

Simplifying Probabilistic Programs Using Computer Algebra^{*}

Jacques Carette (✉)¹ and Chung-chieh Shan²

¹ McMaster University carette@mcmaster.ca

² Indiana University ccshan@indiana.edu

Abstract. We transform probabilistic programs to run more efficiently and read more easily, by composing three semantics-preserving transformations: (1) apply the denotational semantics; (2) improve the resulting integral; then (3) invert the denotational semantics. Whereas step 1 is a straightforward transformation from monadic to continuation-passing style, the rest builds on computer algebra: step 2 reorders and performs integrals, and step 3 represents density functions as differential operators.

Keywords: probabilistic programming, computer algebra, integrating continuous distributions

1 Introduction

The success of machine learning has made it clear that computing with probability distributions is very useful. Given a distribution, we might want a simpler representation of it or to generate random samples from it. These two computations go hand in hand: If we are lucky, we can find a simple representation, maybe even an exact one, that renders all further calculation trivial. If not, we would likely want to generate random samples, a process that can be made dramatically more efficient and accurate by any simplification we manage to perform.

1.1 Contributions

We introduce a way to simplify probabilistic programs: starting with a monadic representation of a distribution (Section 2), we transform it to an integral (Section 4), improve the integral by controlled use of computer algebra (Section 6), then transform back (Section 5). Whereas the denotational semantics that maps probabilistic programs to integrals is straightforward and well known (Section 3),

^{*} Thanks to Mike Kucera, Praveen Narayanan, Natalie Perna, and Robert Zinkov for developing the Hakaru probabilistic programming system, the home of our research. Thanks also to the anonymous referees for many comments, which improved this paper. 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.

we put it to work and show how to turn integrals into simpler programs. In particular, we identify necessary transformations on mathematical and integral expressions that direct an existing computer algebra system (CAS) to simplify the kind of multidimensional integrals that probabilistic programs tend to mean. We also apply computer algebra to convert integrals back to programs robustly. The CAS we use is Maple, but other systems such as Mathematica are similar.

1.2 Three examples

We showcase our approach using three natural (if small) examples.

A discrete distribution Start with the following sampling procedure: choose two real numbers x, y independently from the uniform distribution between 0 and 1, then return whether $x < y$. The corresponding term in our language

$$\text{Bind}(\text{Uniform}(0, 1), x, \text{Bind}(\text{Uniform}(0, 1), y, \text{If}(x < y, \text{Ret}(\text{true}), \text{Ret}(\text{false})))) \quad (1)$$

denotes a measure over Booleans (\mathbb{B}). As explained in Section 3, a measure is equivalent to a linear operator. This equivalence means that we can integrate with respect to this distribution an arbitrary function $h : \mathbb{B} \rightarrow \mathbb{R}^+$, and the result completely determines the denotation of the term (1). That result is the number

$$\int_0^1 \int_0^1 \left(\begin{cases} h(\text{true}) & \text{if } x < y \\ h(\text{false}) & \text{otherwise} \end{cases} \right) dy dx. \quad (2)$$

Off the shelf, a CAS like Maple can simplify this integral expression to $\frac{1}{2} \cdot h(\text{true}) + \frac{1}{2} \cdot h(\text{false})$. From this integral, we read off the simpler term in our language

$$\text{Msum}(\text{Weight}(1/2, \text{Ret}(\text{true})), \text{Weight}(1/2, \text{Ret}(\text{false}))). \quad (3)$$

This term expresses a linear combination of measures, namely the average of the two measures $\text{Ret}(\text{true})$ and $\text{Ret}(\text{false})$. So this term denotes the same measure as (1), but it is more efficient as a sampling procedure: instead of drawing two numbers from a uniform distribution, just flip one fair coin and return either `true` or `false` accordingly. This term is also easier for humans to understand. Indeed, this representation amounts to a compact table that lists the outcomes (`true` and `false`) alongside their exact probabilities (1/2 each). Nothing could be better.

A continuous distribution Let us take two steps on a one-dimensional random walk: first choose x from the Gaussian distribution with mean 0 and standard deviation 1, then choose y from the Gaussian distribution with mean x and standard deviation 1. If we only care about the final point y , then the corresponding term

$$\text{Bind}(\text{Gaussian}(0, 1), x, \text{Bind}(\text{Gaussian}(x, 1), y, \text{Ret}(y))) \quad (4)$$

denotes a measure over \mathbb{R} . In other words, we can integrate with respect to this distribution any function h from \mathbb{R} to \mathbb{R}^+ . The result is the number

$$\int_{-\infty}^{\infty} \frac{\exp(-\frac{x^2}{2})}{\sqrt{2 \cdot \pi}} \cdot \int_{-\infty}^{\infty} \frac{\exp(-\frac{(y-x)^2}{2})}{\sqrt{2 \cdot \pi}} \cdot h(y) dy dx. \quad (5)$$

With Maple's help, we can simplify this integral expression to

$$\int_{-\infty}^{\infty} \frac{\exp\left(-\frac{y^2}{4}\right)}{2 \cdot \sqrt{\pi}} \cdot h(y) dy. \quad (6)$$

From this, we read off the much simpler term in our language $\text{Gaussian}(0, \sqrt{2})$, which expresses the fact that taking two steps on a random walk with standard deviation 1 is equivalent to taking one step on a random walk with standard deviation $\sqrt{2}$. This term denotes the same measure over \mathbb{R} as (4), but it is more efficient as a sampling procedure: instead of drawing two numbers from Gaussian distributions, just draw one number from a third Gaussian. This term is also easier for humans to understand (for instance, easier to visualize as a bell curve).

A conditional distribution What if, instead of caring about the second step y and throwing away x , we *observe* y from an actual walk and want to use that information to *infer* the first step x ? One way to express the observation is to replace the random choice of y in (4) by a *probability density* \mathcal{D} :

$$\text{Bind}(\text{Gaussian}(0, 1), x, \text{Weight}(\mathcal{D}(\text{Gaussian}(x, 1), y), \text{Ret}(x))) \quad (7)$$

This term denotes a measure over \mathbb{R} . As a sampling procedure, it means to draw x from $\text{Gaussian}(0, 1)$ then return x with the *weight* $\mathcal{D}(\text{Gaussian}(x, 1), y)$. This weight, defined in Table 1, is a positive real number that reaches its maximum when $x = y$ and decreases as x moves farther from y .

The integral of h with respect to the distribution (7) is similar to (5),

$$\int_{-\infty}^{\infty} \frac{\exp\left(-\frac{x^2}{2}\right)}{\sqrt{2 \cdot \pi}} \cdot \frac{\exp\left(-\frac{(y-x)^2}{2}\right)}{\sqrt{2 \cdot \pi}} \cdot h(x) dx. \quad (8)$$

From this integral, we read off the term

$$\text{Weight}\left(\frac{\exp\left(-\frac{y^2}{4}\right)}{2 \cdot \sqrt{\pi}}, \text{Gaussian}\left(\frac{y}{2}, \frac{1}{\sqrt{2}}\right)\right) \quad (9)$$

for the *conditional* distribution of x given the observation y . This term tells us that this distribution is proportional to a Gaussian with mean $y/2$ and standard deviation $1/\sqrt{2}$. It denotes the same measure over \mathbb{R} as (7), but it is more efficient as a sampling procedure because every sample of x is returned with the same weight, depending only on y . Statisticians call the equivalence between (7) and (9) a *conjugacy* between the two Gaussians in (7). This simplification underpins the popular *Kalman filter* (Maybeck 1979).

2 Adding measures to a mathematical language

Instead of building our own language from scratch, we start with a CAS whose language already expresses integral calculus. For example, Maple represents the integral (2) internally as the term

$$\text{Int}(\text{Int}(\text{If}(x < y, h(\text{true}), h(\text{false})), y = 0..1), x = 0..1) \quad (10)$$

$$\begin{array}{c}
\frac{a : \mathbb{R} \quad b : \mathbb{R}}{\text{Uniform}(a, b) : \mathbb{MIR}} \quad \frac{\mu : \mathbb{R} \quad \sigma : \mathbb{R}^+}{\text{Gaussian}(\mu, \sigma) : \mathbb{MIR}} \quad \frac{\mu : \mathbb{R} \quad \gamma : \mathbb{R}^+}{\text{Cauchy}(\mu, \gamma) : \mathbb{MIR}} \\
\frac{\nu : \mathbb{R} \quad \mu : \mathbb{R} \quad \gamma : \mathbb{R}^+}{\text{StudentT}(\nu, \mu, \gamma) : \mathbb{MIR}} \quad \frac{\alpha : \mathbb{R}^+ \quad \beta : \mathbb{R}^+}{\text{Beta}(\alpha, \beta) : \mathbb{MIR}^+} \quad \frac{k : \mathbb{R}^+ \quad \theta : \mathbb{R}^+}{\text{Gamma}(k, \theta) : \mathbb{MIR}^+} \\
\begin{array}{c} [x : A] \\ \vdots \end{array} \\
\frac{e : A}{\text{Ret}(e) : \mathbb{MA}} \quad \frac{m : \mathbb{MA} \quad m' : \mathbb{MB}}{\text{Bind}(m, x, m') : \mathbb{MB}} \quad \frac{\forall i \leq n, m_i : \mathbb{MA}}{\text{Msum}(m_1, \dots, m_n) : \mathbb{MA}} \\
\begin{array}{c} [\tilde{h} : A \rightarrow \mathbb{R}^+] \\ \vdots \end{array} \\
\frac{e : \mathbb{R}^+ \quad m : \mathbb{MA}}{\text{Weight}(e, m) : \mathbb{MA}} \quad \frac{e : \mathbb{B} \quad m : \mathbb{MA} \quad m' : \mathbb{MA}}{\text{If}(e, m, m') : \mathbb{MA}} \quad \frac{g : \mathbb{R}^+}{\text{LO}(\tilde{h}, g) : \mathbb{MA}}
\end{array}$$

Fig. 1. Informal typing rules for how we represent measures

using the syntax constructors `Int` `If` `true` `false` `<` `=` `..` (the last three of which are written infix) and the variables x y and h .

We add to this language a handful of constructors that amount to an abstract data type of measures. These constructors are summarized in Figure 1, using informal “typing rules” even though Maple is not statically typed. If A is a type, then \mathbb{MA} is our informal type of measures over A . In particular, the top two rows of Figure 1 show several primitive measures of types \mathbb{MIR} and \mathbb{MIR}^+ , where \mathbb{R} is the type of reals and \mathbb{R}^+ is the type of non-negative reals (including $+\infty$).

Moving on in Figure 1, the constructors `Ret` and `Bind` represent the unit and bind operations of the measure monad. In `Bind`(m, x, m'), the variable x takes scope over the measure term m' . Here we write x as a metavariable that ranges over variables in the syntax; whereas for terms we use a variety of metavariables, including $e, a, b, k, \alpha, \beta, \gamma, \mu, \nu, \sigma, \theta$, and especially m for a measure term, and g for an integral term.

The next two constructors express non-negative linear combinations of measures. The term `Msum`(m_1, \dots, m_n) represents the sum of n measures, so `Msum`() is the zero measure. The term `Weight`(e, m) represents multiplying a measure m by a non-negative scalar factor e .

We use the conditional `If` in measure terms as well as in ordinary expressions denoting numbers. Thus `If`($x < y, \text{Ret}(\text{true}), \text{Ret}(\text{false})) is a measure term, whereas the expression `If`($x < y, h(\text{true}), h(\text{false})) denotes a number. In our Maple implementation, we actually handle the multiway conditional construct `piecewise` (Carette 2007), but we describe only if-then-else in this paper to keep the notation simple.$$

The last constructor is `LO`, short for “linear operator”. In `LO`(\tilde{h}, g), the metavariable \tilde{h} stands for an *integrand* variable, like h in equation (2). As the typing rule indicates, \tilde{h} takes scope over g . We use `LO` to name a measure by specifying how it integrates a function: `LO`(\tilde{h}, g) means the measure m such that the integral of a measurable non-negative function \tilde{h} with respect to m is g .

An operational way to interpret a measure term is to run it as an *importance sampler*, which generates a random outcome along with a weight. Interpreted thus, `Gaussian` means to draw a number from a Gaussian distribution, `Ret` means to produce the given outcome, and `Bind` means to sequence two importance samplers. As equation (3) exemplifies, a non-negative linear combination of measures

$$\text{Msum}(\text{Weight}(e_1, m_1), \dots, \text{Weight}(e_n, m_n)) \quad (11)$$

is an n -way random choice: we choose one of the subterms m_i with probability proportional to e_i , and at the same time multiply the current weight (which starts at 1) by $\sum_{i=1}^n e_i$. In particular, if $m_1 = \dots = m_n = 1/n$ then we choose m_i uniformly and leave the weight unaffected. If an immediate subterm m_i of $\text{Msum}(m_1, \dots, m_n)$ is not built with `Weight`, then we treat m_i like `Weight(1, m_i)`. And if the term `Weight(e, m)` occurs in isolation, then its sampling interpretation is as with $\text{Msum}(\text{Weight}(e, m))$: scale the weight by e and continue with m .

We can also interpret each measure term denotationally—as a measure. To do this, we specify the integral of a function with respect to the measure. Our approach to simplification starts with this definition, so we turn to it next.

3 Connecting abstract and concrete integration

A fundamental result of measure theory is that measures can be viewed in two equivalent ways (Theorems (12) and (13) in Pollard 2001).

On one hand, we can view a measure as a function that maps sets to their sizes in \mathbb{R}^+ , by *measuring* the sets' length or volume.

On this view, the uniform distribution on the unit interval (`Uniform(0, 1)` in our language) is a function that maps the interval $[2/3, 2]$ to the number $1/3$, because $[0, 1] \cap [2/3, 2]$ has length $1/3$. In other words, the probability is $1/3$ for a point randomly drawn from `Uniform(0, 1)` to lie in $[2/3, 2]$. The same function maps the singleton set $\{1/2\}$ to the number 0, which is the probability for a point randomly drawn from `Uniform(0, 1)` to be exactly $1/2$. Similarly, the *Dirac distribution* at $1/2$ (which we write `Ret(1/2)`) is a function that maps each set S to 1 if $1/2 \in S$, and to 0 otherwise. So it maps the interval $[2/3, 1]$ to 0 and the singleton $\{1/2\}$ to 1. After all, if we draw the point $1/2$ deterministically, then the probability that the point lies in $[2/3, 1]$ is 0, and the probability that the point is exactly $1/2$ is 1.

On the other hand, we can view a measure as a function that maps functions to their integrals in \mathbb{R}^+ , with higher-order type $(\dots \rightarrow \mathbb{R}^+) \rightarrow \mathbb{R}^+$, so mathematicians call it an *operator*. On this view, `Uniform(0, 1)` is a function that maps h to $\int_0^1 h(x) dx$. For example, it maps $\lambda x. x$ to $\int_0^1 x dx = 1/2$. In other words, the expected value of a random draw from `Uniform(0, 1)` is $1/2$. We can view the Dirac distribution `Ret(1/2)` as a function as well. It maps h to the number $h(1/2)$. And lo and behold, the expected value of drawing the number $1/2$ deterministically is $1/2$.

It may seem strange to call the application $h(1/2)$ an integral of h , but it satisfies the important properties of integration that we use. To start with, if

h is non-negative, then so is $h(1/2)$, just like the integral $\int_0^1 h(x) dx$. Moreover, application is *linear*: for example, we have

$$(\lambda x. a \cdot h(x) + b \cdot k(x))(1/2) = a \cdot h(1/2) + b \cdot k(1/2) \quad (12)$$

for any functions $h, k : \mathbb{R} \rightarrow \mathbb{R}^+$ and weights $a, b : \mathbb{R}^+$. Thus the application of a measure m as a linear operator to a function h is called the (*abstract*) *integral* of the *integrand* h with respect to m . We represent it by `integrate`(m, h).

These two views are equivalent, as mentioned above. Specifically, the size of a set S with respect to a measure m is the integral `integrate`(m, χ_S) of its *characteristic function* χ_S , defined by $\chi_S(x) = 1$ if $x \in S$ and $\chi_S(x) = 0$ otherwise.

Despite this correspondence, we need to view a measure as an integrator, for three reasons. First, abstract integration is required to define the `Bind` operation of the measure monad (Giry 1982; Ramsey and Pfeffer 2002), and `Bind` is essential in the probabilistic programs we seek to simplify. Second, existing CASes can handle concrete integrals, as demonstrated in the introduction. Third, existing CASes are weak at representing and measuring sets.

The ability to compute integrals is useful to us in the common cases where abstract integration involves concrete integration, like with `Uniform` and `Gaussian`. But to take advantage of this ability, we need to relate concrete integration to abstract integration, which CASes off the shelf do not even represent. For example, Maple has plenty of facilities we covet for simplifying and transforming concrete integrals \int , sums \sum , function applications, and their iterated combinations, but the facilities for integration, for summation, and for function application are separate, and there is no single way to represent, say, a product measure.

Hence, our plan is to simplify measure terms in three steps:

1. View the given measure term as a linear operator, which specifies the abstract integral of an arbitrary integrand.
2. Improve the integral using computer algebra, keeping the integrand arbitrary.
3. Read off a new measure term from the improved linear operator.

A major contribution of this paper is to automate not only step 1 (Section 4) but also step 3 (Section 5), so that the task of simplifying probabilistic programs reduces to the task of improving integrals. This reduction leads to another major contribution of this paper, namely to improve integrals using the building blocks provided by an existing CAS (Section 6).

4 From measure term to linear operator

Step 1 is to convert a given measure term m to a term of the form `LO`(\tilde{h}, g). Recall that the \tilde{h} in `LO`(\tilde{h}, g) represents an arbitrary function to be integrated, so we set \tilde{h} to a fresh name. Then, to preserve the meaning of m , the g in `LO`(\tilde{h}, g) should be the abstract integral of \tilde{h} with respect to m , so we set g to `integrate`(m, \tilde{h}), defined in Figure 2.

Figure 2 defines `integrate`(m, h) by structural induction on measure term m , so our definition serves as a compositional denotational semantics of measure terms,

$$\begin{aligned}
\text{integrate}(m, h) &= \int_{\mathcal{L}(m)}^{\mathcal{U}(m)} \mathcal{D}(m, x) \cdot h(x) dx \\
&\quad \text{where } m \text{ is a primitive measure term and } x \text{ is fresh} \\
\text{integrate}(\text{Ret}(e), h) &= h(e) \\
\text{integrate}(\text{Bind}(m, x, m'), h) &= \text{integrate}(m, \lambda x. \text{integrate}(m', h)) \\
\text{integrate}(\text{Msum}(m, \dots), h) &= \text{integrate}(m, h) + \dots \\
\text{integrate}(\text{Weight}(e, m), h) &= e \cdot \text{integrate}(m, h) \\
\text{integrate}(\text{If}(e, m, m'), h) &= \text{If}(e, \text{integrate}(m, h), \text{integrate}(m', h)) \\
\text{integrate}(\text{LO}(\tilde{h}, g), h) &= g\{\tilde{h} \mapsto h\} \\
\text{integrate}(m, h) &= \text{Integrate}(m, h) \quad \text{otherwise}
\end{aligned}$$

Fig. 2. Translating a measure term m to a linear operator $\text{LO}(\tilde{h}, \text{integrate}(m, \tilde{h}))$, where \tilde{h} is fresh

and a standard one at that. The integrand h is an accumulator argument that may not just be a name \tilde{h} , as can be seen in the right-hand side of the **Bind** case. Actually, h amounts to a continuation, so step 1 amounts to a one-pass transform from monadic to continuation-passing style (CPS) (Hatcliff and Danvy 1994) (more precisely, to continuation-composing style (Danvy and Filinski 1990)).

When m is primitive, **integrate** produces a concrete integral using the standard properties of m listed in Table 1. In this paper, all primitive measures range over real intervals. The properties we use from Table 1 are

- $\mathcal{L}(m)$, the lower bound of the interval, possibly $-\infty$;
- $\mathcal{U}(m)$, the upper bound of the interval, possibly $+\infty$; and
- $\mathcal{D}(m, x)$, the *density* of m with respect to the Lebesgue measure.

Many more measures could be added to Table 1, but some measures cannot be so represented, such as **Ret**(0) and the logistic distribution.

The last line in Figure 2 handles the case where m is a free variable representing an unknown measure. The result is a residual term **Integrate**(m, h).

For the example (1), step 1 produces the term

$$\text{LO} \left(\tilde{h}, \int_0^1 \frac{1}{1-0} \cdot \int_0^1 \frac{1}{1-0} \cdot \left(\begin{cases} h(\text{true}) & \text{if } x < y \\ h(\text{false}) & \text{otherwise} \end{cases} dy dx \right), \quad (13)$$

in which we notate the conditional **If** with a curly left brace. Maple immediately removes the factor $\frac{1}{1-0}$, yielding equation (2).

As equation (13) illustrates, the output of this **integrate** step is patently a linear operator. In other words, just by examining the integral produced syntactically, without any deep reasoning, we can tell that it is linear in the integrand \tilde{h} . Formally, we say that g is *patently linear in \tilde{h}* iff g is generated by the grammar

$$g ::= \tilde{h}(e) \mid g_1 + \dots + g_n \mid \int_a^b g dx \mid e \cdot g \mid \text{If}(e, g, g') \mid \text{Integrate}(m, \lambda x. g) \quad (14)$$

where a, b, e, m do not contain \tilde{h} free, and $x \neq \tilde{h}$. (For simplicity, this grammar omits $\sum_a^b g dx$, and this paper omits measures over \mathbb{Z} .)

Table 1. Defining properties of primitive measures

Measure term m	Bounds		Density	Holonomic representation
	$\mathcal{L}(m)$	$\mathcal{U}(m)$	$\mathcal{D}(m, x)$	$\frac{p_0(x)}{p_1(x)} = -\frac{\frac{d}{dx}\mathcal{D}(m, x)}{\mathcal{D}(m, x)}$
Uniform(a, b)	a	b	$\frac{1}{b-a}$	0
Gaussian(μ, σ)	$-\infty$	$+\infty$	$\frac{\exp(-\frac{1}{2} \cdot (\frac{x-\mu}{\sigma})^2)}{\sqrt{2} \cdot \pi \cdot \sigma}$	$\frac{x-\mu}{\sigma^2}$
Cauchy(μ, γ)	$-\infty$	$+\infty$	$\frac{(1 + (\frac{x-\mu}{\gamma})^2)^{-1}}{\pi \cdot \gamma}$	$\frac{2 \cdot (x-\mu)}{(x-\mu)^2 + \gamma^2}$
StudentT(ν, μ, γ)	$-\infty$	$+\infty$	$\frac{(1 + (\frac{x-\mu}{\gamma})^2)^{-\nu/2}}{\frac{\Gamma(\nu/2)}{\Gamma((\nu+1)/2)} \cdot \sqrt{\pi} \cdot \nu \cdot \gamma}$	$\frac{(\nu+1) \cdot (x-\mu)}{(x-\mu)^2 + \gamma^2 \cdot \nu}$
Beta(α, β)	0	1	$\frac{x^{\alpha-1} \cdot (1-x)^{\beta-1}}{B(\alpha, \beta)}$	$\frac{(\alpha + \beta - 2) \cdot x - (\alpha - 1)}{x \cdot (1-x)}$
Gamma(k, θ)	0	$+\infty$	$\frac{x^{k-1} \cdot \exp(-\frac{x}{\theta})}{\Gamma(k) \cdot \theta^k}$	$\frac{\frac{x}{\theta} + 1 - k}{x}$

Of course, patent linearity entails linearity, but the converse is not the case. For example, the term $\sin(h(a))^2 + \cos(h(a))^2 - 1$ is linear in h (because it is zero) but not patently linear. Still, the right-hand sides in Figure 2 all maintain patent linearity, so step 1 produces a patently linear integral, which we then subject to algebraic manipulations in step 2. As long as our manipulations of the integral preserve patent linearity, the result can be turned back to a measure term. That is the job of step 3, which we present next.

5 From linear operator back to measure term

Figure 3 defines $\text{unintegrate}(\hbar, g, c)$, whose goal is to find a measure term m (built without LO) such that $\text{integrate}(m, \hbar) = g$. If we think of step 1 above as a transform from monadic style to CPS, then step 3 inverts this.

Like integrate , unintegrate proceeds by structural induction on the input term. Most lines of this definition just handle a case in the grammar (14) by inverting a corresponding line defining integrate in Figure 2. The main deviations from this pattern are three, described in the three subsections below.

5.1 Peephole optimizations

As unintegrate builds a measure term, it invokes the smart constructors bind and weight , defined in Figure 3. They are semantically equivalent to Bind and Weight but perform peephole optimizations using algebraic laws: bind uses the right and

$$\begin{aligned}
\text{unintegrate}(\hbar, \int_a^b g dx, c) &= \text{weight}(e_1, \text{bind}(m, x, \text{weight}(e_2, m'))) \\
&\quad \text{where } x \neq \hbar \text{ and } (e, m') = \text{unweight}(\text{unintegrate}(\hbar, g, c \wedge a < x < b)) \\
&\quad (m, e') = \text{recognize}(e, x, a, b, c) \\
&\quad e_1 \cdot e_2 = e' \quad \text{where } e_1 \text{ does not contain } x \text{ free} \\
\text{unintegrate}(\hbar, \hbar(e), c) &= \text{Ret}(e) \\
\text{unintegrate}(\hbar, g + \dots, c) &= \text{Msum}(\text{unintegrate}(\hbar, g, c), \dots) \\
\text{unintegrate}(\hbar, e \cdot g, c) &= \text{weight}(e, \text{unintegrate}(\hbar, g, c)) \\
&\quad \text{where } e \text{ does not contain } \hbar \text{ free} \\
\text{unintegrate}(\hbar, \text{If}(e, g, g'), c) &= \text{If}(e, \text{unintegrate}(\hbar, g, c \wedge e), \text{unintegrate}(\hbar, g', c \wedge \neg e)) \\
\text{unintegrate}(\hbar, \text{Integrate}(m, h), c) &= \text{bind}(m, x, \text{unintegrate}(\hbar, h(x), c)) \quad \text{where } x \text{ is fresh} \\
\text{unintegrate}(\hbar, g, c) &= \text{LO}(\hbar, g) \quad \text{otherwise} \\
\text{bind}(m, x, \text{Ret}(x)) &= m & \text{weight}(1, m) &= m \\
\text{bind}(\text{Ret}(e), x, m) &= m\{x \mapsto e\} & \text{weight}(0, m) &= \text{Msum}() \\
\text{bind}(m, x, m') &= \text{Bind}(m, x, m') & \text{weight}(e, \text{Weight}(e', m)) &= \text{weight}(e \cdot e', m) \\
& \text{otherwise} & \text{weight}(e, m) &= \text{Weight}(e, m) \\
& & & \text{otherwise} \\
\text{unweight}(\text{Weight}(e, m)) &= (e, m) \\
\text{unweight}(\text{Msum}(m_1, \dots, m_n)) &= (e, \text{Msum}(\text{weight}(\frac{1}{e}, m_1), \dots, \text{weight}(\frac{1}{e}, m_n))) \\
& \quad \text{where } e = e_1 + \dots + e_n \text{ and } (e_i, -) = \text{unweight}(m_i) \\
\text{unweight}(m) &= (1, m) \quad \text{otherwise}
\end{aligned}$$

Fig. 3. Translating a linear operator $\text{LO}(\hbar, g)$ to a measure term $\text{unintegrate}(\hbar, g, \text{true})$

left monad identity laws to eliminate Bind , and weight uses linearity to eliminate Weight . We wait until step 3 to perform these optimizations, so that they apply to terms produced by density recognition (Section 5.3).

5.2 Assumption context

As unintegrate traverses the input integral g , it maintains a *context* of assumptions about the variables in g .

For example, when unintegrate processes $\text{If}(e, g, g')$, it conjoins the condition e onto the current context c while processing the then-case g , and it conjoins $\neg e$ onto c while processing the else-case g' . So if e is the condition $x > 0$, then the recursive call to $\text{unintegrate}(\hbar, g, c \wedge e)$ can simplify any occurrence of $\sqrt{x^2}$ to x , and $\text{sgn } x$ to 1. These simplifications are performed by the call to recognize and benefit later calls to recognize , as explained near the end of Section 5.3 below.

Similarly, when unintegrate processes $\int_a^b g dx$, it conjoins $a < x < b$ onto the current context c . (Actually, that recursive call to $\text{unintegrate}(\hbar, g, c \wedge a < x < b)$ is free to ignore any countable number of points in the interval $[a, b]$, but unintegrate does not currently exercise that freedom.)

5.3 Density recognition

When `unintegrate` processes $\int_a^b g \, dx$, it tries to recognize a concrete integral that belongs to a primitive measure. For example, when `unintegrate` processes the concrete integral in equation (2), it recognizes $\int_0^1 \dots dx$ as belonging to `Uniform(0, 1)` and so produces a term of the form `Bind(Uniform(0, 1), x, ...)`. This requires non-trivial computer algebra because there is a gulf between the body of that integral, namely the integral $\int_0^1 (\begin{cases} h(\text{true}) & \text{if } x < y \\ h(\text{false}) & \text{otherwise} \end{cases}) dy$, and $\mathcal{D}(m, x) \cdot h(x)$ at the top of Figure 2, namely the product $\frac{1}{b-a} \cdot h(x)$. These terms do not unify syntactically.

For recognizing primitive measures other than `Uniform`, the gulf is even wider, and syntactic matching is even more futile. For example, how can an *algorithm* recognize that the density in (8) is proportional to a certain Gaussian distribution and read off the measure term (9)? When `unintegrate` processes the integral (8), the recursive call to `unintegrate(h, g, c')` returns the measure term

$$\text{Weight}\left(\frac{\exp(-\frac{x^2}{2})}{\sqrt{2 \cdot \pi}} \cdot \frac{\exp(-\frac{(y-x)^2}{2})}{\sqrt{2 \cdot \pi}}, \text{Ret}(x)\right), \quad (15)$$

and the call to `unweight` (defined in Figure 3) decomposes this term into its two subterms, the weight

$$f(x) = \frac{\exp(-\frac{x^2}{2})}{\sqrt{2 \cdot \pi}} \cdot \frac{\exp(-\frac{(y-x)^2}{2})}{\sqrt{2 \cdot \pi}} \quad (16)$$

and the measure `Ret(x)`. The weight is a function in the integration variable x . We call this function the *target density*. We need an algorithm to recognize that it is proportional to $\mathcal{D}(\text{Gaussian}(\mu, \sigma))$ and not, say, $\mathcal{D}(\text{Cauchy}(\mu, \gamma))$.

Holonomic representation We solve this density recognition problem in one fell swoop for *all* our primitive measures, by treating the target density f as a *holonomic*³ expression (Wilf and Zeilberger 1992; Chyzak and Salvy 1998). A good tutorial on holonomicity is provided by Kauers (2013). In short, it means we find a *homogeneous linear differential equation*

$$p_n(x) \cdot f^{(n)}(x) + \dots + p_1(x) \cdot f'(x) + p_0(x) \cdot f(x) = 0 \quad (17)$$

that defines $f(x)$ up to a constant factor, in which each $p_i(x)$ is a polynomial in x . For example, the density (16) is defined up to a constant factor by

$$1 \cdot f'(x) + (2 \cdot x - y) \cdot f(x) = 0, \quad (18)$$

so $n = 1$, $p_1(x) = 1$, and $p_0(x) = 2 \cdot x - y$.

The general algorithm for finding a differential equation (17) from a closed-form expression $f(x)$ is well established and efficiently implemented by the Maple

³ More precisely, a *D-finite* expression. We speak of holonomicity, which encompasses D-finite functions and *P-recursive* sequences, because we foresee recognizing discrete measures using P-recursive equations. Those are almost as effective as D-finite equations, modulo the issue of *accurate summation* (Abramov and Petkovšek 2005).

function `gfun[holexprtodiffeq]`⁴ (Salvy and Zimmermann 1994). Such differential equations are only determined up to scaling by a polynomial in x . For example, scaling (18) by x yields another differential equation

$$x \cdot f'(x) + (2 \cdot x^2 - y \cdot x) \cdot f(x) = 0, \quad (19)$$

which is also satisfied by (16), as well as by some new but singular solutions (which can be found by Laplace-transform methods). To reduce this degree of freedom and eliminate such spurious solutions, we divide the differential equation by the leading coefficient $p_n(x)$, so that the coefficients become the rational functions $1, p_{n-1}(x)/p_n(x), \dots, p_0(x)/p_n(x)$. We then use the non-leading coefficients to identify the primitive measure.

For all our primitive measures (listed in Table 1), it turns out that $n = 1$, so we have to consider just one ratio $p_0(x)/p_1(x)$, which is $2 \cdot x - y$ for the density (16). In general, when $n = 1$, this ratio is equal to $-f'(x)/f(x)$, but the algorithm that computes this ratio is completely different from differentiating f then dividing by f and hoping that the factors that do not form a rational function cancel out. Rather, the algorithm achieves efficiency and robustness by performing linear algebra on *Ore algebra* elements, which represent linear differential operators and enjoy many closure properties that are efficiently computable (Salvy 2005).

Rational-function matching Thanks to existing efficient algorithms for normalizing rational functions (such as Euclid’s algorithm), we can easily and robustly match rational functions against each other, and thus recognize all our primitive measures in the same way. For example, given that $n = 1$ and the ratio $p_0(x)/p_1(x) = 2 \cdot x - y$ is linear in x for our target density (16), we look through the rightmost column of Table 1 for a rational function whose numerator’s degree in x is one higher than its denominator’s.⁵ The only candidate is $m = \text{Gaussian}(\mu, \sigma)$, and its integration bounds match those of (8), so we equate corresponding coefficients ($2 = 1/\sigma^2$ and $-y = -\mu/\sigma^2$) and solve to get $\mu = y/2$ and $\sigma = 1/\sqrt{2}$. This solving step can be done by calling Maple’s `solve` or by designing a custom matcher for each primitive measure. We prefer the latter as it is more efficient.

The recognize algorithm In general, `unintegrate` in Figure 3 tries to recognize a target density e by invoking `recognize(e, x, a, b, c)`. Here the density expression e may contain the integration variable x free, a and b are the bounds on x , and c is the context. The goal of `recognize` is to robustly satisfy the equation

$$\text{recognize}(\mathcal{D}(m, x) \cdot e', x, \mathcal{L}(m), \mathcal{U}(m), c) = (m, e'), \quad (20)$$

where m is a primitive measure term and e' does not contain x free. To do so, `recognize` proceeds through four steps:

⁴ The version of `gfun` shipped with Maple has several important bugs, fixed in the version at <http://perso.ens-lyon.fr/bruno.salvy/software/the-gfun-package/>

⁵ We only consider the difference between the two degrees, because the numerator and denominator may not be relatively prime, like for $\text{Beta}(1, \beta)$ and $\text{Beta}(\alpha, 1)$.

1. Convert e as a function of x to a holonomic representation.
2. Find the primitive measure m by rational-function matching under the assumption context c . Thus new measures can be added by extending Table 1.
3. Solve for the constant factor e' , by equating the target e and the matched $\mathcal{D}(m, x) \cdot e'$ at certain points x . To choose x , either use initial conditions returned by `gfun[holexprtodiffeq]`, if any, or default to $x \in \{0, 1/2, 1\}$.
4. Simplify e' using the assumption context c . Often the e' returned by one call to $(m, e') = \text{recognize}(e, x, a, b, c)$ is then passed to the next call as part of e , so it helps if we can eliminate tricky constructs such as `sgn` from e' .

If any step above fails, then `recognize` resorts to the following:

- if $-\infty < a < b < \infty$, then return $(\text{Uniform}(a, b), e_1 \cdot (b - a))$;
- otherwise, return $(\text{LO}(\hbar', \int_a^b \hbar'(x) dx), e_1)$, where \hbar' is fresh;

even though e_1 may contain x free. Returning `LO` is our unobtrusive way to admit failure, which is also done in the the catch-all case at the bottom of Figure 3 when g is not patently linear.

This approach to density recognition is general and modular enough so that it took one of us only a couple of hours to add each primitive measure in Table 1, once we had implemented the general infrastructure.

6 Improving linear operators algebraically

The previous two sections established a two-way bridge between measure terms and patently linear expressions. Patently linear expressions are real-valued expressions, and computer algebra gives us many more tools for these than for measure terms. We now put these tools to work.

6.1 What automatic simplification can and cannot do

To start with, every CAS incessantly performs so-called *automatic simplification* on every expression, using linear-time algorithms that are much faster than naive rewrite rules. Automatic simplification ensures that

- addition $+$ and multiplication \cdot are commutative and associative, with identities 0 and 1, so for example $e_1 \cdot (1 \cdot e_2 \cdot g)$ is equivalent to $(e_1 \cdot e_2) \cdot g$;
- repeated terms are collected, so for example $g_1 + (g_1 + g_2)$ becomes $2 \cdot g_1 + g_2$;
- arithmetic on rational numbers happens, so the factor $\frac{1}{1-0}$ in (13) disappears.

Our language of patently linear expressions take advantage of automatic simplification pervasively. In particular, in Figures 3 and 4, we pattern-match against $+$ and \cdot using commutativity and associativity. Hence just by composing steps 1 and 3 (Figures 2 and 3), we already simplify the measure term

$$\text{Weight}(e_1, \text{Weight}(1, \text{Weight}(e_2, \text{Msum}(m_1, \text{Msum}(m_1, m_2))))) \quad (21)$$

to $\text{Weight}(e_1 \cdot e_2, \text{Msum}(\text{Weight}(2, m_1), m_2))$.

This immediate success of automatic simplification might embolden us to feed patently linear expressions to Maple functions such as `value`, which triggers symbolic integration, and `simplify`, whose behaviour is not precisely defined (Moses 1971; Carette 2004). But such optimism would be misplaced. First, CASes tend to be ineffective at multidimensional integrals, which are typical of probabilistic programs. For example, Maple’s `value` can only simplify the double integral (5) (to (6)) if we reverse the order of integration between x and y . On the other hand, even though there is no hope to improve a nested integral of the form

$$\iiint \exp(x \cdot y \cdot z) \cdot h(x, y, z) \, dz \, dy \, dx \quad (22)$$

because h is arbitrary, the amount of time Maple’s `value` takes to return the integral unimproved explodes as the nesting depth increases.

Second, the nebulous notion of simplification that CASes use can wreck havoc on patently linear expressions. For example, consider this probabilistic program, which draws a number from a `Gaussian` then returns its absolute value:

$$\text{Bind}(\text{Gaussian}(0, 1), x, \text{If}(x < 0, \text{Ret}(-x), \text{Ret}(x))) \quad (23)$$

In step 1, we compute the corresponding linear operator

$$\text{LO} \left(h, \int_{-\infty}^{\infty} \frac{\exp(-\frac{x^2}{2})}{\sqrt{2 \cdot \pi}} \cdot \left(\begin{cases} h(-x) & \text{if } x < 0 \\ h(x) & \text{otherwise} \end{cases} \right) dx \right). \quad (24)$$

If we feed this linear operator straight to step 3, then we recover the measure term (23) unchanged. But if we feed it to Maple’s `simplify` or `value`, the Gaussian density becomes duplicated and unrecognized: we get the linear operators

$$\text{LO} \left(h, \int_{-\infty}^{\infty} \left(\begin{cases} \frac{\exp(-\frac{x^2}{2})}{\sqrt{2 \cdot \pi}} \cdot h(-x) & \text{if } x < 0 \\ \frac{\exp(-\frac{x^2}{2})}{\sqrt{2 \cdot \pi}} \cdot h(x) & \text{otherwise} \end{cases} \right) dx \right) \quad (25)$$

$$\text{LO} \left(h, \int_{-\infty}^0 \frac{\exp(-\frac{x^2}{2})}{\sqrt{2 \cdot \pi}} \cdot h(-x) \, dx + \int_0^{\infty} \frac{\exp(-\frac{x^2}{2})}{\sqrt{2 \cdot \pi}} \cdot h(x) \, dx \right), \quad (26)$$

which step 3 cannot express as measure terms without `LO`. Therefore, trying to improve a patently linear expression by feeding it willy-nilly to symbolic integration or simplification can be not only ineffective but even counter-productive.

6.2 Following the grammar of patently linear expressions

We achieve effective improvement by *controlled* application of symbolic integration and algebraic simplification. Informally, our idea is to apply these CAS facilities throughout a patently linear expression *except* where the grammar (14) calls for a patently linear subexpression. Intuitively, these places form the control-flow tree of a probabilistic program. For example, in (24) we withhold the conditional `If` subexpression from CAS facilities, and in (5) we withhold $h(y)$.

$$\begin{aligned}
\text{reduce}(\hbar, \int_a^b g \, dx, c) &= \begin{cases} \text{reduce}(\hbar, g', c) & \text{if no argument to } \hbar \text{ in } g \text{ contains } x \text{ free} \\ & \hbar' \text{ is fresh} \\ & g' = \text{banish}(\text{LO}(\hbar', \int_a^b \hbar'(x) \, dx), x, \hbar, g) \\ & g' \text{ contains fewer } \int \text{ signs than } \int_a^b g \, dx \text{ does} \\ 0 & \text{if } c' \text{ is inconsistent} \\ \text{if}(c'_1, \cdot \int_{a'}^{b'} \text{if}(c'_2, g', 0) \, dx, 0) & \text{if } c'_1 \wedge c'_2 \wedge a' < x < b' \text{ is equivalent to } c' \\ & c'_1 \text{ does not contain } x \text{ free} \end{cases} \\
\text{where } x \neq \hbar \text{ and } \text{Indicator}(c_1) \cdot \dots \cdot \text{Indicator}(c_n) \cdot g' &= \text{reduce}(\hbar, g, c \wedge a < x < b) \\
c' &= c_1 \wedge \dots \wedge c_n \wedge a < x < b \\
\text{reduce}(\hbar, g + \dots, c) &= \text{reduce}(\hbar, g, c) + \dots \\
\text{reduce}(\hbar, e \cdot g, c) &= \text{simplifyIf}(\text{simplify}(e, c)) \cdot \text{reduce}(\hbar, g, c) \\
&\quad \text{where } e \text{ does not contain } \hbar \text{ free} \\
\text{reduce}(\hbar, \text{If}(e, g, g'), c) &= \text{if}(\text{simplify}(e, c), \text{reduce}(\hbar, g, c \wedge e), \text{reduce}(\hbar, g', c \wedge \neg e)) \\
\text{reduce}(\hbar, \text{Integrate}(m, h), c) &= \text{Integrate}(m, \lambda x. \text{reduce}(\hbar, h(x), c)) \quad \text{where } x \text{ is fresh} \\
\text{reduce}(\hbar, g, c) &= \text{simplify}(g, c) \quad \text{otherwise} \\
\text{banish}(m, x, \hbar, g) &= (\text{do the integral in } \text{integrate}(m, \lambda x. 1)) \cdot g \\
&\quad \text{if } g \text{ does not contain } x \text{ free} \\
\text{banish}(m, x, \hbar, \int_a^b g \, dy) &= \int_a^b \text{banish}(m, x, \hbar, g) \, dy \\
&\quad \text{if } a, b \text{ do not contain } x, \hbar \text{ free and } x \neq \hbar \\
&\quad \text{where } m \text{ does not contain } y \text{ free and } x \neq y \\
\text{banish}(m, x, \hbar, g + \dots) &= \text{banish}(m, x, \hbar, g) + \dots \\
\text{banish}(m, x, \hbar, e \cdot g) &= \text{banish}(\text{Bind}(m, x, \text{Weight}(e, \text{Ret}(x))), x, \hbar, g) \\
&\quad \text{where } e \text{ does not contain } \hbar \text{ free} \\
\text{banish}(m, x, \hbar, \text{If}(e, g, g')) &= \begin{cases} \text{banish}(\text{Bind}(m, x, \text{If}(e, \text{Ret}(x), \text{Msum}())), x, \hbar, g) + \\ \text{banish}(\text{Bind}(m, x, \text{If}(e, \text{Msum}(), \text{Ret}(x))), x, \hbar, g') & \text{if } e \text{ contains } x \text{ free} \\ \text{If}(e, \text{banish}(m, x, \hbar, g), \text{banish}(m, x, \hbar, g')) & \text{otherwise} \end{cases} \\
\text{banish}(m, x, \hbar, \text{Integrate}(n, h)) &= \text{Integrate}(n, \lambda y. \text{banish}(m, x, \hbar, h(y))) \\
&\quad \text{if } n \text{ does not contain } x \text{ free} \\
&\quad \text{where } m \text{ does not contain } \hbar \text{ free and } y \text{ is fresh} \\
\text{simplifyIf}(\text{If}(c, e, e')) &= \text{if}(c, \text{simplifyIf}(e), \text{simplifyIf}(e')) \quad \text{if } (c, e, 0) = \text{Indicator}(c) \cdot e \\
\text{simplifyIf}(e) &= e \quad \text{otherwise} \quad \text{if } (c, e, e') = \text{If}(c, e, e') \text{ otherwise}
\end{aligned}$$

Fig. 4. Improving a linear operator $\text{LO}(\hbar, g)$ to $\text{LO}(\hbar, \text{reduce}(\hbar, g, \text{true}))$. Some cases of banish and if are relegated to Figure 5 in our technical report (Carette and Shan 2015).

Our workhorse is the `reduce` function defined in Figure 4. Starting with a linear operator $\text{LO}(\hbar, g)$, the goal of $\text{reduce}(\hbar, g, c)$ is to improve g under the assumption context c (initially `true`, as in Section 5.2). The `reduce` function exercises control over CAS facilities by following the grammar (14). To start with, `reduce` only invokes $\text{simplify}(e, c)$ in three places, namely where the grammar (14) calls for a general e rather than some patently linear g . The call $\text{simplify}(e, c)$ means to simplify the expression e under the context c (using Maple’s `simplify`).

When `reduce` encounters an integral $\int_a^b g dx$, it first checks if the integration variable x is used in any call to the arbitrary integrand \hbar (in other words, if x is used to Return any outcome). If not, then the variable x could be *eliminated* (Dechter 1998) or *integrated out*. An example is the variable x in (5). To perform the integral over x , we push it in past the integral over y , to form the expression

$$\int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} \frac{\exp(-\frac{x^2}{2})}{\sqrt{2 \cdot \pi}} \cdot \frac{\exp(-\frac{(y-x)^2}{2})}{\sqrt{2 \cdot \pi}} dx \right) \cdot h(y) dy, \quad (27)$$

then ask Maple to do just the inner integral over x , which no longer contains any call to h . In general, we push the integral to be eliminated all the way in, so that it becomes the innermost integral and does not contain any call to \hbar , then ask Maple to do just that integral. This pushing and symbolic integration is performed by `banish`. The specification of `banish` is that $\text{LO}(\hbar, \text{banish}(m, x, \hbar, g))$ should mean the same measure as $\text{Bind}(m, x, \text{LO}(\hbar, g))$ but integrate over x innermost rather than outermost. Achieving this specification requires doubling the work when g is an `If`. The upshot of this improvement is that we manage to simplify the measure term (4) to $\text{Gaussian}(0, \sqrt{2})$.

If the integral $\int_a^b g dx$ cannot be eliminated altogether, then `reduce` tries to shrink the integral’s bounds a, b by checking g to see if it looks like $\text{If}(a' < x < b', g', 0)$ or $\text{If}(a' < x < b', 1, 0) \cdot g'$. If it does, then the recursive call to $\text{reduce}(\hbar, g, c \wedge a < x < b)$ would return a product like $\text{Indicator}(a' < x < b') \cdot g'$. The meaning of $\text{Indicator}(\dots)$ is same as $\text{If}(\dots, 1, 0)$, but we introduce a separate constructor `Indicator` for this particular use of `If` to encode a domain restriction that `reduce` should incorporate into the integration bounds. The constructor `Indicator` is called in the appropriate cases by the smart constructor `if`. Any domain restriction so gathered is then used by `reduce` to shrink the integration bounds. The upshot of this improvement is that we manage to simplify the measure terms

$$\text{Bind}(\text{Uniform}(0, 1), x, \text{If}(0 < x < 1/2, \text{Ret}(x), \text{Msum}())) \quad (28)$$

$$\text{Bind}(\text{Uniform}(0, 1), x, \text{Weight}(\text{If}(0 < x < 1/2, 1, 0), \text{Ret}(x))) \quad (29)$$

both to $\text{Weight}(1/2, \text{Uniform}(0, 1/2))$.

Besides these improvements that constitute `reduce`, we have found two other operations on patently linear expressions to be sometimes useful: pushing an integral inward but not all the way in; and reparameterizing an integral so that the integration variable is what’s passed to the arbitrary integrand \hbar . These operations are not always beneficial, so `reduce` does not perform them automatically; rather, human or heuristic guidance is required to invoke them for now.

7 Evaluation

Through our work on probabilistic programming, we have assembled a test suite with 66 input-output pairs of measure terms, including probabilistic models and inference procedures that arise in practice. Our simplifier passes all of these tests. A referee asked if some heuristics had crept in; they had, in droves, in a previous simplifier. We present a rational reconstruction, where the only heuristics are the ones documented in Section 5.

Many tests have the same input and output, because our simplifier should not degrade any term, especially one that is already as simple as can be. As explained in Section 6, this requirement—although basic—is not trivial to satisfy. Indeed, if we only apply step 1 followed by step 3, 42 of our tests still pass.

Our test suite can be seen online⁶ alongside our implementation.⁷ The following ablation results give a sense of our tests’ coverage and variety. If we remove `banish`, 14 tests fail. If we remove assumption contexts, 8 tests fail. If we remove calls to Maple’s `simplify`, 8 tests fail. If we remove shrinking integration bounds, 9 tests fail. If we remove density recognition from step 3, fully 53 tests fail.

8 Related work

Statisticians have long used CASes to analyze the expected values, densities, moments, and other properties of particular distributions (Andrews and Stafford 2000). For modularity, however, we need to express distributions by *composing* building blocks, especially using monadic `Bind`. Researchers handling such distributions have long acknowledged that tasks like exact inference (Dechter 1998), lifted inference (de Salvo Braz et al. 2005), and sampler code generation (Tristan et al. 2014) amount to computer algebra. By bringing an off-the-shelf CAS to bear on these tasks, and by producing output in the same language as the input, we hope to have laid the groundwork for more robustness and reuse. For example, once we defined the densities and density recognition for `Gaussian` and `Gamma`, the handful of conjugacies among those distributions all fall out from the CAS.

Compared to our earlier (unpublished) attempts at simplifying probabilistic programs using a CAS, this version is quite compact: 775 lines of Maple code (including comments and testing infrastructure). More importantly, the design of this version can be explained. Especially effective and robust are our conversion between measure terms and linear operators, and our use of differential equations for density recognition. These techniques pave the way for two pressing tasks:

- Add arrays. This will let us express many more distributions in wide use. However, it is a wide-open problem how to deal with product measures and their linear-operator counterparts, especially when their length is symbolic.
- Recover domain restrictions from more advanced uses of `If`, such as multivariate conditions. Potential solutions include SMT solvers and methods based on normal forms for inequalities.

⁶ <https://github.com/hakaru-dev/hakaru/blob/master/maple/NewSLOTests.mpl>

⁷ <https://github.com/hakaru-dev/hakaru/blob/master/maple/NewSLO.mpl>

References

- Sergei A. Abramov and Marko Petkovšek. Gosper’s algorithm, accurate summation, and the discrete Newton-Leibniz formula. In *Proceedings of the International Symposium on Symbolic and Algebraic Computation (ISSAC)*, pages 5–12. ACM Press, 2005.
- David F. Andrews and James E. H. Stafford. *Symbolic Computation for Statistical Inference*. Oxford University Press, 2000.
- Jacques Carette. Understanding expression simplification. In *Proceedings of the International Symposium on Symbolic and Algebraic Computation (ISSAC)*, pages 72–79. ACM Press, 2004.
- Jacques Carette. A canonical form for piecewise defined functions. In *Proceedings of the International Symposium on Symbolic and Algebraic Computation (ISSAC)*, pages 77–84. ACM Press, 2007.
- Jacques Carette and Chung-chieh Shan. Simplifying probabilistic programs using computer algebra. Technical Report 719, Indiana University, 2015. URL <http://www.cs.indiana.edu/cgi-bin/techreports/TRNNN.cgi?trnum=TR719>.
- Frédéric Chyzak and Bruno Salvy. Non-commutative elimination in Ore algebras proves multivariate holonomic identities. *Journal of Symbolic Computation*, 26(2):187–227, 1998.
- Olivier Danvy and Andrzej Filinski. Abstracting control. In *Proceedings of the Conference on LISP and Functional Programming*, pages 151–160. ACM Press, 1990.
- Rodrigo de Salvo Braz, Eyal Amir, and Dan Roth. Lifted first-order probabilistic inference. In *Proceedings of the 19th International Joint Conference on Artificial Intelligence*, pages 1319–1325, 2005.
- Rina Dechter. Bucket elimination: A unifying framework for probabilistic inference. In Michael I. Jordan, editor, *Learning and Inference in Graphical Models*. Kluwer, Dordrecht, 1998. Paperback: *Learning in Graphical Models*, MIT Press.
- Michèle Giry. A categorical approach to probability theory. In Bernhard Banaschewski, editor, *Categorical Aspects of Topology and Analysis: Proceedings of an International Conference Held at Carleton University, Ottawa, August 11–15, 1981*, pages 68–85. Springer, 1982.
- John Hatcliff and Olivier Danvy. A generic account of continuation-passing styles. In *Proceedings of the 21st Symposium on Principles of Programming Languages (POPL)*, pages 458–471. ACM Press, 1994.
- Manuel Kauers. The holonomic toolkit. In Carsten Schneider and Johannes Blümlein, editors, *Computer Algebra in Quantum Field Theory*, pages 119–144. Springer, 2013.
- Peter S. Maybeck. *Stochastic Models, Estimation, and Control*. Academic Press, 1979.
- Joel Moses. Algebraic simplification: A guide for the perplexed. *Communications of the ACM*, 14(8):548–560, 1971.
- David Pollard. *A User’s Guide to Measure Theoretic Probability*. Cambridge University Press, 2001.

- Norman Ramsey and Avi Pfeffer. Stochastic lambda calculus and monads of probability distributions. In *Proceedings of the 29th Symposium on Principles of Programming Languages (POPL)*, pages 154–165. ACM Press, 2002.
- Bruno Salvy. D-finiteness: Algorithms and applications. In *Proceedings of the International Symposium on Symbolic and Algebraic Computation (ISSAC)*, pages 2–3. ACM Press, 2005.
- Bruno Salvy and Paul Zimmermann. Gfun: a Maple package for the manipulation of generating and holonomic functions in one variable. *ACM Transactions on Mathematical Software*, 20(2):163–177, 1994.
- Jean-Baptiste Tristan, Daniel Huang, Joseph Tassarotti, Adam C. Pockock, Stephen J. Green, and Guy Lewis Steele, Jr. Augur: Data-parallel probabilistic modeling. In *Advances in Neural Information Processing Systems*, pages 2600–2608, 2014.
- Herbert S. Wilf and Doron Zeilberger. An algorithmic proof theory for hypergeometric (ordinary and “q”) multisum/integral identities. *Inventiones mathematicae*, 108:557–633, 1992.