
Equational reasoning for conditioning as disintegration

Chung-chieh Shan
ccshan@post.harvard.edu

Dylan Thurston
dpt@bostoncoop.net

Conditional distributions are widely used for practical inference, even when the condition has zero probability (such as setting a continuous variable to an observed value). This popularity contrasts with the scary pitfalls (such as Borel’s paradox) that beset rigorous treatments of conditioning. In general, conditional expectations may arise that do not correspond to any conditional distribution at all. This gap between theory and practice has made it difficult to automate conditioning calculations soundly. In particular, given a generative model, we want to produce its conditional expectations mechanically, so as to optimize them numerically, simplify them symbolically, and so on.

Disintegration [1] is a rigorous approach to conditioning that covers a wide range of applications yet admits intuitive conditional distributions and allows their ‘guilt-free manipulation’. In the present work, we mechanize this approach by adding a ‘Lebesgue measure’ operation to the usual monadic representation of stochastic experiments [2]. We show how to compute conditional distributions by equating expressions in this representation.

By ‘compute’, we mean producing a symbolic algebraic expression usually and giving up in failure occasionally, not approximating a real number to arbitrary precision. In fact, the latter task is impossible for conditional probabilities in general [3]. As for the former task, a system for computing conditional distributions that guarantees any expression it produces to be correct in some sense (explained below) is useful because we can feed a mathematical expression to existing software such as computer algebra systems, numerical solvers, and optimized code generators.

We use the following stochastic experiment as a running example. It flips a fair coin c to decide whether to scale a uniform random number x by 1 or by 3.

$$\begin{aligned} c &\sim \text{Uniform}[0, 1] \\ s &= \text{if } c < 1/2 \text{ then } 1 \text{ else } 3 \\ x &\sim \text{Uniform}[-1, 1] \\ y &= sx \end{aligned} \tag{1}$$

Suppose we observe y to be 0. The propositions $y = 0$ and $x = 0$ are logically equivalent, so they are one and the same event. When we condition on $x = 0$, the posterior marginal distribution of s is of course concentrated on $s = 1$ and $s = 3$, each with probability $1/2$. However, when we condition on $y = 0$, the posterior should be $s = 1$ with probability $3/4$ and $s = 3$ with probability $1/4$. This example is thus a simplified version of Borel’s paradox.

We can think of this experiment as defining a measure μ . We can then compute the expectation of a ‘test function’ f from c, s, x, y to reals:

$$E(f) = \int f(c, s, x, y) \mu = \int_0^1 \left(\left(\int_{-1}^1 \frac{1}{2} (f(c, s, x, y)|_{y=sx}) dx \right) \Big|_{s=\text{if } c < 1/2 \text{ then } 1 \text{ else } 3} \right) dc. \tag{2}$$

Despite its clumsy notation, this formula corresponds directly to (1): it is an integral over c of a (Dirac delta) integral over s of an integral over x of a (Dirac delta again) integral over y .

If there are any distributions conditional on y for us to speak of, they had better be a family μ_y of measures over c, s, x , parametrized by y . We do not even require these measures to be normalized. The basic idea of disintegration is that we want μ_y to satisfy

$$E(f) = \int \left(\int f(c, s, x, y) \mu_y \right) \nu \tag{3}$$

for all measurable f . The outer integral is with respect to some measure ν over y , which may be normalized like the marginal distribution over y or unnormalized like the standard Lebesgue measure over y . The inner integral is with respect to the measure μ_y over c, s, x , which again may be normalized or not. Note that we want to find μ_y independently of any particular test function f , so that we may call μ_y a conditional distribution, and not just call $\int f(c, s, x, y) \mu_y$ a conditional expectation. That is one reason disintegration is a stronger notion of conditioning than Kolmogorov's, and one reason μ_y does not always exist (and when it exists, is unique only up to an almost-sure equivalence).

The requirement (3) can be understood intuitively in several different ways:

1. If the original probability measure μ can be chopped into thin slices—one for each y —then the slices can be used as unnormalized conditional distributions.
2. If the original stochastic experiment can be rearranged so that y is chosen first (by sampling from ν), without affecting the ultimate measure μ , then we can describe a 'conditional stochastic experiment' by replacing that first step to choose y with a step that clamps y at the given value.
3. A conditional probability density can be defined as the ratio of a joint density to a marginal density, even when the marginal density is not with respect to the Lebesgue measure over y , but with respect to any measure ν . For instance, if ν is the marginal distribution over y , then the marginal density is simply 1.

The second intuition above is especially useful, because it turns out that the 'monadic' representation exemplified by (1) is a convenient format in which to reason about equality of measures. To represent unnormalized measures in the same format, we add two kinds of operations that go beyond describing steps in a stochastic experiment. First, we write ' $x \sim \text{Lebesgue}$ ' to denote the Lebesgue measure over x (which corresponds to the integration $\int_{-\infty}^{\infty} \dots dx$). Second, if δ is a scalar expression then we write 'weight δ ' to scale the integral by the density δ . For example, the expressions

$$x \sim \text{Uniform}[-1, 1] \quad \text{and} \quad x \sim \text{Lebesgue} \quad \text{weight (if } -1 < x < 1 \text{ then } 1/2 \text{ else } 0) \quad (4)$$

are equivalent, because the integrals

$$\int_{-1}^1 \frac{1}{2} f(x) dx \quad \text{and} \quad \int_{-\infty}^{\infty} (\text{if } -1 < x < 1 \text{ then } 1/2 \text{ else } 0) f(x) dx \quad (5)$$

are equal. Also useful is the fact that the expressions

$$y \sim \text{Lebesgue} \quad \text{and} \quad z \sim \text{Lebesgue} \quad (6) \\ z = x + y \quad \quad \quad y = z - x$$

are equivalent. Finally, if $s > 0$ then the expressions

$$x \sim \text{Lebesgue} \quad \text{and} \quad y \sim \text{Lebesgue} \quad (7) \\ y = sx \quad \quad \quad x = y/s \\ \quad \quad \quad \text{weight } 1/s$$

are equivalent.

Now, to condition our running example (1) on y , we reason equationally as follows:

$$\begin{aligned} (1) &= \begin{array}{l} c \sim \text{Uniform}[0, 1] \\ s = \text{if } c < 1/2 \text{ then } 1 \text{ else } 3 \\ x \sim \text{Lebesgue} \\ \text{weight (if } -1 < x < 1 \text{ then } 1/2 \text{ else } 0) \\ y = sx \end{array} \quad \left. \vphantom{\begin{array}{l} c \sim \text{Uniform}[0, 1] \\ s = \text{if } c < 1/2 \text{ then } 1 \text{ else } 3 \\ x \sim \text{Lebesgue} \\ \text{weight (if } -1 < x < 1 \text{ then } 1/2 \text{ else } 0) \\ y = sx \end{array}} \right\} \text{by (4)} \\ &= \begin{array}{l} c \sim \text{Uniform}[0, 1] \\ s = \text{if } c < 1/2 \text{ then } 1 \text{ else } 3 \\ x \sim \text{Lebesgue} \\ y = sx \\ \text{weight (if } -1 < x < 1 \text{ then } 1/2 \text{ else } 0) \end{array} \quad \left. \vphantom{\begin{array}{l} c \sim \text{Uniform}[0, 1] \\ s = \text{if } c < 1/2 \text{ then } 1 \text{ else } 3 \\ x \sim \text{Lebesgue} \\ y = sx \\ \text{weight (if } -1 < x < 1 \text{ then } 1/2 \text{ else } 0) \end{array}} \right\} \text{exchange scalar with (Dirac delta) integral} \end{aligned}$$

- [5] Sooraj Bhat, Ashish Agarwal, Richard Vuduc, and Alexander Gray. A type theory for probability density functions. In *POPL '12: Symposium on Principles of Programming Languages*, January 2012.