

Continuous Probabilistic Programming¹

Steven Holtzen

s.holtzen@northeastern.edu

October 17, 2023

¹ CS7470 Fall 2023: Foundations of Probabilistic Programming.

Some logistics:

- Friday Oct. 20 is Systems day: Please sign up to give a presentation on a system, and be sure to post your slides ahead of time in the Slack
- I will release and adjust the deadline for the continuous language this weekend based on what we've covered. It should be a pretty simple project.
- We will begin reading papers next week. See the instructions and consult the syllabus for which papers to read.

1 CONT: A tiny continuous PPL

Syntax:

```
1 v ::= real | bool
2 τ ::= ℬ | ℝ | Dist(ℝ)
3 e ::= unif | x ← e | x | v | e + e | e < e | e ∨ e | e ∧ e | ¬ e
```

- As usual, our recipe will have two parts: new denotation, new operational semantics. We begin with denotation.
- **Goal:** give a denotation for CONT that maps terms to probability distributions
- The `unif` keyword denotes a *uniform probability distribution on the interval* $[0, 1]$. This means we need to be able to handle probability distributions on the interval.
- For example, we can write the following program:

```
1 x ← unif;
2 return x < 1/2
```

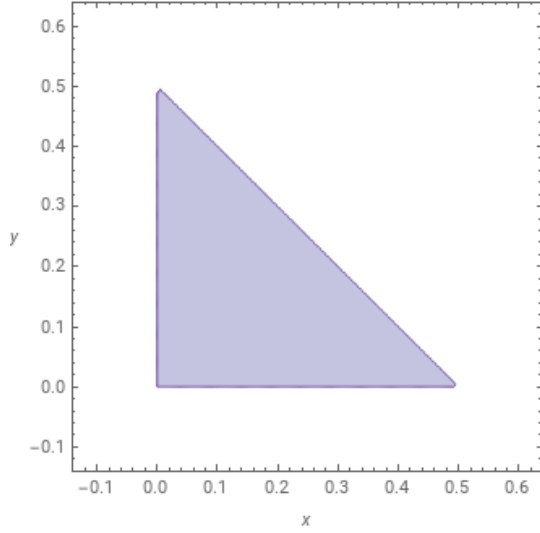
How should we interpret this program? Intuitively, this program should return a distribution:

$$[\text{true} \mapsto 1/2, \text{false} \mapsto 1/2]$$

- What about this more interesting program:

```
1 x ← unif;
2 y ← unif;
3 return x + y < 1/2
```

This is significantly more challenging! The probability that this returns true is equal to the area of this triangle:



Put differently, it is an integral:

$$\int_0^1 \int_0^1 \mathbb{1}(x+y < 1/2) \, dx \, dy = 1/8 \quad (1)$$

Probability distributions on the interval

- **Goal:** Define a probability distribution on the sample space $\Omega = [0, 1]$
- Recall the challenge from last time: we cannot define a distribution naively on points because the probability of any particular real value $\omega \in \Omega$ should be 0.²
- Hence, we need to consider a subset of well-behaved semi-algebra of events on which to define a probability measure. A standard choice are the *collection of all open intervals* contained in $[0, 1]$:

$$\mathcal{I} = \{(l, u) \mid 0 \leq l < u \leq 1\} \cup \{\emptyset, \Omega\}$$

Why is this a good choice of events? Because we know a good way to define their probability:

$$\mu((l, u)) = u - l \quad \mu(\emptyset) = 0 \quad \mu(\Omega) = 1 \quad (2)$$

- The **Borel σ -algebra** is the smallest σ -algebra that contains \mathcal{I} . We denote the Borel σ -algebra as \mathcal{B} .³

For an excellent discussion on the development of this section, see Axler [2020].

² Why is this the case? Assume that there is a probability distribution $\mu : \Omega \rightarrow [0, 1]$ that assigns a non-zero probability to a countably-infinite set of points in $\omega \in \Omega$. Then, $\sum_{\omega \in \Omega} \mu(\omega) = \infty$, which is a contradiction.

³ See Definition 2.27 of Axler [2020] for a discussion of why the notion of a smallest σ -algebra is well-founded.

- The unique probability measure $\lambda : \mathcal{B} \rightarrow [0, 1]$ given by the extension of Eq. (2) to \mathcal{B} is called the **Lebesgue measure** on the interval. For example:

$$\lambda \left(\left(0, \frac{1}{2}\right) \cup \left(\frac{3}{4}, 1\right) \right) = 1/2 + 1/4 = 3/4.$$

- Now we will need a new generalized notion of random variable for this new setting. A random variable will be defined as a well-behaved function out of the sample space:

Definition 1 (Measurable function). *Let $(\Omega, \mathcal{F}, \mu)$ be a probability space, (X, \mathcal{A}) be a measurable space, and $f : \Omega \rightarrow X$ be a function. The function f is called measurable if, for any $A \in \mathcal{A}$, it is the case that $f^{-1}(A) \in \mathcal{F}$.*

Measurable functions $f : [0, 1] \rightarrow \mathbb{R}$ defined on $([0, 1], \mathcal{B}, \lambda)$ are often called *Borel-measurable*. Most functions you have seen are Borel-measurable.⁴

⁴ For instance, all continuous functions are Borel measurable.

- Measurable functions behave like random variables in that they induce probability measures on their co-domains via a *push-forward*:

Definition 2. *Let $(\Omega, \mathcal{F}, \mu)$ be a probability space, (X, \mathcal{A}) be a measurable space, and $f : \Omega \rightarrow X$ be a measurable variable. Then, the push-forward probability measure $\nu : \mathcal{A} \rightarrow [0, 1]$ is defined:*

$$\nu(A) = \mu(f^{-1}(A)). \quad (3)$$

2 Denotation of `CONT`

- These semantics are in the style of *Markov kernels* [Kozen, 1979, Barthe et al., 2020]
- The case for pure terms is relatively simple
- First we give the semantics:

The semantics of $x \leftarrow e_1; e_2$ is the bind operator for the Giry monad [Giry, 2006, Ramsey and Pfeffer, 2002]. This definition involves a *generalized integral* wrt. a base measure.

$$\begin{aligned} \llbracket \text{unif} : \text{Dist}(\mathbb{R}) \rrbracket (A) &= \lambda(A) \\ \llbracket \text{return } e : \text{Dist}(\mathbb{R}) \rrbracket (A) &= \mathbb{1}(\llbracket e \rrbracket \in A) = \delta_{\llbracket e \rrbracket} \\ \llbracket x \leftarrow e_1; e_2 \rrbracket (A) &= \mathbb{E}_{v \sim \llbracket e_1 \rrbracket} \left[\llbracket e_2[v/x] \rrbracket \right] = \int_{\mathbb{R}} \llbracket e_2[v/x] \rrbracket (A) \, d \llbracket e_1 \rrbracket (v) \end{aligned}$$

The *Diract-delta* $\delta_v(A)$ denotes a distribution that assigns a probability of 1 to any event that contains $v \in \Omega$.

- Integration wrt. the Lebesgue base measure is written $\int f(x) \, d\lambda(x)$. It behaves just like Riemann integration (integrals from high school) if f is a suitably well-behaved function (i.e., continuous).
- Example:

$$\begin{aligned}
\llbracket x \leftarrow \text{unif}; \text{return } x + x \rrbracket((0, 1/2)) &= \int_0^1 \llbracket \text{return } v + v \rrbracket(v)((0, 1/2)) \, d\lambda(v) \\
&= \int_0^1 \mathbb{1}(v + v \in (0, 1/2)) \, d\lambda(v) \\
&= \int_0^{1/4} 1 \, d\lambda(v) \\
&= 1/4
\end{aligned}$$

- Integration wrt. the Dirac measure δ_v is defined:

$$\int_{\mathbb{R}} f(x) \, d\delta_v(x) \triangleq f(v). \quad (4)$$

This lets us interpret programs with return:

$$\begin{aligned}
\llbracket x \leftarrow \text{return } 1; \text{return } x + x \rrbracket(A) &= \int_{\mathbb{R}} \llbracket \text{return } v + v \rrbracket(A) \, d\delta_1(v) \\
&= \int_{\mathbb{R}} \delta_{v+v}(A) \, d\delta_1(v) \\
&= \delta_2(A).
\end{aligned}$$

3 Sampling semantics for CONT

- These are very similar the sampling semantics for DISC from last time, except we need a different source of random bits
- We will entropy $\sigma \in [0, 1]^{\mathbb{N}}$, i.e. σ is a countably infinite stream of real values drawn from the interval. We again define π_L and π_R as ways of splitting the entropy into independent sets.
- Then we will define a big-step relation $\sigma \vdash e \Downarrow v$:

$$\begin{array}{c}
v :: \sigma \vdash \text{unif} \Downarrow v \quad \frac{\pi_L(\sigma) \vdash e_1 \Downarrow v_1 \quad \pi_R(\sigma) \vdash e_2[v_1/x] \Downarrow v_2}{\sigma \vdash x \leftarrow e_1; e_2 \Downarrow v_2} \\
\sigma \vdash \text{return } v \Downarrow v
\end{array}$$

Theorem 1. Assume e is a well-typed closed term and Pr is the distribution on the entropy space. Then, for any event $A \in \mathcal{B}$, we have that:

$$\mathbb{E}_{\sigma \sim \text{Pr}} [\mathbb{1}[\text{eval}(e, \sigma) \in A]] = \llbracket e \rrbracket(A). \quad (5)$$

4 Continuous observation

Add a new observe keyword to `CONT`:

```

1 v ::= real | bool
2  $\tau$  ::=  $\mathbb{B}$  |  $\mathbb{R}$  |  $\text{Dist}(\mathbb{R})$ 
3 e ::= unif |  $x \leftarrow e$  |  $x \mid v$  |  $e + e$  |  $e * e$  |  $e < e$  |  $e \vee e$  |  $e \wedge e$  |  $\neg e$ 
4   | observe e; e

```

This keyword is surprisingly powerful. Suppose we want to infer something about the distribution of people’s heights. Suppose we assume that people are between 65 and 73 inches tall, and that people vary from the true value by 1 inch. Then, we observe that someone’s height is between 71 and 73 inches, and want to update prior beliefs about how people’s heights are distributed. We can model this as the following program:

```

1 h  $\leftarrow$  unif;
2 height_dist  $\leftarrow$  (h * 4) + 69
3 eps1  $\leftarrow$  unif;
4 observe 71 < height_dist + eps1 < 73;
5 return height_dist

```

References

- Sheldon Axler. *Measure, integration & real analysis*. Springer Nature, 2020.
- Gilles Barthe, Joost-Pieter Katoen, and Alexandra Silva. *Foundations of probabilistic programming*. Cambridge University Press, 2020.
- Michele Giry. A categorical approach to probability theory. In *Categorical Aspects of Topology and Analysis: Proceedings of an International Conference Held at Carleton University, Ottawa, August 11–15, 1981*, pages 68–85. Springer, 2006.
- Dexter Kozen. Semantics of probabilistic programs. In *20th Annual Symposium on Foundations of Computer Science (sfcs 1979)*, pages 101–114. IEEE, 1979.
- Norman Ramsey and Avi Pfeffer. Stochastic lambda calculus and monads of probability distributions. In *Proceedings of the 29th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 154–165, 2002.