star \lambda y. c(t, y). \quad (17)$$

Often, as in the few examples below, μ is the Lebesgue measure

$$\Lambda = \lambda^*c. \int_{-\infty}^{\infty} c(t) dt. \quad (18)$$

For example [as discussed on 2015-04-15], suppose α is \mathbb{R} , μ is the Lebesgue measure Λ , and ξ is the measure

$$\xi = \lambda^*c. \int_4^7 \frac{c(\sin y, y)}{3} dy, \quad (19)$$

which by the way is the denotation of the program

$$\text{do } \{y \leftarrow \text{uniform } 4 \ 7; \text{ return } (\sin y, y)\}. \quad (20)$$

We let $t = \sin y$ and change the integration variable from y to t . We have $|dt/dy| = |\cos y| = \sqrt{1-t^2}$. Solving for y in terms of t gives the countably infinite number of solutions

$$y = 2 \times \pi \times n + \arcsin t \quad \text{or} \quad y = 2 \times \pi \times n + \pi - \arcsin t \quad (21)$$

where $n \in \mathbb{Z}$. Hence

$$\xi = \lambda^*c. \sum_{n \in \mathbb{Z}} \int_{-1}^1 \sum_{y \in \{2 \times \pi \times n + \arcsin t, 2 \times \pi \times n + \pi - \arcsin t\}, 4 < y < 7} \frac{c(t, y)}{3 \times \sqrt{1-t^2}} dt \quad (22)$$

so by Tonelli's theorem

$$\xi = \lambda^*c. \int_{-1}^1 \sum_{n \in \mathbb{Z}} \sum_{y \in \{2 \times \pi \times n + \arcsin t, 2 \times \pi \times n + \pi - \arcsin t\}, 4 < y < 7} \frac{c(t, y)}{3 \times \sqrt{1-t^2}} dt. \quad (23)$$

Comparing this against (17) shows that the kernel

$$\nu = \lambda t. \lambda^*c'. \text{if } -1 < t < 1$$

$$\text{then } \sum_{n \in \mathbb{Z}} \sum_{y \in \{2 \times \pi \times n + \arcsin t, 2 \times \pi \times n + \pi - \arcsin t\}, 4 < y < 7} \frac{c'(y)}{3 \times \sqrt{1-t^2}} \quad (24)$$

else 0

is a disintegration of ξ (with respect to Λ).

To take another example [as discussed on 2015-04-14], suppose α is \mathbb{R} , μ is the Lebesgue measure Λ , and ξ is some measure of the form

$$\xi = \lambda^*c. \iint w(x, y) \times c(h_1(x, y) + h_2(x, x, y)) dy dx \in \mathbb{M}(\mathbb{R} \times \beta_1 \times \beta_2) \quad (25)$$

in which the integrals over x and y are some measures on β_1 and β_2 , so $\beta = \beta_1 \times \beta_2$. In the expression $h_1(x, y) + h_2(x)$, note that the second term does not depend on y . Suppose now that we have found a family of kernels $v_x \in \mathbb{R} \rightarrow \mathbb{M} \beta_2$ ranging over $x \in \beta_1$, such that for each x , the kernel v_x is a disintegration of

$$\xi_x = \lambda^* c_x. \int w(x, y) \times c_x(h_1(x, y), y) dy \in \mathbb{M}(\mathbb{R} \times \beta_2) \quad (26)$$

with respect to Λ . In other words, suppose now we have

$$\int w(x, y) \times c_x(h_1(x, y), y) dy = \int_{-\infty}^{\infty} v_x t_1 \star \lambda y. c_x(t_1, y) dt_1 \quad (27)$$

for each $x \in \beta_1$ and $c_x \in (\mathbb{R} \times \beta_2) \rightarrow [0, \infty]$. Then for any $x \in \beta_1$ and $c \in (\mathbb{R} \times \beta_1 \times \beta_2) \rightarrow [0, \infty]$, we can let

$$c_x = \lambda(t_1, y). c(t_1 + h_2(x), x, y), \quad (28)$$

so by (27) we have

$$\begin{aligned} \int w(x, y) \times c(h_1(x, y) + h_2(x), x, y) dy \\ = \int_{-\infty}^{\infty} v_x t_1 \star \lambda y. c(t_1 + h_2(x), x, y) dt_1. \end{aligned} \quad (29)$$

Applying this equation in (25) gives

$$\xi \star c = \int \int_{-\infty}^{\infty} v_x t_1 \star \lambda y. c(t_1 + h_2(x), x, y) dt_1 dx. \quad (30)$$

We let $t = t_1 + h_2(x)$ and change the inner integration variable from t_1 to t :

$$\xi \star c = \int \int_{-\infty}^{\infty} v_x(t - h_2(x)) \star \lambda y. c(t, x, y) dt dx. \quad (31)$$

If we could apply Tonelli's theorem (in particular, if we knew the measure on β_1 to be σ -finite), then we would have

$$\xi \star c = \int_{-\infty}^{\infty} \int v_x(t - h_2(x)) \star \lambda y. c(t, x, y) dx dt \quad (32)$$

and so would be able to set

$$v = \lambda t. \lambda^* c'. \int v_x(t - h_2(x)) \star \lambda y. c'(x, y) dx. \quad (33)$$

3. Syntax and semantics

Figure 4 defines the syntax and type system of our language.

If $\Gamma \vdash e : \beta$, where Γ is the type environment $x_1 : \alpha_1, \dots, x_n : \alpha_n$, then the denotation $\llbracket e \rrbracket$ of e is a measurable function [TODO: Emphasize measurability. What about σ -finiteness?] from the measurable space $\prod \Gamma = \alpha_1 \times \dots \times \alpha_n$ to the measurable space β . Figure 5 specifies this denotation $\llbracket e \rrbracket$ by induction on e . We treat each element of $\prod \Gamma$ as a function that maps each variable name x_i to an element of the corresponding space α_i . [TODO: Explain syntax and semantics in tandem?]

(If $\Gamma \vdash h : \Delta$, where Δ is the sequence of types of variables bound by the heap h , then the denotation $\llbracket h \rrbracket$ of h is a measurable function from the measurable space $\prod \Gamma$ to the measurable space $\prod \Delta$. Again we treat each element of these product spaces as a function from variable names.) Extend do -notation to a binary operation that turns a heap and an expression into an expression:

$$\text{do } \{ \square; e \} = e \quad (34)$$

$$\text{do } \{ h; g; e \} = \text{do } \{ h; \text{do } \{ g; e \} \} \quad (35)$$

4. Partial evaluation

Our binding-time analysis is online: the metalanguage is static/earlier-stage evaluation (*Lazy s repr* in the Haskell code); the object language is dynamic/later-stage evaluation (*repr* in the Haskell code).

By convention, we write the names of term constructors in bold. For example, **fst** and **inl** are term constructors, so **fst** (x, y) and x are two different terms in our language. In contrast, we write the names of metalanguage functions in italic. For example, the following equations define the metalanguage function *fst*, which projects a head normal form of type $\alpha \times \beta$ to a term of type α :

$$\begin{aligned} \text{fst } (e_1, e_2) &= e_1 \\ \text{fst } a &= \mathbf{fst } a \end{aligned}$$

(The metavariable a stands for an atomic term.) Informally, we write the signature

$$\text{fst}: [\Gamma; \Delta \vdash \alpha \times \beta] \rightarrow [\Gamma, \Delta \vdash \alpha]$$

in which $[\Gamma; \Delta \vdash \alpha \times \beta]$ means head normal forms of type $\alpha \times \beta$ in type environment Γ, Δ (but only variables in Γ are considered atomic because the variables in Δ are bound in the heap), and $[\Gamma, \Delta \vdash \alpha]$ means terms of type α in type environment Γ, Δ . Hence *fst* (x, y) equals the term x . Analogously, we define the metalanguage function *snd*:

$$\begin{aligned} \text{snd}: [\Gamma; \Delta \vdash \alpha \times \beta] &\rightarrow [\Gamma, \Delta \vdash \beta] \\ \text{snd } (e_1, e_2) &= e_2 \\ \text{snd } a &= \mathbf{snd } a \end{aligned}$$

The metalanguage functions *fst* and *snd* are thus partially evaluating counterparts to the term constructors **fst** and **snd**.

Following this convention, the term constructors that we write in bold include arithmetic operations. For example, **exp** and **+** and **-** are term constructors, so **exp** 0 and **3 + (-2)** and **1** are two different terms in our language. In contrast, when we want to exponentiate the number 0 or add the number 3 to the negation of the number 2 in the metalanguage, we write non-bold $\text{exp } 0$ or $3 + (-2)$, which equals the number 1. Informally, we write the signatures

$$\begin{aligned} + : \mathbb{R} &\rightarrow \mathbb{R} \rightarrow \mathbb{R}, \\ - : \mathbb{R} &\rightarrow \mathbb{R}, \\ \text{exp} : \mathbb{R} &\rightarrow \mathbb{R}. \end{aligned}$$

Moreover, we extend these metalanguage functions from operating on concrete numbers to operating on head normal forms (which include concrete numbers). That is, we define

$$\begin{aligned} + : [\Gamma; \Delta \vdash \mathbb{R}] &\rightarrow [\Gamma; \Delta \vdash \mathbb{R}] \rightarrow [\Gamma; \Delta \vdash \mathbb{R}], \\ - : [\Gamma; \Delta \vdash \mathbb{R}] &\rightarrow [\Gamma; \Delta \vdash \mathbb{R}], \\ \text{exp} : [\Gamma; \Delta \vdash \mathbb{R}] &\rightarrow [\Gamma; \Delta \vdash \mathbb{R}] \end{aligned}$$

by extending the usual operations with the fallback cases

$$\begin{aligned} n_1 + n_2 &= n_1 + n_2 && \text{if } n_1 \text{ or } n_2 \text{ is atomic,} \\ -a &= -a && \text{given } a \text{ is atomic,} \\ \text{exp } a &= \mathbf{exp } a && \text{given } a \text{ is atomic.} \end{aligned}$$

(It is straightforward to add algebraic simplifications such as $n + 0 = n$.) Similarly, we define the metalanguage function

$$< : [\Gamma; \Delta \vdash \mathbb{R}] \rightarrow [\Gamma; \Delta \vdash \mathbb{R}] \rightarrow [\Gamma; \Delta \vdash \mathbb{1} + \mathbb{1}]$$

by extending the usual comparison on concrete numbers

$$\begin{aligned} r_1 < r_2 &= \mathbf{inl } () && \text{if } r_1 \text{ is less than } r_2, \\ r_1 < r_2 &= \mathbf{inr } () && \text{if } r_1 \text{ is greater than or equal to } r_2 \end{aligned}$$

(where $r_1, r_2 \in \mathbb{R}$) with the fallback case

$$n_1 < n_2 = n_1 < n_2 \quad \text{if } n_1 \text{ or } n_2 \text{ is atomic.}$$

For pattern matching on sum types, our language includes the projection constructs **do** {**let inl** $x = e; e'$ } and **do** {**let inr** $x = e; e'$ }.

Types	$\alpha, \beta, \gamma ::= \mathbb{R} \mid \mathbb{1} \mid \alpha \times \beta \mid \alpha + \beta \mid \mathbb{M}\alpha$	Terms	e
Type environments	$\Gamma, \Delta ::= [] \mid \Gamma, x : \alpha$	Variables	x, y
Type judgments	$\Gamma \vdash e : \alpha$	Real numbers	$r \in \mathbb{R}$

$\frac{}{\Gamma, x : \alpha, \Delta \vdash x : \alpha}$	$\frac{r \in \mathbb{R}}{\Gamma \vdash r : \mathbb{R}}$	$\frac{\Gamma \vdash e : \mathbb{R}}{\Gamma \vdash -e : \mathbb{R}}$	$\frac{\Gamma \vdash e : \mathbb{R}}{\Gamma \vdash e^{-1} : \mathbb{R}}$	$\frac{\Gamma \vdash e : \mathbb{R}}{\Gamma \vdash \mathbf{exp} e : \mathbb{R}}$	$\frac{\Gamma \vdash e : \mathbb{R}}{\Gamma \vdash \mathbf{log} e : \mathbb{R}}$	$\frac{\Gamma \vdash e_1 : \mathbb{R} \quad \Gamma \vdash e_2 : \mathbb{R}}{\Gamma \vdash e_1 + e_2 : \mathbb{R}}$	$\frac{\Gamma \vdash e_1 : \mathbb{R} \quad \Gamma \vdash e_2 : \mathbb{R}}{\Gamma \vdash e_1 \times e_2 : \mathbb{R}}$
$\frac{\Gamma \vdash e_1 : \mathbb{R} \quad \Gamma \vdash e_2 : \mathbb{R}}{\Gamma \vdash e_1 < e_2 : \mathbb{1} + \mathbb{1}}$	$\frac{}{\Gamma \vdash () : \mathbb{1}}$	$\frac{\Gamma \vdash e_1 : \alpha \quad \Gamma \vdash e_2 : \beta}{\Gamma \vdash (e_1, e_2) : \alpha \times \beta}$	$\frac{\Gamma \vdash e : \alpha \times \beta}{\Gamma \vdash \mathbf{fst} e : \alpha}$	$\frac{\Gamma \vdash e : \alpha \times \beta}{\Gamma \vdash \mathbf{snd} e : \beta}$	$\frac{\Gamma \vdash e : \alpha}{\Gamma \vdash \mathbf{inl} e : \alpha + \beta}$	$\frac{\Gamma \vdash e : \beta}{\Gamma \vdash \mathbf{inr} e : \alpha + \beta}$	
$\frac{}{\Gamma \vdash \mathbf{lebesgue} : \mathbb{M}\mathbb{R}}$	$\frac{\Gamma \vdash e : \alpha}{\Gamma \vdash \mathbf{return} e : \mathbb{M}\alpha}$	$\frac{\Gamma \vdash e : \alpha + \beta \quad \Gamma, x : \alpha \vdash e' : \mathbb{M}\gamma}{\Gamma \vdash \mathbf{do} \{\mathbf{let} \mathbf{inl} x = e; e'\} : \mathbb{M}\gamma}$	$\frac{\Gamma \vdash e : \alpha + \beta \quad \Gamma, x : \beta \vdash e' : \mathbb{M}\gamma}{\Gamma \vdash \mathbf{do} \{\mathbf{let} \mathbf{inr} x = e; e'\} : \mathbb{M}\gamma}$	$\frac{\Gamma \vdash e : \mathbb{M}\alpha \quad \Gamma, x : \alpha \vdash e' : \mathbb{M}\beta}{\Gamma \vdash \mathbf{do} \{x \leftarrow e; e'\} : \mathbb{M}\beta}$	$\frac{\Gamma \vdash e : \mathbb{R} \quad \Gamma \vdash e' : \mathbb{M}\alpha}{\Gamma \vdash \mathbf{do} \{\mathbf{factor} e; e'\} : \mathbb{M}\alpha}$	$\frac{}{\Gamma \vdash \mathbf{mzero} : \mathbb{M}\alpha}$	$\frac{\Gamma \vdash e_1 : \mathbb{M}\alpha \quad \Gamma \vdash e_2 : \mathbb{M}\alpha}{\Gamma \vdash \mathbf{mplus} e_1 e_2 : \mathbb{M}\alpha}$

Bindings (guards)	$g ::= \mathbf{let} \mathbf{inl} x = e \mid \mathbf{let} \mathbf{inr} x = e \mid x \leftarrow e \mid \mathbf{factor} e$
Heaps	$h ::= [] \mid h; g$
Head normal forms	$n, m ::= a \mid r \mid () \mid (e_1, e_2) \mid \mathbf{inl} e \mid \mathbf{inr} e \mid \mathbf{lebesgue} \mid \mathbf{return} e \mid \mathbf{do} \{g; e\} \mid \mathbf{mzero} \mid \mathbf{mplus} e_1 e_2$
Atomic terms	$a ::= x$ (not bound in the heap) $\mid -a \mid a^{-1} \mid \mathbf{exp} a \mid \mathbf{log} a \mid a + n \mid n + a \mid a \times n \mid n \times a \mid a < n \mid n < a \mid \mathbf{fst} a \mid \mathbf{snd} a$

Figure 4. The syntax and type system of our language

In these constructs, the body expression e' must be of measure type. Given $e : \alpha + \beta$, these constructs test whether the result of e is an **inl** value or an **inr** value. If e is an **inl** value, say, then $\mathbf{do} \{\mathbf{let} \mathbf{inl} x = e; e'\}$ binds $x : \alpha$ in e' , whereas $\mathbf{do} \{\mathbf{let} \mathbf{inr} x = e; e'\}$ simply fails. In measure-theoretic terms, this construct converts a measure on α into a measure on $\alpha + \beta$, via the inclusion map **inl**.

To project from sum types in our partial evaluator, we define a pair of metalanguage functions *outl* and *outr*. Conceptually, *outl* projects a head normal form of type $\alpha + \beta$ to a term of type α . But because it is not known whether an atomic term is **inl**, we need to define *outl* in continuation-passing style [4, 6, 12]:

$$\begin{aligned} \mathit{outl}: [\Gamma; \Delta \vdash \alpha + \beta] &\rightarrow ([\Gamma; \Delta \vdash \alpha] \rightarrow \langle \Gamma \vdash \Delta \rangle \rightarrow [\Gamma; \vdash \mathbb{M}\gamma]) \\ &\rightarrow \langle \Gamma \vdash \Delta \rangle \rightarrow [\Gamma; \vdash \mathbb{M}\gamma] \\ \mathit{outl} a \quad c h &= \mathbf{do} \{\mathbf{let} \mathbf{inl} x = a; c x h\} \\ \mathit{outl} (\mathbf{inl} e) c h &= c e h \\ \mathit{outl} (\mathbf{inr} e) c h &= \mathbf{mzero} \end{aligned}$$

In the type above, $\langle \Gamma \vdash \Delta \rangle$ means heaps that bind the variables in Δ using the variables in Γ . The definition of

$$\mathit{outr}: [\Gamma; \Delta \vdash \alpha + \beta] \rightarrow ([\Gamma; \Delta \vdash \beta] \rightarrow \langle \Gamma \vdash \Delta \rangle \rightarrow [\Gamma; \vdash \mathbb{M}\gamma]) \rightarrow \langle \Gamma \vdash \Delta \rangle \rightarrow [\Gamma; \vdash \mathbb{M}\gamma]$$

is analogous.

Prove by mutual induction: a head normal form of type \mathbb{R} does not use any variable bound in the heap; an atomic term of any type does not use any variable bound in the heap. So those terms can be strengthened (i.e., have variables bound in the heap removed from their type environments: $[\Gamma; \Delta \vdash \mathbb{R}] \rightarrow [\Gamma; \vdash \mathbb{R}]$).

Weakening is admissible (i.e., any term can have variables added to its type environment: $\forall \Gamma' \geq \Gamma. \forall \Delta' \geq \Delta. ([\Gamma; \Delta \vdash \alpha] \rightarrow [\Gamma'; \Delta' \vdash \alpha], [\Gamma; \Delta \vdash \alpha] \rightarrow [\Gamma'; \Delta' \vdash \alpha])$).

$$\begin{aligned} \mathit{abs}: [\Gamma; \Delta \vdash \mathbb{R}] &\rightarrow (\forall \Gamma' \geq \Gamma. [\Gamma'; \Delta \vdash \mathbb{R}] \rightarrow \langle \Gamma' \vdash \Delta \rangle \rightarrow [\Gamma'; \vdash \mathbb{M}\gamma]) \\ &\rightarrow \langle \Gamma \vdash \Delta \rangle \rightarrow [\Gamma; \vdash \mathbb{M}\gamma] \\ \mathit{abs} r c h &= c \mid r \mid h \quad \text{given } r \in \mathbb{R} \\ \mathit{abs} a c h &= \mathbf{do} \{x \leftarrow \mathbf{mplus} (\mathbf{do} \{\mathbf{let} \mathbf{inl} _ = a < 0; \mathbf{return} (-a)\}) \\ &\quad (\mathbf{do} \{\mathbf{let} \mathbf{inr} _ = a < 0; \mathbf{return} a\}); \\ &\quad c x h\} \quad \text{given } a \text{ is atomic} \end{aligned}$$

Call-by-need PE/supercompilation [3, 11, 14] though our side effect is commutative yet non-idempotent.

To preserve sharing in lazy evaluation, use a heap; for partial evaluation, the heap leaves some locations unbound [8].

Example—uniform distribution on $[0, 1]$:

$$\mathbf{random} = \mathbf{do} \{x \leftarrow \mathbf{lebesgue}; \quad (36) \\ \mathbf{let} \mathbf{inl} _ = 0 < x; \\ \mathbf{let} \mathbf{inl} _ = x < 1; \\ \mathbf{return} x\}$$

For disintegration to succeed on this simple example, it is essential that our constructs for discriminating **inl** from **inr** do not force evaluation of the scrutinee! Follow-up example: sum of uniform distributions on $[0, 1]$ and $[2, 3]$, to show how we handle **mplus**, which is where one invocation of partial evaluation turns into two.

Acknowledgments

Thanks to Jacques Carette, Ryan Culpepper, Praveen Narayanan, Norman Ramsey, Mitchell Wand, and Robert Zinkov for helpful comments and discussions.

This research was supported by DARPA grant FA8750-14-2-0007, NSF grant CNS-0723054, Lilly Endowment, Inc. (through its support for the Indiana University Pervasive Technology Institute), and the Indiana METACyt Initiative. The Indiana METACyt Initiative at IU is also supported in part by Lilly Endowment, Inc.

$$\begin{aligned}
\llbracket x \rrbracket \rho &= \rho(x) \\
\llbracket r \rrbracket \rho &= r \\
\llbracket -e \rrbracket \rho &= -\llbracket e \rrbracket \rho \\
\llbracket e^{-1} \rrbracket \rho &= (\llbracket e \rrbracket \rho)^{-1} \\
\llbracket \mathbf{exp} \ e \rrbracket \rho &= \exp(\llbracket e \rrbracket \rho) \\
\llbracket \mathbf{log} \ e \rrbracket \rho &= \log(\llbracket e \rrbracket \rho) \\
\llbracket e_1 + e_2 \rrbracket \rho &= \llbracket e_1 \rrbracket \rho + \llbracket e_2 \rrbracket \rho \\
\llbracket e_1 \times e_2 \rrbracket \rho &= \llbracket e_1 \rrbracket \rho \times \llbracket e_2 \rrbracket \rho \\
\llbracket e_1 < e_2 \rrbracket \rho &= \mathbf{true} && \text{if } \llbracket e_1 \rrbracket \rho < \llbracket e_2 \rrbracket \rho \\
\llbracket e_1 < e_2 \rrbracket \rho &= \mathbf{false} && \text{if } \llbracket e_1 \rrbracket \rho \geq \llbracket e_2 \rrbracket \rho \\
\llbracket () \rrbracket \rho &= () \\
\llbracket (e_1, e_2) \rrbracket \rho &= (\llbracket e_1 \rrbracket \rho, \llbracket e_2 \rrbracket \rho) \\
\llbracket \mathbf{fst} \ e \rrbracket \rho &= x && \text{if } \llbracket e \rrbracket \rho = (x, y) \\
\llbracket \mathbf{snd} \ e \rrbracket \rho &= y && \text{if } \llbracket e \rrbracket \rho = (x, y) \\
\llbracket \mathbf{inl} \ e \rrbracket \rho &= \mathbf{inl} (\llbracket e \rrbracket \rho) \\
\llbracket \mathbf{inr} \ e \rrbracket \rho &= \mathbf{inr} (\llbracket e \rrbracket \rho) \\
\llbracket \mathbf{lebesgue} \rrbracket \rho &= \lambda^* c. \int_{-\infty}^{\infty} c(x) dx \\
\llbracket \mathbf{return} \ e \rrbracket \rho &= \lambda^* c. c(\llbracket e \rrbracket \rho) \\
\llbracket \mathbf{do} \ \{\mathbf{let} \ \mathbf{inl} \ x = e; e'\} \rrbracket \rho &= \llbracket e' \rrbracket (\rho \{x \mapsto a\}) && \text{if } \llbracket e \rrbracket \rho = \mathbf{inl} \ a \\
\llbracket \mathbf{do} \ \{\mathbf{let} \ \mathbf{inl} \ x = e; e'\} \rrbracket \rho &= \lambda^* c. 0 && \text{if } \llbracket e \rrbracket \rho = \mathbf{inr} \ b \\
\llbracket \mathbf{do} \ \{\mathbf{let} \ \mathbf{inr} \ x = e; e'\} \rrbracket \rho &= \llbracket e' \rrbracket (\rho \{x \mapsto b\}) && \text{if } \llbracket e \rrbracket \rho = \mathbf{inr} \ b \\
\llbracket \mathbf{do} \ \{\mathbf{let} \ \mathbf{inr} \ x = e; e'\} \rrbracket \rho &= \lambda^* c. 0 && \text{if } \llbracket e \rrbracket \rho = \mathbf{inl} \ a \\
\llbracket \mathbf{do} \ \{x \leftarrow e; e'\} \rrbracket \rho &= \lambda^* c. \llbracket e \rrbracket \rho * \lambda a. \llbracket e' \rrbracket (\rho \{x \mapsto a\}) * c \\
\llbracket \mathbf{do} \ \{\mathbf{factor} \ e; e'\} \rrbracket \rho &= \lambda^* c. (\llbracket e \rrbracket \rho) \times (\llbracket e' \rrbracket \rho * c) \\
\llbracket \mathbf{mzero} \rrbracket \rho &= \lambda^* c. 0 \\
\llbracket \mathbf{mplus} \ e_1 \ e_2 \rrbracket \rho &= \lambda^* c. (\llbracket e_1 \rrbracket \rho * c) + (\llbracket e_2 \rrbracket \rho * c)
\end{aligned}$$

Figure 5. The denotational semantics of our language

References

- [1] S. Bhat, A. Agarwal, R. Vuduc, and A. Gray. A type theory for probability density functions. In *POPL '12: Conference Record of the Annual ACM Symposium on Principles of Programming Languages*, pages 545–556, New York, Jan. 2012. ACM Press.
- [2] S. Bhat, J. Borgström, A. D. Gordon, and C. V. Russo. Deriving probability density functions from probabilistic functional programs. In N. Piterman and S. A. Smolka, editors, *Proceedings of TACAS 2013: 19th International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, number 7795 in Lecture Notes in Computer Science, pages 508–522, Berlin, 16–24 Mar. 2013. Springer. ISBN 978-3-642-36741-0.
- [3] M. Bolingbroke and S. Peyton Jones. Supercompilation by evaluation. In J. Gibbons, editor, *Proceedings of the 3rd ACM SIGPLAN Symposium on Haskell*, pages 135–146, New York, 2010. ACM Press. ISBN 978-1-4503-0252-4.
- [4] A. Bondorf. Improving binding times without explicit CPS-conversion. In W. D. Clinger, editor, *Proceedings of the 1992 ACM Conference on Lisp and Functional Programming*, volume V(1) of *Lisp Pointers*, pages 1–10, New York, 22–24 June 1992. ACM Press.
- [5] J. T. Chang and D. Pollard. Conditioning as disintegration. *Statistica Neerlandica*, 51(3):287–317, Nov. 1997. URL <http://www.stat.yale.edu/jtc5/papers/ConditioningAsDisintegration.pdf>.
- [6] O. Danvy, K. Malmkjær, and J. Palsberg. Eta-expansion does The Trick. *ACM Transactions on Programming Languages and Systems*, 18(6):730–751, 1996.
- [7] B. de Finetti. *Theory of Probability: A Critical Introductory Treatment*, volume 1. Wiley, New York, 1974. Translated from *Teoria Delle Probabilità*, 1970.
- [8] S. Fischer, J. Silva, S. Tamarit, and G. Vidal. Preserving sharing in the partial evaluation of lazy functional programs. In A. King, editor, *Revised Selected Papers from LOPSTR 2007: 17th International Symposium on Logic-Based Program Synthesis and Transformation*, number 4915 in Lecture Notes in Computer Science, pages 74–89, Berlin, 2008. Springer. ISBN 978-3-540-78768-6.
- [9] M. Giry. A categorical approach to probability theory. In B. Banaschewski, editor, *Categorical Aspects of Topology and Analysis: Proceedings of an International Conference Held at Carleton University, Ottawa, August 11–15, 1981*, number 915 in Lecture Notes in Mathematics, pages 68–85, Berlin, 1982. Springer.
- [10] A. H. Jazwinski. *Stochastic Processes and Filtering Theory*. Academic Press, San Diego, CA, 1970.
- [11] J. Jørgensen. Generating a compiler for a lazy language by partial evaluation. In *POPL '92: Conference Record of the Annual ACM Symposium on Principles of Programming Languages*, pages 258–268, New York, 1992. ACM Press. ISBN 0-89791-453-8.
- [12] J. L. Lawall and O. Danvy. Continuation-based partial evaluation. In *Proceedings of the 1994 ACM Conference on Lisp and Functional Programming*, pages 227–238, New York, 1994. ACM Press. URL <http://www.diku.dk/users/julia/lawall-danvy-lfp94.ps.gz>.
- [13] D. J. C. MacKay. Introduction to Monte Carlo methods. In M. I. Jordan, editor, *Learning and Inference in Graphical Models*. Kluwer, Dordrecht, 1998. URL <http://www.cs.toronto.edu/mackay/abstracts/erice.html> <http://www.inference.phy.cam.ac.uk/mackay/BayesMC.html>. Paperback: *Learning in Graphical Models*, MIT Press.
- [14] N. Mitchell. Rethinking supercompilation. In P. Hudak and S. Weirich, editors, *ICFP '10: Proceedings of the ACM International Conference on Functional Programming*, pages 309–320, New York, 2010. ACM Press. ISBN 978-1-60558-794-3.
- [15] S. Park, F. Pfenning, and S. Thrun. A probabilistic language based upon sampling functions. <http://www.postech.ac.kr/gla/>, Aug. 2006. URL <http://www.postech.ac.kr/gla/paper/toplas06.ps> <http://www.postech.ac.kr/gla/paper/toplas06.pdf>. Submitted.
- [16] D. Pollard. *A User's Guide to Measure Theoretic Probability*. Cambridge University Press, Cambridge, 2001. URL <http://www.stat.yale.edu/pollard/Books/UGMTP/>.
- [17] A. Stuhlmüller and N. D. Goodman. Reasoning about reasoning by nested conditioning: Modeling theory of mind with probabilistic programs. *Cognitive Systems Research*, 28:80–99, 2014.

For all $h : \langle \Gamma \vdash \Delta \rangle, m : \lfloor \Gamma; \Delta \vdash \mathbb{M} \alpha \rfloor, e' : \lceil \Gamma, \Delta, x: \alpha \vdash \mathbb{M} \gamma \rceil,$	$\llbracket \text{fwdExec } m (\lambda n. \lambda h'. \mathbf{do} \{h'; e' \{x \mapsto n\}\}) h \rrbracket = \llbracket \mathbf{do} \{h; x \rightsquigarrow m; e'\} \rrbracket.$
For all $h : \langle \Gamma \vdash \Delta \rangle, e : \lceil \Gamma, \Delta \vdash \alpha \rceil, e' : \lceil \Gamma, \Delta, x: \alpha \vdash \mathbb{M} \gamma \rceil,$	$\llbracket \text{fwdEval } e (\lambda n. \lambda h'. \mathbf{do} \{h'; e' \{x \mapsto n\}\}) h \rrbracket = \llbracket \mathbf{do} \{h; e' \{x \mapsto e\}\} \rrbracket.$
For all $h : \langle \Gamma \vdash \Delta \rangle, m : \lfloor \Gamma; \Delta \vdash \mathbb{M} \mathbb{R} \rfloor, e' : \lceil \Gamma, \Delta, x: \mathbb{R} \vdash \mathbb{M} \gamma \rceil,$	$\llbracket \mathbf{do} \{x \rightsquigarrow \mathbf{lebesgue}; \text{bwdExec } m x (\lambda h'. \mathbf{do} \{h'; e'\}) h \rrbracket = \llbracket \mathbf{do} \{h; x \rightsquigarrow m; e'\} \rrbracket$
and for all $n : \lfloor \Gamma; \Delta \vdash \mathbb{R} \rfloor,$	$\llbracket (\text{bwdExec } m x (\lambda h'. \mathbf{do} \{h'; e'\}) h) \{x \mapsto n\} \rrbracket = \llbracket \text{bwdExec } m n (\lambda h'. \mathbf{do} \{h'; e'\}) h \rrbracket.$
For all $h : \langle \Gamma \vdash \Delta \rangle, e : \lceil \Gamma, \Delta \vdash \mathbb{R} \rceil, e' : \lceil \Gamma, \Delta, x: \mathbb{R} \vdash \mathbb{M} \gamma \rceil,$	$\llbracket \mathbf{do} \{x \rightsquigarrow \mathbf{lebesgue}; \text{bwdEval } e x (\lambda h'. \mathbf{do} \{h'; e'\}) h \rrbracket = \llbracket \mathbf{do} \{h; e' \{x \mapsto e\}\} \rrbracket$
and for all $n : \lfloor \Gamma; \Delta \vdash \mathbb{R} \rfloor,$	$\llbracket (\text{bwdEval } e x (\lambda h'. \mathbf{do} \{h'; e'\}) h) \{x \mapsto n\} \rrbracket = \llbracket \text{bwdEval } e n (\lambda h'. \mathbf{do} \{h'; e'\}) h \rrbracket.$

Figure 6. The specification of our lazy partial evaluator and disintegrator

Lazy partial evaluation to head normal form

$fwdExec: [\Gamma; \Delta \vdash \mathbb{M} \alpha] \rightarrow (\forall \Gamma' \geq \Gamma. \forall \Delta' \geq \Delta. [\Gamma'; \Delta' \vdash \alpha] \rightarrow \langle \Gamma' \vdash \Delta' \rangle \rightarrow [\Gamma'; \vdash \mathbb{M} \gamma]) \rightarrow \langle \Gamma \vdash \Delta \rangle \rightarrow [\Gamma; \vdash \mathbb{M} \gamma]$	
$fwdExec \ a \quad c \ h$	$= \mathbf{do} \{x \leftarrow a; \quad c \ x \ h\}$ given a is atomic
$fwdExec \ \mathbf{lebesgue} \quad c \ h$	$= \mathbf{do} \{x \leftarrow \mathbf{lebesgue}; \ c \ x \ h\}$
$fwdExec \ (\mathbf{return} \ e) \quad c \ h$	$= fwdEval \ e \ c \ h$
$fwdExec \ (\mathbf{do} \{g; e\}) \quad c \ h$	$= fwdEval \ e \ (\lambda m. fwdExec \ m \ c) \ (h; g)$
$fwdExec \ \mathbf{mzero} \quad c \ h$	$= \mathbf{mzero}$
$fwdExec \ (\mathbf{mplus} \ e_1 \ e_2) \quad c \ h$	$= \mathbf{mplus} \ (fwdEval \ e_1 \ (\lambda m. fwdExec \ m \ c) \ h) \ (fwdEval \ e_2 \ (\lambda m. fwdExec \ m \ c) \ h)$
$fwdEval: [\Gamma; \Delta \vdash \alpha] \rightarrow (\forall \Gamma' \geq \Gamma. \forall \Delta' \geq \Delta. [\Gamma'; \Delta' \vdash \alpha] \rightarrow \langle \Gamma' \vdash \Delta' \rangle \rightarrow [\Gamma'; \vdash \mathbb{M} \gamma]) \rightarrow \langle \Gamma \vdash \Delta \rangle \rightarrow [\Gamma; \vdash \mathbb{M} \gamma]$	
$fwdEval \ n \quad c \ h$	$= c \ n \ h$ given n is in head normal form
$fwdEval \ (\mathbf{fst} \ e_0) \quad c \ h$	$= fwdEval \ e_0 \ (\lambda n_0. fwdEval \ (fst \ n_0) \ c) \ h$ unless e_0 is atomic
$fwdEval \ (\mathbf{snd} \ e_0) \quad c \ h$	$= fwdEval \ e_0 \ (\lambda n_0. fwdEval \ (snd \ n_0) \ c) \ h$ unless e_0 is atomic
$fwdEval \ (-e_0) \quad c \ h$	$= fwdEval \ e_0 \ (\lambda n_0. c \ (-n_0)) \ h$
$fwdEval \ (e_0^{-1}) \quad c \ h$	$= fwdEval \ e_0 \ (\lambda n_0. c \ (n_0^{-1})) \ h$
$fwdEval \ (\mathbf{exp} \ e_0) \quad c \ h$	$= fwdEval \ e_0 \ (\lambda n_0. c \ (\exp n_0)) \ h$
$fwdEval \ (\mathbf{log} \ e_0) \quad c \ h$	$= fwdEval \ e_0 \ (\lambda n_0. c \ (\log n_0)) \ h$
$fwdEval \ (e_1 + e_2) \quad c \ h$	$= fwdEval \ e_1 \ (\lambda n_1. fwdEval \ e_2 \ (\lambda n_2. c \ (n_1 + n_2))) \ h$
$fwdEval \ (e_1 \times e_2) \quad c \ h$	$= fwdEval \ e_1 \ (\lambda n_1. fwdEval \ e_2 \ (\lambda n_2. c \ (n_1 \times n_2))) \ h$
$fwdEval \ (e_1 < e_2) \quad c \ h$	$= fwdEval \ e_1 \ (\lambda n_1. fwdEval \ e_2 \ (\lambda n_2. c \ (n_1 < n_2))) \ h$
$fwdEval \ x \quad c \ (h_1; x \leftarrow e; h_2)$	$= fwdEval \ e \quad (\lambda m. fwdExec \ m \ (\lambda n. \lambda h'_1. c \ n \ (h'_1; x \leftarrow \mathbf{return} \ n; h_2))) \ h_1$
$fwdEval \ x \quad c \ (h_1; \mathbf{let} \ \mathbf{inl} \ x = e_0; h_2)$	$= fwdEval \ e_0 \ (\lambda n_0. \mathit{outl} \ n_0 \ (\lambda e. fwdEval \ e \ (\lambda n. \lambda h'_1. c \ n \ (h'_1; x \leftarrow \mathbf{return} \ n; h_2)))) \ h_1$
$fwdEval \ x \quad c \ (h_1; \mathbf{let} \ \mathbf{inr} \ x = e_0; h_2)$	$= fwdEval \ e_0 \ (\lambda n_0. \mathit{outr} \ n_0 \ (\lambda e. fwdEval \ e \ (\lambda n. \lambda h'_1. c \ n \ (h'_1; x \leftarrow \mathbf{return} \ n; h_2)))) \ h_1$

Disintegration from head normal form

$bwdExec: [\Gamma; \Delta \vdash \mathbb{M} \mathbb{R}] \rightarrow [\Gamma; \Delta \vdash \mathbb{R}] \rightarrow (\forall \Gamma' \geq \Gamma. \forall \Delta' \geq \Delta. \langle \Gamma' \vdash \Delta' \rangle \rightarrow [\Gamma'; \vdash \mathbb{M} \gamma]) \rightarrow \langle \Gamma \vdash \Delta \rangle \rightarrow [\Gamma; \vdash \mathbb{M} \gamma]$	
$bwdExec \ \mathbf{lebesgue} \quad n \ c \ h$	$= c \ h$
$bwdExec \ (\mathbf{return} \ e) \quad n \ c \ h$	$= bwdEval \ e \ n \ c \ h$
$bwdExec \ (\mathbf{do} \{g; e\}) \quad n \ c \ h$	$= fwdEval \ e \ (\lambda m. bwdExec \ m \ n \ c) \ (h; g)$
$bwdExec \ \mathbf{mzero} \quad n \ c \ h$	$= \mathbf{mzero}$
$bwdExec \ (\mathbf{mplus} \ e_1 \ e_2) \quad n \ c \ h$	$= \mathbf{mplus} \ (fwdEval \ e_1 \ (\lambda m. bwdExec \ m \ n \ c) \ h) \ (fwdEval \ e_2 \ (\lambda m. bwdExec \ m \ n \ c) \ h)$
$bwdEval: [\Gamma; \Delta \vdash \mathbb{R}] \rightarrow [\Gamma; \Delta \vdash \mathbb{R}] \rightarrow (\forall \Gamma' \geq \Gamma. \forall \Delta' \geq \Delta. \langle \Gamma' \vdash \Delta' \rangle \rightarrow [\Gamma'; \vdash \mathbb{M} \gamma]) \rightarrow \langle \Gamma \vdash \Delta \rangle \rightarrow [\Gamma; \vdash \mathbb{M} \gamma]$	
$bwdEval \ (\mathbf{fst} \ e_0) \quad n \ c \ h$	$= fwdEval \ e_0 \ (\lambda n_0. bwdEval \ (fst \ n_0) \ n \ c) \ h$ unless e_0 is atomic
$bwdEval \ (\mathbf{snd} \ e_0) \quad n \ c \ h$	$= fwdEval \ e_0 \ (\lambda n_0. bwdEval \ (snd \ n_0) \ n \ c) \ h$ unless e_0 is atomic
$bwdEval \ (-e_0) \quad n \ c \ h$	$= bwdEval \ e_0 \ (-n) \ c \ h$
$bwdEval \ (e_0^{-1}) \quad n \ c \ h$	$= \mathbf{do} \{\mathbf{factor} \ (n \times n)^{-1}; \ bwdEval \ e_0 \ n^{-1} \ c \ h\}$
$bwdEval \ (\mathbf{exp} \ e_0) \quad n \ c \ h$	$= \mathit{outl} \ (0 < n) \ (\lambda _ . \lambda h'. \mathbf{do} \{\mathbf{factor} \ n^{-1}; \ bwdEval \ e_0 \ (\log n) \ c \ h'\}) \ h$
$bwdEval \ (\mathbf{log} \ e_0) \quad n \ c \ h$	$= \mathbf{do} \{\mathbf{factor} \ (\exp n); \ bwdEval \ e_0 \ (\exp n) \ c \ h\}$
$bwdEval \ (e_1 + e_2) \quad n \ c \ h$	$= fwdEval \ e_1 \ (\lambda n_1. bwdEval \ e_2 \ (n + (-n_1)) \ c) \ h$ $\sqcup \ fwdEval \ e_2 \ (\lambda n_2. bwdEval \ e_1 \ (n + (-n_2)) \ c) \ h$
$bwdEval \ (e_1 \times e_2) \quad n \ c \ h$	$= fwdEval \ e_1 \ (\lambda n_1. \mathit{abs} \ n_1 \ (\lambda n'_1. \lambda h'. \mathbf{do} \{\mathbf{factor} \ n_1'^{-1}; \ bwdEval \ e_2 \ (n \times n_1^{-1}) \ c \ h'\})) \ h$ $\sqcup \ fwdEval \ e_2 \ (\lambda n_2. \mathit{abs} \ n_2 \ (\lambda n'_2. \lambda h'. \mathbf{do} \{\mathbf{factor} \ n_2'^{-1}; \ bwdEval \ e_1 \ (n \times n_2^{-1}) \ c \ h'\})) \ h$
$bwdEval \ x \quad n \ c \ (h_1; x \leftarrow e; h_2)$	$= fwdEval \ e \quad (\lambda m. bwdExec \ m \ n \ (\lambda h'_1. c \ (h'_1; x \leftarrow \mathbf{return} \ n; h_2))) \ h_1$
$bwdEval \ x \quad n \ c \ (h_1; \mathbf{let} \ \mathbf{inl} \ x = e_0; h_2)$	$= fwdEval \ e_0 \ (\lambda n_0. \mathit{outl} \ n_0 \ (\lambda e. bwdEval \ e \ n \ (\lambda h'_1. c \ (h'_1; x \leftarrow \mathbf{return} \ n; h_2)))) \ h_1$
$bwdEval \ x \quad n \ c \ (h_1; \mathbf{let} \ \mathbf{inr} \ x = e_0; h_2)$	$= fwdEval \ e_0 \ (\lambda n_0. \mathit{outr} \ n_0 \ (\lambda e. bwdEval \ e \ n \ (\lambda h'_1. c \ (h'_1; x \leftarrow \mathbf{return} \ n; h_2)))) \ h_1$

[TODO: uncurry heap argument?]

Figure 7. The implementation of our lazy partial evaluator and disintegrator