

Trustworthy, Useful Languages for Probabilistic Modeling and Inference

Neil Toronto

A dissertation submitted to the faculty of
Brigham Young University
in partial fulfillment of the requirements for the degree of
Doctor of Philosophy

Jay McCarthy, Chair
Kevin Seppi
Chris Grant
Eric Mercer
Dan Olsen

Department of Computer Science

Brigham Young University

February 2014

Copyright © 2014 Neil Toronto

All Rights Reserved

ABSTRACT

Trustworthy, Useful Languages for Probabilistic Modeling and Inference

Neil Toronto

Department of Computer Science, BYU

Doctor of Philosophy

The ideals of exact modeling, and of putting off approximations as long as possible, make Bayesian practice both successful and difficult. Automating it has been taking the form of languages for modeling probabilistic processes and implementations that answer questions about them.

Unfortunately, very few of these languages have a mathematical specification. This makes them difficult to trust: there is no way to distinguish between a bug and a feature, and there is no standard by which to prove optimizations correct. Further, because the languages are based on the incomplete theories of probability typically used in Bayesian practice, they place seemingly artificial restrictions on legal programs and queries, such as disallowing unbounded recursion and allowing only simple equality constraints.

We prove it is possible to make trustworthy probabilistic languages for Bayesian practice by using functional programming theory to define them mathematically and prove them correct. The specifications interpret programs using the most complete known theory of probability, so they place no restrictions on programs and allow arbitrary conditions with arbitrarily low probability.

We demonstrate that the results are useful by implementing the languages, and using them to model and answer questions about typical probabilistic processes. We also model and answer questions about processes that are either difficult or impossible to reason about precisely using Bayesian mathematical tools.

Keywords: Bayesian, Probability, Domain-Specific Languages, Functional Programming, Semantics

To my wife, Amy.

ACKNOWLEDGMENTS

I must first thank my parents for bringing me up right.

I must secondly thank my parents again, for encouraging my curiosity and intellectual development. Only now do I hear stories of the difficulties they had finding age-appropriate computer accessories and programming books for my Christmas presents—in the mid-1980s, when such things were scarce and few were made for children.

I must thirdly thank my parents yet again, for trying to instill in me a love of hard work. It eventually took. I assure them it was worth it.

My PhD advisor, Jay McCarthy, has been stellar in every respect. He has shown me how to do programming language research with rigor, taught me how to structure my communication, and offered excellent insights and suggestions on research that crosses into a field he initially knew little about. Without Jay's guidance, there would be little to distinguish this work from other work in probabilistic languages. Together, we have done it in a way that is demonstrably the most complete and reliable.

I cannot confine my thanks only to Jay's professional instruction. He has also ensured that I could take care of my family without going into debt, and ensured that I have the connections necessary to secure a good career. The wait for “when Daddy gets a real job” has been bearable, with no ever-increasing financial load and with good prospects to look forward to. The positive repercussions will be nearly endless. I am eternally grateful.

My wife has constantly supported my academic pursuits for 11 years. This is truly marvellous in a way that words are too narrow to contain.

I thank my children, for praying for me every day, and for being sweet, loving, talented, and so very interesting.

Lastly, I thank my God and Savior, whose justice and mercy makes our existence meaningful and beautiful.

Table of Contents

List of Figures	ix
1 Thesis	3
1.1 Introduction	3
1.2 Terms	4
1.3 Proof and Supporting Evidence	5
1.4 Exposition Transition System	6
2 Background	9
2.1 Bayesian Practice	10
2.2 Functional Programming Theory	22
3 Computing in Cantor’s Paradise With λ_{ZFC}	41
3.1 Motivation	41
3.2 Language Tower and Terminology	43
3.3 Metalanguage: First-Order Set Theory	44
3.4 λ_{ZFC} ’s Grammar	49
3.5 λ_{ZFC} ’s Big-Step Reduction Semantics	53
3.6 Syntactic Sugar and a Small Set Library	56
3.7 Example: The Reals From the Rationals	58
3.8 Example: Computable Real Limits	60
3.9 Related Work	64
3.10 Conclusions	66

4	Using λ_{ZFC}	67
4.1	Computations and Values	67
4.2	Auxiliary Type Systems	68
4.3	Using ZFC Values and Theorems	69
4.4	Internal Equality and External Equivalence	70
4.5	Additional Functions and Syntactic Forms	71
4.6	Extensional Functions	71
5	Countable Models and Implementation	73
5.1	Introduction	73
5.2	The Expression Language	75
5.3	The Query Language	78
5.4	Conditional Queries	82
5.5	The Statement Language	84
5.6	Why Separate Statements and Queries?	91
5.7	Conclusions	92
6	Interlude: Uncountable Outcomes and Recursion	93
7	Preimage Computation Theory: Running Programs Backwards	95
7.1	Introduction	95
7.2	Arrows and First-Order Semantics	98
7.3	The Bottom Arrow	105
7.4	Deriving the Mapping Arrow	106
7.5	Lazy Preimage Mappings	112
7.6	Deriving the Preimage Arrow	117
7.7	Preimages Under Partial, Probabilistic Functions	122
7.8	Output Probabilities and Measurability	134
7.9	Approximating Semantics	135

7.10	Implementations	144
7.11	Conclusions	145
8	Preimage Computation Implementation	147
8.1	Set Representation	147
8.2	Preimages Under Real Functions	147
8.3	Primitive Implementation	159
8.4	Partitioned Sampling	159
9	Measurability Theorems	175
9.1	Basic Definitions	175
9.2	Measurable Pure Computations	176
9.3	Measurable Probabilistic Computations	180
9.4	Measurable Projections	183
10	Sampling Theorems	185
10.1	Basic Definitions	185
10.2	Sampling Proofs	194
11	Related Work	199
11.1	Implementations	199
11.2	Semantics	200
11.3	Techniques	201
11.4	Somewhat Related Work	202
12	Conclusions and Future Work	203
12.1	Conclusions	203
12.2	Future Work	204
	References	207

List of Figures

1.1	A transition system for reading this dissertation	7
2.1	Computing probabilities using the standard normal density function	15
2.2	Joint density model plot	16
2.3	Bayes' law for densities, in pictures	19
2.4	Big-step operational semantics example	26
2.5	Implementation of the big-step semantics	28
2.6	Big-step operational semantics with nondeterminism	28
2.7	Denotational semantics and implementation	30
2.8	Denotational semantics with nondeterminism	31
2.9	Categorical semantics with nondeterminism	34
2.10	Abstract semantics with nondeterminism	38
3.1	Definition of λ_{ZFC}^-	49
3.2	Semantic function $\mathcal{F}[\![\cdot]\!]$	50
3.3	λ_{ZFC} 's grammar	53
3.4	λ_{ZFC} 's semantics	54
5.1	Random variable expression semantics	77
5.2	Implementation of $\mathcal{R}[\![\cdot]\!]$	78
5.3	Implementation of finite approximation and distribution queries	82
5.4	State monad functions for queries and statements	87
5.5	Theory extension and query semantic functions	88

7.1	First-order semantics	103
7.2	Bottom arrow definitions	105
7.3	Additional mapping operations	107
7.4	Mapping arrow definitions	108
7.5	Lazy preimage mappings	113
7.6	Comparison of arrows used as target categories	118
7.7	Preimage arrow definitions	119
7.8	Associative store arrow transformer	125
7.9	A random source $\omega \in \Omega$	126
7.10	Specific preimage arrow lifts	135
7.11	A finite model of a rectangular subset of Ω	137
7.12	Implementable, approximating arrows	140
8.1	148
8.2	153
8.3	154
8.4	162

“I think you’re begging the question,” said Haydock, “and I can see looming ahead one of those terrible exercises in probability where six men have white hats and six men have black hats and you have to work it out by mathematics how likely it is that the hats will get mixed up and in what proportion. If you start thinking about things like that, you would go round the bend. Let me assure you of that!”

Agatha Christie, *The Mirror Crack’d*

Chapter 1

Thesis

This branch of mathematics [Probability] is the only one, I believe, in which good writers frequently get results which are entirely erroneous.

Charles S. Peirce

1.1 Introduction

Probability is notorious for being counterintuitive. Ask anyone who is wrestling with the birthday paradox, the Monty Hall problem, or the Bayesian phenomenon *explaining away*.

Probability's habit of violating intuition makes any automation of probabilistic reasoning helpful. In Bayesian statistics, automation has been taking the form of modeling languages for probabilistic processes. The languages' implementations compute answers to questions about the processes under constraints.

Probabilistic languages should have mathematical specifications. The reason is simple: if a probabilistic language is implemented to be faithful to its maker's intuitions instead of a specification, it is almost certainly faulty.

Unfortunately, there are currently few probabilistic languages that have mathematical specifications. This state of affairs is partly responsible for another: until now, every probabilistic language that can support Bayesian practice is artificially limited in what it can express. Most commonly, probabilistic languages disallow unbounded loops and recursion, allow only discrete or continuous distributions, and restrict constraints to the form $X = c$.

The thesis statement is essentially that these states of affairs need not continue.

Thesis Statement: Functional programming theory and measure-theoretic probability provide a solid foundation for trustworthy, useful languages for constructive probabilistic modeling and inference.

1.2 Terms

To **model** something is to make it into a model of a theory, by developing the theory. For example, physicists model gravity by developing theories of gravitation; the physical phenomenon is a model of the theories. Likewise, Bayesians model probabilistic processes by developing probabilistic theories for which the physical processes are models. When there are mathematical models of theories, the mathematical models can be used to predict the physical models' behavior and discover their properties.

Bayesians write theories in many ways. One is to write them **constructively**: in such a way that the theory contains enough information to directly construct one of its mathematical models. This is often regarded as the ideal way to write them.

Inference means answering questions about theories. In this context, it implies **conditioning**: constraining the model in a way that preserves certain relative probabilities.

Measure-theoretic probability [27] is the most successful theory of probability in precision, maturity, and explanatory power. It was first developed in the early 1900s to formalize intuitive ideas about probability, to unify notions of discrete and continuous random variables, and to settle paradoxes that arise from incorrectly applying intuition to infinities.

Functional programming theory is used to give mathematically precise meaning to programs and to give rules for executing them. In it, the λ -calculus serves as a model of computation and as a minimal language in which to reason by substitution.

A **trustworthy** language has a mathematical meaning called a **semantics**. By defining a language mathematically, it is possible to prove theorems about it, which apply to all faithful implementations. Further, if a language implementation computes something unexpected, its semantics provides a way to determine whether its behavior is correct.

We generally think of languages as being **useful** when they save time by automating calculations. Languages are also useful when they allow us to express ideas naturally and reason about them precisely, or provide abstraction mechanisms so we can express ideas and reason about them at high levels.

1.3 Proof and Supporting Evidence

All but usefulness in the thesis can be proved, and for usefulness, we give evidence. To prove and demonstrate the thesis, we define two semantics, prove them correct, implement approximations of them, and test the implementations.

The first semantics is an initial investigation into our general approach: transforming Bayesian theories into λ -calculus terms that build exact measure-theoretic models of the theories, and then changing the transformation to build approximate models to carry out computations. To keep the investigation simple, we restrict theories to countable probability distributions and finitely many statements.

We ensure the first language is trustworthy by deriving its semantics from an idealized expected meaning of Bayesian theories, and computing answers to queries from approximate models in a way that converges to the correct answers according to the exact models. We demonstrate the language is useful by implementing the approximating transformation, and encoding theories and running queries that are difficult to model directly without it.

The second semantics handles uncountable probability distributions and recursion by transforming a first-order functional language with probabilistic choice into λ -calculus terms that build models. Again, we change the transformation to build approximate models and use them to carry out computations.

We show the language is trustworthy by proving

- Exact queries always terminate with correct answers (Theorem 7.51).
- All probabilistic programs have sensible output distributions, regardless of nontermination (Theorems 7.52 and 7.53).

- The approximations are sound, always terminate, and have other desirable properties (Theorems 7.58 through 7.61).
- Answers computed using the approximations correctly converge (Theorem 8.28).

Further, Theorems 7.52 and 7.53 apply to any probabilistic programming language that can be transformed into ours. Because ours is Turing-equivalent (with a random oracle) and is easy to extend with uncomputable operations such as real limits and decidable equality, this includes all probabilistic programming languages to date, and likely almost all future probabilistic programming languages.

We demonstrate the second language is useful by implementing the approximating transformation and encoding some typical Bayesian theories and running queries. In all of our tests, the theory encodings are straightforward and the queries are efficient.

To demonstrate further usefulness, we encode theories and run queries that are impossible to reason about precisely using typical Bayesian mathematical tools. One example draws inferences from a correctly modeled thermometer. Another is a simple, direct theory of light transport and a query that together carry out stochastic ray tracing.

1.4 Exposition Transition System

While this work is designed to make sense when read straight through, readers may skip some depending on their goals.

In principle, answers to questions such as “Which chapters should I read if I am only interested in implementations of probabilistic languages with conditioning and recursion?” can be answered using a dependency graph. However, the graph would be a mess of arrows: for example, everything after Chapter 2 (Background) depends on Chapter 2; similarly for Chapter 4 (Using λ_{ZFC}). Figure 1.1 shows an alternative: a transition system on chapters for which following any path (with backtracking permitted) guarantees a reader will not miss out on prerequisites.

Chapter 2 gives the necessary background in Bayesian practice and functional pro-

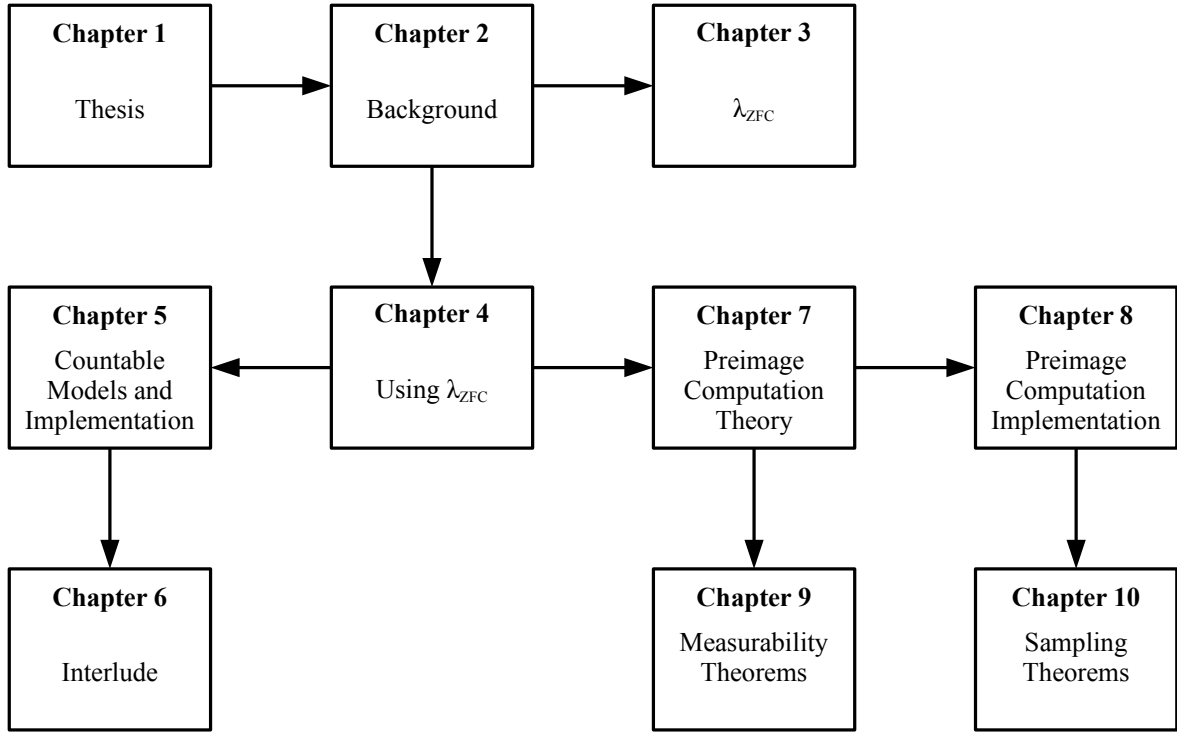


Figure 1.1: A transition system showing possible paths through this dissertation.

gramming theory, and motivates using measure theory.

The two semantics mentioned in the preceding section transform programs into λ -calculus terms. This target language has three requirements that are unusual for a λ -calculus: it must be able to represent infinite objects and operations on them, it must have nonterminating programs, and measure-theoretic theorems must apply directly to its terms. Before this work, such a λ -calculus did not exist. Chapter 3 defines one, λ_{ZFC} , with the precision necessary to carry out proofs with it.

While this precision is necessary for doing the rest of our work and verifying it, such precision is not necessary for understanding it. Readers who are not verifying our work may therefore skip from Chapter 2 to Chapter 4, which gives an overview of λ_{ZFC} and its relationship with contemporary mathematics, gives examples of use, and defines some common terminology and functions.

Chapter 5 defines a semantics for Bayesian notation restricted to countable probability distributions and finitely many statements. Chapter 6 explains why its specific way of transforming notation into models does not extend easily to theories with recursion, which motivates a slight change in tactics.

Following the new tactics, Chapter 7 defines a semantics for a probabilistic language with uncountable distributions, recursion, and arbitrary probabilistic conditions. Chapter 8 gives details that should be common to all implementations, and details specific to ours.

Chapter 9 contains proofs of theorems critical to correctness, but whose inclusion in Chapter 7 would interrupt the narrative flow too much for readers unfamiliar with measure theory. Chapter 10 is similar, but contains proofs of theorems from Chapter 8. While familiarity with measure theory is helpful while reading these two chapters, it is not strictly necessary: both explain the necessary concepts, and import enough definitions and lemmas from other sources to verify the proofs.

Chapter 2

Background

This work brings together Bayesian practice and functional programming theory, and bridges them using measure-theoretic probability.

In this chapter, we attempt to give enough background to motivate and explain Bayesian practice and functional programming theory. We only *motivate* measure-theoretic probability here. Further on, we introduce, explain and use parts as we need them.

It is difficult to find two areas in computer science as different as Bayesian practice and functional programming theory. In Bayesian practice, we find deeply held belief that unknowns can and should be *quantified* (by probabilities), reasoning by probabilistic inference, willingness to accept many kinds of approximations, and common notation that is—to put it kindly—flexible. In functional programming theory, we find deeply held belief that unknowns should be *qualified* (usually universally), reasoning by logical inference, little tolerance for unsound approximations even if they converge, and common notation that is—to put it kindly—almost precise enough to feed a compiler.

There is one common trait that makes bridging both areas even conceivable. While Bayesians model processes and functional programming researchers model languages, both approach their tasks methodically, and both create theories in which every entity they want to reason about is represented explicitly. If something important is going on behind the scenes—whether a hidden Markov process or mutating the program’s store—it is brought to the fore and fully characterized. In both areas, the extra time and cognitive burden are considered worthwhile payment for reliable artifacts and repeatable results.

It is this trait, explicit representation, that makes it possible to automate Bayesian inference, and again this trait that makes it possible to prove the automation correct.

2.1 Bayesian Practice

From Bayesian practice, the requisite background knowledge includes probability mass functions, probability densities, queries and manipulation rules, Bayesian modeling, and Bayesian inference. We assume readers know arithmetic, some set theory, functions, and the basic ideas behind integration.

2.1.1 Discrete Probability and Joint Distribution Models

In a probabilistic model of a real-world process, distinguished identifiers called **random variables** denote random values. These are regarded as free variables, but with additional information that quantifies the likelihoods of every *combination* of their values, or **observable outcomes**. The additional information is completely characterized by a function called a **joint distribution**.

For example, suppose $X, Y \in \{h, t\}$ are random variables that each represent the outcome of a coin toss. Further, let the joint distribution $p_{X,Y} : \{h, t\} \times \{h, t\} \rightarrow [0, 1]$ quantify the likelihood of every possible combination of observable outcomes by defining

$$p_{X,Y} = \left[(h, h) \mapsto \frac{1}{4}, (h, t) \mapsto \frac{1}{4}, (t, h) \mapsto \frac{1}{6}, (t, t) \mapsto \frac{1}{3} \right] \quad (2.1)$$

(The subscript “ X, Y ” is just part of the function name, and has no special meaning.) This is a **probability mass function**: its outputs sum to 1.

The probability that $X = t$ and $Y = h$ is $p_{X,Y}(t, h) = \frac{1}{6}$. Another way to write “the probability that $X = t$ and $Y = h$ ” is with a **probability query**:

$$\Pr[X = t \wedge Y = h] = p_{X,Y}(t, h) \quad (2.2)$$

A probability query always implicitly refers to some ambient joint distribution. In general,

the result of a probability query $\Pr[e]$ is the sum of probabilities of observable outcomes for which the proposition e is true. Conjunctions are often separated by a comma; e.g. $\Pr[X = x, Y = y]$ means $\Pr[X = x \wedge Y = y]$.

Probability queries have manipulation rules. One is that random variables may be “summed out” to consider the probabilities of the values of others independently. For example, to consider just the probabilities of values of X , we may sum out Y :

$$\Pr[X = x] = \sum_{y \in \{h, t\}} \Pr[X = x, Y = y] \quad (2.3)$$

According to this rule, X represents the outcome of a fair coin toss (independent of Y):

$$\begin{aligned} \Pr[X = h] &= \Pr[X = h, Y = h] + \Pr[X = h, Y = t] = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \\ \Pr[X = t] &= \Pr[X = t, Y = h] + \Pr[X = t, Y = t] = \frac{1}{6} + \frac{1}{3} = \frac{1}{2} \end{aligned} \quad (2.4)$$

Another manipulation rule allows fixing the value of one random variable to consider the probabilities of the values of others *dependently*. For example, if x is fixed and $\Pr[X = x] > 0$, then the probability that $Y = y$ is

$$\Pr[Y = y \mid X = x] = \frac{\Pr[X = x, Y = y]}{\Pr[X = x]} \quad (2.5)$$

This **conditional probability** query is read “the probability that $Y = y$ given $X = x$.” Using this rule, we can determine that, if X is known to be t , then Y represents a coin toss that is not fair:

$$\begin{aligned} \Pr[Y = h \mid X = t] &= \frac{\Pr[X = t, Y = h]}{\Pr[X = t]} = \frac{\frac{1}{6}}{\frac{1}{2}} = \frac{1}{3} \\ \Pr[Y = t \mid X = t] &= \frac{\Pr[X = t, Y = t]}{\Pr[X = t]} = \frac{\frac{1}{3}}{\frac{1}{2}} = \frac{2}{3} \end{aligned} \quad (2.6)$$

We could similarly show $\Pr[Y = h \mid X = h] = \frac{1}{2}$ and $\Pr[Y = t \mid X = h] = \frac{1}{2}$, so when X is known to be h , Y represents a fair coin toss.

To avoid the condition $\Pr[X = x] > 0$, the preceeding rule is often written

$$\begin{aligned}\Pr[X = x, Y = y] &= \Pr[X = x] \cdot \Pr[Y = y \mid X = x] \\ &= \Pr[Y = y] \cdot \Pr[X = x \mid Y = y]\end{aligned}\tag{2.7}$$

In this form, it is called the **chain rule**.

As a function of x , $\Pr[X = x]$ is a probability mass function, but over just X instead of X and Y together. With any fixed x , as a function of y , $\Pr[Y = y \mid X = x]$ is also a probability mass function. Most Bayesian models are constructed by reifying these queries as functions called (respectively) **distributions** and **conditional distributions**, and using the chain rule to build a joint distribution. For the present example,

$$\begin{aligned}p_X &= \left[h \mapsto \frac{1}{2}, t \mapsto \frac{1}{2} \right] \\ p_{Y|X}(y \mid x) &= \begin{cases} x = h & \left[h \mapsto \frac{1}{2}, t \mapsto \frac{1}{2} \right](y) \\ x = t & \left[h \mapsto \frac{1}{3}, t \mapsto \frac{2}{3} \right](y) \end{cases} \\ p_{X,Y}(x, y) &= p_X(x) \cdot p_{Y|X}(y \mid x)\end{aligned}\tag{2.8}$$

In the conditional distribution $p_{Y|X}$, the “ $Y|X$ ” subscript is simply part of the function name, and in applying it, $(y \mid x)$ is just another way to write the arguments (y, x) .¹ The syntax simply connotes that we expect $\Pr[Y = y \mid X = x] = p_{Y|X}(y \mid x)$.

This model can be more compactly specified by a constructive theory about X and Y , which states only the properties that a joint distribution model must satisfy:

$$\begin{aligned}X &\sim \left[h \mapsto \frac{1}{2}, t \mapsto \frac{1}{2} \right] \\ Y &\sim \begin{cases} X = h & \left[h \mapsto \frac{1}{2}, t \mapsto \frac{1}{2} \right] \\ X = t & \left[h \mapsto \frac{1}{3}, t \mapsto \frac{2}{3} \right] \end{cases}\end{aligned}\tag{2.9}$$

Here, $X \sim e$ is read “ X is distributed e .” In this leaner form, it is perhaps easier to understand

¹It is common to leave off subscripts such as $Y|X$ and use the form of application to distinguish the different ps ; e.g. $p(x)$ means $p_X(x)$ and $p(y \mid x)$ means $p_{Y|X}(y \mid x)$. Doing so is helpful when there are many random variables, and it is usually unambiguous, but the practice often confuses and frustrates newcomers.

the process being modeled, which is

1. Toss a coin and call its outcome X .
2. If $X = h$, toss a fair coin and call its outcome Y .
3. If $X = t$, toss a biased coin with heads probability $\frac{1}{3}$ and call its outcome Y .

It is usually easy to manually translate constructive theories into programs that sample random variable values.

By combining a conditional probability query and the chain rule (or using the chain rule twice), we get **Bayes' law**: if $\Pr[Y = y] > 0$, then

$$\Pr[X = x \mid Y = y] = \frac{\Pr[X = x] \cdot \Pr[Y = y \mid X = x]}{\Pr[Y = y]} \quad (2.10)$$

This is different than $\Pr[Y = y \mid X = x]$. This time, we are interested in the conditional probability that $X = x$ given we know $Y = y$ for some fixed y .

For example, suppose we did not observe X , but were allowed to observe $Y = t$. Given that we know the second coin toss is tails, the probability that $X = t$ is

$$\begin{aligned} \Pr[X = t \mid Y = t] &= \frac{\Pr[X = t] \cdot \Pr[Y = t \mid X = t]}{\Pr[Y = t]} \\ &= \frac{\frac{1}{2} \cdot \frac{2}{3}}{\sum_{x \in \{h, t\}} \Pr[X = x, Y = t]} \\ &= \frac{\frac{1}{2} \cdot \frac{2}{3}}{\frac{1}{4} + \frac{1}{3}} = \frac{\frac{1}{3}}{\frac{7}{12}} = \frac{4}{7} \end{aligned} \quad (2.11)$$

which is greater than $\Pr[X = t] = \frac{1}{2}$. Similarly, $\Pr[X = h \mid Y = t] = \frac{3}{7}$, which is less than $\Pr[X = h] = \frac{1}{2}$. Observing the effects of the first coin toss—even random effects—allows us to draw stronger conclusions about the first coin toss. In this case, we know it is more likely to be tails.

Using Bayes' law to draw probabilistic conclusions about probabilistic processes given observed probabilistic effects is called **Bayesian inference**.

It is easy for probability newcomers with logical background to think of conditional queries as a fuzzy kind of logical implication. A short example demonstrates that doing so

leads to faulty intuition. To compute the probability that $Y = t \implies X = t$, we apply logical rules to the query until it is a disjunction of distinct observable outcomes, and add their probabilities:

$$\begin{aligned}
\Pr[Y = t \implies X = t] &= \Pr[\neg(Y = t) \vee X = t] \\
&= \Pr[Y = h \vee X = t] \\
&= \Pr[(Y = h \wedge X = t) \vee (Y = h \wedge X = h) \vee (Y = t \wedge X = t)] \\
&= \Pr[Y = h \wedge X = t] + \Pr[Y = h \wedge X = h] + \Pr[Y = t \wedge X = t] \\
&= \frac{1}{6} + \frac{1}{4} + \frac{1}{3} = \frac{3}{4}
\end{aligned} \tag{2.12}$$

This is clearly not $\Pr[X = t | Y = t] = \frac{4}{7}$. In a similar fashion, $\Pr[Y = t \implies X = h] = \frac{2}{3}$, which is not $\Pr[X = h | Y = t] = \frac{3}{7}$.

In conditioning on $Y = t$, we did not consider any outcomes in which $Y \neq t$: we restricted the possible outcomes to those for which $Y = t$ and renormalized the probabilities of $X = h$ and $X = t$. On the other hand, in computing the probability that $Y = t \implies X = t$, we had to consider *all* of the outcomes in which $Y \neq t$, and only *one* outcome in which $Y = t$.

2.1.2 Probability Densities and Density Models

Probability mass functions cannot quantify the likelihoods of uncountably many observable outcomes, such as when $X \in \mathbb{R}$. In these cases, the distributions, conditional distributions, and joint distributions are specified using **probability density functions**: functions over outcomes that *integrate* to 1 instead of sum to 1.²

For example, this probability density function defines the **standard normal distribution**, or bell curve:

$$f_N(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) \tag{2.13}$$

²For readers familiar with measure theory: we use the word *density* for densities with respect to Lebesgue measure, and *mass* for densities with respect to counting measure. We call all other densities *derivatives*.

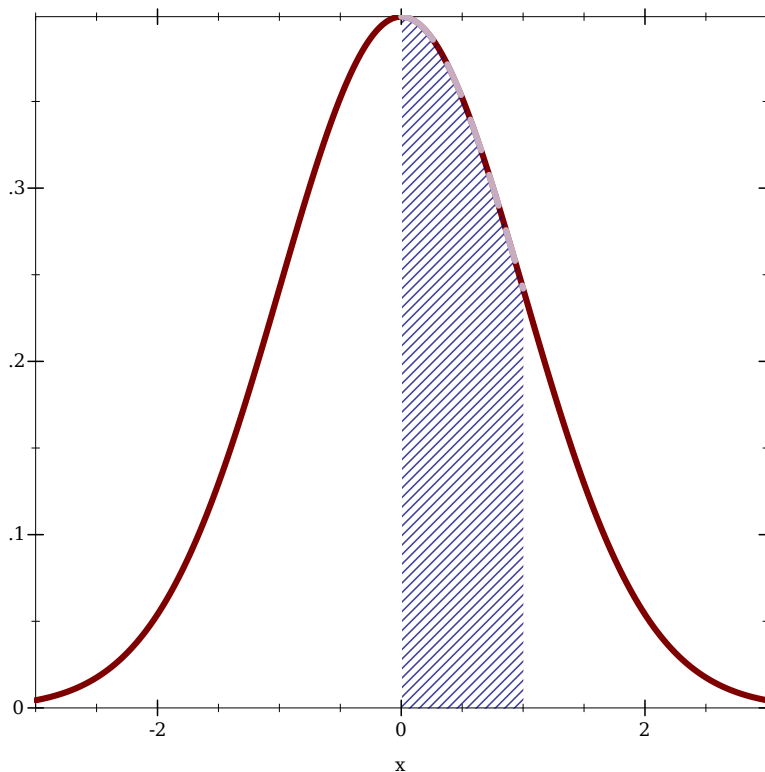


Figure 2.1: Integrating under the standard normal density to compute $\Pr[X \in (0, 1)] \approx 0.34$.

If X has a standard normal distribution, then the probability that $X \in (0, 1)$ is

$$\Pr[X \in (0, 1)] = \int_0^1 f_N(x) dx \approx 0.3413447460685 \quad (2.14)$$

Figure 2.1 plots this density and illustrates integrating under it to compute $\Pr[X \in (0, 1)]$.

When probabilities are computed by integrating density functions, sets of outcomes may have positive probability, but every *single* outcome has zero probability. For any random variable $X \in \mathbb{R}$, outcome $x \in \mathbb{R}$, and density function $f_X : \mathbb{R} \rightarrow [0, +\infty)$,

$$\Pr[X = x] = \int_x^x f_X(x) dx = (f_X(x) - 0) \cdot (x - x) = f_X(x) \cdot 0 = 0 \quad (2.15)$$

As a consequence, interval endpoints do not matter; i.e. $\Pr[X \in [a, b]] = \Pr[X \in (a, b)]$. We discuss other, more difficult consequences further on.

The normal distribution can be extended to a **distribution family** by parameterizing

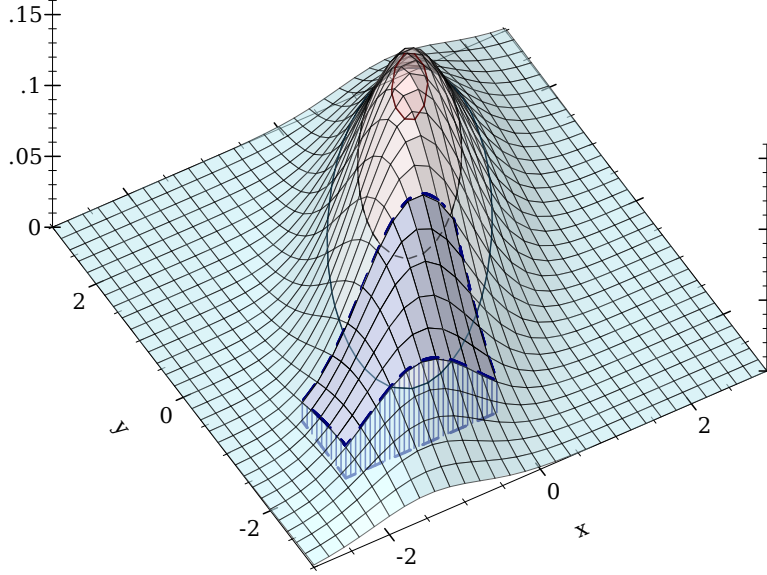


Figure 2.2: The joint density model $f_{X,Y}$ constructed from the density f_X and the conditional density $f_{Y|X}$. Integrating under $f_{X,Y}$ on the set $(-2, 0) \times (-2, -1)$ computes $\Pr[X \in (-2, 0), Y \in (-2, -1)]$.

it on a *mean* μ and a *standard deviation* σ :

$$f_N(x | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \quad (2.16)$$

(Again, the application syntax $(x | \mu, \sigma)$ simply means (x, μ, σ) .) Using parameterized distributions, we can define a **joint density model** of a probabilistic process involving two random variables $X, Y \in \mathbb{R}$:

$$\begin{aligned} f_X(x) &= f_N(x | 0, 1) \\ f_{Y|X}(y | x) &= f_N(y | x, 1) \\ f_{X,Y}(x, y) &= f_X(x) \cdot f_{Y|X}(y | x) \end{aligned} \quad (2.17)$$

Integrating under $f_{X,Y}$ computes probability queries:

$$\Pr[X \in (a, b), Y \in (c, d)] = \int_a^b \int_c^d f_{X,Y}(x, y) dy dx \quad (2.18)$$

Figure 2.2 illustrates the joint density model $f_{X,Y}$ and computing a probability query.

As with discrete models, density models can be specified by constructive theories:

$$\begin{aligned} X &\sim \text{Normal}(0, 1) \\ Y &\sim \text{Normal}(X, 1) \end{aligned} \quad (2.19)$$

The probabilistic process being modeled is

1. Choose X according to the standard normal distribution.
2. Choose Y according to the normal distribution with mean X , standard deviation 1.

Again, it is usually easy to manually translate such theories into programs that sample random variable values.

When all single outcomes have zero probability, interpreting theories in terms of expected conditional probability queries is difficult. (In fact, doing so requires measure theory.) Fortunately, densities have rules analogous to rules for manipulating probability queries, which allow practitioners to derive joint density models from theories and compute a restricted class of conditional queries. Instances of the two most important rules are

$$\begin{aligned} f_X(x) &= \int_{-\infty}^{+\infty} f_{X,Y}(x, y) dy \\ f_{X,Y}(x, y) &= f_X(x) \cdot f_{Y|X}(y|x) = f_Y(y) \cdot f_{X|Y}(x|y) \end{aligned} \quad (2.20)$$

The second rule is the **chain rule for densities**, and it justifies constructing the joint density $f_{X,Y}$ from the density f_X and the conditional density $f_{Y|X}$.

The rules for densities can be used to derive **Bayes' law for densities**: if $f_Y(y) > 0$,

then, starting from the right-hand side of the chain rule,

$$\begin{aligned}
f_Y(y) \cdot f_{X|Y}(x|y) &= f_X(x) \cdot f_{Y|X}(y|x) \\
f_{X|Y}(x|y) &= \frac{f_X(x) \cdot f_{Y|X}(y|x)}{f_Y(y)} \\
&= \frac{f_X(x) \cdot f_{Y|X}(y|x)}{\int_{-\infty}^{+\infty} f_{X,Y}(x,y) dx} \\
&= \frac{f_X(x) \cdot f_{Y|X}(y|x)}{\int_{-\infty}^{+\infty} f_X(x) \cdot f_{Y|X}(y|x) dx}
\end{aligned} \tag{2.21}$$

The last form is conveniently in terms of f_X and $f_{Y|X}$, which we have on-hand.

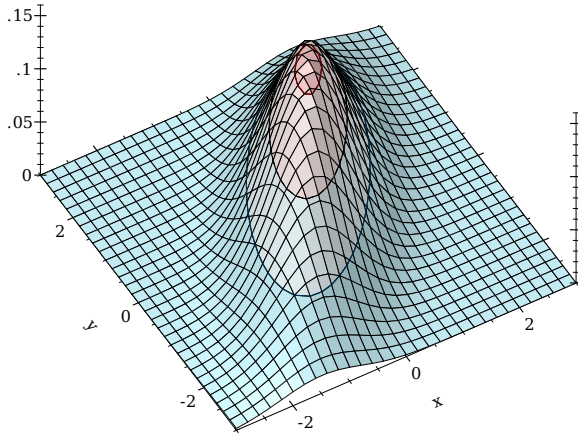
Using Bayes' law for densities, we can draw conclusions about X given knowledge about Y . For example, suppose we want to know the distribution of X given $Y = 2$, as a density. A quite lengthy derivation finally results in

$$\begin{aligned}
f_{X|Y}(x|2) &= \frac{f_X(x) \cdot f_{Y|X}(2|x)}{\int_{-\infty}^{+\infty} f_X(x) \cdot f_{Y|X}(2|x) dx} \\
&= \frac{f_N(x|0,1) \cdot f_N(2|x,1)}{\int_{-\infty}^{+\infty} f_N(x|0,1) \cdot f_N(2|x,1) dx} \\
&\dots \\
&= f_N(x|1, \sqrt{\frac{1}{2}})
\end{aligned} \tag{2.22}$$

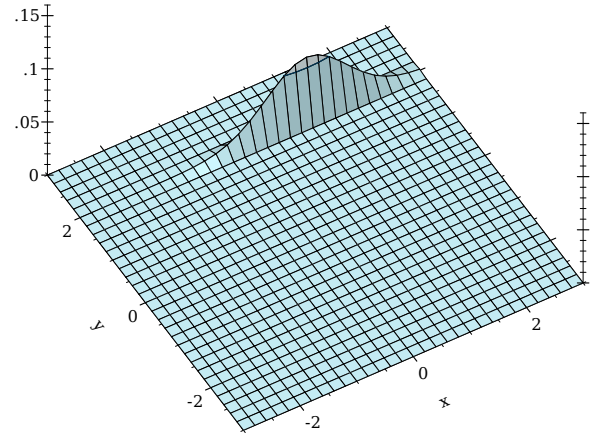
We can answer conditional probability queries such as “the probability that $X \in (a, b)$ given $Y = 2$ ” by integrating $f_{X|Y}(x|2) = f_N(x|1, \sqrt{\frac{1}{2}})$:

$$\Pr[X \in (a, b) | Y = 2] = \int_a^b f_N(x|1, \sqrt{\frac{1}{2}}) dx \tag{2.23}$$

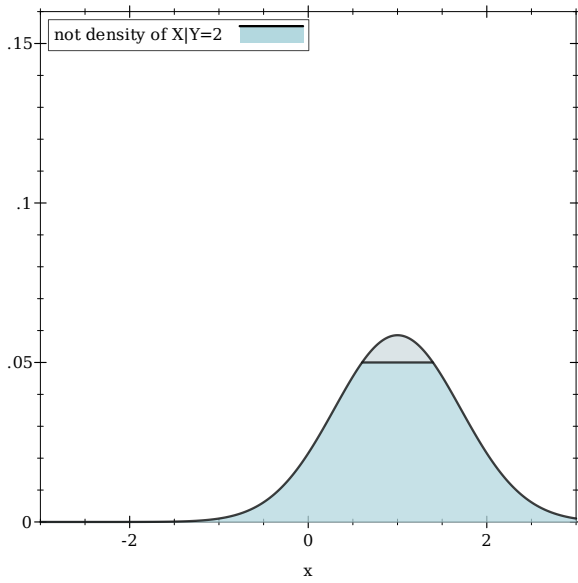
While using Bayes' law for densities is difficult, it is easy to visualize, at least in two dimensions. We may think of using it as happening in three steps: restrict, project, then renormalize. Figure 2.3 illustrates them. First, we restrict the joint density model (Figure 2.3a) to the subset of $\mathbb{R} \times \mathbb{R}$ where $y = 2$ (Figure 2.3b). Second, because the resulting model integrates to zero and therefore all answers to probability queries using it would be zero,



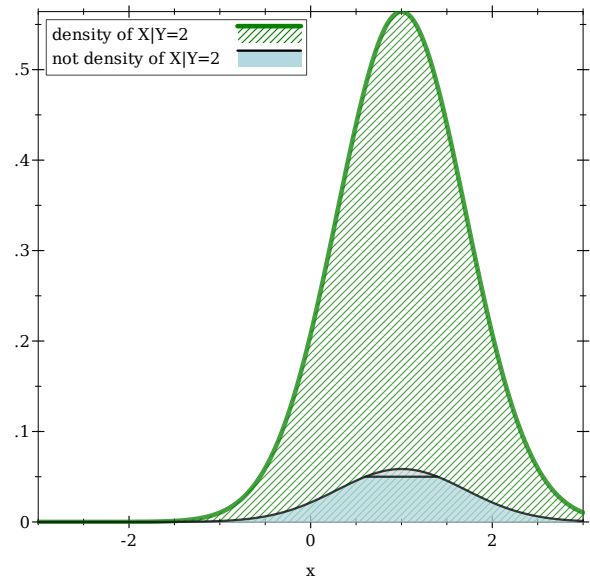
(a) The original joint density model.



(b) Restricting the model to the subset of its domain where $y = 2$. The probability of the subset is zero.



(c) Projecting the restricted model onto the x axis results in a density that integrates to a nonzero constant that is less than 1.



(d) Normalizing the restricted, projected model results in a probability density that characterizes the distribution of X given $Y = 2$.

Figure 2.3: Bayes' law for densities, in pictures.

we project it onto the x axis (Figure 2.3c), on which it integrates to a positive value. Third, because the total probability is now less than 1, we renormalize the restricted, projected model by dividing its output by its area, and obtain a probability density (Figure 2.3d).

Using Bayes' law for densities is not only often technically challenging, but in general there are no closed-form solutions. In such cases, practitioners turn to Monte Carlo integration, or sampling, to answer conditional probability queries.

2.1.3 Motivating Measure-Theoretic Models

It is easy to define a random variable whose distribution cannot be characterized by a density. Suppose we have a thermometer whose output is not quite correct, but is usually within 1 degree. We could model the error with a normal distribution with standard deviation 1. If the thermometer cannot show a number greater than 100 and we want to model that fact, assuming the correct temperature is 99, we could write a theory about it like this:

$$\begin{aligned} T' &\sim \text{Normal}(99, 1) \\ T &= \min(T', 100) \end{aligned} \tag{2.24}$$

so that T represents the thermometer's output. Because $T' \geq 100$ if and only if $T = 100$, $\Pr[T' \geq 100] = \Pr[T = 100] > 0$. But we know that if T has a density, $\Pr[T = 100] = 0$.

Even though T has no density, it must have a distribution, because $T \in (a, b)$ for any interval (a, b) has a sensible probability, which we can compute by integrating the density of T' up to 100 and adding $\Pr[T' \geq 100]$ if the interval happens to contain 100:

$$\Pr[T \in (a, b)] = \int_{\min(a, 100)}^{\min(b, 100)} f_N(x | 99, 1) dx + \begin{cases} \Pr[T' \geq 100] & \text{if } 100 \in (a, b) \\ 0 & \text{if } 100 \notin (a, b) \end{cases} \tag{2.25}$$

Further, it is easy to write a program that samples T values, and we expect the outputs of programs with probabilistic choice to have well-defined probability distributions.

Another way to define a distribution that cannot be characterized by a density is

to restrict a joint density model with a non-axial condition. For example, with the present theory and density model for $X, Y \in \mathbb{R}$, suppose we want to answer questions like

$$\Pr[X \in (a, b), Y \in (c, d) \mid \sqrt{X^2 + Y^2} = 1] \quad (2.26)$$

which restricts the model to a unit circle. To do so with a density, we need $f_{X,Y|\sqrt{X^2+Y^2}}(x, y \mid z)$. But for any z , a joint density for X, Y restricted to $\sqrt{X^2 + Y^2} = z$ cannot exist because the area of that circular set is zero.

In the preceding examples, on a part of the distribution's domain that has length or area 0 (i.e. the single outcome 100 or the set $\{x \in \mathbb{R}, y \in \mathbb{R} \mid \sqrt{x^2 + y^2} = z\}$), the probability of that set is nonzero. Integrating on such a set can yield only 0, so the distributions cannot be defined by densities.

There are many more sensible distributions that cannot be defined by densities, such as the distribution of a random variable with outcomes in $\mathbb{R} \cup \mathbb{R}^2$ or $\mathbb{R}^{\mathbb{N}}$. Not only can sets of such outcomes have sensible probabilities, but as in the thermometer example, *it is easy to define random variables with these distributions*.

Measure theory's answer to this shocking lack of densities is to define probability distributions using **probability measures**, which map *sets to probabilities* instead of *values to instantaneous changes* in probability. The measure defining T 's distribution directly answers queries such as $\Pr[T \in (a, b)]$. There is also a measure defined on sets of $\mathbb{R} \times \mathbb{R}$ that directly answers queries such as $\Pr[X \in (a, b), Y \in (c, d) \mid \sqrt{X^2 + Y^2} = 1]$. Probability measures exist on spaces with varying and infinite dimension.

Chapter 5 introduces measure theory by using it to interpret discrete Bayesian theories mechanically. Chapter 7 gives a measure-theoretic interpretation of probabilistic programs, which can encode Bayesian theories that do not have density models.

2.2 Functional Programming Theory

From functional programming theory, the requisite background knowledge includes the λ -calculus, big-step operational semantics, denotational semantics, categorical semantics, and abstract interpretation. We assume readers know basic computer science theory, including propositional logic, relations, functions, proof by induction, context-free grammars, and nondeterminism.

2.2.1 λ -Calculus

The following grammar defines a set of variable names X and a language E (a set of terms) called the **pure λ -calculus**.

$$\begin{aligned} e &::= x \mid e \ e \mid \lambda x. e \\ x &::= [\text{variable names}] \end{aligned} \tag{2.27}$$

Terms $\lambda x. e$ are unnamed functions of one argument, terms x refer to function arguments, and terms $e_1 \ e_2$ apply e_1 to e_2 (i.e. “call” function e_1 with argument e_2). For readers unused to the λ -calculus but familiar with other mathematical languages, perhaps the most difficult thing to get used to is that juxtaposition means application instead of multiplication.

As in most mathematical languages, parentheses are optional. Lambda terms greedily enclose their bodies in implicit parentheses, so $\lambda x. \lambda y. e$ (with some assumed-meaningful function body e) is the same term as $\lambda x. (\lambda y. e)$: a function that receives an x and returns a function of y in which x is available. Application is left-associative, so $e_1 \ e_2 \ e_3$ is the same term as $(e_1 \ e_2) \ e_3$. Here, $e_1 \ e_2$ returns a function, which is applied to e_3 .

This duality makes it easy to write two-argument functions using nested lambdas, and apply them using sequences of arguments. For example, $(\lambda x. \lambda y. e) \ e_x \ e_y$ defines a function of two arguments and applies it to e_x and e_y .

Like its younger brother the Turing machine, the pure λ -calculus is a universal model of computation. Also like the Turing machine, it would be quite painful to program with it.

Unlike the Turing machine, it is easy to get something practical with a few extensions such as pairs and numbers, and a few primitive functions to operate on them.

In most programming languages, implementation details define the meaning of function application. It typically involves a jump from one machine address to another, and if the function returns, a jump back. However, in the pure λ -calculus and its extensions, there are no jumps or machine addresses. Function application is defined entirely in terms of substitution, as in algebra. For example, suppose *hypot* is a term in a λ -calculus extended with real numbers, defined by

$$\textit{hypot} = \lambda x. \lambda y. \sqrt{x^2 + y^2} \quad (2.28)$$

Applying *hypot* eliminates lambdas by substituting their formal arguments with the supplied actual arguments:

$$\begin{aligned} \textit{hypot} \ 3 \ 4 &= \left(\lambda x. \lambda y. \sqrt{x^2 + y^2} \right) \ 3 \ 4 \\ &= \left(\left(\lambda x. \left(\lambda y. \sqrt{x^2 + y^2} \right) \right) \ 3 \right) \ 4 \\ &\equiv \left(\lambda y. \sqrt{3^2 + y^2} \right) \ 4 \\ &\equiv \sqrt{3^2 + 4^2} \end{aligned} \quad (2.29)$$

The two equivalences at the end of (2.29) are called β -**reductions**, or just **reductions**. We would expect $\sqrt{3^2 + 4^2}$ to further reduce to 5.

Computer implementations of an extended λ -calculus, such as the programming language Racket, necessarily use jumps and machine addresses to implement function application. However, the meanings of their programs are defined mathematically as the results of carrying out reductions. It is therefore possible to reason about programs algebraically and inductively, without having to consider complicating machine details.

It is sometimes convenient to define a λ -calculus whose variables refer to function arguments by *number* instead of by *name*. Such numeric references are called **De Bruijn**³

³Typically pronounced “deh brOIN,” and named after Dutch mathematician Nicolaas de Bruijn.

indexes. One form of the pure λ -calculus with De Bruijn indexes is

$$\begin{aligned} e &::= \mathbf{env} \ n \mid e \ e \mid \lambda. e \\ n &::= 0 \mid 1 \mid 2 \mid \dots \end{aligned} \tag{2.30}$$

where a “variable” term $\mathbf{env} \ 0$ refers to the innermost lambda’s argument.

Suppose we define *hypot* as a term in a λ -calculus with De Bruijn indexes, extended with real numbers:

$$\mathit{hypot} = \lambda. \lambda. \sqrt{(\mathbf{env} \ 1)^2 + (\mathbf{env} \ 0)^2} \tag{2.31}$$

Here, $\mathbf{env} \ 1$ (which was previously x) refers to the outer lambda’s argument and $\mathbf{env} \ 0$ refers to the inner lambda’s argument. Reducing an application of *hypot* proceeds this way:

$$\begin{aligned} \mathit{hypot} \ 3 \ 4 &= \left(\lambda. \lambda. \sqrt{(\mathbf{env} \ 1)^2 + (\mathbf{env} \ 0)^2} \right) \ 3 \ 4 \\ &\equiv \left(\lambda. \sqrt{3^2 + (\mathbf{env} \ 0)^2} \right) \ 4 \\ &\equiv \sqrt{3^2 + 4^2} \end{aligned} \tag{2.32}$$

So far, we have been taking a certain evaluation order for granted when computing reductions. To highlight an ambiguity, consider this lambda term, which returns 0 given any argument:

$$\mathit{zero} = \lambda x. 0 \tag{2.33}$$

Suppose $1 / 0$ does not reduce to any value, as in algebra. Should *zero* $(1 / 0)$ reduce to 0, or likewise not reduce? In other words, should we accept this reduction:

$$\begin{aligned} \mathit{zero} \ (1 / 0) &= (\lambda x. 0) \ (1 / 0) \\ &\equiv 0 \end{aligned} \tag{2.34}$$

or should we require function arguments to reduce before substituting them? Always reducing function arguments first is **call-by-value** reduction, and substituting without reducing arguments is **call-by-name**. Both policies have their place, but we mostly use call-by-value

reduction, in which *zero* ($1 / 0$) does not reduce.

2.2.2 Big-Step Operational Semantics

Instead of describing evaluation order using English phrases with scattered mathematical terms, we could instead give our λ -calculus a **semantics**: a precise mathematical definition of the meaning of its terms. To specify evaluation order and other operational aspects specifically, we would typically give it an **operational semantics**.

An operational semantics is defined by a **reduction relation**, which relates program terms to other program terms. There are two main kinds of operational semantics:

- **Small-step**, specified by a subset of $E \times E$, where E is the set of program terms.
- **Big-step**, specified by a subset of $E \times V$, where E is the set of program terms and $V \subseteq E$ is the set of irreducible program values (e.g. the number 4, the pair $\langle 10, 23 \rangle$).

For example, suppose we have a lambda term

$$inc = \lambda x. x + 1 \tag{2.35}$$

A small-step semantics would typically “stop” after a function application. If “ \Rightarrow ” is a small-step reduction relation, then $(inc\ 4) \Rightarrow (4 + 1)$ should be true, and also $(4 + 1) \Rightarrow 5$, so we can conclude $inc\ 4$ reduces to 5 in two **small steps**. On the other hand, a big-step semantics cannot “stop” after most function applications. If “ \Downarrow ” is a big-step reduction relation, then we cannot expect $(inc\ 4) \Downarrow (4 + 1)$ because by any reasonable definition, the term $4 + 1$ is an *expression* but not a *value*. We should expect, however, that $(inc\ 4) \Downarrow 5$ is true; i.e. inc reduces to 5 in one **big step**.

As with call-by-name and call-by-value, small-step and big-step semantics both have their place. However, as Chapter 3 contains the only operational semantics in this dissertation and it is a big-step semantics, we concentrate on big-step in this overview.

Figure 2.4 defines a language and its semantics by giving a grammar and a big-step reduction relation “ \Downarrow ”. The language is even simpler than the pure λ -calculus: its terms

$$\begin{array}{c}
\frac{e ::= v \mid \mathbf{add} \ e \ e}{v ::= 0 \mid 1 \mid 2 \mid \dots} \\
\text{(a) A grammar to define sets } E \text{ and } V
\end{array}
\qquad
\begin{array}{c}
\frac{}{v \Downarrow v} \text{ (val)} \qquad \frac{e_1 \Downarrow v_1 \quad e_2 \Downarrow v_2}{(\mathbf{add} \ e_1 \ e_2) \Downarrow (v_1 + v_2)} \text{ (add)} \\
\text{(b) Reduction rules to define } \Downarrow \subseteq E \times V
\end{array}$$

Figure 2.4: A big-step operational semantics for a simple addition language.

simply represent adding concrete numbers. The relation “ \Downarrow ” is defined by **reduction rules** in the form

$$\frac{\text{premise}_1 \quad \text{premise}_2 \quad \dots}{\text{conclusion}} \quad \text{(name)} \tag{2.36}$$

Grammar nonterminals are implicitly universally quantified, premises are implicitly conjuncted, and the rule is interpreted as an implication. For example, the (add) rule in Figure 2.4b means “for all $e_1, e_2 \in E$ and $v_1, v_2 \in V$, if e_1 reduces to v_1 and e_2 reduces to v_2 , then $\mathbf{add} \ e_1 \ e_2$ reduces to $v_1 + v_2$.” The (val) rule means “for all $v \in V$, v reduces to v ” or equivalently, “for all $v \in V$, *true* implies v reduces to v .”

The reduction relation “ \Downarrow ” is defined as the *smallest* subset of $E \times V$ for which the reduction rules hold. Defining it as the smallest subset precludes unintended conclusions such as $4 \Downarrow 5$, which are not otherwise precluded by interpreting the rules as implications. Equivalently, it restricts “ \Downarrow ” to conclusions that are provable from the reduction rules.

Reduction rules can be used directly to build **derivation trees**, which represent both computation steps and proofs of conclusions. For example, suppose we want to use the reduction rules in Figure 2.4b to compute the value of $\mathbf{add} \ (\mathbf{add} \ 4 \ 5) \ 90$. We start by writing it as a conclusion without premises:

$$\overline{(\mathbf{add} \ (\mathbf{add} \ 4 \ 5) \ 90) \Downarrow v_1} \tag{2.37}$$

There is only one rule (add) with a matching conclusion, so we add its premises, renaming variables as appropriate:

$$\frac{\overline{(\mathbf{add} \ 4 \ 5) \Downarrow v_2} \quad \overline{90 \Downarrow v_3}}{\overline{(\mathbf{add} \ (\mathbf{add} \ 4 \ 5) \ 90) \Downarrow v_1}} \tag{2.38}$$

There is only one rule (val) matching the conclusion $90 \Downarrow v_3$, and it has no premises. We thus only add premises for the (add) rule matching (add 4 5):

$$\frac{\frac{4 \Downarrow v_4 \quad 5 \Downarrow v_5}{(\text{add } 4 \ 5) \Downarrow v_2} \quad 90 \Downarrow v_3}{(\text{add } (\text{add } 4 \ 5) \ 90) \Downarrow v_1} \quad (2.39)$$

It is easy to find values of v_3 , v_4 and v_5 that make the leaf premises true, so we substitute them and recursively fill in the conclusions:

$$\frac{\frac{4 \Downarrow 4 \quad 5 \Downarrow 5}{(\text{add } 4 \ 5) \Downarrow v_2} \quad 90 \Downarrow v_3}{(\text{add } (\text{add } 4 \ 5) \ 90) \Downarrow v_1} \Longrightarrow \frac{\frac{4 \Downarrow 4 \quad 5 \Downarrow 5}{(\text{add } 4 \ 5) \Downarrow 9} \quad 90 \Downarrow 90}{(\text{add } (\text{add } 4 \ 5) \ 90) \Downarrow v_1} \Longrightarrow \frac{\frac{4 \Downarrow 4 \quad 5 \Downarrow 5}{(\text{add } 4 \ 5) \Downarrow 9} \quad 90 \Downarrow 90}{(\text{add } (\text{add } 4 \ 5) \ 90) \Downarrow 99} \quad (2.40)$$

Thus, the rightmost derivation tree in (2.40) is a proof that $(\text{add } (\text{add } 4 \ 5) \ 90) \Downarrow 99$.

In most cases, reduction relations can be mathematically constructed by iterating a function that uses the reduction rules to add more conclusions given known premises. A fixpoint is reachable in countably many iterations, and as a consequence, derivation trees are always finite. On the other hand, Chapter 3 defines a λ -calculus in which the iterating function must be applied uncountably many times to reach a fixpoint, and as a consequence, its derivation trees may be infinite. Despite this minor difference in size, the basic principles behind the reduction relation’s construction and use are the same.

If a big-step reduction relation “ \Downarrow ” relates each left-hand side term to exactly one right-hand side term, it is a total function, or $\Downarrow : E \rightarrow V$. If it relates each left-hand side term to *at most* one right-hand side term, it is a partial function, or $\Downarrow : E \rightarrowtail V$. In either case, if its derivation trees are finite, it can be implemented as a recursive function.

Figure 2.5 gives a Racket implementation of “ \Downarrow ” in Figure 2.4b. (We say Racket is the implementation’s **host language**.) The implementation defines a structure type `add` to model `add` expressions, and uses Racket’s built-in big integers to model V . Computation recursively interprets expressions, and proceeds similarly to the derivation tree construction in (2.37) through (2.40). As an example of use, at DrRacket’s Read-Eval-Print Loop (REPL),

```

(define value? exact-nonnegative-integer?)
(struct add (e1 e2))

(define (interp e)
  (match e
    [(? value? v) v]
    [(add e1 e2) (define v1 (interp e1))
                  (define v2 (interp e2))
                  (+ v1 v2)]))

```

Figure 2.5: Racket implementation of the semantics defined in Figure 2.4.

$ \begin{array}{l} e ::= v \mid \text{add } e \ e \mid \text{choose } e \ e \\ v ::= 0 \mid 1 \mid 2 \mid \dots \end{array} $	$ \begin{array}{c} \frac{}{v \Downarrow v} \text{ (val)} \quad \frac{e_1 \Downarrow v_1 \quad e_2 \Downarrow v_2}{(\text{add } e_1 \ e_2) \Downarrow (v_1 + v_2)} \text{ (add)} \\ \frac{e_1 \Downarrow v_1}{(\text{choose } e_1 \ e_2) \Downarrow v_1} \text{ (left)} \quad \frac{e_2 \Downarrow v_2}{(\text{choose } e_1 \ e_2) \Downarrow v_2} \text{ (right)} \end{array} $
(a) A grammar to define sets E and V	(b) Reduction rules to define $\Downarrow \subseteq E \times V$

Figure 2.6: Big-step operational semantics for a language with nondeterministic choice.

we get

```

> (interp (add (add 4 5) 90))
99

```

as expected.

Figure 2.6 extends the present example language with nondeterministic choice, which results in a reduction relation that is *not* a function. The culprits are the new rules (left) and (right), which can both match the same conclusion. For example, suppose we want to use the reduction rules in Figure 2.6b to compute the value of `add (choose 4 5) 90`. We start as before, by writing it as a conclusion without premises:

$$\overline{(\text{add } (\text{choose } 4 \ 5) \ 90) \Downarrow v_1} \tag{2.41}$$

We match the conclusion to the (add) rule and add its premises:

$$\frac{\overline{(\text{choose } 4 \ 5) \Downarrow v_2} \quad \overline{90 \Downarrow v_3}}{(\text{add } (\text{choose } 4 \ 5) \ 90) \Downarrow v_1} \quad (2.42)$$

Again, there is only one rule (val) matching the conclusion $90 \Downarrow v_3$, and it has no premises. For the conclusion $(\text{choose } 4 \ 5) \Downarrow v_2$, however, we may choose either (left) or (right), leading to two different derivation trees:

$$\frac{\overline{4 \Downarrow v_4} \quad \overline{(\text{choose } 4 \ 5) \Downarrow v_2} \quad \overline{90 \Downarrow v_3}}{(\text{add } (\text{choose } 4 \ 5) \ 90) \Downarrow v_1} \quad \frac{\overline{5 \Downarrow v_5} \quad \overline{(\text{choose } 4 \ 5) \Downarrow v_2} \quad \overline{90 \Downarrow v_3}}{(\text{add } (\text{choose } 4 \ 5) \ 90) \Downarrow v_1} \quad (2.43)$$

After replacing v_3 , v_4 and v_5 with the only values that make the leaf premises true and recursively filling in the conclusions, we would find that both $(\text{add } (\text{choose } 4 \ 5) \ 90) \Downarrow 94$ and $(\text{add } (\text{choose } 4 \ 5) \ 90) \Downarrow 95$ are true, and would have derivation trees to prove these facts.

An implementation of a nondeterministic semantics would be correct if, for every interpretation of a term e that produced value v , $e \Downarrow v$ were a valid conclusion. For **choose** 4 5, for example, a correct implementation may always choose 4, always choose 5, choose randomly, choose the number that gives the best or worst outcome according to some objective function, or always choose 4 on weekends or during the fall equinox. Its choice is simply not modeled by the semantics.

Suppose we wanted to compute results for every possible combination of nondeterministic choices. We could define a big-step relation $\Downarrow : E \rightarrow \mathcal{P} V$, which returns (when used as a function) a set of values, and implement an interpreter for it. However, we are saving that example for the next section.

2.2.3 Denotational Semantics

A **denotational semantics** is defined by a deterministic **semantic function** from language terms to values *in another language*. The other language is called the **metalanguage** or **target language**, and is often an axiomatic logic such as first-order set theory (i.e. ordinary

$\llbracket \cdot \rrbracket : E \rightarrow \mathbb{N}$	<code>(define-syntax compile</code>
$\llbracket v \rrbracket = v$	<code> (syntax-rules (add)</code>
$\llbracket \text{add } e_1 \ e_2 \rrbracket = \llbracket e_1 \rrbracket + \llbracket e_2 \rrbracket$	<code> [(_ (add e1 e2)) (+ (compile e1)</code>
	<code> (compile e2))]</code>
	<code> [(_ v) v]))</code>
(a) A semantic function for the addition language	(b) An implementation of the semantic function as a syntax transformer

Figure 2.7: A denotational semantics and its implementation.

mathematics).

Figure 2.7a defines a denotational semantics for the addition language without **choose** by defining a semantic function $\llbracket \cdot \rrbracket : E \rightarrow \mathbb{N}$. The double square brackets are simply a different application syntax: they connote nothing mathematically, but serve as a visual cue to read applications of the semantic function as “the meaning of” or “the denotation of.” For example, the meaning of `add (add 4 5) 90` is

$$\begin{aligned}
 \llbracket \text{add (add 4 5) 90} \rrbracket &= \llbracket \text{add 4 5} \rrbracket + \llbracket 90 \rrbracket \\
 &= (\llbracket 4 \rrbracket + \llbracket 5 \rrbracket) + \llbracket 90 \rrbracket \\
 &= (4 + 5) + 90 \\
 &= 99
 \end{aligned}
 \tag{2.44}$$

The semantic function is **compositional**: it gives meaning to terms *by combining the meanings of their direct subterms*. Compositionality allows most proofs of program properties to be done by structural induction, as we will demonstrate shortly.

When the results of applying $\llbracket \cdot \rrbracket$ are computable, because it is compositional, it is often easy to implement it as local syntax transformation or compilation. Figure 2.7b shows a Racket implementation of $\llbracket \cdot \rrbracket$ as a transformation from meaningless parenthetical syntax (an `add` function does not exist) to runnable Racket syntax. The syntax transformer is barely more than a transcription of the semantic function’s definition, with a little extra code to signal to Racket that it is to be applied to the syntax of expressions before compiling or

$$\begin{aligned}
\llbracket \cdot \rrbracket &: E \rightarrow \mathcal{P} \mathbb{N} \\
\llbracket v \rrbracket &= \{v\} \\
\llbracket \text{add } e_1 \ e_2 \rrbracket &= \{v_1 + v_2 \mid v_1 \in \llbracket e_1 \rrbracket, v_2 \in \llbracket e_2 \rrbracket\} \\
\llbracket \text{choose } e_1 \ e_2 \rrbracket &= \llbracket e_1 \rrbracket \cup \llbracket e_2 \rrbracket
\end{aligned}$$

Figure 2.8: A denotational semantics for the addition language with nondeterministic choice.

evaluating them (i.e. `define-syntax` instead of `define`) and to identify the symbol `add` as terminal symbol.

The results of compilation seem to be equivalent to the results of interpretation:

```
> (compile (add (add 4 5) 90))
99
```

but the REPL does not show transformed syntax. Fortunately, `expand-syntax` can show it:

```
> (expand-syntax #'(add (add 4 5) 90))
#'+ (+ (+ 4 5) 90))
```

Figure 2.8 defines a compositional function $\llbracket \cdot \rrbracket : E \rightarrow \mathcal{P} \mathbb{N}$, which transforms the addition language with nondeterministic choice into sets of natural numbers. For example, the meaning of 4 is $\{4\}$, the meaning of `choose 4 5` is $\{4\} \cup \{5\} = \{4, 5\}$, and the meaning of `add (choose 4 5) 90` is

$$\begin{aligned}
\llbracket \text{add (choose 4 5) 90} \rrbracket &= \{v_1 + v_2 \mid v_1 \in \llbracket \text{choose 4 5} \rrbracket, v_2 \in \llbracket 90 \rrbracket\} \\
&= \{v_1 + v_2 \mid v_1 \in (\llbracket 4 \rrbracket \cup \llbracket 5 \rrbracket), v_2 \in \llbracket 90 \rrbracket\} \\
&= \{v_1 + v_2 \mid v_1 \in (\{4\} \cup \{5\}), v_2 \in \{90\}\} \\
&= \{v_1 + v_2 \mid v_1 \in \{4, 5\}, v_2 \in \{90\}\} \\
&= \{4 + 90, 5 + 90\} \\
&= \{94, 95\}
\end{aligned} \tag{2.45}$$

We know that under “ \Downarrow ,” `add (choose 4 5) 90` reduces to both 94 and 95, so it appears $\llbracket \cdot \rrbracket$ is correct. It would be nice to know whether it is *always* correct. The following theorem states

correctness precisely in terms of “ \Downarrow ,” and critically uses $\llbracket \cdot \rrbracket$ ’s compositionality in a proof by induction on the structure of e .

Theorem 2.1 (correctness). *For all $v \in V$ and $e \in E$, $v \in \llbracket e \rrbracket \iff e \Downarrow v$.*

Proof. Let $v \in V$ and $e \in E$. The proof is by induction on the structure of e .

Base case $e \in V$. If $e = v$, then $v \in \llbracket e \rrbracket = \{e\} = \{v\}$ by definition of $\llbracket \cdot \rrbracket$, and $e \Downarrow v$ by the (val) rule. Similarly, if $e \neq v$, then $v \notin \llbracket e \rrbracket$, and not $e \Downarrow v$.

Inductive case $e = \mathbf{add} \ e_1 \ e_2$ for some $e_1 \in E$ and $e_2 \in E$.

Suppose $v \in \llbracket \mathbf{add} \ e_1 \ e_2 \rrbracket$. By definition of $\llbracket \cdot \rrbracket$, there exist $v_1 \in \llbracket e_1 \rrbracket$, $v_2 \in \llbracket e_2 \rrbracket$ such that $v = v_1 + v_2$. By the inductive hypothesis, $e_1 \Downarrow v_1$ and $e_2 \Downarrow v_2$. By the (add) rule, $(\mathbf{add} \ e_1 \ e_2) \Downarrow v$.

Conversely, if $(\mathbf{add} \ e_1 \ e_2) \Downarrow v$, by (add), there exist v_1, v_2 such that $e_1 \Downarrow v_1$, $e_2 \Downarrow v_2$ and $v = v_1 + v_2$. By hypothesis, $v_1 \in \llbracket e_1 \rrbracket$ and $v_2 \in \llbracket e_2 \rrbracket$. By definition of $\llbracket \cdot \rrbracket$, $v \in \llbracket \mathbf{add} \ e_1 \ e_2 \rrbracket$.

Proof of the inductive case $e = \mathbf{choose} \ e_1 \ e_2$ is similar to the preceding, though each “ \iff ” direction has in inner case for nondeterministic choice. \square

Now that we know $\llbracket \cdot \rrbracket$ is correct, we can regard any implementation of it as an implementation of “ \Downarrow ” as well. In general, it is easy to transfer theorems about “ \Downarrow ” to a correct $\llbracket \cdot \rrbracket$, and vice-versa.

What if we wanted to represent nondeterministic choices using lists instead of sets, or model a different computational effect, such as mutation or probabilistic choice? We could define a different semantic function for each model, but there is a more elegant way.

2.2.4 Categorical Semantics

When computer scientists from any area want to extend a fixed process without having to repeat themselves more than necessary, they **abstract**: they decouple the desired varying part from the fixed process, and parameterize the previously fixed process on the varying part. This characterizes modular, object-oriented, functional, and even semantic abstraction.

To abstract a denotational semantics, we parameterize its semantic function on the meaning it produces. The parameter takes the form of a **category**.⁴ In semantics, the category is comprised of a collection of objects called **computations** (i.e. possible program meanings) and operations on them called **combinators**.

The appropriate category for the addition language with **choose** contains sets of numbers as computations, or $\mathcal{P} \mathbb{N}$, and operations on them. While there are many possible collections of combinators, one kind of collection that functional programmers and theorists have found very useful are **monads**.⁵ The **set monad** operates on set-valued computations and is defined by these two combinators:

$$\begin{aligned} \text{return}_{\text{set}} v &= \{v\} \\ \text{bind}_{\text{set}} A f &= \bigcup_{v \in A} f v \end{aligned} \tag{2.46}$$

Evidently, from a semantic function $\llbracket \cdot \rrbracket_a$ parameterized on a monad a we should expect $\llbracket v \rrbracket_{\text{set}} = \text{return}_{\text{set}} v = \{v\}$. How to use bind_{set} is less clear, however. It apparently applies f to the objects in set A to yield a set for each, and collects these sets' members in a big union. Turning the set comprehension in the definition of $\llbracket \text{add } e_1 \ e_2 \rrbracket$ into an indexed union (as in bind_{set}) makes its use clearer:

$$\begin{aligned} \llbracket \text{add } e_1 \ e_2 \rrbracket &= \{v_1 + v_2 \mid v_1 \in \llbracket e_1 \rrbracket, v_2 \in \llbracket e_2 \rrbracket\} \\ &= \bigcup_{v_1 \in \llbracket e_1 \rrbracket} \bigcup_{v_2 \in \llbracket e_2 \rrbracket} \{v_1 + v_2\} \\ &\equiv \bigcup_{v_1 \in \llbracket e_1 \rrbracket} \bigcup_{v_2 \in \llbracket e_2 \rrbracket} \text{return}_{\text{set}} (v_1 + v_2) \\ &\equiv \bigcup_{v_1 \in \llbracket e_1 \rrbracket} \text{bind}_{\text{set}} \llbracket e_2 \rrbracket (\lambda v_2. \text{return}_{\text{set}} (v_1 + v_2)) \\ &\equiv \text{bind}_{\text{set}} \llbracket e_1 \rrbracket (\lambda v_1. \text{bind}_{\text{set}} \llbracket e_2 \rrbracket (\lambda v_2. \text{return}_{\text{set}} (v_1 + v_2))) \end{aligned} \tag{2.47}$$

Thus, we expect $\llbracket \text{add } e_1 \ e_2 \rrbracket_{\text{set}} = \text{bind}_{\text{set}} \llbracket e_1 \rrbracket_{\text{set}} (\lambda v_1. \text{bind}_{\text{set}} \llbracket e_2 \rrbracket_{\text{set}} (\lambda v_2. \text{return}_{\text{set}} (v_1 + v_2)))$.

⁴The word “category” comes from category theory, an alternative axiomatization of mathematics. Fortunately, little knowledge of category theory is necessary to define or understand categorical semantics.

⁵Strictly speaking, in category theory, they are *strong* monads.

$$\begin{aligned}
\llbracket \cdot \rrbracket_a &: E \rightarrow M_a \mathbb{N} \\
\llbracket v \rrbracket_a &= \text{return}_a v \\
\llbracket \text{add } e_1 \ e_2 \rrbracket_a &= \text{bind}_a \llbracket e_1 \rrbracket_a (\lambda v_1. \text{bind}_a \llbracket e_2 \rrbracket_a (\lambda v_2. \text{return}_a (v_1 + v_2))) \\
\llbracket \text{choose } e_1 \ e_2 \rrbracket_a &= \text{merge}_a \llbracket e_1 \rrbracket_a \llbracket e_2 \rrbracket_a
\end{aligned}$$

Figure 2.9: A categorical semantics for the addition language with nondeterministic choice.

Finally, we need to extend the set monad with an operation for **choose** expressions. We define

$$\text{merge}_{\text{set}} A_1 A_2 = A_1 \cup A_2 \quad (2.48)$$

so that $\llbracket \text{choose } e_1 \ e_2 \rrbracket_{\text{set}} = \text{merge}_{\text{set}} \llbracket e_1 \rrbracket_{\text{set}} \llbracket e_2 \rrbracket_{\text{set}}$.

In Figure 2.9, guided by our expectations for $\llbracket \cdot \rrbracket_{\text{set}}$, we define a categorical semantics for the addition language with **choose**, by defining a semantic function $\llbracket \cdot \rrbracket_a$ parameterized on a target monad a . The parameterized function M_a returns the monad's computations. If $M_{\text{set}} X = \mathcal{P} X$, then $\llbracket \cdot \rrbracket_{\text{set}} : E \rightarrow M_{\text{set}} \mathbb{N}$ is equivalent to $\llbracket \cdot \rrbracket : E \rightarrow \mathcal{P} \mathbb{N}$ as defined in Figure 2.8, as expected.

Because Figure 2.9 does not refer to sets or set operations, it is abstract enough to interpret programs as many different kinds of computations. For example, let $M_{\text{list}} X = [X]$, where $[X]$ denotes all the lists of X , and define the **list monad** extended with *merge* by

$$\begin{aligned}
\text{return}_{\text{list}} v &= [v] \\
\text{bind}_{\text{list}} vs f &= \text{concat} (\text{map } f \text{ } vs) \\
\text{merge}_{\text{list}} vs_1 \ vs_2 &= \text{append } vs_1 \ vs_2
\end{aligned} \quad (2.49)$$

Here, $[v]$ is a list containing just v , *map* applies a function to every element in a list and returns the list of results, and *concat* : $[[X]] \rightarrow [X]$ appends the elements in a list of lists. Now $\llbracket \cdot \rrbracket_{\text{list}} : E \rightarrow [\mathbb{N}]$ models nondeterminism with lists of numbers instead of sets of numbers.

For example, the meaning of `choose 4 5` as a list of nondeterministic choices is

$$\begin{aligned}
\llbracket \text{choose } 4 \ 5 \rrbracket_{list} &= \text{merge}_{list} \llbracket 4 \rrbracket_{list} \llbracket 5 \rrbracket_{list} \\
&= \text{merge}_{list} (\text{return}_{list} \ 4) (\text{return}_{list} \ 5) \\
&\equiv \text{merge}_{list} [4] [5] \\
&\equiv \text{append} [4] [5] \\
&\equiv [4, 5]
\end{aligned} \tag{2.50}$$

The meaning of `add (choose 4 5) (choose 4 5)` is thus

$$\begin{aligned}
&\llbracket \text{add (choose 4 5) (choose 4 5)} \rrbracket_{list} \\
&\equiv \text{bind}_{list} [4, 5] (\lambda v_1. \text{bind}_{list} [4, 5] (\lambda v_2. \text{return}_{list} (v_1 + v_2))) \\
&\equiv \text{concat} (\text{map} (\lambda v_1. \text{bind}_{list} [4, 5] (\lambda v_2. \text{return}_{list} (v_1 + v_2))) [4, 5]) \\
&\equiv \text{concat} [\text{bind}_{list} [4, 5] (\lambda v_2. \text{return}_{list} (4 + v_2)), \\
&\quad \text{bind}_{list} [4, 5] (\lambda v_2. \text{return}_{list} (5 + v_2))] \\
&\equiv \text{concat} [\text{concat} (\text{map} (\lambda v_2. \text{return}_{list} (4 + v_2)) [4, 5]), \\
&\quad \text{concat} (\text{map} (\lambda v_2. \text{return}_{list} (5 + v_2)) [4, 5])] \\
&\equiv \text{concat} [\text{concat} [\text{return}_{list} (4 + 4), \text{return}_{list} (4 + 5)], \\
&\quad \text{concat} [\text{return}_{list} (5 + 4), \text{return}_{list} (5 + 5)]] \\
&\equiv \text{concat} [\text{concat} [[8], [9]], \text{concat} [[9], [10]]] \\
&\equiv \text{concat} [[8, 9], [9, 10]] \\
&\equiv [8, 9, 9, 10]
\end{aligned} \tag{2.51}$$

In contrast, $\llbracket \text{add (choose 4 5) (choose 4 5)} \rrbracket_{set} \equiv \{8, 9, 10\}$.

The semantic function $\llbracket \cdot \rrbracket_a$ can be parameterized not just on the set and list monads, but any monad a for which merge_a can be sensibly defined. This includes monads for any kind of nondeterminism (e.g. all possibilities, angelic/demonic, random, probabilistic) with

any kind of encoding for nondeterministic values (e.g. sets, lists, worst/best choices, execution paths, random values, probability distributions). It also includes monads that combine nondeterminism with other effects, such as input/output or backtracking search. Abstracting has granted the desired flexibility.

As evidenced by the long derivations in (2.50) and (2.51), like most other abstractions, semantic abstraction increases complexity in return for its flexibility and generalization. There are many ways to deal with this, including inferring the behavior of effects from computation types, and classifying effectful behaviors as belonging to different categories. The programming language Haskell benefits greatly from categorical semantics by using them to hide the encodings of effects, which, being an implementation of an effect-free λ -calculus, it cannot compute directly, by design. Its primary way to deal with the increase in complexity is to use just one built-in, standard semantic function that targets any monad, which transforms syntax that many Haskell programmers find (or learn to find) intuitive.

Besides increasing complexity, abstraction affects the semantics in another way that we have only hinted at by using “ \equiv ” instead of “ $=$ ” in some of our equations: *it no longer targets first-order set theory*. Instead, the semantic function $\llbracket \cdot \rrbracket_a$ targets a λ -calculus.

Targeting a λ -calculus restricts a denotational semantics to be *directly implementable* as a syntax transformer. This restriction is generally regarded as good, because it makes the proof of the direct implementation’s correctness trivial. However, we want to define semantic functions for Bayesian notation, which often denotes uncountable things such as probability distributions over \mathbb{R} . The entire reason for the work in Chapter 3 is to define a λ -calculus with a semantics that gives meaning to operations on infinite values of any size, so that we can define categorical semantics for probabilistic languages in Chapter 5 and Chapter 7.

Categorical abstraction has affected the semantics in a third way. Compare the rule for **add** in Figure 2.7a with the corresponding rule in Figure 2.9:

$$\begin{aligned} \llbracket \mathbf{add} \ e_1 \ e_2 \rrbracket &= \{v_1 + v_2 \mid v_1 \in \llbracket e_1 \rrbracket, v_2 \in \llbracket e_2 \rrbracket\} \\ \llbracket \mathbf{add} \ e_1 \ e_2 \rrbracket_a &= \mathit{bind}_a \ \llbracket e_1 \rrbracket_a \ (\lambda v_1. \mathit{bind}_a \ \llbracket e_2 \rrbracket_a \ (\lambda v_2. \mathit{return}_a \ (v_1 + v_2))) \end{aligned} \tag{2.52}$$

Because $\llbracket \text{add } e_1 \ e_2 \rrbracket$ does not specify the order of evaluating $\llbracket e_1 \rrbracket$ and $\llbracket e_2 \rrbracket$, an implementation is free to choose the order, evaluate them in parallel, or let the host language decide. On the other hand, $\llbracket e_1 \rrbracket_{list}$ *must* be evaluated first, because changing the evaluation order changes the results:

$$\begin{aligned} bind_{list} [4, 5] (\lambda v_1. bind_{list} [1, 2, 3] (\lambda v_2. return_{list} (v_1 + v_2))) &\equiv [5, 6, 7, 6, 7, 8] \\ bind_{list} [1, 2, 3] (\lambda v_1. bind_{list} [4, 5] (\lambda v_2. return_{list} (v_1 + v_2))) &\equiv [5, 6, 6, 7, 7, 8] \end{aligned} \quad (2.53)$$

In general, parameterizing a semantics on a monad allows certain monads to impose a total order on evaluation, regardless of the host language’s evaluation order.

The combinators in a category must obey certain laws. For example, to define a monad, $return_a$ and $bind_a$ must obey these laws:

$$\begin{aligned} bind_a (return_a x) f &\equiv f x && \text{left identity} \\ bind_a m return_a &\equiv m && \text{right identity} \\ bind_a (bind_a m f) g &\equiv bind_a m (\lambda x. bind_a (f x) g) && \text{associativity} \end{aligned} \quad (2.54)$$

It is not necessary for readers to understand these laws deeply, just that they exist, are expected to hold, are occasionally useful, and that we interpret them a little more broadly than is typical. In particular, “ \equiv ” is almost always understood to be the default equivalence for the λ -calculus in which the combinators are defined. When programming in Haskell, this helpfully ensures that using the laws to transform programs maintains program equivalence.

When defining categorical semantics, however, there is no reason for “ \equiv ” to be defined so narrowly. In fact, it is often useful to define equivalence per-category. For example, we might say that two lists are equivalent when the sets of their elements are equal. In Chapter 3’s *limit monad*, computations are infinite sequences, and are equivalent when they converge to the same value. Chapter 7 defines a notion of equivalence for each of the categories it uses to interpret probabilistic programs.

Two other kinds of categories besides monads are useful targets for categorical semantics: **idioms** and **arrows**. Each kind of category has its own combinators, orderings,

$$\begin{aligned}
\mathcal{L}[\![\cdot]\!] &: E \rightarrow \mathbb{N} \\
\mathcal{L}[\![v]\!] &= 1 \\
\mathcal{L}[\![\text{add } e_1 \ e_2]\!] &= \mathcal{L}[\![e_1]\!] \cdot \mathcal{L}[\![e_2]\!] \\
\mathcal{L}[\![\text{choose } e_1 \ e_2]\!] &= \mathcal{L}[\![e_1]\!] + \mathcal{L}[\![e_2]\!]
\end{aligned}$$

Figure 2.10: An abstract semantics for the addition language with nondeterministic choice.

and laws. We do not review them here because they are not as well-known in functional programming theory as monads, so the chapters that use them also review them.

2.2.5 Abstract Interpretation

When we want to discover something about *every* evaluation of a program, we might do it with **abstract interpretation**: evaluating by operating on just the properties of terms instead of their actual values. Equivalently, we can think of an abstract interpretation as operating on sets of values for which those properties hold. The properties or sets of values are called **abstract values**. The actual values are called **concrete values**.

Perhaps the most common example of abstract interpretation is type checking. In this case, the abstract values are types, which represent properties such as “is a number” or “is a function from lists to natural numbers.” During abstract interpretation, expressions are not evaluated on concrete values, but are checked to determine whether they preserve the properties that concrete values should have.

As with concrete interpretation, abstract interpretation is specified by a semantics. As with concrete semantics, any of a language’s abstract semantics can be defined using rules or semantic functions. Type systems and type checkers are typically defined by rules with premises and conclusions similar to reduction rules. Because Chapter 7 defines an abstract interpretation using a semantic function, we give a small example of that approach here.

Figure 2.10 defines $\mathcal{L}[\![\cdot]\!]$, which defines an abstract semantics for the addition language with **choose**. (The prefix \mathcal{L} means nothing mathematically; it simply differentiates this semantic function from the others we have defined.) The abstract values are the lengths of

lists or cardinalities of sets; i.e. natural numbers. The abstract meaning of a term is a *bound* on the number of nondeterministic values computed. For example, the abstract meaning of `add (choose 4 5) (choose 4 5)` is

$$\begin{aligned}
\mathcal{L}[\text{add (choose 4 5) (choose 4 5)}] &= \mathcal{L}[\text{choose 4 5}] \cdot \mathcal{L}[\text{choose 4 5}] \\
&= (\mathcal{L}[4] + \mathcal{L}[5]) \cdot (\mathcal{L}[4] + \mathcal{L}[5]) \\
&= (1 + 1) \cdot (1 + 1) \\
&= 4
\end{aligned} \tag{2.55}$$

Indeed, $\llbracket \text{add (choose 4 5) (choose 4 5)} \rrbracket_{\text{set}} \equiv \{8, 9, 10\}$, which is no more than 4 values.

This example demonstrates a pervasive fact about abstract semantics: almost every abstract semantics trades precision to get efficiency, tractability, or even computability. Certainly $|\{8, 9, 10\}| \neq 4$.

Usually, we need abstract interpretations to be **sound**, which roughly means that the abstract values are always a conservative approximation. (There is a way to formalize this notion using Galois connections, but that brings in more complexity than we need.) When abstraction interpretations must be sound, the abstract semantics must have a soundness theorem relating it to a concrete semantics, such as the following.

Theorem 2.2 ($\mathcal{L}[\cdot]$ soundness). *For all $e \in E$, $|\llbracket e \rrbracket_{\text{set}}| \leq \mathcal{L}[e]$.*

Proof. By structural induction on e . □

A soundness theorem sometimes suggests how abstraction interpretations might be used. For a type system, soundness implies that accepted programs never compute concrete values with the wrong type, so operations on concrete values may be specialized in ways that would otherwise be unsafe or incorrect. (A child class's methods may be inlined, for example.) By Theorem 2.2, we can use $\mathcal{L}[\cdot]$ to determine how much space to preallocate for results in a less direct but faster implementation of $\llbracket \cdot \rrbracket_{\text{set}}$, and we will never allocate too little.

Sometimes the abstraction is both sound and precise, as $\mathcal{L}[\cdot]$ is with respect to $\llbracket \cdot \rrbracket_{\text{list}}$.

Theorem 2.3 ($\mathcal{L}[\![\cdot]\!]$ soundness and precision). *For all $e \in E$, $\text{length } \llbracket e \rrbracket_{\text{list}} = \mathcal{L}[\![e]\!]$.*

Proof. By structural induction on e . □

Having both soundness and precision is unusual, however.

Abstract interpretation is often used for program analysis in which determining precise properties is only semidecidable or is undecidable. An example is determining which functions are applied at every application site in a program written in a λ -calculus. In these cases, another concrete semantics is created, whose concrete interpretations—if they could be evaluated on a computer—would collect the necessary information. An abstract semantics is then created, whose abstract interpretations are computable, and which overapproximate the necessary information.

Such analyses are sometimes said to “embrace the infinite.” In this work, we must do the same to interpret Bayesian notation—but instead embrace the *uncountably* infinite. Doing so with a concrete categorical semantics requires a powerful λ -calculus.

Chapter 3

Computing in Cantor’s Paradise With λ_{ZFC}

This chapter is derived from work published at the 11th *International Symposium on Functional and Logic Programming (FLOPS)*, 2012.

No one shall expel us from the Paradise that Cantor has created.

David Hilbert

3.1 Motivation

Georg Cantor first proved some of the surprising consequences of assuming infinite sets exist. David Hilbert passionately defended Cantor’s set theory as a mathematical foundation, coining the term “Cantor’s Paradise” to describe the universe of transfinite sets in which most mathematics now takes place.

The calculations done in Cantor’s Paradise range from computable to unimaginably uncomputable. Still, its inhabitants increasingly use computers to answer questions. We want to make domain-specific languages (DSLs) for writing these questions, with implementations that compute exact and approximate answers.

Such a DSL should have two meanings: an exact mathematical semantics, and an approximate computational one. A traditional, denotational approach is to give the exact as a transformation to first-order set theory, and because set theory is unlike any intended implementation language, the approximate as a transformation to a λ -calculus. However,

deriving approximations while switching target languages is rife with opportunities to commit errors.

A more certain way is to define the exact semantics in a proof assistant like HOL [30] or Coq [8], prove theorems, and extract programs. The type systems confer an advantage: if the right theorems are proved, the programs are certainly correct.

Unfortunately, reformulating and re-proving theorems in such an exacting way causes significant delays. For example, half of Joe Hurd’s 2002 dissertation on probabilistic algorithms [22] is devoted to formalizing early-1900s measure theory in HOL. Our work in Bayesian inference would require at least three times as much formalization, even given the work we could build on.

Some middle ground is clearly needed: something between the traditional, error-prone way and the slow, absolutely certain way.

Instead of using a typed, higher-order logic, suppose we defined, in first-order set theory, an untyped λ -calculus that contained infinite sets and operations on them. We could interpret DSL terms exactly as uncomputable programs in this λ -calculus. But instead of redoing a century of work to extract programs that compute approximations, we could directly reuse first-order theorems to derive them from the uncomputable programs.

Conversely, set theory, which lacks lambdas and general recursion, is an awkward target language for a semantics that is intended to be implemented. Suppose we extended set theory with untyped lambdas (as objects, not quantifiers). We could still interpret DSL terms as operations on infinite objects. But instead of leaping from infinite sets and operations on them to implementations, we could replace those operations with computable approximations piece at a time.

If we had a λ -calculus with infinite sets as values, we could approach computability from above in a principled way, gradually changing programs for Cantor’s Paradise until they can be implemented in Church’s Purgatory.

We define that λ -calculus, λ_{ZFC} , and a call-by-value, big-step reduction semantics. To

show that it is expressive enough, we code up the real numbers, arithmetic and limits, following standard analysis. To show that it simplifies language design, we define the uncomputable limit monad in λ_{ZFC} , and derive a computable, directly implementable replacement monad by applying standard topological theorems. When certain proof obligations are met, the outputs of programs that use the computable monad converge to the same values as the outputs of programs that use the uncomputable monad.

Readers interested only in probabilistic programming languages may skip to Chapter 4, which reviews this chapter’s highlights, without missing important prerequisites.

3.2 Language Tower and Terminology

λ_{ZFC} ’s metalanguage is **first-order set theory**: first-order logic with equality extended with ZFC, or the Zermelo Fraenkel axioms and Choice (equivalently well-ordering). We also assume the existence of an inaccessible cardinal. Section 3.3 reviews the axioms, from which we will derive λ_{ZFC} ’s primitives.

To help ensure λ_{ZFC} ’s definition conservatively extends set theory, we encode its terms as sets. For example, ordered pairs of sets x and y are encoded as $\langle x, y \rangle = \{\{x\}, \{x, y\}\}$, and $\langle t_{\mathcal{P}}, \mathbb{R} \rangle = \{\{t_{\mathcal{P}}\}, \{t_{\mathcal{P}}, \mathbb{R}\}\}$ encodes the expression that applies the powerset operator to \mathbb{R} .

λ_{ZFC} ’s semantics reduces terms to terms; e.g. $\langle t_{\mathcal{P}}, \mathbb{R} \rangle$ reduces to the actual powerset of \mathbb{R} . Thus, λ_{ZFC} contains infinite terms. Infinitary languages are useful and definable: the infinitary λ -calculus [25] is an example, and Aczel’s broadly used work [3] on inductive sets treats infinite inference rules explicitly.

For convenience, we define a language λ_{ZFC}^- of finite terms and a function $\mathcal{F}[\![\cdot]\!]$ from λ_{ZFC}^- to λ_{ZFC} . We can then write $\mathcal{P} \mathbb{R}$, meaning $\mathcal{F}[\![\mathcal{P} \mathbb{R}]\!] = \langle t_{\mathcal{P}}, \mathbb{R} \rangle$.

Semantic functions like $\mathcal{F}[\![\cdot]\!]$ and the interpretation of BNF grammars are defined in set theory’s metalanguage, or the *meta*-metalanguage. Distinguishing metalanguages helps avoid paradoxes of definition such as Berry’s paradox, which are particularly easy to stumble onto when dealing with infinities.

We write λ_{ZFC}^- terms in **sans serif** font, and the metalanguage and meta-metalanguage in *math font*. We write common keywords in **bold** and invented keywords in ***bold italics***. We abbreviate proofs for space.

3.3 Metalanguage: First-Order Set Theory

We assume readers are familiar with classical first-order logic with equality and its inference rules, but not set theory. Hrbacek and Jech [20] is a fine introduction.

Set theory extends classical first-order logic with equality, which distinguishes between truth-valued formulas ϕ and object-valued terms x . Set theory allows only sets as objects, and quantifiers like “ \forall ” may range only over sets.

We define predicates and functions using “ $:=$ ”; e.g. $\text{nand}(\phi_1, \phi_2) := \neg(\phi_1 \wedge \phi_2)$. They must be nonrecursive so they can be exhaustively applied. Such definitions are **conservative extensions**: they do not prove more theorems.

To develop set theory, we make **proper extensions**, which prove more theorems, by adding symbols and axioms to first-order logic. For example, we first add “ \emptyset ” and “ \in ”, and the **empty set axiom** $\forall x. x \notin \emptyset$.

We use “ $:\equiv$ ” to define syntax; e.g. $\forall x \in A. P(x) :\equiv \forall x. (x \in A \Rightarrow P(x))$, where predicate application $P(x)$ represents a formula that may depend on x . We allow recursion in meta-metalanguage definitions if substitution terminates, so $\forall x_1 x_2 \dots x_n. \phi :\equiv \forall x_1. \forall x_2 \dots x_n. \phi$ can bind any number of names.

We already have Axiom 0 (empty set). Now for the rest.

Axiom 1 (extensionality). Define $A \subseteq B := \forall x \in A. x \in B$ and assume $A = B$ if A and B mutually are subsets; i.e. assume $\forall A B. (A \subseteq B \wedge B \subseteq A \Rightarrow A = B)$. \square

The converse follows from substituting A for B or B for A .

Axiom 2 (foundation). Define $A \not\cap B := \forall x. (x \in A \Rightarrow x \notin B)$ (“ A and B are disjoint”) and assume $\forall A. (A = \emptyset) \vee \exists x \in A. x \not\cap A$. \square

Foundation implies that the following nondeterministic procedure always terminates: If input $A = \emptyset$, return A ; otherwise restart with any $A' \in A$. Thus, sets are roots of trees in which every upward path is unbounded but finite. Foundation is analogous to “all data constructors are strict.”

Axiom 3 (powerset). Add “ \mathcal{P} ” and assume $\forall A x. (x \in \mathcal{P}(A) \iff x \subseteq A)$. \square

A **hereditarily finite** set is finite and has only hereditarily finite members. Each such set first appears in some $\mathcal{P}(\mathcal{P}(\dots\mathcal{P}(\emptyset)\dots))$. For example, after $\{x, \dots\}$ (literal set syntax) is defined, $\{\emptyset\} \in \mathcal{P}(\mathcal{P}(\emptyset))$. $\{\mathbb{R}\}$ is not hereditarily finite.

Axiom 4 (union). Add “ \bigcup ” (“big” union) and assume arbitrary unions of sets of sets exist; i.e. $\forall A x. (x \in \bigcup A \iff \exists y. x \in y \wedge y \in A)$. \square

For example, after $\{x, \dots\}$ is defined, $\bigcup\{\{x, y\}, \{y, z\}\} = \{x, y, z\}$. Also, because all objects are sets, “ \bigcup ” can extract the object in a singleton set: if $A = \{x\}$, then $x = \bigcup A$.

Axiom 5 (replacement schema). A binary predicate R can act as a function if it relates each x to exactly one y ; i.e. $\forall x \in A. \exists! y. R(x, y)$, where “ $\exists!$ ” means unique existence (read “there exists exactly one”). We cannot quantify over predicates in first-order logic, but we can assume, for each such definable R , that $\forall y. (y \in \{y' \mid x \in A \wedge R(x, y')\} \iff \exists x \in A. R(x, y))$. Roughly, treating R as a function, if R ’s domain is a set, its image (range) is also a set. \square

An **axiom schema** represents countably many axioms. If $R(n, m) \iff m = n + 1$, for example, then there is an instance of Axiom 5 for $R(n, m)$.

It is not hard to show by Axiom 5 that (after \mathbb{N} is defined) $\{m \mid n \in \mathbb{N} \wedge R(n, m)\}$ increments every natural number, yielding the set of positive naturals. But the syntax is cumbersome, so we define $\{F(x) \mid x \in A\} \equiv \{y \mid x \in A \wedge y = F(x)\}$, analogous to **map F A**, for **functional replacement**. Now the more familiar $\{n + 1 \mid n \in \mathbb{N}\}$ is the positive naturals.

It might seem replacement should be *defined* functionally, but predicates allow powerful nonconstructivism. Suppose $Q(y)$ for exactly one y . The **description operator**

$$\iota y. Q(y) \equiv \bigcup\{y \mid x \in \mathcal{P}(\emptyset) \wedge Q(y)\} \quad (3.1)$$

finds “the y such that $Q(y)$.”

From the six axioms so far, we can define $A \cup B$ (binary union), $\{x, \dots\}$ (literal finite sets), $\langle x, y, z, \dots \rangle$ (ordered pairs and lists), $\{x \in A \mid Q(x)\}$ (bounded selection), $A \setminus B$ (relative complement), $\cap A$ (“big” intersection), $\bigcup_{x \in A} F(x)$ (indexed union), $A \times B$ (cartesian product), and $A \rightarrow B$ (total function spaces). For details, we recommend Paulson’s remarkably lucid development in HOL [43].

3.3.1 The Gateway to Cantor’s Paradise: Infinity

From the six axioms so far, we cannot construct a set that is closed under unboundedly many operations, such as the language of a recursive grammar.

Example 3.1 (interpreting a grammar). We want to interpret $z ::= \emptyset \mid \langle \emptyset, z \rangle$. It should mean the least fixpoint of a function F_z , which, given a subset of z ’s language, returns a larger subset. To define F_z , replace “ \mid ” with “ \cup ”, the terminal \emptyset with $\{\emptyset\}$, and the rule $\langle \emptyset, z \rangle$ with functional replacement:

$$F_z(Z) := \{\emptyset\} \cup \{\langle \emptyset, z \rangle \mid z \in Z\} \quad (3.2)$$

We could define $Z(0) := \emptyset$, then $Z(1) := F_z(Z(0)) = \{\emptyset, \langle \emptyset, \emptyset \rangle\}$, then $Z(2) = F_z(Z(1)) = \{\emptyset, \langle \emptyset, \emptyset \rangle, \langle \emptyset, \emptyset, \emptyset \rangle\}$, and so on. The language should be the union of all the $Z(n)$, but we cannot construct it without a set of all n . \diamond

We follow Von Neumann, defining $0 := \emptyset$ as the **first ordinal number** and $s(n) := n \cup \{n\}$ to generate **successor ordinals**. Then $1 := s(0) = \{0\}$, $2 := s(1) = \{0, 1\}$, and $3 := s(2) = \{0, 1, 2\}$, and so on, so that every ordinal is defined as the set of its predecessors. The set of such numbers is the language of $n ::= 0 \mid s(n)$, which should be the least fixpoint of $F_n(N) := \{0\} \cup \{s(n) \mid n \in N\}$, similar to (3.2). Before we can prove this set exists, we must assume *some* fixpoint exists.

Axiom 6 (infinity). $\exists I. I = F_n(I)$. \square

I is a bounding set, so it may contain more than just finite ordinals. But F_n is monotone in I , so by the Knaster-Tarski theorem (suitably restricted [42]),

$$\omega := \bigcap \{N \subseteq I \mid N = F_n(N)\} \quad (3.3)$$

is the least fixpoint of F_n : the finite ordinals, a model of the natural numbers.

Example 3.2 (interpreting a grammar). We build the language defined by $z ::= \emptyset \mid \langle \emptyset, z \rangle$ recursively:

$$\begin{aligned} Z(0) &= \emptyset \\ Z(s(n)) &= F_z(Z(n)), \quad n \in \omega \\ Z(\omega) &= \bigcup_{n \in \omega} Z(n) \end{aligned} \quad (3.4)$$

By induction, $Z(n)$ exists for every $n \in \omega$; therefore $Z(\omega)$ exists, so (3.4) is a conservative extension of set theory. It is not hard to prove (also by induction) that $Z(\omega)$ is the set of all finite lists of \emptyset , and that it is the least fixpoint of F_z . \diamond

Similarly to building the language $Z(\omega)$ of z in (3.4), we can build the set $\mathcal{V}(\omega)$ of all hereditarily finite sets (see Axiom 3) by iterating \mathcal{P} instead of F_z :

$$\begin{aligned} \mathcal{V}(0) &= \emptyset \\ \mathcal{V}(s(n)) &= \mathcal{P}(\mathcal{V}(n)), \quad n \in \omega \\ \mathcal{V}(\omega) &= \bigcup_{n \in \omega} \mathcal{V}(n) \end{aligned} \quad (3.5)$$

The set ω is not just a model of the natural numbers. It is also a number itself: the **first countable ordinal**. Indeed, ω is strikingly similar to every finite ordinal in two ways. First, it is defined as the set of its predecessors. Second, it has a successor $s(\omega) = \omega \cup \{\omega\}$. (Imagine it as $\{0, 1, 2, \dots, \omega\}$.) However, unlike finite, nonzero ordinals, ω has no *immediate* predecessor—it is a **limit ordinal**.

Defining more limit ordinals allows iterating \mathcal{P} further. It is not hard to build $\omega + \omega$,

ω^2 and ω^ω as least fixpoints. The **Von Neumann hierarchy** generalizes (3.5):

$$\begin{aligned}\mathcal{V}(0) &= \emptyset \\ \mathcal{V}(s(\alpha)) &= \mathcal{P}(V(\alpha)), \text{ ordinal } \alpha \\ \mathcal{V}(\beta) &= \bigcup_{\alpha \in \beta} \mathcal{V}(\alpha), \text{ limit ordinal } \beta\end{aligned}\tag{3.6}$$

It is a theorem of ZFC that every set first appears in $\mathcal{V}(\alpha)$ for some ordinal α .

Equations (3.4,3.5,3.6) demonstrate **transfinite recursion**, set theory's **unfold**: defining a function V on ordinals, with $V(\beta)$ in terms of $V(\alpha)$ for every $\alpha \in \beta$.

3.3.2 Every Set Can Be Sequenced: Well-Ordering

A **sequence** is a total function from an ordinal to a codomain; e.g. $f \in 3 \rightarrow A$ is a length-3 sequence of A 's elements. (An ordinal is comprised of its predecessors, so $3 = \{0, 1, 2\}$.) A **well-order** of A is a bijective sequence of A 's elements.

Axiom 7 (well-ordering). Suppose the predicate *Ord* identifies ordinals and $B \leftrightarrow A$ is the set of all bijective mappings from B to A . Assume $\forall A. \exists \alpha f. \text{Ord}(\alpha) \wedge f \in \alpha \leftrightarrow A$; i.e. every set can be well-ordered. \square

Because the bijective sequence f is not unique, a well-ordering primitive could make λ_{ZFC} 's semantics nondeterministic. Fortunately, the existence of a cardinality operator is equivalent to well-ordering [52], so we will give λ_{ZFC} a cardinality primitive.

The **cardinality** of a set A is the smallest ordinal that can be put in bijection with A . Formally, if F is the set of A 's well-orderings, then $|A| = \bigcap \{\text{domain}(f) \mid f \in F\}$.

3.3.3 Infinity's Infinity: An Inaccessible Cardinal

The set $\mathcal{V}(\omega)$ of hereditarily finite sets is closed under powerset, union, replacement (with predicates restricted to $\mathcal{V}(\omega)$), and cardinality. It is also **transitive**: if $A \in \mathcal{V}(\omega)$, then $x \in \mathcal{V}(\omega)$ for all $x \in A$. These closure properties make it a **Grothendieck universe**: a set that acts like a set of all sets.

$$\begin{aligned}
e &::= n \mid v \mid e e \mid \text{if } e e e \mid e \in e \mid \bigcup e \mid \text{take } e \mid \mathcal{P} e \mid \text{image } e e \mid \text{card } e \\
v &::= \text{false} \mid \text{true} \mid \lambda. e \mid \emptyset \mid \omega \\
n &::= 0 \mid 1 \mid 2 \mid \dots
\end{aligned}$$

Figure 3.1: The definition of λ_{ZFC}^- , which represents countably many λ_{ZFC} terms.

λ_{ZFC} 's values should contain ω and be closed under its primitives. But a Grothendieck universe containing ω cannot be proved from the typical axioms. If it exists, it must be equal to $\mathcal{V}(\kappa)$ for some **inaccessible cardinal** κ .

Axiom 8 (inaccessible cardinal). Suppose $GU(V)$ if and only if V is a Grothendieck universe. Add “ κ ” and assume $\text{Ord}(\kappa) \wedge (\kappa > \omega) \wedge GU(\mathcal{V}(\kappa))$. \square

We call the sets in $\mathcal{V}(\kappa)$ **hereditarily accessible**.

Inaccessible cardinals are not usually assumed but are widely believed consistent. Set theorists regard them as no more dangerous than ω . Interpreting category theory with small and large categories, second-order set theory, or CIC in first-order set theory requires at least one inaccessible cardinal [6, 53, 55].

Constructing a set $A \notin \mathcal{V}(\kappa)$ requires assuming κ or an equivalent, so $\mathcal{V}(\kappa)$ easily contains most mathematics. In fact, most can be modeled well within $\mathcal{V}(2^\omega)$; e.g. the model of \mathbb{R} we define in Section 3.7 is in $\mathcal{V}(\omega + 11)$. Besides, if λ_{ZFC} needed to contain large cardinals, we could always assume even larger ones.

3.4 λ_{ZFC} 's Grammar

We define λ_{ZFC} 's terms in three steps. First, we define λ_{ZFC}^- , a language of finite terms with primitives that correspond with the ZFC axioms. Second, we encode these terms as sets. Third, guided by the first two steps, we define λ_{ZFC} by defining its terms, most of which are infinite, as sets in $\mathcal{V}(\kappa)$.

Figure 3.1 shows λ_{ZFC}^- 's grammar. Expressions e are typical: variables, values, appli-

$$\begin{array}{lcl}
\text{Distinct } t_{\text{var}}, t_{\text{app}}, t_{\text{if}}, t_{\in}, t_{\cup}, t_{\text{take}}, t_{\mathcal{P}}, t_{\text{image}}, t_{\text{card}}, t_{\text{set}}, t_{\text{atom}}, t_{\lambda}, t_{\text{false}}, t_{\text{true}} \\
\mathcal{F}[[n]] & := & \langle t_{\text{var}}, n \rangle & \mathcal{F}[[\emptyset]] & := & \text{set}(\emptyset) & \mathcal{F}[[\omega]] & := & \text{set}(\omega) \\
\mathcal{F}[[e_f e_x]] & := & \langle t_{\text{app}}, \mathcal{F}[[e_f]], \mathcal{F}[[e_x]] \rangle & \mathcal{F}[[\text{false}]] & := & a_{\text{false}} & a_{\text{false}} & := & \langle t_{\text{atom}}, t_{\text{false}} \rangle \\
\mathcal{F}[[e_x \in e_A]] & := & \langle t_{\in}, \mathcal{F}[[e_x]], \mathcal{F}[[e_A]] \rangle & \mathcal{F}[[\text{true}]] & := & a_{\text{true}} & a_{\text{true}} & := & \langle t_{\text{atom}}, t_{\text{true}} \rangle \\
& \dots & & \text{set}(A) & = & \langle t_{\text{set}}, \{\text{set}(x) \mid x \in A\} \rangle
\end{array}$$

Figure 3.2: The semantic function $\mathcal{F}[[\cdot]]$ from λ_{ZFC}^- terms to λ_{ZFC} terms.

cation, if, and domain-specific primitives, for membership, union, extraction (**take**), powerset, functional replacement (**image**), and cardinality. Values v are also typical: booleans and lambdas, and the domain-specific constants \emptyset and ω .

In set theory, $\cup \{A\} = A$ holds for all A , so \cup can extract the element from a singleton. In λ_{ZFC} , the encoding of $\cup \{A\}$ reduces to A only if A is an encoded set. Therefore, the primitives must include **take**, which extracts A from $\{A\}$. In particular, extracting a lambda from an ordered pair requires **take**.

We use De Bruijn indexes with 0 referring to the innermost binding. Because we will define λ_{ZFC} terms as well-founded sets, by Axiom 2, countably many indexes is sufficient for λ_{ZFC} as well as λ_{ZFC}^- .

Figure 3.2 shows part of the meta-metalanguage function $\mathcal{F}[[\cdot]]$ that encodes λ_{ZFC}^- terms as λ_{ZFC} terms. It distinguishes sorts of terms in the standard way, by pairing them with tags; e.g. if t_{set} is the “set” tag, then $\langle t_{\text{set}}, \emptyset \rangle$ encodes \emptyset .

To recursively tag sets, we add the axiom $\text{set}(A) = \langle t_{\text{set}}, \{\text{set}(x) \mid x \in A\} \rangle$. The **well-founded recursion theorem** proves that for all A , $\text{set}(A)$ exists, so this axiom is a conservative extension. The actual proof is tedious, but in short, set is structurally recursive. Now $\text{set}(\emptyset) = \langle t_{\text{set}}, \emptyset \rangle$ and $\text{set}(\omega)$ encodes ω .

3.4.1 An Infinite Set Rule For Finite BNF Grammars

There is no sensible reduction relation for λ_{ZFC}^- . (For example, $\mathcal{P} \emptyset$ cannot correctly reduce to a value because no value in λ_{ZFC}^- corresponds to $\{\emptyset\}$.) The easiest way to ensure a reduction relation exists for λ_{ZFC} is to include encodings of all the sets in $\mathcal{V}(\kappa)$ as values.

To define λ_{ZFC} 's terms, we first extend BNF with a set rule: $\{y^{*\alpha}\}$, where α is a cardinal number. Roughly, it means sets comprised of no more than α terms from the language of y . Formally, it means $\mathcal{P}_{<}(Y, \alpha)$, where Y is a subset of y 's language generated while building a least fixpoint, and the bounded powerset operation is defined by

$$\mathcal{P}_{<}(Y, \alpha) := \{x \in \mathcal{P}(Y) \mid |x| < \alpha\} \quad (3.7)$$

meaning $\mathcal{P}_{<}(Y, \alpha)$ returns all subsets of Y with cardinality less than α .

Example 3.3 (finite sets). The grammar $h ::= \{h^{*\omega}\}$ should represent all hereditarily finite sets, or $\mathcal{V}(\omega)$. Intuitively, the single rule for h should be equivalent to countably many rules $h ::= \{\} \mid \{h\} \mid \{h, h\} \mid \{h, h, h\} \mid \dots$.

Its language is the least fixpoint of $F_h(H) := \mathcal{P}_{<}(H, \omega)$. Further on, we will prove that F_h 's least fixpoint is $\mathcal{V}(\omega)$ using a general theorem. \diamond

Example 3.4 (accessible sets). The language of $a ::= \{a^{*\kappa}\}$ is the least fixpoint of $F_a(A) := \mathcal{P}_{<}(A, \kappa)$, which should be $\mathcal{V}(\kappa)$. \diamond

The following theorem schemas will make it easy to find least fixpoints.

Theorem 3.5. *Let F be a unary function. Define V by transfinite recursion:*

$$\begin{aligned} V(0) &= \emptyset \\ V(s(\alpha)) &= F(V(\alpha)) \\ V(\beta) &= \bigcup_{\alpha \in \beta} V(\alpha), \text{ limit ordinal } \beta \end{aligned} \quad (3.8)$$

*Let γ be an ordinal. If F is monotone on $V(\gamma)$, V is monotone on γ , and $V(\gamma)$ is a fixpoint of F , then $V(\gamma)$ is also the **least** fixpoint of F .*

Proof. By induction: successor case by monotonicity; limit by a property of \bigcup . \square

All the F s we define are monotone. In particular, the interpretations of $\{y^{*\alpha}\}$ rules are monotone because \mathcal{P} is monotone. Further, all the F s we define give rise to a monotone V . Grammar terminals “seed” every iteration with singleton sets, and $\{y^{*\alpha}\}$ rules seed every iteration with \emptyset .

From here on, we write F^α instead of $V(\alpha)$ to mean α iterations of F .

Theorem 3.6. *Suppose a grammar with $\{y^{*\alpha}\}$ rules and iterating function F . The language of the grammar, F ’s least fixpoint, is F^γ , where γ is a regular cardinal not less than any α .*

Proof. Fixpoint by Aczel [3, Theorem 1.3.4]; least fixpoint by Theorem 3.5. \square

Example 3.7 (finite sets). Because ω is regular, by Theorem 3.6, F_h ’s least fixpoint is F_h^ω . Further, $F_h(H) = \mathcal{P}(H)$ for all hereditarily finite H , and $\mathcal{V}(\omega)$ is closed under \mathcal{P} , so $F_h^\omega = \mathcal{V}(\omega)$, the set of all hereditarily finite sets. \diamond

Example 3.8 (accessible sets). By a similar argument, F_a ’s least fixpoint is $F_a^\kappa = \mathcal{V}(\kappa)$, the set of all hereditarily accessible sets. \diamond

Example 3.9 (encoded accessible sets). The language of $v ::= \langle t_{\text{set}}, \{v^{*\kappa}\} \rangle$ is comprised of the *encodings* of all the hereditarily accessible sets. \diamond

3.4.2 The Grammar of Infinite, Encoded Terms

There are three main differences between λ_{ZFC} ’s grammar in Fig. 3.3 and λ_{ZFC}^- ’s grammar in Fig. 3.1. First, λ_{ZFC} ’s grammar defines a language of terms that are already encoded as sets. Second, instead of the symbols \emptyset and ω , it includes, as values, encoded sets of values. Most of these value terms are infinite, such as the encoding of ω . Third, it includes encoded sets of *expressions*.

$$\begin{aligned}
e &::= n \mid v \mid \langle t_{\text{app}}, e, e \rangle \mid \langle t_{\text{if}}, e, e, e \rangle \mid \langle t_{\in}, e, e \rangle \mid \langle t_{\cup}, e \rangle \mid \langle t_{\text{take}}, e \rangle \mid \langle t_{\mathcal{P}}, e \rangle \mid \\
&\quad \langle t_{\text{image}}, e, e \rangle \mid \langle t_{\text{card}}, e \rangle \mid \langle t_{\text{set}}, \{e^{*\kappa}\} \rangle \\
v &::= a_{\text{false}} \mid a_{\text{true}} \mid \langle t_{\lambda}, e \rangle \mid \langle t_{\text{set}}, \{v^{*\kappa}\} \rangle \\
n &::= \langle t_{\text{var}}, 0 \rangle \mid \langle t_{\text{var}}, 1 \rangle \mid \dots
\end{aligned}$$

Figure 3.3: λ_{ZFC} 's grammar. Here, $\{e^{*\kappa}\}$ means sets comprised of no more than κ terms from the language of e .

The language of n is $N := \{\langle t_{\text{var}}, i \rangle \mid i \in \omega\}$. The rules for e and v are mutually recursive. Interpreted, but leaving out some of e 's rules, they are

$$\begin{aligned}
F_e(E, V) &:= N \cup V \cup \{\langle t_{\text{app}}, e_f, e_x \rangle \mid \langle e_f, e_x \rangle \in E \times E\} \cup \dots \cup \{\langle t_{\text{set}}, e \rangle \mid e \in \mathcal{P}_{<}(E, \kappa)\} \\
F_v(E, V) &:= \{a_{\text{false}}, a_{\text{true}}\} \cup \{\langle t_{\lambda}, e \rangle \mid e \in E\} \cup \{\langle t_{\text{set}}, v \rangle \mid v \in \mathcal{P}_{<}(V, \kappa)\}
\end{aligned} \tag{3.9}$$

To use Theorem 3.6, we need to iterate a single function. Note that the language pair $\langle E, V \rangle = \langle \{e, \dots\}, \{v, \dots\} \rangle$ is isomorphic to the single set of tagged terms $EV = \{\langle 0, e \rangle, \dots, \langle 1, v \rangle, \dots\}$. Binary **disjoint union**, denoted $E \sqcup V$, creates such sets. We define F_{ev} by $F_{ev}(E \sqcup V) = F_e(E, V) \sqcup F_v(E, V)$. By Theorem 3.6, its least fixpoint is F_{ev}^{κ} , so we define E and V by $E \sqcup V = F_{ev}^{\kappa}$.

To make well-founded substitution easy, we will use capturing substitution, which does not capture when used on closed terms. Let $Cl(e)$ indicate whether a term is closed—this is structurally recursive. Then $E' := \{e \in E \mid Cl(e)\}$ and $V' := \{v \in V \mid Cl(v)\}$ contain only closed terms. Lastly, we define $\lambda_{\text{ZFC}} := E'$.

3.5 λ_{ZFC} 's Big-Step Reduction Semantics

We distinguish sets from other expressions using E_{set} and V_{set} , which merely check tags. We also lift set constructors to operate on encoded sets. For example, for cardinality, $\widehat{C}(v_A) := \text{set}(|\text{snd}(v_A)|)$ extracts the tagged set from v_A , applies $|\cdot|$, and recursively tags the

$$\frac{}{v \Downarrow v} \text{ (val)} \quad \frac{e_f \Downarrow \langle t_\lambda, e_y \rangle \quad e_x \Downarrow v_x \quad e_y[0 \setminus v_x] \Downarrow v_y}{\langle t_{\text{app}}, e_f, e_x \rangle \Downarrow v_y} \text{ (ap)} \quad \frac{e_c \Downarrow a_{\text{true}} \quad e_t \Downarrow v_t \quad e_c \Downarrow a_{\text{false}} \quad e_f \Downarrow v_f}{\langle t_{\text{if}}, e_c, e_t, e_f \rangle \Downarrow v_t} \text{ (if)}$$

(a) Standard call-by-value reduction rules

$$\frac{e_A \Downarrow v_A \quad V_{\text{set}}(v_A) \quad e_x \Downarrow v_x \quad v_x \in \text{snd}(v_A)}{\langle t_{\in}, e_x, e_A \rangle \Downarrow a_{\text{true}}} \quad \frac{e_A \Downarrow v_A \quad V_{\text{set}}(v_A) \quad e_x \Downarrow v_x \quad v_x \notin \text{snd}(v_A)}{\langle t_{\in}, e_x, e_A \rangle \Downarrow a_{\text{false}}} \text{ (in)}$$

$$\frac{e_A \Downarrow v_A \quad V_{\text{set}}(v_A) \quad \forall v_x \in \text{snd}(v_A). V_{\text{set}}(v_x)}{\langle t_{\cup}, e_A \rangle \Downarrow \widehat{U}(v_A)} \text{ (union)} \quad \frac{e_A \Downarrow v_A \quad V_{\text{set}}(v_A)}{\langle t_{\mathcal{P}}, e_A \rangle \Downarrow \widehat{\mathcal{P}}(v_A)} \text{ (pow)}$$

$$\frac{e_A \Downarrow v_A \quad V_{\text{set}}(v_A) \quad e_f \Downarrow \langle t_\lambda, e_y \rangle \quad \widehat{I}(\langle t_\lambda, e_y \rangle, v_A) \Downarrow v_y}{\langle t_{\text{image}}, e_f, e_A \rangle \Downarrow v_y} \text{ (image)} \quad \frac{e_A \Downarrow v_A \quad V_{\text{set}}(v_A)}{\langle t_{\text{card}}, e_A \rangle \Downarrow \widehat{C}(v_A)} \text{ (card)}$$

$$\frac{E_{\text{set}}(e_A) \quad \forall e_x \in \text{snd}(e_A). \exists v_x. e_x \Downarrow v_x}{e_A \Downarrow \langle t_{\text{set}}, \{v_x \mid e_x \in \text{snd}(e_A) \wedge e_x \Downarrow v_x\} \rangle} \text{ (set)} \quad \frac{e_A \Downarrow \langle t_{\text{set}}, \{v_x\} \rangle}{\langle t_{\text{take}}, e_A \rangle \Downarrow v_x} \text{ (take)}$$

(b) λ_{ZFC} -specific rules

Figure 3.4: Reduction rules defining λ_{ZFC} 's big-step, call-by-value semantics.

resulting cardinal number. The rest are

$$\begin{aligned} \widehat{\mathcal{P}}(v_A) &:= \langle t_{\text{set}}, \{ \langle t_{\text{set}}, v_x \rangle \mid v_x \in \mathcal{P}(\text{snd}(v_A)) \} \rangle \\ \widehat{U}(v_A) &:= \langle t_{\text{set}}, \cup \{ \text{snd}(v_x) \mid v_x \in \text{snd}(v_A) \} \rangle \\ \widehat{I}(v_f, v_A) &:= \langle t_{\text{set}}, \{ \langle t_{\text{app}}, v_f, v_x \rangle \mid v_x \in \text{snd}(v_A) \} \rangle \end{aligned} \tag{3.10}$$

All but \widehat{I} return values. Sets returned by \widehat{I} are intended to be reduced further.

We use $e[n \setminus v]$ for De Bruijn substitution. Because e and v are closed, it is easy to define it using simple structural recursion on terms; it is thus conservative.

Figure 3.4 shows the reduction rules that define the reduction relation “ \Downarrow ”. Figure 3.4a has standard call-by-value rules: values reduce to themselves, and applications reduce by substitution. Figure 3.4b has the λ_{ZFC} -specific rules. Most simply use V_{set} to check tags before applying a lifted operator. The (image) rule replaces each value v_x in the set v_A with an application, generating a set expression, and the (set) rule reduces all the terms inside a set expression.

To define “ \Downarrow ” as a least fixpoint, we adapt Aczel’s treatment [3]. We first define a bounding set for “ \Downarrow ” using closed terms, or $\mathcal{U} := E' \times V'$, so that $\Downarrow \subseteq \mathcal{U}$.

The rules in Fig. 3.4 can be used to define a predicate $D(R, \langle e, v \rangle)$. This predicate indicates whether some reduction rule, after replacing every “ \Downarrow ” in its premise with the approximation R , derives the conclusion $e \Downarrow v$.¹ Using D , we define a function that derives new conclusions from the known conclusions in R :

$$F_{\Downarrow}(R) := \{c \in \mathcal{U} \mid D(R, c)\} \quad (3.11)$$

For example, $F_{\Downarrow}(\emptyset) = \{\langle v, v \rangle \mid v \in V\}$, by the (val) rule. $F_{\Downarrow}(F_{\Downarrow}(\emptyset))$ includes all pairs of non-value expressions and the values they reduce to in one derivation, as well as $\{\langle v, v \rangle \mid v \in V\}$. Generally, (val) ensures that iterating F_{\Downarrow} is monotone.

For F_{\Downarrow} itself to be nonmonotone, for some $R \subseteq R' \subseteq \mathcal{U}$, there would have to be a conclusion $c \in F_{\Downarrow}(R)$ that is not in $F_{\Downarrow}(R')$. In other words, having more known conclusions could falsify a premise. None of the rules in Fig. 3.4 can do so.

Because F_{\Downarrow} is monotone and iterating it is monotone, we can define $\Downarrow := F_{\Downarrow}^{\gamma}$ for some ordinal γ . If λ_{ZFC} had only finite terms, $\gamma = \omega$ iterations would reach a fixpoint. But a simple countable term shows why “ \Downarrow ” cannot be F_{\Downarrow}^{ω} .

Example 3.10 (countably infinite term). If s is the successor function in λ_{ZFC} , the term $t := \langle t_{\text{set}}, \{0, \langle t_{\text{app}}, s, 0 \rangle, \langle t_{\text{app}}, s, \langle t_{\text{app}}, s, 0 \rangle \rangle, \dots \} \rangle$ should reduce to $\text{set}(\omega)$. The (set) rule’s premises require each of t ’s subterms to reduce—using at least F_{\Downarrow}^{ω} because each subterm requires a finite, unbounded number of (ap) derivations. Though $F_{\Downarrow}^{s(\omega)}$ reduces t , for larger terms, we must iterate F_{\Downarrow} much further. \diamond

Theorem 3.11. $\Downarrow := F_{\Downarrow}^{\kappa}$ is the least fixpoint of F_{\Downarrow} .

Proof. Fixpoint by Aczel [3, Theorem 1.3.4]; least fixpoint by Theorem 3.5. \square

Lastly, ZFC theorems that do not depend on κ can be applied to λ_{ZFC} terms.

¹ D is definable in first-order logic, but its definition does not aid understanding much.

Theorem 3.12. λ_{ZFC} 's set values and $\langle t_{\in}, \cdot, \cdot \rangle$ are a model of ZFC- κ .

Proof. $\mathcal{V}(\kappa)$, a model of ZFC- κ , is isomorphic to $v ::= \langle t_{\text{set}}, \{v^{*\kappa}\} \rangle$. □

3.6 Syntactic Sugar and a Small Set Library

From here on, we write only λ_{ZFC}^- terms, assume $\mathcal{F}[\cdot]$ is applied, and no longer distinguish λ_{ZFC}^- from λ_{ZFC} .

We use names instead of De Bruijn indexes and assume names are converted. We get alpha equivalence for free; for example, $\lambda x. x = \langle t_{\lambda}, \langle t_{\text{var}}, 0 \rangle \rangle = \lambda y. y$.

λ_{ZFC} does not contain terms with free variables. To get around this technical limitation, we assume free variables are metalanguage names for closed terms.

We allow the primitives (\in) , \bigcup , **take**, \mathcal{P} , **image** and **card** to be used as if they were functions. Enclosing infix operators in parenthesis refers to them as functions, as in (\in) . We partially apply infix functions using Haskell-like sectioning rules, so $(x \in)$ means $\lambda A. x \in A$ and $(\in A)$ means $\lambda x. x \in A$.

We define first-order objects using “ $:=$ ”, as in $0 := \emptyset$, and syntax with “ \equiv ”, as in $\lambda x_1 x_2 \dots x_n. e \equiv \lambda x_1. \lambda x_2 \dots x_n. e$ to automatically curry. Function definitions expand to lambdas (using fixpoint combinators for recursion); for example, $x = y := x \in \{y\}$ and $(=) := \lambda x y. x \in \{y\}$ equivalently define $(=)$ in terms of (\in) . We destructure pairs implicitly in binding patterns, as in $\lambda \langle x, y \rangle. f \ x \ y$.

To do anything useful, we need a small set library. The definitions are similar to the metalanguage definitions we omitted in Section 3.3, and we similarly elide most of the λ_{ZFC} definitions. However, some deserve special mention.

Because λ_{ZFC} has only *functional* replacement, we cannot define unbounded “ \forall ” and

“ \exists ”. But we can define bounded quantifiers in terms of bounded selection, or

$$\begin{aligned}\text{select } f \ A &:= \bigcup (\text{image } (\lambda x. \text{if } (f \ x) \ \{x\} \ \emptyset) \ A) \\ \forall x \in e_A. e_f &\equiv (\text{select } (\lambda x. e_f) \ e_A) = e_A \\ \exists x \in e_A. e_f &\equiv (\text{select } (\lambda x. e_f) \ e_A) = \emptyset\end{aligned}\tag{3.12}$$

We also define a bounded description operator using the **filter**-like **select**:

$$\iota x \in e_A. e_f \equiv \text{take } (\text{select } (\lambda x. e_f) \ e_A)\tag{3.13}$$

Note $\iota x \in e_A. e_f$ reduces only if $e_f \Downarrow \text{true}$ for exactly one $x \in e_A$.

Unlike in first-order logic, converting a predicate to an object in λ_{ZFC} requires a bounding set as well as unique existence. For example, if

$$\langle e_x, e_y \rangle \equiv \{\{e_x\}, \{e_x, e_y\}\}\tag{3.14}$$

defines ordered pairs, then

$$\text{fst } p := \iota x \in (\bigcup p). \exists y \in (\bigcup p). p = \langle x, y \rangle\tag{3.15}$$

takes the first element by identifying it in the bounding set $\bigcup p$ using a predicate.

The **set monad** simulates nondeterministic choice. We define it by

$$\begin{aligned}\text{return}_{\text{set}} \ a &:= \{a\} \\ \text{bind}_{\text{set}} \ A \ f &:= \bigcup (\text{image } f \ A)\end{aligned}\tag{3.16}$$

Using $\text{bind } m \ f = \text{join } (\text{lift } f \ m)$, evidently $\text{lift}_{\text{set}} := \text{image}$ and $\text{join}_{\text{set}} := \bigcup$. The proofs of the monad laws follow the proofs for the list monad. We also define

$$\{x \in e_A\}. e_f \equiv \text{bind}_{\text{set}} (\lambda x. e_f) \ e_A\tag{3.17}$$

read “choose x in e_A , then e_f .” For example, binary cartesian product is defined by

$$A \times B := \{x \in A\}. \{y \in B\}. \text{return}_{\text{set}} \langle x, y \rangle \quad (3.18)$$

Every $f \in A \rightarrow B$ is shaped $f = \{\langle x_1, y_1 \rangle, \langle x_2, y_2 \rangle, \dots\}$ and is total on A . To distinguish these hash tables from lambdas, we call them **mappings**. Mappings can be applied by $\text{ap } f \ x := \iota y \in (\text{range } f). \langle x, y \rangle \in f$, but we write just $f \ x$. We define

$$\text{mapping } f \ A := \text{image } (\lambda x. \langle x, f \ x \rangle) \ A \quad (3.19)$$

to convert a lambda to a mapping or to restrict a mapping to A . We usually use

$$\lambda x \in e_A. e_y := \text{mapping } (\lambda x. e_y) \ e_A \quad (3.20)$$

to define mappings.

A sequence of A is a mapping $xs \in \alpha \rightarrow A$ for some ordinal α . For example, $ns := \lambda n \in \omega. n$ is a countable sequence in $\omega \rightarrow \omega$ of increasing finite ordinals. We assume useful sequence functions like **map**, **map2** and **drop** are defined.

3.7 Example: The Reals From the Rationals

Here, we demonstrate that λ_{ZFC} is computationally powerful enough to construct the real numbers. For a clear, well-motivated, rigorous treatment in first-order set theory without lambdas, we recommend Abbott’s excellent introductory text [2].

Assume we have a model $\mathbb{Q}, +_{\mathbb{Q}}, -_{\mathbb{Q}}, \times_{\mathbb{Q}}, \div_{\mathbb{Q}}$ of the rationals and rational arithmetic.² To get the reals, we close the rationals under countable limits.

We represent limits of rationals with sequences in $\omega \rightarrow \mathbb{Q}$. To select only the converging ones, we must define what convergence means. We start with convergence to zero

²Though the λ_{ZFC} development of \mathbb{Q} is short and elegant, it does not fit in this paper.

and equivalence. Given \mathbb{Q}^+ , ' $<_{\mathbb{Q}}$ ' and $|\cdot|_{\mathbb{Q}}$, define

$$\begin{aligned}\text{conv-zero?}_R \text{ xs} &:= \forall \varepsilon \in \mathbb{Q}^+. \exists N \in \omega. \forall n \in \omega. (N \in n \Rightarrow |\text{xs } n|_{\mathbb{Q}} <_{\mathbb{Q}} \varepsilon) \\ \text{xs} =_R \text{ ys} &:= \text{conv-zero?}_R (\text{map2 } (-_{\mathbb{Q}}) \text{ xs ys})\end{aligned}\tag{3.21}$$

So a sequence $\text{xs} \in \omega \rightarrow \mathbb{Q}$ converges to zero if, for any positive ε , there is some index N after which all xs are smaller than ε . Two sequences are equivalent ($=_R$) if their pointwise difference converges to zero.

We should be able to drop finitely many elements from a converging sequence without changing its limit. Therefore, a sequence of rationals converges to *something* when it is equivalent to all of its suffixes. We thus define an equivalent to the Cauchy convergence test, and use it to select the converging sequences:

$$\begin{aligned}\text{conv?}_R \text{ xs} &:= \forall n \in \omega. \text{xs} =_R (\text{drop } n \text{ xs}) \\ R &:= \text{select conv?}_R (\omega \rightarrow \mathbb{Q})\end{aligned}\tag{3.22}$$

But R (equipped with the equivalence relation $=_R$) is not the real numbers as they are normally defined: converging sequences in R may be equivalent but not equal. To decide real equality using λ_{ZFC} 's "=", we partition R into disjoint sets of equivalent sequences—we make a **quotient space**. Thus,

$$\begin{aligned}\text{quotient } A (=_A) &:= \text{image } (\lambda x. \text{select } (=_A x) A) A \\ \mathbb{R} &:= \text{quotient } R (=_R)\end{aligned}\tag{3.23}$$

defines the reals with extensional equality.

To define real arithmetic, we must lift rational arithmetic to sequences and then to sets of sequences. The `map2` function lifts, say, $+_{\mathbb{Q}}$ to sequences, as in $(+_R) := \text{map2 } (+_{\mathbb{Q}})$. To lift $+_R$ to sets of sequences, note that sets of sequences are models of nondeterministic sequences, suggesting the set monad. We define $\text{lift2}_{\text{set}} f A B := \{a \in A\}. \{b \in B\}. \text{return}_{\text{set}} (f a b)$ to lift two-argument functions to the set monad. Now $(+) := \text{lift2}_{\text{set}} (+_R)$, and similarly for the other operators.

Using $\text{lift2}_{\text{set}}$ is atypical, so we prove that $A + B \in \mathbb{R}$ when $A \in \mathbb{R}$ and $B \in \mathbb{R}$, and similarly for the other operators. It follows from the fact that the rational operators lifted to sequences are surjective morphisms, and this theorem:

Theorem 3.13. *Suppose $=_X$ is an equivalence relation on X , and define its quotient $\mathbb{X} := \text{quotient } X (=_X)$. If op is surjective on X and a binary morphism for $=_X$, then $(\text{lift2}_{\text{set}} \text{ op } A \ B) \in \mathbb{X}$ for all $A \in \mathbb{X}$ and $B \in \mathbb{X}$.*

Proof. Reduce to an equality. Case “ \subseteq ” by morphism; case “ \supseteq ” by surjection. \square

Now for real limits. If \mathbb{R}^+ , ‘ $<$ ’, and $|\cdot|$ are defined, we can define $\text{conv-zero?}_{\mathbb{R}}$, which is like (3.21) but operates on real sequences $\text{xs} \in \omega \rightarrow \mathbb{R}$. We then define $\text{limit}_{\mathbb{R}} \text{ xs} := \iota y \in \mathbb{R}. \text{conv-zero?}_{\mathbb{R}} (\text{map } (-y) \text{ xs})$ to calculate their limits.

From here, it is not difficult to treat \mathbb{Q} and \mathbb{R} uniformly by redefining $\mathbb{Q} \subset \mathbb{R}$.

3.8 Example: Computable Real Limits

Exact real computation has been around since Turing’s seminal paper [51]. The novelty here is how we do it. We define the **limit monad** in λ_{ZFC} for expressing calculations involving limits, with bind_{lim} defined in terms of a general limit. We then derive a limit-free, computable replacement $\text{bind}'_{\text{lim}}$. Replacing bind_{lim} with $\text{bind}'_{\text{lim}}$ in a λ_{ZFC} term incurs proof obligations. If they can be met, the computable λ_{ZFC} term has the same limit as the original, uncomputable term.

In other words, entirely in λ_{ZFC} , we define uncomputable things, and gradually turn them into computable, directly implementable approximations.

The proof obligations are related to topological theorems [37] that we will import as lemmas. By Theorem 3.12, we are allowed to use them directly.

At this point, it is helpful to have a simple, informal type system, which we can easily add to the untyped λ_{ZFC} . $A \Rightarrow B$ is a lambda or mapping type. $A \rightarrow B$ is the set of total mappings from A to B . A set is a membership proposition.

3.8.1 The Limit Monad

We first need a universe \mathbb{U} of values that is closed under sequencing; i.e. if $A \subset \mathbb{U}$ then so is $\omega \rightarrow A$. Define \mathbb{U} as the language of $u ::= \mathbb{R} \mid \omega \rightarrow u$. A complete product metric $\delta : \mathbb{U} \Rightarrow \mathbb{U} \Rightarrow \mathbb{R}$ exists; therefore, a function $\text{limit} : (\omega \rightarrow \mathbb{U}) \Rightarrow \mathbb{U}$ similar to $\text{limit}_{\mathbb{R}}$ exists that calculates limits. These are all λ_{ZFC} -definable.

The limit monad's computations are of type $\omega \rightarrow \mathbb{U}$. The type does not imply convergence, which must be proved separately. Its `run` function is `limit`.

Example 3.14 (infinite series). Define `partial-sums` : $(\omega \rightarrow \mathbb{R}) \Rightarrow (\omega \rightarrow \mathbb{R})$ first by `partial-sums' xs := $\lambda n.$ if (n = 0) (xs 0) ((xs n) + (partial-sums' xs (n - 1)))`. (The sequence is recursively defined, so we cannot use $\lambda n \in \omega. e$ to immediately create it.) Then convert it to a mapping: `partial-sums xs := mapping (partial-sums' xs) ω` .

Now $\sum_{n \in \omega} e \equiv \text{limit} (\text{partial-sums } \lambda n \in \omega. e)$, or the limit of partial sums. Even if `xs` converges, `partial-sums xs` may not; e.g. if `xs` = $\lambda n \in \omega. \frac{1}{n+1}$. \diamond

The limit monad's `returnlim` : $\mathbb{U} \Rightarrow (\omega \rightarrow \mathbb{U})$ creates constant sequences, and its `bindlim` : $(\omega \rightarrow \mathbb{U}) \Rightarrow (\mathbb{U} \Rightarrow (\omega \rightarrow \mathbb{U})) \Rightarrow (\omega \rightarrow \mathbb{U})$ simply takes a limit:

$$\begin{aligned} \text{return}_{\text{lim}} x &:= \lambda n \in \omega. x \\ \text{bind}_{\text{lim}} xs f &:= f (\text{limit } xs) \end{aligned} \tag{3.24}$$

The left identity and associativity monad laws hold using “=” for equivalence. However, right identity holds only in the limit, so we define equivalence by $xs =_{\text{lim}} ys := \text{limit } xs = \text{limit } ys$.

Example 3.15 (lifting). Define `liftlim` f xs := `bindlim xs $\lambda x.$ returnlim (f x)`, as is typical. Substituting `bindlim` and reducing reveals that `f (limit xs) = limit (liftlim f xs)`. That is, using `liftlim` pulls `limit` out of `f`'s argument. \diamond

Example 3.16 (exponential). The Taylor series expansion of the exponential function is `exp-seq` : $\mathbb{R} \Rightarrow (\omega \rightarrow \mathbb{R})$, defined by `exp-seq x := partial-sums $\lambda n \in \omega. \frac{x^n}{n!}$` . It always converges,

so $\text{limit} (\text{exp-seq } x) = \sum_{n \in \omega} \frac{x^n}{n} = \text{exp } x$ for $x \in \mathbb{R}$. To exponentiate converging sequences, define $\text{exp}_{\text{lim}} xs := \text{bind}_{\text{lim}} xs \text{ exp-seq}$. \diamond

3.8.2 The Computable Limit Monad

We derive the computable limit monad in two steps. In the first, longest step, we replace the limit monad's defining functions with those that do not use `limit`. But computations will still have type $\omega \rightarrow \mathbb{U}$, whose inhabitants are not directly implementable, so in the second step, we give them a lambda type.

We define $\text{return}'_{\text{lim}} := \text{return}_{\text{lim}}$. A drop-in, limit-free replacement for bind_{lim} does not exist, but there is one that incurs three proof obligations. Without imposing rigid constraints on using bind_{lim} , we cannot meet them automatically. But we can separate them by factoring bind_{lim} into lift_{lim} and join_{lim} .

Limit-Free Lift. Substituting to get $\text{lift}_{\text{lim}} f xs = \text{return}_{\text{lim}} (f (\text{limit } xs))$ exposes the use of `limit`. Removing it requires continuity and definedness.

Lemma 3.17 (continuity in metric spaces). *Let $f : A \Rightarrow B$ with A a metric space. Then f is continuous at $x \in A$ if and only if for all $xs \in \omega \rightarrow A$ for which $\text{limit } xs = x$ and f is defined on all elements of xs , $f (\text{limit } xs) = \text{limit} (\text{map } f xs)$.*

So if $f : \mathbb{U} \Rightarrow \mathbb{U}$ is continuous at $\text{limit } xs$, and f is defined on all xs , then

$$\begin{aligned} \text{limit} (\text{lift}_{\text{lim}} f xs) &= \text{limit} (\text{return}_{\text{lim}} (f (\text{limit } xs))) \\ &= \text{limit} (\text{return}_{\text{lim}} (\text{limit} (\text{map } f xs))) \\ &= \text{limit} (\text{map } f xs) \end{aligned} \tag{3.25}$$

Thus, $\text{lift}_{\text{lim}} f xs =_{\text{lim}} \text{map } f xs$, so $\text{lift}'_{\text{lim}} f xs := \text{map } f xs$. Using $\text{lift}_{\text{lim}} f xs$ instead of $\text{lift}'_{\text{lim}} f xs$ requires f to be continuous at $\text{limit } xs$ and defined on all xs .

Limit-Free Join. Using $\text{join } m = \text{bind } m \lambda x. x$ results in $\text{join}_{\text{lim}} = \text{limit}$. Removing `limit` might seem hopeless—until we distribute it pointwise over xss .

Lemma 3.18 (limits of sequences). *Let $f \in \omega \rightarrow \omega \rightarrow A$, where $\omega \rightarrow A$ has a product topology. Then $\text{limit } f = \lambda n \in \omega. \text{limit } (\text{flip } f \ n)$, where $\text{flip } f \ x \ y := f \ y \ x$.*

A countable product metric defines a product topology, so $\text{join}_{\text{lim}} \text{ xss} := \lambda n \in \omega. \text{limit } (\text{flip } \text{xss } n)$. Now we can remove limit by restricting join_{lim} 's input.

Definition 3.19 (uniform convergence). *A sequence $f \in \omega \rightarrow \omega \rightarrow \mathbb{U}$ converges **uniformly** if $\forall \varepsilon \in \mathbb{R}^+. \exists N \in \omega. \forall n, m > N. (\delta (f \ n \ m) (\text{limit } (f \ n))) < \varepsilon$.*

Lemma 3.20 (collapsing limits). *If $f \in \omega \rightarrow \omega \rightarrow \mathbb{U}$ converges uniformly, and $r, s : \omega \Rightarrow \omega$ increase, then $\text{limit } \lambda n \in \omega. \text{limit } (f \ n) = \text{limit } \lambda n \in \omega. f \ (r \ n) \ (s \ n)$.*

So if $\text{flip } \text{xss}$ converges uniformly, then

$$\begin{aligned} \text{limit } (\text{join}_{\text{lim}} \text{ xss}) &= \text{limit } \lambda n \in \omega. \text{limit } (\text{flip } \text{xss } n) \\ &= \text{limit } \lambda n \in \omega. \text{flip } \text{xss } (r \ n) \ (s \ n) \end{aligned} \tag{3.26}$$

We define $\text{join}'_{\text{lim}} : (\omega \rightarrow \omega \rightarrow \mathbb{U}) \Rightarrow (\omega \rightarrow \mathbb{U})$ by $\text{join}'_{\text{lim}} \text{ xss} := \lambda n \in \omega. \text{xss } n \ n$. Replacing $\text{join}_{\text{lim}} \text{ xss}$ with $\text{join}'_{\text{lim}} \text{ xss}$ requires that $\text{flip } \text{xss}$ converge uniformly.

Limit-Free Bind. Define $\text{bind}'_{\text{lim}} \text{ xs } f := \text{join}'_{\text{lim}} (\text{lift}'_{\text{lim}} f \ \text{xs})$. It inherits obligations to prove that f is continuous at $\text{limit } \text{xs}$ and defined on all xs , and to prove that $\text{flip } (\text{map } f \ \text{xs})$ converges uniformly.

Example 3.21 (exponential cont.). Define exp'_{lim} by replacing bind_{lim} by $\text{bind}'_{\text{lim}}$ in exp_{lim} , so $\text{exp}'_{\text{lim}} \text{ xs} := \text{bind}'_{\text{lim}} \text{ xs } \text{exp-seq}$. We now meet the proof obligations.

Lemma 3.22. *Let $f : A \Rightarrow (\omega \rightarrow B)$. If $\omega \rightarrow B$ has a product topology, then f is continuous if and only if $(\text{flip } f) \ n$ is continuous for every $n \in \omega$.*

We have a product topology, so for the first obligation, pointwise continuity is enough. Let $g := \text{flip } \text{exp-seq}$. Every $g \ n$ is a finite polynomial, and thus continuous. The second obligation, that exp-seq is defined on all xs , is obvious. The third, that $\text{flip } (\text{map } \text{exp-seq } \text{xs})$ converges uniformly, can be proved using the Weierstrass M test [2, Theorem 6.4.5]. \diamond

Example 3.23 (π). The definition of $\text{arctan}_{\text{lim}}$ is like exp_{lim} 's. Defining $\text{arctan}'_{\text{lim}}$, including proving correctness, is like defining exp'_{lim} . To compute π , we use

$$\pi_{\text{lim}} := ((\text{return}_{\text{lim}} 16) \times_{\text{lim}} (\text{arctan}_{\text{lim}} (\text{return}_{\text{lim}} \frac{1}{5}))) -_{\text{lim}} ((\text{return}_{\text{lim}} 4) \times_{\text{lim}} (\text{arctan}_{\text{lim}} (\text{return}_{\text{lim}} \frac{1}{239}))) \quad (3.27)$$

where $(\cdot)_{\text{lim}}$ are lifted arithmetic operators. Because (3.27) does not directly use bind_{lim} , defining the limit-free π'_{lim} imposes no proof obligations. \diamond

In general, using functions defined in terms of $\text{bind}'_{\text{lim}}$ requires little more work than using functions on finite values. The implicit limits are pulled outward and collapse on their own, hidden within monadic computations.

Computable Sequences. Lambdas are the simplest model of $\omega \rightarrow \mathbb{U}$. After manipulating some terms, we define the final, computable limit monad by $\text{return}'_{\text{lim}} x := \lambda n. x$ and $\text{bind}'_{\text{lim}} xs f := \lambda n. f (xs n) n$. Computations have type $\omega \Rightarrow \mathbb{U}'$, where \mathbb{U}' contains countable sequences of rationals.

Implementation. We have transliterated $\text{return}'_{\text{lim}}$, $\text{bind}'_{\text{lim}}$, exp'_{lim} , $\text{arctan}'_{\text{lim}}$ and π'_{lim} into Racket [16], using its built-in models of ω and \mathbb{Q} . Even without optimizations, π'_{lim} 141 yields a rational approximation in a few milliseconds that is correct to 200 digits. More importantly, exp'_{lim} , $\text{arctan}'_{\text{lim}}$ and π'_{lim} are almost identical to their counterparts in the uncomputable limit monad, and meet their proof obligations. The code is clean, short, correct and reasonably fast, and resides in a directory named `flops2012` at <https://github.com/ntoronto/plt-stuff/>.

3.9 Related Work

O'Connor's completion monad [39] is quite similar to the limit monad. Both operate on general metric spaces and compute to arbitrary precision. O'Connor starts with computable

approximations and completes them using a monad. Implementing it in Coq took five months. It is certainly correct.

We start with a monad for exact values and define a computable replacement. It was two weeks from conception to implementation. Between directly using well-known theorems, and deriving the computable monad from the uncomputable monad without switching languages, we are as certain as we can be without mechanically verifying it. We have found our middle ground.

Higher-order logics such as HOL [30], CIC [8], MT [7] (Map Theory) and EFL* [15] continue Church’s programme to found mathematics on the λ -calculus. Like λ_{ZFC} , interpreting them in set theory seems to require a slightly stronger theory than plain ZFC. HOL and CIC ensure consistency using types, and use the Curry-Howard correspondence to extract programs.

MT and EFL* are more like λ_{ZFC} in that they are untyped. MT ensures consistency partly by making nontermination a truth value, and EFL* partly by tagging propositions. Both support classical reasoning. MT and EFL* are interpreted in set theory using a straightforward extension of Scott-style denotational semantics to κ -sized domains, while λ_{ZFC} is interpreted in set theory using a straightforward extension of operational semantics to κ -sized relations.

The key difference between λ_{ZFC} and these higher-order logics is that λ_{ZFC} is not a logic. It is a programming language with infinite terms, which by design includes a transitive model of set theory (Theorem 3.12). Therefore, ZFC theorems can be applied to its set-valued terms with only trivial interpretation, whereas the interpretation it takes to apply ZFC theorems to lambda terms that represent sets in MT or EFL* can be nontrivial. Applying a ZFC theorem in HOL or CIC requires re-proving it to the satisfaction of a type checker.

The infinitary λ -calculus [25] has “infinitely deep” terms. Although it exists for investigating laziness, cyclic data, and undefinedness in finitary languages, it is possible to encode uncomputable mathematics in it. In λ_{ZFC} , such up-front encodings are unnecessary.

Hypercomputation [40] describes many Turing machine extensions, including completion of transfinite computations. Much of the research is devoted to discovering the properties of computation in physically plausible extensions. While λ_{ZFC} might offer a civilized way to program such machines, we do not think of our work as hypercomputation, but as approaching computability from above.

3.10 Conclusions

We defined λ_{ZFC} , which can express essentially anything constructible in contemporary mathematics, in a way that makes it compatible with existing first-order theorems. We demonstrated that it makes deriving computational meaning easier by defining the limit monad in it, deriving a computable replacement, and computing real numbers to arbitrary accuracy with acceptable speed.

Now that we have a suitably expressive target language for exact and approximating categorical semantics, we can get back to defining languages for Bayesian modeling and inference. But more generally, we no longer have to hold back when a set-theoretic construction could be defined elegantly with untyped lambdas or recursion, or generalized precisely with higher-order functions. If we can derive a computable replacement, we might help someone in Cantor’s Paradise compute the apparently uncomputable.

Chapter 4

Using λ_{ZFC}

The previous chapter defined λ_{ZFC} , an untyped, call-by-value, operational λ -calculus. It is designed for deriving implementable programs from programs that carry out infinite computations. We will mostly use it as a target language for categorical semantics.

There are two reasonably accurate, short characterizations of λ_{ZFC} . First, it can be regarded as contemporary mathematics (Zermelo-Fraenkel set theory with the axiom of Choice, or ZFC) with well-defined lambdas. Second, it can be regarded as a pure functional programming language with infinite sets. The previous chapter defines λ_{ZFC} in a way that makes these characterizations absolutely precise.

Fortunately, understanding and writing λ_{ZFC} code, and knowing how to prove λ_{ZFC} code correct, requires much less detail. We review the important details here.

4.1 Computations and Values

In λ_{ZFC} , essentially every set is a value, as well as every lambda and every set of lambdas. For example, these are all λ_{ZFC} values:

$$\begin{aligned} &\{1, 2, 3\} \\ &\{(\lambda a. a), (\lambda b. b + 1), (\lambda c. \{c, c + 1\})\} \\ &\mathbb{N}, \mathbb{Q}, \mathbb{R}, \mathbb{R} \times \mathbb{R}, \mathbb{R}^{\mathbb{N}}, \mathbb{R}^{\mathbb{R}} \end{aligned} \tag{4.1}$$

(We generally write λ_{ZFC} terms in **sans serif**.) All primitive operations on values, including operations on infinite sets, are assumed to complete instantly if they terminate.

Nonterminating λ_{ZFC} programs are similar to nonterminating programs in any other call-by-value λ -calculus. For example, a function that does not terminate on any input because of runaway recursion (i.e. an infinite loop) is

$$\text{count-from } n := \text{count-from } (n + 1) \quad (4.2)$$

We could say that `count-from 0` does not terminate because it attempts “infinitely deep” computation. This limitation is necessary; for example, it prevents λ_{ZFC} from having a program that solves its own halting problem, which would make it inconsistent.

Terminating, infinite computations are “infinitely wide,” as in either of these equivalent expressions:

$$\text{image } (\lambda n. n + 1) \mathbb{N} \quad \{n + 1 \mid n \in \mathbb{N}\} \quad (4.3)$$

Both yield the set of all positive natural numbers. It is usually fine to think of terminating, infinite computations as being run in parallel.

As in ZFC, in λ_{ZFC} , all algebraic data structures are encoded as sets; e.g. the pair $\langle x, y \rangle$ can be encoded as $\{\{x\}, \{x, y\}\}$. Only the *existence* of encodings into sets is important, as it means data structures inherit a defining characteristic of sets: strictness. More precisely, in every data structure, every path between the root and a leaf must have finite length. Less precisely, as with computations, values may be “infinitely wide,” such as \mathbb{R} , but not “infinitely deep,” such as infinite trees and infinite lists made from nested pairs.

4.2 Auxiliary Type Systems

Though λ_{ZFC} is untyped, it often helps to define an auxiliary type system. When we use a type system, we use a manually checked, polymorphic one characterized by these rules:

- A free type variable is universally quantified; if uppercase, it denotes a set.
- A set denotes a member of that set.
- $x \Rightarrow y$ denotes a partial function.

- $\langle x, y \rangle$ denotes a pair of values with types x and y .
- $\text{Set } x$ denotes a set with members of type x .

Because the type $\text{Set } X$ denotes the same values as the set $\mathcal{P} X$ (i.e. subsets of the set X) we regard them as equivalent. Similarly, $\langle X, Y \rangle$ is equivalent to $X \times Y$.

All function arrows are right-associative. Recall that in a λ -calculus, application is left-associative. This duality makes writing multi-argument function types easy. For example, $f : \mathbb{N} \Rightarrow (\mathbb{N} \Rightarrow \mathbb{N})$ denotes that function f returns a $\mathbb{N} \Rightarrow \mathbb{N}$ function. If $m : \mathbb{N}$ and $n : \mathbb{N}$ (equivalently $m \in \mathbb{N}$ and $n \in \mathbb{N}$), it can be applied twice using $(f \ m) \ n$. Alternatively, it can be applied using $f \ m \ n$, and its type may be written $f : \mathbb{N} \Rightarrow \mathbb{N} \Rightarrow \mathbb{N}$.

Other examples of types are those of the λ_{ZFC} primitives powerset $\mathcal{P} : \text{Set } x \Rightarrow \text{Set } (\text{Set } x)$, its left inverse, big union $\bigcup : \text{Set } (\text{Set } x) \Rightarrow \text{Set } x$, and the map-like primitive $\text{image} : (x \Rightarrow y) \Rightarrow \text{Set } x \Rightarrow \text{Set } y$.

It is often helpful to create type aliases. For example, to avoid repeatedly writing “ $\bigcup \{\perp\}$ ” we might define

$$X \rightsquigarrow_{\perp} Y ::= X \Rightarrow Y \cup \{\perp\} \tag{4.4}$$

so that $f : X \rightsquigarrow_{\perp} Y$ and $f : X \Rightarrow Y \cup \{\perp\}$ are equivalent type-level statements.

4.3 Using ZFC Values and Theorems

Almost everything definable in ZFC can be defined by a finite λ_{ZFC} program. The previous chapter, for example, defined the real numbers, arithmetic, and limits. The only ZFC values that cannot are those that *must* be defined **nonconstructively**: by proving existence and uniqueness, without giving a bound (no matter how loose) on the length or cardinality of the value. Mathematicians avoid such nonconstructive definitions, and most would consider that definition of “nonconstructive” too liberal.¹

¹Constructivists would object to allowing the law of excluded middle, which is derivable from λ_{ZFC} ’s if, and almost everyone else would object to allowing choice functions.

Almost all known ZFC theorems apply to λ_{ZFC} 's set values without alteration.² Proofs about λ_{ZFC} 's set values apply directly to ZFC sets.³

We often import well-known ZFC theorems as lemmas; for example:

Lemma 4.1 (set equality is extensional). *For all $A : \text{Set } x$ and $B : \text{Set } x$, $A = B$ if and only if $A \subseteq B$ and $B \subseteq A$.*

Or, $A = B$ if and only if they contain the same members.

4.4 Internal Equality and External Equivalence

Any λ_{ZFC} term e used as a truth statement means “ e reduces to **true**” or “ e evalutes to **true**.” Therefore, the terms $(\lambda a. a) 1$ and 1 are (externally) unequal, but $(\lambda a. a) 1 = 1$.

Because of the way λ_{ZFC} 's lambda terms are defined, lambda equality is alpha equivalence, or equivalence up to renaming identifiers. For example, $(\lambda a. a) = (\lambda b. b)$, but not $(\lambda a. 2) = (\lambda a. 1 + 1)$.

If $e_1 = e_2$, then e_1 and e_2 both terminate, and substituting one for the other in an expression does not change its value. Substitution is also safe if both e_1 and e_2 do not terminate, leading to a coarser notion of equivalence.

Definition 4.2 (observational equivalence). *For terms e_1 and e_2 , $e_1 \equiv e_2$ when $e_1 = e_2$, or both e_1 and e_2 do not terminate.*

It might seem helpful to define basic equivalence even more coarsely, so that we can say $\lambda a. 2$ is equivalent to $\lambda a. 1 + 1$. However, we want internal equality and basic external equivalence to be similar, and we want to be able to extend “ \equiv ” with type-specific rules.

²The only exceptions are theorems that rely critically on the existence of an inaccessible cardinal.

³Assuming the existence of an inaccessible cardinal, which is a modest assumption, as $\text{ZFC} + \kappa$ is a smaller theory than Coq's [6].

4.5 Additional Functions and Syntactic Forms

We use heavily sugared syntax, with automatic currying (including primitive applications, so `image fst` means $\lambda A. \text{image fst } A$), binding forms such as indexed unions $\bigcup_{x \in e_A} e$, destructuring binds as in `swap` $\langle a, b \rangle := \langle b, a \rangle$, and comprehensions like $\{a \in A \mid a \in B\}$. We assume logical operators, bounded quantifiers, and typical set operations are defined. To refer to binary operators as values, we enclose them in parentheses, as in (\in) and (\subseteq) .

A less typical set operation we use is disjoint union:

$$\begin{aligned} (\uplus) : \text{Set } x \Rightarrow \text{Set } x \Rightarrow \text{Set } x \\ A \uplus B := \text{if } (A \cap B = \emptyset) (A \cup B) (\text{take } \emptyset) \end{aligned} \tag{4.5}$$

The primitive `take` : $\text{Set } x \Rightarrow x$ returns the element in a singleton set, and does not terminate when applied to a non-singleton set. Thus, $A \uplus B$ terminates only when A and B are disjoint.

Operator precedence is the same as in ordinary mathematics; e.g. $a + b \cdot c = a + (b \cdot c)$. Application has the highest precedence, so $f \ a + g \ b = (f \ a) + (g \ b)$.

4.6 Extensional Functions

In mathematics, logic, and computer science, there are two general classes of functions:

- **Extensional**: functions whose equality, like that of sets, is determined only by their external properties, and not by how they are defined or constructed.
- **Intensional**: functions whose equality is determined only by their internal properties, or by how they are defined or constructed.

In λ_{ZFC} , lambda equality is decided by comparing body expressions, so lambdas are intensional.

In ZFC and λ_{ZFC} , function extensionality is achieved by encoding functions as sets of input-output pairs. For example, the increment function for the natural numbers is $\{\langle 0, 1 \rangle, \langle 1, 2 \rangle, \langle 2, 3 \rangle, \dots\}$. (It is fine to think of such encodings as infinite hash tables.) We call these encodings **mappings**. We use **function** to mean either a lambda or a mapping, and

use adjacency (e.g. $(f\ a)$ or $f\ a$) to apply either kind.

Syntax for constructing unnamed mappings is defined by

$$\lambda x_a \in e_A. e_b := \text{mapping } (\lambda x_a. e_b) e_A \quad (4.6)$$

$$\begin{aligned} \text{mapping} : (X \Rightarrow Y) &\Rightarrow \text{Set } X \Rightarrow (X \rightarrow Y) \\ \text{mapping } f\ A &:= \text{image } (\lambda a. \langle a, f\ a \rangle) A \end{aligned} \quad (4.7)$$

For symmetry with partial functions $x \Rightarrow y$, **mapping** returns a member of the set $X \rightarrow Y$ of all *partial* mappings from X to Y ; i.e. if $g : X \rightarrow Y$, then g 's domain may be a subset of X . Two common partial mapping operations we use in the next chapter are

$$\begin{aligned} \text{domain} : (X \rightarrow Y) &\Rightarrow \text{Set } X \\ \text{domain } g &:= \text{image fst } g \end{aligned} \quad (4.8)$$

$$\begin{aligned} \text{preimage} : (X \rightarrow Y) &\Rightarrow \text{Set } Y \Rightarrow \text{Set } X \\ \text{preimage } g\ B &:= \{a \in \text{domain } g \mid g\ a \in B\} \end{aligned} \quad (4.9)$$

The **preimage** function finds g 's inputs whose corresponding outputs are in B .

The set $J \rightarrow X$ contains all the *total* mappings from J to X ; equivalently, all the vectors of X indexed by J , which may be infinite. We use infinite vectors of type $J \rightarrow [0, 1]$ in Chapter 7 as infinite sources of uniformly random numbers.

In short, in addition to lambdas in λ_{ZFC} , we have every necessary mathematical object at our disposal.

Chapter 5

Countable Models and Implementation

This chapter is derived from work published at the 22nd *Symposium on Implementation and Application of Functional Languages (IFL)*, 2010.

An approximate answer to the right question is worth a good deal more than the exact answer to an approximate problem.

John W. Tukey

5.1 Introduction

Bayesians write theories without regard to whether future calculations are closed-form or tractable. They are loath to make simplifying assumptions. (If answering questions about some probabilistic process involves an unsolvable integral, so be it.) When they must approximate, they often create two theories: an “ideal” theory first, and a second that approximates it.

Because they create theories without regard to future calculations, they usually must accept approximate answers to queries about them. Typically, they adapt algorithms that compute converging approximations in programming languages they are familiar with. The process is tedious and error-prone, and involves much performance tuning and manual optimization. It is by far the most time-consuming part of their work—and also the most automatable part.

They follow this process to adhere to an overriding philosophy: an approximate answer to the right question is worth more than an exact answer to an approximate question. Thus, they put off approximating as long as possible.

We must also adhere to this philosophy because Bayesian practitioners are unlikely to use a language that requires them to approximate early, or that approximates earlier than they would. We have found that a good way to put the philosophy into practice in language design is to create two semantics: an “ideal,” or *exact* semantics first, and a converging, *approximating* semantics.

5.1.1 Approach

Measure-theoretic probability is the most successful theory of probability in precision, maturity, and explanatory power. In particular, it is believed to explain every Bayesian theory. We therefore define the exact semantics as a transformation from Bayesian notation to measure-theoretic calculations.

Measure theory treats finite, countably infinite, and uncountably infinite probabilistic outcomes uniformly, but with significant complexity. Though there are relatively few important Bayesian models that require countably many outcomes but not uncountably many, in our preliminary work, we deal with only countable sets. This choice avoids most of measure theory’s complexity while retaining its functional structure, and still requires approximation.

For three syntactic categories of Bayesian notation, we

1. Manually interpret an unambiguous subclass of common notation.
2. Mechanize the interpretation with a semantic function.
3. If necessary, create an approximation and prove convergence.
4. Implement the approximation in Racket [16].

This approach is most effective if the target language can express measure-theoretic calculations and is similar to Racket in structure and semantics. We therefore use λ_{ZFC} .

The Bayesian notation we interpret falls into these syntactic categories:

- **Expressions**, which have no side effects, interpreted by $\mathcal{R}[\![\cdot]\!]$.
- **Statements**, which create side effects, interpreted by $\mathcal{M}[\![\cdot]\!]$.
- **Queries**, which observe side effects, interpreted by $\mathbf{P}[\![\cdot]\!]$ and $\mathbf{D}[\![\cdot]\!]$.

We write Bayesian notation in *italics*, Racket in **fixed width**, common keywords in **bold** and invented keywords in ***bold italics***. We omit proofs for space.

5.2 The Expression Language

5.2.1 Background Theory: Random Variables

Most practitioners of probability understand random variables as free variables whose values have ambient probabilities. But measure-theoretic probability defines a **random variable** X as a total mapping

$$X : \Omega \rightarrow S_X \tag{5.1}$$

where Ω and S_X are sets called **sample spaces**, with elements called **outcomes**. Random variables define and limit what is observable about any outcome $\omega \in \Omega$, so we call outcomes in S_X ***observable outcomes***.

Example 5.1. Suppose we want to encode, as a random variable E , the act of observing whether the outcome of a die roll is even or odd.

A complicated way is to define Ω as the possible states of the universe. $E : \Omega \rightarrow \{\text{even}, \text{odd}\}$ must simulate the universe until the die is still, and then recognize the outcome. Hopefully, the probability that $E \omega = \text{even}$ is close to $\frac{1}{2}$.

A tractable way defines $\Omega := \{1, 2, 3, 4, 5, 6\}$ and $E : \Omega \rightarrow \{\text{even}, \text{odd}\}$ so that $E \omega = \text{even}$ if $\omega \in \{2, 4, 6\}$, otherwise **odd**. The probability that $E \omega = \text{even}$ is the sum of probabilities of every even $\omega \in \Omega$, or $\frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{1}{2}$.

If we are interested in observing only evenness, we can define $\Omega := \{\text{even}, \text{odd}\}$, each with probability $\frac{1}{2}$, and $E \omega := \omega$. ◇

Random variables enable a kind of probabilistic abstraction. The example does it twice. The first makes calculating the probability that $\mathbf{E} \omega = \mathbf{even}$ tractable. The second is an optimization. In fact, redefining Ω , the random variables, and the probabilities of outcomes—without changing the probabilities of *observable* outcomes—is the essence of measure-theoretic optimization.

Defining random variables as functions is also a good factorization: it separates nondeterminism from assigning probabilities. It allows us to interpret expressions involving random variables without considering probabilities at all.

5.2.2 Interpreting Random Variable Expressions As Computations

When random variables are regarded as free variables, arithmetic with random variables is no different from deterministic arithmetic. Measure-theoretic probability uses the same notation, but regards it as implicit pointwise lifting (as in vector arithmetic). For example, if $A, B, C : \Omega \rightarrow \mathbb{R}$ are random variables, $C := A + B$ means $C \omega := (A \omega) + (B \omega)$, and $B := 4 + A$ means $B \omega := 4 + (A \omega)$.

Because we use λ_{ZFC} , we can extend the class of random variables from $\Omega \rightarrow S_X$ to $\Omega \Rightarrow S_X$. Including lambdas as well as mappings makes it easy to interpret unnamed random variables: $4 + A$, or in prefix form $((+) 4 A)$, means $\lambda\omega. ((+) 4 (A \omega))$. Lifting constants allows us to interpret expressions uniformly: if we interpret $(+)$ as $\text{Plus} := \lambda\omega. (+)$ and 4 as $\text{Four} := \lambda\omega. 4$, then $((+) 4 A)$ means

$$\lambda\omega. ((\text{Plus } \omega) (\text{Four } \omega) (A \omega)) \quad (5.2)$$

We abstract lifting and application with these combinators:

$$\begin{aligned} \text{pure}_r c &:= \lambda\omega. c \\ \text{ap}_r^* F X_1 \dots X_n &:= \lambda\omega. ((F \omega) (X_1 \omega) \dots (X_n \omega)) \end{aligned} \quad (5.3)$$

$$\begin{aligned}
\mathcal{R}[\![X]\!] &::= X \\
\mathcal{R}[\![x]\!] &::= \text{pure}_{\text{rv}} x \\
\mathcal{R}[\![v]\!] &::= \text{pure}_{\text{rv}} v \\
\mathcal{R}[\![e_f \ e_1 \ \dots \ e_n]\!] &::= \text{ap}_{\text{rv}}^* \ \mathcal{R}[\![e_f]\!] \ \mathcal{R}[\![e_1]\!] \ \dots \ \mathcal{R}[\![e_n]\!] \\
\mathcal{R}[\![\lambda x_1 \dots x_n. e]\!] &::= \lambda \omega. \lambda x_1 \dots x_n. (\mathcal{R}[\![e]\!] \ \omega) \\
\text{pure}_{\text{rv}} c &::= \lambda \omega. c \\
\text{ap}_{\text{rv}}^* F \ X_1 \ \dots \ X_n &::= \lambda \omega. ((F \ \omega) (X_1 \ \omega) \ \dots \ (X_n \ \omega))
\end{aligned}$$

Figure 5.1: Random variable expression semantics. The source and target language are both λ_{ZFC} . Conditionals and primitive operators are trivial special cases of application.

In terms of pure_{rv} and ap_{rv}^* , $4 + A$ means

$$\begin{aligned}
\text{ap}_{\text{rv}}^* (\text{pure}_{\text{rv}} (+)) (\text{pure}_{\text{rv}} 4) A &\equiv \text{ap}_{\text{rv}}^* (\lambda \omega. (+)) (\lambda \omega. 4) A \\
&\equiv \lambda \omega. ((\lambda \omega. (+)) \ \omega) ((\lambda \omega. 4) \ \omega) (A \ \omega) \\
&\equiv \lambda \omega. (+) 4 (A \ \omega) \\
&= \lambda \omega. 4 + (A \ \omega)
\end{aligned} \tag{5.4}$$

as desired. These combinators define an **idiom** [35], which is like a monad but can impose a partial order on computations. The *random variable idiom* instantiates the environment idiom with the type constructor $\text{I } \mathbf{a} ::= \Omega \Rightarrow \mathbf{a}$ for some Ω .

$\mathcal{R}[\![\cdot]\!]$ (Figure 5.1), the semantic function that interprets random variable expressions, targets this idiom. It does mechanically what we have done manually, and additionally interprets lambdas. For simplicity, it follows probability convention by assuming single uppercase letters are random variables. Figure 5.1 assumes syntactic sugar has been replaced; e.g. that application is in prefix form.

$\mathcal{R}[\![\cdot]\!]$ may return lambdas that do not terminate when applied to an ω . For now, we assume they terminate for all $\omega \in \Omega$. (Chapter 7 deals with nonterminating programs.)

We will be able to recover mappings using the **mapping** function, which, given a domain, converts a lambda or mapping to a mapping, as in `mapping $\mathcal{R}[\![4 + A]\!] \ \Omega$` .

```

(define-syntax (RV/kernel stx)
  (syntax-parse stx
    [(_ Xs:ids e:expr)
     (syntax-parse #'e #:literal-sets (kernel-literals)
       [X:id #:when (free-id-in? #'Xs #'X) #'X]
       [x:id #'(pure x)]
       [(quote c) #'(pure (quote c))]
       [(%#plain-app e ...) #'(ap* (RV/kernel Xs e) ...)]
       ....))]))

(define-syntax (RV stx)
  (syntax-parse stx
    [(_ Xs:ids e:expr)
     #'(RV/kernel Xs #,(local-expand #'e 'expression empty))]))

```

Figure 5.2: A fragment of our implementation of $\mathcal{R}[\![\cdot]\!]$ in Racket.

5.2.3 Implementation in Racket

Figure 5.2 shows `RV` and a snippet of `RV/kernel`, the macros that implement $\mathcal{R}[\![\cdot]\!]$. `RV` fully expands expressions into Racket’s kernel language, allowing `RV/kernel` to transform any pure Racket expression into a random variable. Both use Racket’s new `syntax-parse` library [14]. `RV/kernel` raises a syntax error on `set!`, but there is no way to disallow applying functions that have effects.

Rather than differentiate between kinds of identifiers, `RV` takes a list of known random variable identifiers as an additional argument. It wraps other identifiers with `pure`, allowing arbitrary Racket values to be random variables.

5.3 The Query Language

It is best to regard statements in Bayesian theories as specifications for the results of later observations. We therefore interpret queries before interpreting statements. First, however, we must define the state objects that queries observe.

5.3.1 Background Theory: Probability Spaces

In practice, functions called **distributions** assign probabilities or probability densities to observable outcomes. Practitioners state distributions for certain random variables, and then calculate the distributions of others.

Measure-theoretic probability generalizes assigning probabilities and densities using **probability measures**, which assign probabilities to *sets* of outcomes. There are typically no special random variables: all random variable distributions are calculated from one global probability measure.

It is generally not possible to assign meaningful probabilities to all subsets of a sample space Ω —except when Ω is countable. We thus deal here with **discrete probability measures** $P : \text{Set } \Omega \rightarrow [0, 1]$, where Ω is countable. Any discrete probability measure is uniquely determined by its value on singleton sets, or by a **probability mass function** $p : \Omega \rightarrow [0, 1]$. It is easy to convert p to a probability measure:

$$\text{sum } p \text{ } A := \sum_{\omega \in A} p \ \omega \quad (5.5)$$

Then $P = \text{sum } p$. Converting the other direction is also easy: $p \ e = P \ \{e\}$.

A **discrete probability space** $\langle \Omega, p \rangle$ embodies all probabilistic nondeterminism introduced by theory statements. It is fine to think of Ω as the set of all possible states of a write-once memory, with p assigning a probability to each state.

5.3.2 Background Theory: Queries

Any probability can be calculated from $\langle \Omega, p \rangle$. For example, suppose we want to calculate, as in Example 5.1, the probability of an even die outcome. We must apply P to the correct subset of Ω . Suppose $\Omega := \{1, 2, 3, 4, 5, 6\}$ and that $p := [1, 2, 3, 4, 5, 6 \rightarrow \frac{1}{6}]$ determines P . The probability that E outputs **even** is

$$P \ \{\omega \in \Omega \mid E \ \omega = \text{even}\} = P \ \{2, 4, 6\} = \text{sum } p \ \{2, 4, 6\} = \frac{1}{2} \quad (5.6)$$

This is a **probability query**.

Alternatively, we could use a **distribution query** to calculate E 's distribution P_E , and then apply it to $\{\text{even}\}$. Measure-theoretic probability elegantly defines P_E as $P \circ (\text{preimage } E)$, but for now we do not need a measure. We only need the probability mass function $p_E : \{\text{even}, \text{odd}\} \rightarrow [0, 1]$, defined by $p_E \ e = \text{sum } p \ (\text{preimage } E \ \{e\})$. Applying it to even yields

$$p_E \ \text{even} = \text{sum } p \ (\text{preimage } E \ \{\text{even}\}) = \text{sum } p \ \{2, 4, 6\} = \frac{1}{2} \quad (5.7)$$

More abstractly, we can calculate discrete distribution queries using

$$\begin{aligned} \text{dist } X \ \langle \Omega, p \rangle &:= \text{let } S_X := \text{image } X \ \Omega \\ &\text{in } \lambda x \in S_X. \text{sum } p \ (\text{preimage } (\text{mapping } X \ \Omega) \ \{x\}) \end{aligned} \quad (5.8)$$

Now $p_E = \text{dist } E \ \langle \Omega, p \rangle$. Recall that the special syntax $\lambda x \in e_A. e$ creates an unnamed mapping with domain e_A , and **mapping** $X \ \Omega$ converts X , which may be a lambda, to a mapping with domain Ω , on which preimages are well-defined.

5.3.3 Interpreting Query Notation

When random variables are regarded as free variables, special notation $\text{Pr}[\cdot]$ replaces applying the probability measure P and sets become propositions. For example, a common way to write “the probability of an even die outcome” in practice is $\text{Pr}[E = \text{even}]$.

The semantic function $\mathcal{R}[\![\cdot]\!]$ turns propositions about random variables into predicates on Ω . The set corresponding to the proposition is the preimage of $\{\text{true}\}$. For the proposition $E = \text{even}$, for example, it is $\text{preimage } (\text{mapping } \mathcal{R}[\![E = \text{even}]\!] \ \Omega) \ \{\text{true}\}$. In general,

$$\text{sum } p \ (\text{preimage } (\text{mapping } \mathcal{R}[\![e]\!] \ \Omega) \ \{\text{true}\}) = \text{dist } \mathcal{R}[\![e]\!] \ \langle \Omega, p \rangle \ \text{true} \quad (5.9)$$

calculates $\text{Pr}[e]$ when e is a proposition; i.e. when $\mathcal{R}[\![e]\!] : \Omega \Rightarrow \{\text{true}, \text{false}\}$.

Although probability queries have common notation, there seems to be no common notation that denotes distributions *per se*. The typical workarounds are to write implicit formulas like $\text{Pr}[E = e]$ and to give distributions suggestive names like p_E . Some theorists

use $\mathcal{L}[\cdot]$, with \mathcal{L} for *law*, an obscure synonym of *distribution*. We define $\mathbf{D}[\![\cdot]\!]$ in place of $\mathcal{L}[\cdot]$. Then $\mathbf{D}[\![E]\!]$ denotes E 's distribution.

Though we could define semantic functions $\mathbf{P}[\![\cdot]\!]$ and $\mathbf{D}[\![\cdot]\!]$ right now, we are putting them off until after interpreting statements.

5.3.4 Approximating Queries

Probabilities are real numbers. They remain real in the approximating semantics; we use floating-point approximation and exact rationals in the implementation.

Arbitrary countable sets are not finitely representable. In the approximating semantics, we restrict Ω to recursively enumerable sets. The implementation encodes them as lazy lists. We trust users to not create “sets” with duplicates.

When A is infinite, $\text{sum } p \ A$ is an infinite series. With A represented by a lazy list, it is easy to compute a converging approximation—but then approximate answers to distribution queries sum to values less than 1. Instead, we approximate Ω and normalize p , which makes the sum finite and the distributions proper.

Suppose $\langle \omega_1, \omega_2, \dots \rangle$ is an enumeration of Ω . Let $z \in \mathbb{N}^+$ be the length of the prefix $\Omega_z := \{\omega_1, \dots, \omega_z\}$ and define $p_z : \Omega \rightarrow [0, 1]$ by $p_z \ \omega = (p \ \omega) / (\text{sum } p \ \Omega_z)$ if $\omega \in \Omega_z$; otherwise 0. Then p_z converges to p pointwise. We define $\text{finitize } \langle \Omega, p \rangle := \langle \Omega_z, p_z \rangle$ with $z \in \mathbb{N}$ as a free variable.

5.3.5 Implementation in Racket

Figure 5.3 shows the implementations of `finitize` and `dist` in Racket. The free variable z appears as a *parameter* `appx-z`: a variable with static scope but dynamic extent. The `cotake` procedure returns the prefix of a lazy list as a finite list.

To implement `dist`, we need to represent mappings in Racket. The applicable struct type `mapping` represents lazy mappings with possibly infinite domains. A `mapping` named `f` can be applied with `(f x)`. We do not ensure `x` is in the domain because checking is

```

(struct mapping (domain proc)
  #:property prop:procedure (λ (f x) ((mapping-proc f) x)))

(struct fmapping (default hash)
  #:property prop:procedure
  (λ (f x) (hash-ref (fmapping-hash f) x (fmapping-default f))))

(define appx-z (make-parameter +inf.0))
(define (finitize ps)
  (match-let* ([ (mapping Ω P) ps]
               [Ωn (cotake Ω (appx-z))]
               [qn (apply + (map P Ωn))])
    (mapping Ωn (λ (ω) (/ (P ω) qn)))))

(define ((dist X) ps)
  (match-define (mapping Ω P) ps)
  (fmapping 0 (for/fold ([h (hash)]) ([ω (in-list Ω)])
    (hash-set h (X ω) (+ (P ω) (hash-ref h (X ω) 0))))))

```

Figure 5.3: Implementation of finite approximation and distribution queries in Racket.

semidecidable and nontermination is a terrible error message. For distributions, checking is not important; the observable domain is.

However, we do not want `dist` to return lazy mappings. Doing so is inefficient: every application of the mapping would filter Ω . Further, `dist` always receives a `finitized` probability space. We therefore define `fmapping` for mappings that are constant on all but a finite set. For these values, `dist` builds a hash table by computing the probabilities of all preimages in one pass through Ω .

We do use `mapping`, but only for probability spaces and stated distributions.

5.4 Conditional Queries

For Bayesian practitioners, the most meaningful queries are **conditional** queries: those *conditioned on*, or *given*, some random variable’s value. (For example, the probability an email is spam given it contains words like “madam,” or the distribution over suspects given security footage.) A probabilistic language without conditional queries is of little more use to

them than a general-purpose language with a **random** primitive.

Measure-theoretic conditional probability is too involved to accurately summarize here. When \mathbf{P} is discrete, however, the conditional probability of set \mathbf{A} given set \mathbf{B} (i.e. asserting that $\omega \in \mathbf{B}$), simplifies to

$$\Pr[\mathbf{A} \mid \mathbf{B}] = \mathbf{P}(\mathbf{A} \cap \mathbf{B}) / \mathbf{P} \mathbf{B} \quad (5.10)$$

In theory and practice, $\Pr[\cdot \mid \cdot]$ is special notation. As with $\Pr[\cdot]$, practitioners apply $\Pr[\cdot \mid \cdot]$ to propositions, and define it with $\Pr[e_A \mid e_B] := \Pr[e_A \wedge e_B] / \Pr[e_B]$.

Example 5.2. Extend Example 5.1 with random variable $\mathbf{L} : \Omega \rightarrow \{\text{low}, \text{high}\}$ defined by $\mathbf{L} \omega = \text{if } (\omega \leq 3) \text{ low high}$. The probability that $E = \text{even}$ given $L = \text{low}$ is

$$\Pr[E = \text{even} \mid L = \text{low}] = \frac{\Pr[E = \text{even} \wedge L = \text{low}]}{\Pr[L = \text{low}]} = \frac{\sum_{\omega \in \{2\}} \mathbf{P} \omega}{\sum_{\omega \in \{1,2,3\}} \mathbf{P} \omega} = \frac{\frac{1}{6}}{\frac{1}{2}} = \frac{1}{3} \quad (5.11)$$

Similarly, $\Pr[E = \text{odd} \mid L = \text{low}] = \frac{2}{3}$. Less precisely, there are proportionally fewer even outcomes when $L = \text{low}$. \diamond

Conditional *distribution* queries ask how one random variable's output influences the distribution of another. As with unconditional distribution queries, practitioners work around a lack of common notation. For example, they might write the distribution of \mathbf{E} given \mathbf{L} as $\Pr[E = e \mid L = l]$ or $p_{E|L}$.

It is tempting to define $\mathbf{P}[\cdot \mid \cdot]$ in terms of $\mathbf{P}[\cdot]$, and $\mathbf{D}[\cdot \mid \cdot]$ in terms of $\mathbf{D}[\cdot]$. However, defining conditioning as an operation on probability spaces instead of on queries is more flexible. The following abstraction returns a discrete probability space in which Ω is restricted to the subset where random variable \mathbf{Y} returns y :

$$\begin{aligned} \text{cond } \mathbf{Y} \ y \ \langle \Omega, \mathbf{p} \rangle &:= \text{let } \Omega' := \text{preimage } (\text{mapping } \mathbf{Y} \ \Omega) \ \{y\} \\ &\quad \mathbf{p}' := \lambda \omega \in \Omega'. (\mathbf{p} \ \omega) / (\text{sum } \mathbf{p} \ \Omega') \\ &\quad \text{in } \langle \Omega', \mathbf{p}' \rangle \end{aligned} \quad (5.12)$$

Then $\Pr[E = \text{even} \mid L = \text{low}]$ means $\text{dist } \mathbf{E} \ (\text{cond } \mathbf{L} \ \text{low} \ \langle \Omega, \mathbf{p} \rangle) \ \text{even}$.

We approximate **cond** by applying **finitize** to the probability space first. Its implementation uses finite list procedures instead of set operators.

5.5 The Statement Language

Random variables influence each other through global probability spaces. However, because practitioners regard random variables as free variables instead of as functions of a probability space, they state facts about random variable distributions instead of facts about probability spaces. Though they call such collections of statements *models*,¹ to us they are **probabilistic theories**. A **model** is a probability space and random variables that imply the stated facts.

Discrete **conditional theories** can always be written to conform to

$$t_i ::= X_i \sim e_i; \ t_{i+1} \mid X_i := e_i; \ t_{i+1} \mid e_a = e_b; \ t_{i+1} \mid \epsilon \quad (5.13)$$

Further, they can always be made **well-formed**: an e_j may refer to some X_i only when $j > i$ (i.e. no circular bindings). We start by interpreting the most common kind of Bayesian theories, which contain only distribution statements.

5.5.1 Interpreting Common Conditional Theories

Example 5.3. Suppose we want to know only whether a die outcome is even or odd, high or low. If L 's distribution is $p_L := [\text{low}, \text{high} \mapsto \frac{1}{2}]$, then E 's distribution depends on L 's output.

Define $p_{E|L} : S_L \rightarrow S_E \rightarrow [0, 1]$ by $p_{E|L} \text{ low} = [\text{even} \mapsto \frac{1}{3}, \text{odd} \mapsto \frac{2}{3}]$ and $p_{E|L} \text{ high} = [\text{even} \mapsto \frac{2}{3}, \text{odd} \mapsto \frac{1}{3}]$.² The conditional theory could be written

$$L \sim p_L; \ E \sim (p_{E|L} \ L) \quad (5.14)$$

If L is a measure-theoretic random variable, $(p_{E|L} \ L)$ does not typecheck: $L : \Omega \rightarrow S_L$ is clearly not in S_L . The *intent* is that $p_{E|L}$ specifies how E 's distribution depends on L . \diamond

¹In the colloquial sense, probably to emphasize their essential incompleteness.

²Usually, $P_{E|L} : S_E \times S_L \rightarrow [0, 1]$. We reorder and curry to simplify interpretation.

We can regard $L \sim p_L$ as a constraint on models: $\text{dist } L \langle \Omega, \mathbf{p} \rangle$ must be \mathbf{p}_L for every model $\langle \Omega, \mathbf{p}, L \rangle$. Similarly, $E \sim (p_{E|L} L)$ means E 's conditional distribution is $\mathbf{p}_{E|L}$. We have been using the model $\Omega := \{1, 2, 3, 4, 5, 6\}$, $\mathbf{p} := [1, 2, 3, 4, 5, 6 \mapsto \frac{1}{6}]$, and the obvious E and L . It is not hard to verify that this is also a model:

$$\begin{aligned}
\Omega &:= \{\text{low}, \text{high}\} \times \{\text{even}, \text{odd}\} \\
L \langle \omega_1, \omega_2 \rangle &:= \omega_1 \\
E \langle \omega_1, \omega_2 \rangle &:= \omega_2 \\
\mathbf{p} &:= [\langle \text{low}, \text{even} \rangle, \langle \text{high}, \text{odd} \rangle \mapsto \frac{1}{6}, \langle \text{low}, \text{odd} \rangle, \langle \text{high}, \text{even} \rangle \mapsto \frac{2}{6}]
\end{aligned} \tag{5.15}$$

The construction of Ω , L and E in (5.15) clearly generalizes, but \mathbf{p} is trickier. Fully justifying the generalization (including that it meets implicit independence assumptions that we have not mentioned) is rather tedious, so we do not do it here. But, for the present example, it is not hard to check these facts:

$$\begin{aligned}
\mathbf{p} &= \lambda \omega \in \Omega. (\mathbf{p}_L (L \omega)) \cdot (\mathbf{p}_{E|L} (L \omega) (E \omega)) \\
&= \text{mapping } \mathcal{R}[(p_L L) \cdot ((p_{E|L} L) E)] \Omega
\end{aligned} \tag{5.16}$$

Let $K_L := \mathcal{R}[p_L]$ and $K_E := \mathcal{R}[p_{E|L} L]$, which interpret (5.14)'s statements' right-hand sides. Then $\mathbf{p} = \text{mapping } \mathcal{R}[(K_L L) \cdot (K_E E)] \Omega$. This can be generalized.

Definition 5.4 (discrete product model). *Given a well-formed, discrete conditional theory $X_1 \sim e_1; \dots; X_n \sim e_n$, let $K_i : \Omega \Rightarrow S_i \rightarrow [0, 1]$, defined by $K_i = \mathcal{R}[e_i]$ for each $1 \leq i \leq n$. The **discrete product model** of the theory is*

$$\begin{aligned}
\Omega &:= S_1 \times \dots \times S_n \\
X_i \langle \omega_1, \dots, \omega_i, \dots, \omega_n \rangle &:= \omega_i \quad (1 \leq i \leq n) \\
\mathbf{p} &:= \text{mapping } \mathcal{R}[(K_1 X_1) \cdot \dots \cdot (K_n X_n)] \Omega
\end{aligned} \tag{5.17}$$

Theorem 5.5 (semantic intent). *The discrete product model induces the stated conditional distributions and meets implicit independence assumptions.*

When writing distribution statements, practitioners tend to apply first-order distributions to simple random variables. But the discrete product model allows any λ_{ZFC} term e_i whose interpretation is a discrete **transition kernel** $\mathcal{R}[\![e_i]\!] : \Omega \Rightarrow S_i \rightarrow [0, 1]$. In measure theory, transition kernels are used to build **product spaces** such as $\langle \Omega, \mathbf{p} \rangle$. Thus, $\mathcal{R}[\![\cdot]\!]$ links Bayesian practice to measure theory and represents an increase in expressive power in specifying distributions, by turning properly typed λ_{ZFC} terms into precisely what measure theory requires.

5.5.2 Interpreting Statements as Monadic Computations

Some conditional theories state more than just distributions [34, 50]. Interpreting theories with different kinds of statements requires recursive, rather than whole-theory, interpretation. Fortunately, well-formedness amounts to lexical scope, making it straightforward to interpret statements as monadic computations.

We use the state monad with probability-space-valued state: computations are functions from probability spaces to probability spaces paired with a statement-specific value. The probability space monad's **return** and **bind** are defined as

$$\begin{aligned} \text{return}_{\text{ps}} \times \langle \Omega, \mathbf{p} \rangle &:= \langle \Omega, \mathbf{p}, \mathbf{x} \rangle \\ \text{bind}_{\text{ps}} \, m \, f \, \langle \Omega, \mathbf{p} \rangle &:= \text{let } \langle \Omega', \mathbf{p}', \mathbf{x} \rangle := m \, \langle \Omega, \mathbf{p} \rangle \\ &\quad \text{in } f \times \langle \Omega', \mathbf{p}' \rangle \end{aligned} \tag{5.18}$$

Figure 5.4 shows the additional dist_{ps} , cond_{ps} and $\text{extend}_{\text{ps}}$. The first two simply reimplement **dist** and **cond**. But $\text{extend}_{\text{ps}}$, which interprets statements, needs more explanation.

According to (5.17), interpreting $X_i \sim e_i$ results in $\Omega_i = \Omega_{i-1} \times S_i$, with S_i extracted from $K_i : \Omega_{i-1} \Rightarrow S_i \rightarrow [0, 1]$. A more precise type for K_i is the dependent type $(\omega : \Omega_{i-1}) \Rightarrow (S'_i \, \omega) \rightarrow [0, 1]$, which reveals a complication. To extract S_i , we first must extract the random variable $S'_i : \Omega_{i-1} \rightarrow \text{Set } S_i$. So let $S'_i \, \omega = \text{domain } (K_i \, \omega)$; then $S_i = \bigcup (\text{image } S'_i \, \Omega_{i-1})$.

But this makes query implementation inefficient: if the union has little overlap or is disjoint, \mathbf{p} will assign 0 to most ω . In more general terms, we actually have a *dependent*

$$\begin{aligned}
\text{dist}_{\text{ps}} X \langle \Omega, p \rangle &:= \text{let } S_X := \text{image } X \Omega \\
&\quad p_X := \lambda x \in S_X. \text{sum } p \text{ (preimage (mapping } X \Omega) \{x\})} \\
&\text{in } \langle \Omega, p, p_X \rangle \\
\text{cond}_{\text{ps}} Y y \langle \Omega, p \rangle &:= \text{let } \Omega' := \text{preimage (mapping } Y \Omega) \{y\} \\
&\quad p' := \lambda \omega \in \Omega'. (p \ \omega) / (\text{sum } p \ \Omega') \\
&\text{in } \langle \Omega', p', _ \rangle \\
\text{extend}_{\text{ps}} K \langle \Omega, p \rangle &:= \text{let } S' \ \omega := \text{domain } (K \ \omega) \\
&\quad \Omega' := (\omega \in \Omega) \times (S' \ \omega) \\
&\quad X \ \omega := \omega_j \quad (\text{where } j \text{ is the length of any } \omega \in \Omega) \\
&\quad p' := \text{mapping } \mathcal{R}[\![p \cdot (K \ X)]\!] \ \Omega' \\
&\text{in } \langle \Omega', p', X \rangle \\
\text{empty}_{\text{ps}} &:= \langle \{\langle \rangle\}, \lambda \omega \in \{\langle \rangle\}. 1 \rangle \\
\text{run}_{\text{ps}} m &:= \text{let } \langle \Omega, p, x \rangle := m \text{ empty}_{\text{ps}} \\
&\quad \text{in } x
\end{aligned}$$

Figure 5.4: State monad functions that represent queries and statements. The state is probability-space-valued.

cartesian product $(\omega \in \Omega_{i-1}) \times (S'_i \ \omega)$, a generalization of the cartesian product.³ To extend Ω , $\text{extend}_{\text{ps}}$ calculates this product instead.

Dependent cartesian products are elegantly expressed using the set monad:

$$\begin{aligned}
\text{return}_{\text{set}} a &:= \{a\} \\
\text{bind}_{\text{set}} A f &:= \bigcup (\text{image } f \ A)
\end{aligned} \tag{5.19}$$

Then $(a \in A) \times (B \ a) = \text{bind}_{\text{set}} A \ \lambda a. \text{bind}_{\text{set}} (B \ a) \ \lambda b. \text{return}_{\text{set}} \langle a, b \rangle$.

Figure 5.5 defines $\mathcal{M}[\![\cdot]\!]$, which interprets conditional theories containing definition, distribution, and conditioning statements as probability space monad computations. After it exhausts the statements, it returns the random variables. Returning their names as well would be an obfuscating complication, which we avoid by implicitly extracting them from the theory

³The dependent cartesian product also generalizes disjoint union to arbitrary index sets. It is often called a *dependent sum* and denoted $\Sigma a : A. (B \ a)$.

$$\begin{aligned}
\mathcal{M}[[X_i := e_i; t_{i+1}]] &::= \text{bind}_{\text{ps}} (\text{return}_{\text{ps}} \mathcal{R}[[e_i]]) \lambda X_i. \mathcal{M}[[t_{i+1}]] \\
\mathcal{M}[[X_i \sim e_i; t_{i+1}]] &::= \text{bind}_{\text{ps}} (\text{extend}_{\text{ps}} \mathcal{R}[[e_i]]) \lambda X_i. \mathcal{M}[[t_{i+1}]] \\
\mathcal{M}[[e_a = e_b; t_{i+1}]] &::= \text{bind}_{\text{ps}} (\text{cond}_{\text{ps}} \mathcal{R}[[e_a]] \mathcal{R}[[e_b]]) \lambda _ . \mathcal{M}[[t_{i+1}]] \\
\mathcal{M}[[\epsilon]] &::= \text{return}_{\text{ps}} \langle X_1, \dots, X_n \rangle \\
\\
\mathbf{D}[[e]] \text{ m} &::= \text{run}_{\text{ps}} (\text{bind}_{\text{ps}} \text{ m} \lambda \langle X_1, \dots, X_n \rangle . \text{dist}_{\text{ps}} \mathcal{R}[[e]]) \\
\mathbf{D}[[e_X | e_Y]] \text{ m} &::= \lambda y. \mathbf{D}[[e_X]] (\text{bind}_{\text{ps}} \text{ m} \lambda \langle X_1, \dots, X_n \rangle . \mathcal{M}[[e_Y = y]]) \\
\mathbf{P}[[e]] \text{ m} &::= \mathbf{D}[[e]] \text{ m true} \\
\mathbf{P}[[e_A | e_B]] \text{ m} &::= \mathbf{D}[[e_A | e_B]] \text{ m true true}
\end{aligned}$$

Figure 5.5: The conditional theory and query semantic functions.

before interpretation. (However, the implementation explicitly extracts and returns names.) Figure 5.5 also defines semantic functions for queries. $\mathbf{D}[[e]]$ expands to a distribution-valued computation and runs it with a probability space with the single outcome $\langle \rangle$. $\mathbf{D}[[e_X | e_Y]]$ conditions the probability space and hands off to $\mathbf{D}[[e_X]]$. $\mathbf{P}[[\cdot]]$ is defined in terms of $\mathbf{D}[[\cdot]]$.

5.5.3 Approximating Models and Queries

We compute dependent cartesian products of sets represented by lazy lists in a way similar to enumerating $\mathbb{N} \times \mathbb{N}$. (It cannot be done with a monad as in the exact semantics, but we do not need it to.) The approximating versions of dist_{ps} and cond_{ps} apply `finitize` to the probability space.

5.5.4 Implementation in Racket

$\mathcal{M}[[\cdot]]$'s implementation is `MDL`. Like `RV`, it passes random variable identifiers; unlike `RV`, `MDL` accumulates them. For example, `(MDL [] ([X ~ Px]))` expands to

$$([X] (\text{bind/ps} (\text{extend/ps} (\text{RV} [] \text{Px})) (\lambda (X) (\text{ret/ps} (\text{list } X)))))$$

where `[X]` is the updated list of identifiers and the rest is a model computation.

We store theories in transformer bindings so queries can expand them later. For example, `(define-model die-roll [L ~ P1] [E ~ (Pe/1 L)])` expands to

```
(define-syntax die-roll #'(MDL [] ([L ~ P1] [E ~ (Pe/1 L)])))
```

The macro `with-model` introduces a scope in which a theory's variables are visible. For example, `(with-model die-roll (Dist L E))` looks up `die-roll` and expands it into its identifiers and computation. Using the identifiers as lambda arguments, `Dist` (the implementation of `D[·]`) builds a query computation as in Figure 5.5, and runs it with `(mapping (list empty) (λ (ω) 1))`, the empty probability space.

Using these identifiers would break hygiene, except that `Dist` replaces the lambda arguments' lexical context. This puts the theory's exported identifiers in scope, even when the theory and query are defined in separate modules. Because queries can access only the exported identifiers, it is safe.

Aside from passing identifiers and monkeying with hygiene, the macros are almost transcribed from the semantic functions.

5.5.5 Examples

Consider a conditional distribution with the first-order definition

```
(define (Geometric p)
  (mapping N1 (λ (n) (* p (expt (- 1 p) (- n 1))))))
```

where `N1` is a lazy list of natural numbers starting at `1`. Nahin gives a delightfully morbid use for `Geometric` in his book of probability puzzles [38].

Two idiots duel with one gun. They put only one bullet in it, and take turns spinning the chamber and firing at each other. They know that if they each take one shot at a time, player one usually wins. Therefore, player one takes one shot, and after that, the next player takes one more shot than the previous player, spinning the chamber before each shot. How probable is player two's demise?

The distribution over the number of shots when the gun fires is `(Geometric 1/6)`. Using this procedure to determine whether player one fires shot `n`:

```
(define (p1-fires? n [shots 1])
  (cond [(n . <= . 0) #f]
        [else (not (p1-fires? (- n shots) (add1 shots))))]))
```

we compute the probability that player one wins with

```
(with-model (model [winning-shot ~ (Geometric 1/6)])
  (Pr (p1-fires? winning-shot)))
```

Nahin computes 0.5239191275550995247919843—25 decimal digits—with custom MATLAB code. At $\text{appx-z} \geq 321$, our solution computes the same digits. (Though it appends the digits 9..., so Nahin should have rounded up.) Implementing it took about five minutes. But the problem is not Bayesian.

This is: suppose player one slyly suggests a single coin flip to determine whether they spin the chamber before each shot. You do not see the duel, but learn that player two won. What is the probability they spun the chamber?

Suppose that the well-known `Bernoulli` and discrete `Uniform` conditional distributions are defined. Using these first-order conditional distributions and Racket’s `cond`, we can state a fairly direct theory of the duel:

```
(define-model half-idiot-duel
  [spin? ~ (Bernoulli 1/2)]
  [winning-shot ~ (cond [spin? (Geometric 1/6)]
                        [else (Uniform 1 6)])])
```

Then `(Pr spin? (not (p1-fires? winning-shot)))` converges to about 0.588.

Bayesian practitioners would normally create a new first-order conditional distribution `WinningShot`, and then state `[winning-shot ~ (WinningShot spin?)]`. Most would *like* to state something more direct—such as the above theory, which plainly shows how `spin?`’s value affects `winning-shot`’s distribution. However, without a semantics, they cannot be sure that using the value of a `cond` (or of any `if`-like expression) as a distribution is well-defined. That `winning-shot` has a *different range* for each value of `spin?` makes things more uncertain.

As specified by $\mathcal{R}[\![\cdot]\!]$, our implementation interprets `(cond ...)` above as a stochastic transition kernel. As specified by $\mathcal{M}[\![\cdot]\!]$, it builds the probability space using dependent cartesian products. Thus, the direct theory really is well-defined.

5.6 Why Separate Statements and Queries?

Whether queries should be allowed inside theories is a decision with subtle effects.

Theories are sets of facts. Well-formedness imposes a partial order, but every linearization should be interpreted equivalently. Thus, we can determine whether two kinds of statements can coexist in theories by determining whether they can be exchanged without changing the interpretation. This is equivalent to determining whether the corresponding monad functions commute.

The following definitions suppose a conditional theory $t_1; \dots; t_n$ in which exchanging some t_i and t_{i+1} (where $i < n$) is well-formed. Applying semantic functions in the definitions yields definitions that are independent of syntax but difficult to read, so we give the syntactic versions.

Definition 5.6 (commutativity). *We say that t_i and t_{i+1} **commute** when*

$$\mathcal{M}[\![t_1; \dots; t_i; t_{i+1}; \dots; t_n]\!] \langle \Omega_0, \mathbf{p}_0 \rangle = \mathcal{M}[\![t_1; \dots; t_{i+1}; t_i; \dots; t_n]\!] \langle \Omega_0, \mathbf{p}_0 \rangle.$$

Unfortunately, this notion of commutativity is usually too strong: distribution statements could never commute with each other. We need a weaker test than equality, based on *observable* outcomes.

Definition 5.7 (equivalence in distribution). *Suppose X_1, \dots, X_k are defined in t_1, \dots, t_n . Let $\mathbf{m} := \mathcal{M}[\![t_1, \dots, t_n]\!]$, and \mathbf{m}' be a (usually different) probability space monad computation. We write $\mathbf{m} \equiv_{\mathbf{D}} \mathbf{m}'$ and call \mathbf{m} and \mathbf{m}' **equivalent in distribution** when $\mathbf{D}[\![X_1, \dots, X_k]\!] \mathbf{m} = \mathbf{D}[\![X_1, \dots, X_k]\!] \mathbf{m}'$.*

The following theorem says $\equiv_{\mathbf{D}}$ is like observational equivalence with query contexts:

Theorem 5.8 (context). $\mathbf{D}[\![e_X | e_Y]\!] \mathbf{m} = \mathbf{D}[\![e_X | e_Y]\!] \mathbf{m}'$ for all random variables $\mathcal{R}[\![e_X]\!]$ and $\mathcal{R}[\![e_Y]\!]$ if and only if $\mathbf{m} \equiv_{\mathbf{D}} \mathbf{m}'$.

Definition 5.9 (commutativity in distribution). We say t_i and t_{i+1} commute *in distribution* when $\mathcal{M}[\![t_1; \dots; t_i; t_{i+1}; \dots; t_n]\!] \equiv_{\mathbf{D}} \mathcal{M}[\![t_1; \dots; t_{i+1}; t_i; \dots; t_n]\!]$.

Theorem 5.10. The following table summarizes commutativity of cond_{ps} , dist_{ps} and $\text{extend}_{\text{ps}}$ in the probability space monad:

cond_{ps}	$=$		
$\text{extend}_{\text{ps}}$	$=$	$\equiv_{\mathbf{D}}$	
dist_{ps}	$\neq_{\mathbf{D}}$	$=$	$=$
	cond_{ps}	$\text{extend}_{\text{ps}}$	dist_{ps}

By Thm. 5.10, if we are to maintain the idea that theories are sets of facts, we cannot allow both conditioning and query statements.

5.7 Conclusions

For discrete Bayesian theories, we explained a large subclass of notation as measure-theoretic calculations by transformation into λ_{ZFC} . There is now at least one precisely defined set of expressions that denote discrete conditional distributions in conditional theories, and it is very large and expressive. We gave a converging approximating semantics and implemented it in Racket.

We could have interpreted notation as first-order set theory, in which measure theory is developed. Defining the exact semantics compositionally would have been difficult, and deriving an implementation from the semantics would have involved much hand-waving. By targeting λ_{ZFC} instead, the path from notation to exact meaning to approximation to implementation is clear.

Chapter 6

Interlude: Uncountable Outcomes and Recursion

Now that we are satisfied that using λ_{ZFC} as a target language for categorical semantics of constructive theories and queries works, we turn our attention to uncountable sample spaces and theories with general recursion.

It seems that, having followed measure-theoretic structure so far, the extension to uncountable sample spaces should be fairly smooth. Discrete probability spaces, probability mass functions, summation, conditioning, and discrete transition kernels all have uncountable analogues that can be composed in much the same way. The probability space monad thus has an uncountable analogue that can be used as a target for a categorical semantics for Bayesian notation. There are two difficulties, however.

The first difficulty is practical. As with the discrete probability space monad, we would like to *derive* an implementable semantics by approximating the target category. Unfortunately, this is complicated by the fact that the uncountable computations have large cardinalities. For example, a general probability space on \mathbb{R} is defined as a triple $\langle \mathbb{R}, \Sigma, \mathbf{P} \rangle$, where Σ is a subset of $\mathcal{P} \mathbb{R}$ and $\mathbf{P} : \Sigma \rightarrow [0, 1]$.

The second difficulty is theoretical. Suppose we define the following recursive function in a language with probabilistic choice, which counts the number of times `random < p`:

$$\text{geometric } p := \text{if } (\text{random} < p) \ 0 \ (1 + \text{geometric } p) \quad (6.1)$$

To interpret `geometric p` using the uncountable probability space monad, we must interpret both branches of the `if` as probability spaces and merge them. Unfortunately, doing so naïvely

results in nontermination, as **geometric p** is applied in the *else* branch at every recurrence. Dealing with nontermination requires complicated fixpoint constructions and limits, which, with uncountable probability spaces, would put us on the frontier of research in mathematics instead of in computer science.

Fortunately, we can take a hint from measure-theoretic probability's general approach to infinite processes: define them with respect to a canonical, infinite-dimensional probability space, and encode branching and other complexities into the random variables. The next chapter takes this approach by interpreting whole programs as random variables in which every **random** expression indexes an infinite, random tree.

Chapter 7

Preimage Computation Theory: Running Programs Backwards

I am so in favor of the actual infinite that instead of admitting that Nature abhors it, as is commonly said, I hold that Nature makes frequent use of it everywhere, in order to show more effectively the perfections of its Author.

Georg Cantor

7.1 Introduction

Measure-theoretic probability [27] is widely believed to be able to define every reasonable distribution, including distributions arising from discontinuous transformations and distributions on infinite spaces. It mainly does this by *assigning probabilities to sets*. Functions that do so are **probability measures**.

If a probability measure P assigns probabilities to subsets of X and $g : X \rightarrow Y$, then the distribution over subsets of Y is

$$\Pr[B] = P(\text{preimage } g B) \quad (7.1)$$

where $\text{preimage } g B = \{a \in \text{domain } g \mid g a \in B\}$ is the subset of X for which g yields a value in B . It is well-defined for any g and B .

Measure-theoretic probability supports any kind of condition. If $\Pr[B] > 0$, the probability of $B' \subseteq Y$ given $B \subseteq Y$ is

$$\Pr[B' \mid B] = \Pr[B' \cap B] / \Pr[B] \quad (7.2)$$

If $\Pr[B] = 0$, conditional probabilities can be calculated as the limit of $\Pr[B' \mid B_n]$ for certain positive-probability $B_1 \supseteq B_2 \supseteq B_3 \supseteq \dots$ whose intersection is B [46]. For example, if $Y = \mathbb{R} \times \mathbb{R}$, the distribution of $\langle x, y \rangle \in Y$ given $x + y = 0$ can be calculated using the descending sequence $B_n = \{\langle x, y \rangle \in Y \mid |x + y| < 2^{-n}\}$.

Only special families of **measurable** sets can be assigned probabilities. Proving measurability, taking limits, and other complications tend to make measure-theoretic probability less attractive, even though it is strictly more powerful.

7.1.1 Measure-Theoretic Semantics

Most purely functional languages allow only nontermination as a side effect, and not probabilistic choice. Programmers therefore encode probabilistic programs as functions from random sources to outputs. Monads and other categorical classes such as idioms (i.e. applicative functors) can make doing so easier [22, 48].

It seems this approach should make it easy to interpret probabilistic programs measure-theoretically. For a probabilistic program $g : X \rightarrow Y$, the probability measure on output sets $B \subseteq Y$ should be defined by preimages of B under g and the probability measure on X . Unfortunately, it is difficult to turn this simple-sounding idea into a compositional semantics, for the following reasons.

1. Preimages are definable only for functions with observable domains, which excludes lambdas.
2. If subsets of X and Y must be measurable, taking preimages under g must preserve measurability (we say g itself is measurable). Proving the conditions under which this is true is difficult, especially if g may not terminate.
3. It is difficult to define useful probability measures for arbitrary spaces of measurable functions [5].

Implementing a language based on such a semantics is complicated because

4. Contemporary mathematics is unlike any implementation's host language.

5. It requires running Turing-equivalent programs backward, efficiently, on possibly uncountable sets of outputs.

We address 1 and 4 by developing our semantics in λ_{ZFC} [49], a λ -calculus with infinite sets, and both extensional and intensional functions. We address 5 by deriving and implementing a *conservative approximation* of the semantics.

There seems to be no way to simplify difficulty 2, so we work through a proof of measurability. The outcome is worth it: all probabilistic programs are measurable, regardless of the inputs on which they do not terminate. This includes uncomputable programs; for example, those that contain real equality tests and limits. We believe this result is the first of its kind, and is general enough to apply to almost all past and future work on probabilistic programming languages. To maintain the flow of this chapter, we put it off until Chapter 9.

For difficulty 3, we have discovered that the “first-orderness” of arrows [21] is a perfect fit for the “first-orderness” of measure theory.

7.1.2 Arrow Solution Overview

Using arrows, we define an *exact* semantics and an *approximating* semantics. The exact semantics includes

- A semantic function which, like the arrow calculus semantic function [32], transforms first-order programs into the computations of an arbitrary arrow.
- Arrows for evaluating expressions in different ways.

This commutative diagram describes the relationships among the six arrows used to define the exact semantics:

$$\begin{array}{ccccc}
 X \rightsquigarrow_{\perp} Y & \xrightarrow{\text{lift}_{\text{map}}} & X \rightsquigarrow_{\text{map}} Y & \xrightarrow{\text{lift}_{\text{pre}}} & X \rightsquigarrow_{\text{pre}} Y \\
 \eta_{\perp^*} \downarrow & & \downarrow \eta_{\text{map}^*} & & \downarrow \eta_{\text{pre}^*} \\
 X \rightsquigarrow_{\perp^*} Y & \xrightarrow{\text{lift}_{\text{map}^*}} & X \rightsquigarrow_{\text{map}^*} Y & \xrightarrow{\text{lift}_{\text{pre}^*}} & X \rightsquigarrow_{\text{pre}^*} Y
 \end{array} \tag{7.3}$$

At the top-left, $X \rightsquigarrow_{\perp} Y$ computations (or “bottom arrow computations”) are intensional functions that may raise errors (i.e. return \perp , which is read “bottom”). From bottom arrow

computations, the lift_{map} combinator produces $X \rightsquigarrow_{\text{map}} Y$ computations, which create equivalent extensional functions, or mappings. From mapping arrow computations, the lift_{pre} combinator produces $X \rightsquigarrow_{\text{pre}} Y$ computations, which compute preimages.

Instances of arrows in the bottom row are like those in the top, except they thread an infinite store of random values, and can be constructed to always terminate.

Most of our correctness theorems rely on proofs that every combinator in (7.3) is a homomorphism; for example, that lift_{map} distributes over all bottom arrow combinators.

The approximating semantics uses the same semantic function, but its arrows $X \rightsquigarrow_{\text{pre}} Y$ and $X \rightsquigarrow_{\text{pre}^*}' Y$ compute conservative approximations. Given a library for representing and operating on rectangular sets, it is directly implementable.

7.2 Arrows and First-Order Semantics

Like monads [54] and idioms [35], arrows [21] thread effects through computations in a way that imposes structure. But arrow computations are always

- Function-like: An arrow computation of type $x \rightsquigarrow_a y$ must behave like a corresponding function of type $x \Rightarrow y$ (in a sense we explain shortly).
- First-order: There is no way to derive a computation $\text{app}_a : \langle x \rightsquigarrow_a y, x \rangle \rightsquigarrow_a y$ from the arrow a 's minimal definition, so it is not possible for an arrow computation to apply another arrow computation.

The first property makes arrows a good fit for a compositional translation from expressions to pure functions that operate on random sources. The second property makes arrows a good fit for a measure-theoretic semantics in particular, as it is difficult to define useful measurable sets of functions that make app 's corresponding function measurable [5].

7.2.1 Alternative Arrow Definitions and Laws

To make applying measure-theoretic theorems easier, and to simplify interpreting let-calculus expressions as arrow computations, we do not give typical minimal arrow definitions. For

each arrow a , instead of first_a , we define $(\&\&\&_a)$. This combinator is typically called **fanout**, but its use will be clearer if we call it **pairing**. One way to strengthen an arrow a is to define an additional combinator left_a , which can be used to choose an arrow computation based on the result of another. Again, we define a different combinator, ifte_a (“if-then-else”).

In a nonstrict λ -calculus, defining a choice combinator allows writing recursive functions using nothing but arrow combinators and lifted, pure functions. However, a strict λ -calculus needs an extra combinator **lazy** for deferring conditional branches. For example, define the **function arrow** with choice, in which $x \rightsquigarrow y ::= x \Rightarrow y$:

$$\begin{array}{ll}
\text{arr } f & := f & \text{lift} \\
f_1 >>> f_2 & := \lambda a. f_2 (f_1 a) & \text{composition} \\
f_1 \&\&\& f_2 & := \lambda a. \langle f_1 a, f_2 a \rangle & \text{pairing} \\
\text{ifte } f_1 f_2 f_3 & := \lambda a. \text{if } (f_1 a) (f_2 a) (f_3 a) & \text{if-then-else} \\
\text{lazy } f & := \lambda a. f \ 0 \ a & \text{laziness}
\end{array} \tag{7.4}$$

and try to define the following recursive function:

$$\begin{array}{ll}
\text{halt-on-true} : \text{Bool} \rightsquigarrow \text{Bool} & (\text{i.e. } \text{halt-on-true} : \text{Bool} \Rightarrow \text{Bool}) \\
\text{halt-on-true} & := \text{ifte } (\text{arr id}) (\text{arr id}) \text{halt-on-true} \\
& \equiv \text{ifte id id } (\text{ifte } (\text{arr id}) (\text{arr id}) \text{halt-on-true}) \\
& \equiv \text{ifte id id } (\text{ifte id id } (\text{ifte } (\text{arr id}) (\text{arr id}) \text{halt-on-true}))
\end{array} \tag{7.5}$$

In a strict λ -calculus, the defining expression does not terminate. But the following is

well-defined in λ_{ZFC} , and loops only when applied to **false**:

$$\begin{aligned}
\text{halt-on-true} &:= \text{ifte } (\text{arr id}) (\text{arr id}) (\text{lazy } \lambda 0. \text{halt-on-true}) \\
&\equiv \text{ifte id id } (\lambda a. (\lambda 0. \text{halt-on-true}) 0 a) \\
&\equiv \lambda a. \text{if } (\text{id } a) (\text{id } a) ((\lambda a. (\lambda 0. \text{halt-on-true}) 0 a) a) \\
&\equiv \lambda a. \text{if } a a ((\lambda a. \text{halt-on-true } a) a) \\
&\equiv \lambda a. \text{if } a a (\text{halt-on-true } a)
\end{aligned} \tag{7.6}$$

All of our arrows are arrows with choice and **lazy**, so we simply call them arrows.

Definition 7.1 (arrow). *Let $1 := \{0\}$ (Section 3.3.1). A binary type constructor (\rightsquigarrow_a) and*

$$\begin{aligned}
\text{arr}_a : (x \Rightarrow y) &\Rightarrow (x \rightsquigarrow_a y) && \text{lift} \\
(>>>_a) : (x \rightsquigarrow_a y) &\Rightarrow (y \rightsquigarrow_a z) \Rightarrow (x \rightsquigarrow_a z) && \text{composition} \\
(\&\&\&_a) : (x \rightsquigarrow_a y) &\Rightarrow (x \rightsquigarrow_a z) \Rightarrow (x \rightsquigarrow_a \langle y, z \rangle) && \text{pairing} \\
\text{ifte}_a : (x \rightsquigarrow_a \text{Bool}) &\Rightarrow (x \rightsquigarrow_a y) \Rightarrow (x \rightsquigarrow_a y) \Rightarrow (x \rightsquigarrow_a y) && \text{if-then-else} \\
\text{lazy}_a : (1 \Rightarrow (x \rightsquigarrow_a y)) &\Rightarrow (x \rightsquigarrow_a y) && \text{laziness}
\end{aligned} \tag{7.7}$$

define an **arrow** if certain monoid, homomorphism, and structural laws hold.

The arrow homomorphism laws can be put in terms of more general homomorphism properties that deal with distributing an arrow-to-arrow lift, which we use extensively to prove correctness.

Definition 7.2 (arrow homomorphism). *A function $\text{lift}_b : (x \rightsquigarrow_a y) \Rightarrow (x \rightsquigarrow_b y)$ is an **arrow homomorphism** from arrow **a** to arrow **b** if the following distributive laws hold for*

appropriately typed f , f_1 , f_2 and f_3 :

$$\text{lift}_b (\text{arr}_a f) \equiv \text{arr}_b f \quad (7.8)$$

$$\text{lift}_b (f_1 \ggg_a f_2) \equiv (\text{lift}_b f_1) \ggg_b (\text{lift}_b f_2) \quad (7.9)$$

$$\text{lift}_b (f_1 \&\&\&_a f_2) \equiv (\text{lift}_b f_1) \&\&\&_b (\text{lift}_b f_2) \quad (7.10)$$

$$\text{lift}_b (\text{ifte}_a f_1 f_2 f_3) \equiv \text{ifte}_b (\text{lift}_b f_1) (\text{lift}_b f_2) (\text{lift}_b f_3) \quad (7.11)$$

$$\text{lift}_b (\text{lazy}_a f) \equiv \text{lazy}_b \lambda 0. \text{lift}_b (f \ 0) \quad (7.12)$$

The arrow homomorphism laws state that $\text{arr}_a : (x \Rightarrow y) \Rightarrow (x \rightsquigarrow_a y)$ must be a homomorphism from the function arrow (7.4) to arrow a . Roughly, arrow computations that do not use additional combinators can be transformed into arr_a applied to a pure computation. They must be *function-like*.

Only a few of the other arrow laws play a role in our semantics and its correctness. We need associativity of (\ggg_a) and a pair extraction law:

$$(f_1 \ggg_a f_2) \ggg_a f_3 \equiv f_1 \ggg_a (f_2 \ggg_a f_3) \quad (7.13)$$

$$(\text{arr}_a f_1 \&\&\&_a f_2) \ggg_a \text{arr}_a \text{snd} \equiv f_2 \quad (7.14)$$

and distribution of pure computations over effectful:

$$\text{arr}_a f_1 \ggg_a (f_2 \&\&\&_a f_3) \equiv (\text{arr}_a f_1 \ggg_a f_2) \&\&\&_a (\text{arr}_a f_1 \ggg_a f_3) \quad (7.15)$$

$$\begin{aligned} \text{arr}_a f_1 \ggg_a \text{ifte}_a f_2 f_3 f_4 &\equiv \text{ifte}_a (\text{arr}_a f_1 \ggg_a f_2) \\ &\quad (\text{arr}_a f_1 \ggg_a f_3) \\ &\quad (\text{arr}_a f_1 \ggg_a f_4) \end{aligned} \quad (7.16)$$

$$\text{arr}_a f_1 \ggg_a \text{lazy}_a f_2 \equiv \text{lazy}_a \lambda 0. \text{arr}_a f_1 \ggg_a f_2 \ 0 \quad (7.17)$$

Equivalence between different arrow representations is usually proved in a strongly normalizing λ -calculus [31, 32], in which every function is free of effects, including nontermination. Such a λ -calculus has no need for lazy_a , so we could not derive (7.17) from existing arrow laws. We follow Hughes's reasoning [21] for the original arrow laws: it is a function-like

property (i.e. it holds for the function arrow), and it cannot not lose, reorder or duplicate effects.

The pair extraction law (7.14), which *can* be derived from existing arrow laws, is a more problematic, in nonstrict λ -calculi as well as λ_{ZFC} . If f_1 does not always terminate, using (7.14) to transform a computation can turn a nonterminating expression into a terminating one, or vice-versa. We could require f_1 in the pair extraction law to always terminate. Instead, we require every argument to arr_a to terminate, which simplifies more proofs.

Rather than prove each arrow law for each arrow, we prove arrows are *epimorphic* to arrows for which the laws are known to hold. (Isomorphism is sufficient but not necessary.)

Definition 7.3 (arrow epimorphism). *An arrow homomorphism $\text{lift}_b : (x \rightsquigarrow_a y) \Rightarrow (x \rightsquigarrow_b y)$ that has a right inverse is an **arrow epimorphism** from a to b .*

Theorem 7.4 (epimorphism implies arrow laws). *If $\text{lift}_b : (x \rightsquigarrow_a y) \Rightarrow (x \rightsquigarrow_b y)$ is an arrow epimorphism and the combinators of a define an arrow, then the combinators of b define an arrow.*

Proof. Let lift_b^{-1} be lift_b 's right inverse. For the pair extraction law (7.14),

$$\begin{aligned}
& (\text{arr}_b f_1 \&\&_b f_2) \gg\gg\gg_b \text{arr}_b \text{snd} & (7.18) \\
& \equiv (\text{lift}_b (\text{arr}_a f_1) \&\&_b (\text{lift}_b (\text{lift}_b^{-1} f_2))) \gg\gg\gg_b \text{lift}_b (\text{arr}_a \text{snd}) & \text{Rewrite with } \text{lift}_b \\
& \equiv \text{lift}_b (\text{arr}_a f_1 \&\&_a \text{lift}_b^{-1} f_2) \gg\gg\gg_b \text{lift}_b (\text{arr}_a \text{snd}) & \text{Homomorphism (7.10)} \\
& \equiv \text{lift}_b ((\text{arr}_a f_1 \&\&_a \text{lift}_b^{-1} f_2) \gg\gg\gg_a \text{arr}_a \text{snd}) & \text{Homomorphism (7.9)} \\
& \equiv \text{lift}_b (\text{lift}_b^{-1} f_2) & \text{Pair extraction (7.14)} \\
& \equiv f_2 & \text{Right inverse}
\end{aligned}$$

The proofs for every other law are similar. □

From here on, we assume all expressions are well-defined. (The arrow \mathbf{a} will be clear from context.) Well-definedness does not guarantee that *running* an interpretation terminates. It just simplifies statements about expressions, such as the following theorem, on which most of our semantic correctness results rely.

Theorem 7.6 (homomorphisms distribute over expressions). *Let $\text{lift}_b : (x \rightsquigarrow_a y) \Rightarrow (x \rightsquigarrow_b y)$ be an arrow homomorphism. For all e , $\llbracket e \rrbracket_b \equiv \text{lift}_b \llbracket e \rrbracket_a$.*

Proof. By structural induction. Base cases proceed by expansion and using $\text{arr}_b \equiv \text{lift}_b \circ \text{arr}_a$ (7.8). For example, for constants:

$$\begin{aligned}
\llbracket v \rrbracket_b &\equiv \text{arr}_b (\text{const } v) && \text{Def of } \llbracket \cdot \rrbracket_b && (7.19) \\
&\equiv \text{lift}_b (\text{arr}_a (\text{const } v)) && \text{Homomorphism (7.8)} \\
&\equiv \text{lift}_b \llbracket v \rrbracket_a && \text{Def of } \llbracket \cdot \rrbracket_a
\end{aligned}$$

Inductive cases proceed by expansion, applying the inductive hypothesis on subterms, and applying distributive laws (7.9)–(7.12). For example, for pairing:

$$\begin{aligned}
\llbracket \langle e_1, e_2 \rangle \rrbracket_b &\equiv \llbracket e_1 \rrbracket_b \&\&_b \llbracket e_2 \rrbracket_b && \text{Def of } \llbracket \cdot \rrbracket_b && (7.20) \\
&\equiv (\text{lift}_b \llbracket e_1 \rrbracket_a) \&\&_b (\text{lift}_b \llbracket e_2 \rrbracket_a) && \text{Ind hypothesis} \\
&\equiv \text{lift}_b (\llbracket e_1 \rrbracket_a \&\&_a \llbracket e_2 \rrbracket_a) && \text{Homomorphism (7.10)} \\
&\equiv \text{lift}_b \llbracket \langle e_1, e_2 \rangle \rrbracket_a && \text{Def of } \llbracket \cdot \rrbracket_a
\end{aligned}$$

It is not hard to check the remaining cases. □

If we assume lift_b defines correct behavior for arrow \mathbf{b} in terms of arrow \mathbf{a} , and prove that lift_b is a homomorphism, then by Theorem 7.6, $\llbracket \cdot \rrbracket_b$ is correct.

$$\begin{array}{ll}
X \rightsquigarrow_{\perp} Y ::= X \Rightarrow Y_{\perp} & \text{ifte}_{\perp} : (X \rightsquigarrow_{\perp} \text{Bool}) \Rightarrow (X \rightsquigarrow_{\perp} Y) \Rightarrow (X \rightsquigarrow_{\perp} Y) \Rightarrow (X \rightsquigarrow_{\perp} Y) \\
\text{arr}_{\perp} : (X \Rightarrow Y) \Rightarrow (X \rightsquigarrow_{\perp} Y) & \text{ifte}_{\perp} f_1 f_2 f_3 a := \text{case } f_1 a \\
\text{arr}_{\perp} f := f & \quad \text{true} \longrightarrow f_2 a \\
& \quad \text{false} \longrightarrow f_3 a \\
& \quad \perp \longrightarrow \perp \\
(\ggg_{\perp}) : (X \rightsquigarrow_{\perp} Y) \Rightarrow (Y \rightsquigarrow_{\perp} Z) \Rightarrow (X \rightsquigarrow_{\perp} Z) & \text{lazy}_{\perp} : (1 \Rightarrow (X \rightsquigarrow_{\perp} Y)) \Rightarrow (X \rightsquigarrow_{\perp} Y) \\
(f_1 \ggg_{\perp} f_2) a := \text{if } (f_1 a = \perp) \perp (f_2 (f_1 a)) & \text{lazy}_{\perp} f a := f \ 0 \ a \\
(\lll_{\perp}) : (X \rightsquigarrow_{\perp} Y_1) \Rightarrow (X \rightsquigarrow_{\perp} Y_2) \Rightarrow (X \rightsquigarrow_{\perp} \langle Y_1, Y_2 \rangle) & \\
(f_1 \lll_{\perp} f_2) a := \text{let } b_1 := f_1 a & \\
& \quad b_2 := f_2 a \\
& \quad \text{in if } (b_1 = \perp \text{ or } b_2 = \perp) \perp \langle b_1, b_2 \rangle
\end{array}$$

Figure 7.2: Bottom arrow definitions.

7.3 The Bottom Arrow

Using the diagram in (7.3) as a sort of map, we start in the upper-left corner:

$$\begin{array}{ccccc}
X \rightsquigarrow_{\perp} Y & \xrightarrow{\text{lift}_{\text{map}}} & X \rightsquigarrow_{\text{map}} Y & \xrightarrow{\text{lift}_{\text{pre}}} & X \rightsquigarrow_{\text{pre}} Y \\
\eta_{\perp}^* \downarrow & & \downarrow \eta_{\text{map}}^* & & \downarrow \eta_{\text{pre}}^* \\
X \rightsquigarrow_{\perp}^* Y & \xrightarrow{\text{lift}_{\text{map}}^*} & X \rightsquigarrow_{\text{map}}^* Y & \xrightarrow{\text{lift}_{\text{pre}}^*} & X \rightsquigarrow_{\text{pre}}^* Y
\end{array} \tag{7.21}$$

Through Section 7.6, we move across the top to $X \rightsquigarrow_{\text{pre}} Y$.

To use Theorem 7.6 to prove correct the interpretations of expressions as preimage arrow computations, we need the preimage arrow to be homomorphic to a simpler arrow with easily understood behavior. The function arrow (7.4) is an obvious candidate. However, we will need to explicitly handle nontermination as an error value, so we need a slightly more complicated arrow.

Figure 7.2 defines the **bottom arrow**. Its computations have type $X \rightsquigarrow_{\perp} Y ::= X \Rightarrow Y_{\perp}$, where $Y_{\perp} ::= Y \cup \{\perp\}$ and \perp is a distinguished error value. The type Bool_{\perp} , for example, denotes the members of $\text{Bool} \cup \{\perp\} = \{\text{true}, \text{false}, \perp\}$.

To prove the arrow laws, we need a coarser notion of equivalence.

Definition 7.7 (bottom arrow equivalence). *Two computations $f_1 : X \rightsquigarrow_{\perp} Y$ and $f_2 : X \rightsquigarrow_{\perp} Y$*

are equivalent, or $f_1 \equiv f_2$, when $f_1 a \equiv f_2 a$ for all $a \in X$.

Theorem 7.8. arr_\perp , $(\&\&\&)_\perp$, $(\gg\gg)_\perp$, ifte_\perp and lazy_\perp define an arrow.

Proof. The bottom arrow is epimorphic to (in fact, isomorphic to) the maybe monad's Kleisli arrow. \square

7.4 Deriving the Mapping Arrow

Computing preimages requires an observable domain, which lambdas do not have. Further, theorems about functions in set theory tend to be about mappings, not about lambdas that may raise errors. As in intermediate step, then, we need an arrow whose computations produce mappings or are mappings themselves.

It is tempting to try to make the mapping arrow's computations mapping-valued; i.e. $X \rightsquigarrow_{\text{map}} Y ::= X \multimap Y$. Unfortunately, we could not define $\text{arr}_{\text{map}} : (X \Rightarrow Y) \Rightarrow (X \multimap Y)$: to define a mapping, we need a domain, but lambdas' domains are unobservable.

To parameterize mapping arrow computations on a domain, we define the *mapping arrow* computation type as

$$X \overset{\sim}{\underset{\text{man}}{\rightrightarrows}} Y ::= \text{Set } X \Rightarrow (X \multimap Y) \quad (7.22)$$

The absence of \perp in $\mathbf{Set} \ X \Rightarrow (X \multimap Y)$, and the fact that type parameters X and Y denote sets, will make it easier to apply well-known theorems from measure theory, which know nothing of lambda types and propagating error values.

To use Theorem 7.6 to prove that expressions interpreted using $\llbracket \cdot \rrbracket_{\text{map}}$ behave correctly with respect to $\llbracket \cdot \rrbracket_{\perp}$, we need to define correctness using a lift from the bottom arrow to the mapping arrow. It is helpful to have a standalone function `domainl` that computes the subset

$\text{range} : (X \multimap Y) \Rightarrow \text{Set } Y$	$\langle \cdot, \cdot \rangle_{\text{map}} : (X \multimap Y_1) \Rightarrow (X \multimap Y_2) \Rightarrow (X \multimap Y_1 \times Y_2)$
$\text{range } g := \text{image } \text{snd } g$	$\langle g_1, g_2 \rangle_{\text{map}} := \text{let } A := \text{domain } g_1 \cap \text{domain } g_2$ in $\lambda a \in A. \langle g_1 a, g_2 a \rangle$
$(\circ_{\text{map}}) : (Y \multimap Z) \Rightarrow (X \multimap Y) \Rightarrow (X \multimap Z)$	$(\uplus_{\text{map}}) : (X \multimap Y) \Rightarrow (X \multimap Y) \Rightarrow (X \multimap Y)$
$g_2 \circ_{\text{map}} g_1 := \text{let } A := \text{preimage } g_1 (\text{domain } g_2)$ in $\lambda a \in A. g_2 (g_1 a)$	$g_1 \uplus_{\text{map}} g_2 := \text{let } A := \text{domain } g_1 \uplus \text{domain } g_2$ in $\lambda a \in A. \text{if } (a \in \text{domain } g_1) (g_1 a) (g_2 a)$

Figure 7.3: Additional operations on partial mappings.

of A on which f does not return \perp . We define that first, and then define lift_{map} in terms of it:

$$\text{domain}_{\perp} : (X \rightsquigarrow_{\perp} Y) \Rightarrow \text{Set } X \Rightarrow \text{Set } X \quad (7.23)$$

$$\text{domain}_{\perp} f A := \{a \in A \mid f a \neq \perp\}$$

$$\text{lift}_{\text{map}} : (X \rightsquigarrow_{\perp} Y) \Rightarrow (X \rightsquigarrow_{\text{map}} Y) \quad (7.24)$$

$$\text{lift}_{\text{map}} f A := \text{mapping } f (\text{domain}_{\perp} f A)$$

So $\text{lift}_{\text{map}} f A$ is like $\text{mapping } f A$, except the domain does not contain inputs that produce errors—a good notion of correctness.

If lift_{map} is to be a homomorphism, mapping arrow computation equivalence needs to be more extensional.

Definition 7.9 (mapping arrow equivalence). *Two computations $g_1 : X \rightsquigarrow_{\text{map}} Y$ and $g_2 : X \rightsquigarrow_{\text{map}} Y$ are equivalent, or $g_1 \equiv g_2$, when $g_1 A \equiv g_2 A$ for all $A \subseteq X$.*

Clearly $\text{arr}_{\text{map}} := \text{lift}_{\text{map}} \circ \text{arr}_{\perp}$ meets the first homomorphism law (7.8). The remainder of this section derives $(\&\&\&_{\text{map}})$, $(\>\>\>_{\text{map}})$, ifte_{map} and lazy_{map} from bottom arrow combinators, in a way that ensures lift_{map} is an arrow homomorphism. Figure 7.3 defines the additional necessary mapping operations range , composition, pairing, and disjoint union, and Figure 7.4 contains the resulting mapping arrow combinators.

$X_{\text{map}} \rightsquigarrow Y ::= \text{Set } X \Rightarrow (X \multimap Y)$	$\text{ifte}_{\text{map}} : (X_{\text{map}} \rightsquigarrow \text{Bool}) \Rightarrow (X_{\text{map}} \rightsquigarrow Y) \Rightarrow (X_{\text{map}} \rightsquigarrow Y) \Rightarrow (X_{\text{map}} \rightsquigarrow Y)$
$\text{arr}_{\text{map}} : (X \Rightarrow Y) \Rightarrow (X_{\text{map}} \rightsquigarrow Y)$	$\text{ifte}_{\text{map}} \ g_1 \ g_2 \ g_3 \ A := \text{let } g'_1 := g_1 \ A$
$\text{arr}_{\text{map}} := \text{lift}_{\text{map}} \circ \text{arr}_{\perp}$	$\quad g'_2 := g_2 \ (\text{preimage } g'_1 \ \{\text{true}\})$
	$\quad g'_3 := g_3 \ (\text{preimage } g'_1 \ \{\text{false}\})$
	$\quad \text{in } g'_2 \uplus_{\text{map}} g'_3$
$(\ggg_{\text{map}}) : (X_{\text{map}} \rightsquigarrow Y) \Rightarrow (Y_{\text{map}} \rightsquigarrow Z) \Rightarrow (X_{\text{map}} \rightsquigarrow Z)$	$\text{lazy}_{\text{map}} : (1 \Rightarrow (X_{\text{map}} \rightsquigarrow Y)) \Rightarrow (X_{\text{map}} \rightsquigarrow Y)$
$(g_1 \ggg_{\text{map}} g_2) \ A := \text{let } g'_1 := g_1 \ A$	$\text{lazy}_{\text{map}} \ g \ A := \text{if } (A = \emptyset) \ \emptyset \ (g \ 0 \ A)$
$\quad g'_2 := g_2 \ (\text{range } g'_1)$	
$\quad \text{in } g'_2 \circ_{\text{map}} g'_1$	
$(\&\&\&_{\text{map}}) : (X_{\text{map}} \rightsquigarrow Y_1) \Rightarrow (X_{\text{map}} \rightsquigarrow Y_2) \Rightarrow (X_{\text{map}} \rightsquigarrow (Y_1, Y_2))$	$\text{lift}_{\text{map}} : (X \rightsquigarrow_{\perp} Y) \Rightarrow (X_{\text{map}} \rightsquigarrow Y)$
$(g_1 \&\&\&_{\text{map}} g_2) \ A := \langle g_1 \ A, g_2 \ A \rangle_{\text{map}}$	$\text{lift}_{\text{map}} \ f \ A := \{\langle a, b \rangle \in \text{mapping } f \ A \mid b \neq \perp\}$

Figure 7.4: Mapping arrow definitions.

7.4.1 Composition

Starting with the left side of (7.9), we expand definitions, simplify f by restricting it to a set for which $f_1 \ a \neq \perp$:

$$\begin{aligned}
& \text{lift}_{\text{map}} (f_1 \ggg_{\perp} f_2) \ A && (7.25) \\
& \equiv \text{let } f := \lambda a. \text{if } (f_1 \ a = \perp) \ \perp \ (f_2 \ (f_1 \ a)) && \text{Def of } \text{lift}_{\text{map}}, (\ggg_{\perp}) \\
& \quad A' := \text{domain}_{\perp} f \ A \\
& \quad \text{in } \text{mapping } f \ A' \\
& \equiv \text{let } f := \lambda a. f_2 \ (f_1 \ a) && \text{Simplify } f \\
& \quad A' := \text{domain}_{\perp} f \ (\text{domain}_{\perp} f_1 \ A) \\
& \quad \text{in } \text{mapping } f \ A' \\
& \equiv \text{let } A' := \{a \in \text{domain}_{\perp} f_1 \ A \mid f_2 \ (f_1 \ a) \neq \perp\} && \text{Def of } \text{domain}_{\perp}, \text{ mapping} \\
& \quad \text{in } \lambda a \in A'. f_2 \ (f_1 \ a)
\end{aligned}$$

We finish by converting bottom arrow computations to the mapping arrow and rewriting in terms of mapping composition (\circ_{map}):

$$\begin{aligned}
& \equiv \text{let } g_1 := \text{lift}_{\text{map}} f_1 \ A && \text{Rewrite with } \text{lift}_{\text{map}} \\
& \quad A' := \text{preimage } g_1 \ (\text{domain}_{\perp} f_2 \ (\text{range } g_1)) \\
& \quad \text{in } \lambda a \in A'. f_2 \ (g_1 \ a)
\end{aligned}$$

$$\begin{aligned}
&\equiv \text{let } g_1 := \text{lift}_{\text{map}} f_1 A && \text{Rewrite with } \text{lift}_{\text{map}} \\
&\quad g_2 := \text{lift}_{\text{map}} f_2 (\text{range } g_1) \\
&\quad A' := \text{preimage } g_1 (\text{domain } g_2) \\
&\text{in } \lambda a \in A'. g_2 (g_1 a)
\end{aligned}$$

$$\begin{aligned}
&\equiv \text{let } g_1 := \text{lift}_{\text{map}} f_1 A && \text{Rewrite with } (\circ_{\text{map}}) \\
&\quad g_2 := \text{lift}_{\text{map}} f_2 (\text{range } g_1) \\
&\text{in } g_2 \circ_{\text{map}} g_1
\end{aligned}$$

Substituting g_1 for $\text{lift}_{\text{map}} f_1$ and g_2 for $\text{lift}_{\text{map}} f_2$ gives a definition for (\ggg_{map}) (Figure 7.4) for which (7.9) holds.

7.4.2 Pairing

Starting with the left side of (7.10), we expand definitions, and simplify f by restricting it to a set for which $f_1 a \neq \perp$ and $f_2 a \neq \perp$:

$$\text{lift}_{\text{map}} (f_1 \&\&\&_{\perp} f_2) A \tag{7.26}$$

$$\begin{aligned}
&\equiv \text{let } f := \lambda a. \text{let } b_1 := f_1 a && \text{Def of } \text{lift}_{\text{map}}, (\&\&\&_{\perp}) \\
&\quad b_2 := f_2 a \\
&\quad \text{in if } (b_1 = \perp \text{ or } b_2 = \perp) \perp \langle b_1, b_2 \rangle \\
&\quad A' := \text{domain}_{\perp} f A \\
&\text{in mapping } f A'
\end{aligned}$$

$$\begin{aligned}
&\equiv \text{let } f := \lambda a. \langle f_1 a, f_2 a \rangle && \text{Simplify } f \\
&\quad A' := \text{domain}_{\perp} f_1 A \cap \text{domain}_{\perp} f_2 A \\
&\text{in mapping } f A'
\end{aligned}$$

$$\begin{aligned}
&\equiv \text{let } A' := \text{domain}_{\perp} f_1 A \cap \text{domain}_{\perp} f_2 A && \text{Def of mapping} \\
&\text{in } \lambda a \in A'. \langle f_1 a, f_2 a \rangle
\end{aligned}$$

We finish by converting bottom arrow computations to the mapping arrow and rewriting in terms of $\langle \cdot, \cdot \rangle_{\text{map}}$:

$$\begin{aligned}
&\equiv \text{let } g_1 := \text{lift}_{\text{map}} f_1 A && \text{Rewrite with } \text{lift}_{\text{map}} \\
&\quad g_2 := \text{lift}_{\text{map}} f_2 A \\
&\quad A' := \text{domain } g_1 \cap \text{domain } g_2 \\
&\quad \text{in } \lambda a \in A'. \langle g_1 a, g_2 a \rangle \\
&\equiv \langle \text{lift}_{\text{map}} f_1 A, \text{lift}_{\text{map}} f_2 A \rangle_{\text{map}} && \text{Rewrite with } \langle \cdot, \cdot \rangle_{\text{map}}
\end{aligned}$$

Substituting g_1 for $\text{lift}_{\text{map}} f_1$ and g_2 for $\text{lift}_{\text{map}} f_2$ gives a definition for $(\&\&\&_{\text{map}})$ (Figure 7.4) for which (7.10) holds.

7.4.3 Conditional

Starting with the left side of (7.11), we expand definitions, and simplify f by restricting it to a domain for which $f_1 a \neq \perp$:

$$\begin{aligned}
&\text{lift}_{\text{map}} (\text{ifte}_{\perp} f_1 f_2 f_3) A && (7.27) \\
&\equiv \text{let } f := \lambda a. \text{case } f_1 a && \text{Def of } \text{lift}_{\text{map}}, \text{ifte}_{\perp} \\
&\quad \text{true} \longrightarrow f_2 a \\
&\quad \text{false} \longrightarrow f_3 a \\
&\quad \perp \longrightarrow \perp \\
&\quad A' := \text{domain}_{\perp} f A \\
&\quad \text{in mapping } f A' \\
&\equiv \text{let } f := \lambda a. \text{if } (f_1 a) (f_2 a) (f_3 a) && \text{Simplify } f \\
&\quad g_1 := \text{mapping } f_1 (\text{domain}_{\perp} f_1 A) \\
&\quad A_2 := \text{preimage } g_1 \{\text{true}\} \\
&\quad A_3 := \text{preimage } g_1 \{\text{false}\} \\
&\quad A' := \text{domain}_{\perp} f_2 A_2 \uplus \text{domain}_{\perp} f_3 A_3 \\
&\quad \text{in mapping } f A' \\
&\equiv \text{let } g_1 := \text{mapping } f_1 (\text{domain}_{\perp} f_1 A) && \text{Def of mapping} \\
&\quad A_2 := \text{preimage } g_1 \{\text{true}\} \\
&\quad A_3 := \text{preimage } g_1 \{\text{false}\} \\
&\quad A' := \text{domain}_{\perp} f_2 A_2 \uplus \text{domain}_{\perp} f_3 A_3 \\
&\quad \text{in } \lambda a \in A'. \text{if } (f_1 a) (f_2 a) (f_3 a)
\end{aligned}$$

We finish by converting bottom arrow computations to the mapping arrow and rewriting in terms of (\uplus_{map}) :

$$\begin{aligned}
&\equiv \text{let } g_1 := \text{lift}_{\text{map}} f_1 A && \text{Rewrite with } \text{lift}_{\text{map}} \\
&\quad g_2 := \text{lift}_{\text{map}} f_2 (\text{preimage } g_1 \{ \text{true} \}) \\
&\quad g_3 := \text{lift}_{\text{map}} f_3 (\text{preimage } g_1 \{ \text{false} \}) \\
&\quad A' := \text{domain } g_2 \uplus \text{domain } g_3 \\
&\quad \text{in } \lambda a \in A'. \text{if } (a \in \text{domain } g_2) (g_2 a) (g_3 a) \\
\\
&\equiv \text{let } g_1 := \text{lift}_{\text{map}} f_1 A && \text{Rewrite with } (\uplus_{\text{map}}) \\
&\quad g_2 := \text{lift}_{\text{map}} f_2 (\text{preimage } g_1 \{ \text{true} \}) \\
&\quad g_3 := \text{lift}_{\text{map}} f_3 (\text{preimage } g_1 \{ \text{false} \}) \\
&\quad \text{in } g_2 \uplus_{\text{map}} g_3
\end{aligned}$$

Substituting g_1 for $\text{lift}_{\text{map}} f_1$, g_2 for $\text{lift}_{\text{map}} f_2$, and g_3 for $\text{lift}_{\text{map}} f_3$ gives a definition for ifte_{map} (Figure 7.4) for which (7.11) holds.

7.4.4 Laziness

Starting with the left side of (7.12), we expand definitions:

$$\begin{aligned}
\text{lift}_{\text{map}} (\text{lazy}_{\perp} f) A &\equiv \text{let } A' := \text{domain}_{\perp} (\lambda a. f \ 0 \ a) A && (7.28) \\
&\quad \text{in } \text{mapping } (\lambda a. f \ 0 \ a) A'
\end{aligned}$$

It appears we need an η rule to continue, which λ_{ZFC} does not have (i.e. $\lambda x. e \ x \not\equiv e$ because e may not terminate). Fortunately, we can use weaker facts. If $A \neq \emptyset$, then $\text{domain}_{\perp} (\lambda a. f \ 0 \ a) A \equiv \text{domain}_{\perp} (f \ 0) A$. Further, it terminates if and only if $\text{mapping } (f \ 0) A'$ terminates. Therefore, if $A \neq \emptyset$, we can replace $\lambda a. f \ 0 \ a$ with $f \ 0$. If $A = \emptyset$, then $\text{lift}_{\text{map}} (\text{lazy}_{\perp} f) A = \emptyset$ (the empty mapping), so

$$\begin{aligned}
\text{lift}_{\text{map}} (\text{lazy}_{\perp} f) A &\equiv \text{if } (A = \emptyset) \ \emptyset \ (\text{mapping } (f \ 0) (\text{domain}_{\perp} (f \ 0) A)) && (7.29) \\
&\equiv \text{if } (A = \emptyset) \ \emptyset \ (\text{lift}_{\text{map}} (f \ 0) A)
\end{aligned}$$

Substituting $g \ 0$ for $\text{lift}_{\text{map}} (f \ 0)$ gives a lazy_{map} (Figure 7.4) for which (7.12) holds.

7.4.5 Correctness

Theorem 7.10 (mapping arrow correctness). lift_{map} is a homomorphism.

Proof. By construction. □

Corollary 7.11 (semantic correctness). For all e , $\llbracket e \rrbracket_{\text{map}} \equiv \text{lift}_{\text{map}} \llbracket e \rrbracket_{\perp}$.

Without restrictions, mapping arrow computations can be quite unruly. For example, the following computation is well-typed, but returns the identity mapping on **Bool** when applied to an empty domain, and the empty mapping when applied to any other domain:

$$\begin{aligned} \text{nonmonotone} &: \text{Bool} \rightsquigarrow_{\text{map}} \text{Bool} \\ \text{nonmonotone } A &:= \text{if } (A = \emptyset) (\lambda a \in \text{Bool}. a) \emptyset \end{aligned} \tag{7.30}$$

It would be nice if we could be sure that every $X \rightsquigarrow_{\text{map}} Y$ is not only monotone, but acts as if it returned restricted mappings. The following equivalent property is easier to state, and makes proving the arrow laws simple.

Definition 7.12 (mapping arrow law). Let $g : X \rightsquigarrow_{\text{map}} Y$. If there exists an $f : X \rightsquigarrow_{\perp} Y$ such that $g \equiv \text{lift}_{\text{map}} f$, then g obeys the **mapping arrow law**.

By homomorphism of lift_{map} , mapping arrow combinators preserve this law. It is therefore safe to assume that the mapping arrow law holds for all $g : X \rightsquigarrow_{\text{map}} Y$.

Theorem 7.13. lift_{map} is an arrow epimorphism.

Proof. Follows from Theorem 7.10 and restriction of $X \rightsquigarrow_{\text{map}} Y$ to instances for which the mapping arrow law (Definition 7.12) holds. □

Corollary 7.14. arr_{map} , $(\&\&\&_{\text{map}})$, (\ggg_{map}) , ifte_{map} and lazy_{map} define an arrow.

7.5 Lazy Preimage Mappings

On a computer, we do not often have the luxury of testing each function input to see whether it belongs to a preimage set. Even for finite domains, doing so is often intractable.

$$\begin{array}{ll}
X \xrightarrow{\text{pre}} Y ::= \langle \text{Set } Y, \text{Set } Y \Rightarrow \text{Set } X \rangle & \langle \cdot, \cdot \rangle_{\text{pre}} : (X \xrightarrow{\text{pre}} Y_1) \Rightarrow (X \xrightarrow{\text{pre}} Y_2) \Rightarrow (X \xrightarrow{\text{pre}} Y_1 \times Y_2) \\
\text{pre} : (X \xrightarrow{\text{map}} Y) \Rightarrow (X \xrightarrow{\text{pre}} Y) & \langle \langle Y'_1, p_1 \rangle, \langle Y'_2, p_2 \rangle \rangle_{\text{pre}} := \text{let } Y' := Y'_1 \times Y'_2 \\
& \quad p := \lambda B. \bigcup_{\langle b_1, b_2 \rangle \in B} p_1 \{b_1\} \cap p_2 \{b_2\} \\
\text{pre } g := \langle \text{range } g, \lambda B. \text{preimage } g \ B \rangle & \text{in } \langle Y', p \rangle \\
\\
\text{ap}_{\text{pre}} : (X \xrightarrow{\text{pre}} Y) \Rightarrow \text{Set } Y \Rightarrow \text{Set } X & (\circ_{\text{pre}}) : (Y \xrightarrow{\text{pre}} Z) \Rightarrow (X \xrightarrow{\text{pre}} Y) \Rightarrow (X \xrightarrow{\text{pre}} Z) \\
\text{ap}_{\text{pre}} \langle Y', p \rangle B := p (B \cap Y') & \langle Z', p_2 \rangle \circ_{\text{pre}} h_1 := \langle Z', \lambda C. \text{ap}_{\text{pre}} h_1 (p_2 C) \rangle \\
\\
\text{domain}_{\text{pre}} : (X \xrightarrow{\text{pre}} Y) \Rightarrow \text{Set } X & (\uplus_{\text{pre}}) : (X \xrightarrow{\text{pre}} Y) \Rightarrow (X \xrightarrow{\text{pre}} Y) \Rightarrow (X \xrightarrow{\text{pre}} Y) \\
\text{domain}_{\text{pre}} \langle Y', p \rangle := p Y' & h_1 \uplus_{\text{pre}} h_2 := \text{let } Y' := \text{range}_{\text{pre}} h_1 \cup \text{range}_{\text{pre}} h_2 \\
& \quad p := \lambda B. \text{ap}_{\text{pre}} h_1 B \uplus \text{ap}_{\text{pre}} h_2 B \\
\text{range}_{\text{pre}} : (X \xrightarrow{\text{pre}} Y) \Rightarrow \text{Set } Y & \text{in } \langle Y', p \rangle \\
\text{range}_{\text{pre}} \langle Y', p \rangle := Y' &
\end{array}$$

Figure 7.5: Lazy preimage mappings and operations.

If we wish to compute with infinite sets in the language implementation, we will need an abstraction that makes it easy to replace computation on points with computation on sets whose representations allow efficient operations. Therefore, in the preimage arrow, we confine computation on points to instances of

$$X \xrightarrow{\text{pre}} Y ::= \langle \text{Set } Y, \text{Set } Y \Rightarrow \text{Set } X \rangle \quad (7.31)$$

with the intention to replace $X \xrightarrow{\text{pre}} Y$ instances with an approximation further on. Like a mapping, an $X \xrightarrow{\text{pre}} Y$ has an observable domain—but computing the input-output pairs is delayed. We therefore call these *lazy preimage mappings*.

Converting a mapping to a lazy preimage mapping requires constructing a delayed application of `preimage`:

$$\begin{array}{l}
\text{pre} : (X \rightarrow Y) \Rightarrow (X \xrightarrow{\text{pre}} Y) \\
\text{pre } g := \langle \text{range } g, \lambda B. \text{preimage } g \ B \rangle
\end{array} \quad (7.32)$$

To apply a preimage mapping to some B , we intersect B with its range and apply the preimage

function:

$$\begin{aligned} \mathbf{ap}_{\text{pre}} : (X \xrightarrow{\text{pre}} Y) &\Rightarrow \text{Set } Y \Rightarrow \text{Set } X \\ \mathbf{ap}_{\text{pre}} \langle Y', p \rangle B &:= p (B \cap Y') \end{aligned} \tag{7.33}$$

Preimage arrow correctness depends on this fact: that using \mathbf{ap}_{pre} to compute preimages is the same as computing them from a mapping using $\mathbf{preimage}$.

Lemma 7.15. *Let $g : X \rightarrow Y$. For all $B \subseteq Y$ and Y' such that $\text{range } g \subseteq Y' \subseteq Y$, $\mathbf{preimage } g (B \cap Y') = \mathbf{preimage } g B$.*

Theorem 7.16 (\mathbf{ap}_{pre} computes preimages). *Let $g : X \rightarrow Y$. For all $B \subseteq Y$, $\mathbf{ap}_{\text{pre}} (\mathbf{pre } g) B = \mathbf{preimage } g B$.*

Proof. Expand definitions and apply Lemma 7.15 with $Y' = \text{range } g$. □

Figure 7.5 defines more operations on preimage mappings, including pairing, composition, and disjoint union operations corresponding to the mapping operations in Figure 7.3. To prove them correct, we need preimage mappings to be equivalent when they compute the same preimages.

Definition 7.17 (preimage mapping equivalence). *$h_1 : X \xrightarrow{\text{pre}} Y$ and $h_2 : X \xrightarrow{\text{pre}} Y$ are equivalent, or $h_1 \equiv h_2$, when $\mathbf{ap}_{\text{pre}} h_1 B \equiv \mathbf{ap}_{\text{pre}} h_2 B$ for all $B \subseteq Y$.*

Similarly to proving arrows correct, we prove the operations in Figure 7.5 are correct by proving that \mathbf{pre} is a homomorphism (though not an arrow homomorphism): it distributes over mapping operations to yield preimage mapping operations. The remainder of this section states these distributive properties as theorems and proves them. We will use these theorems to derive the preimage arrow from the mapping arrow.

7.5.1 Composition

To prove \mathbf{pre} distributes over mapping composition, we can make more or less direct use of the fact that $\mathbf{preimage}$ distributes over mapping composition.

Lemma 7.18 (preimage distributes over (\circ_{map})). *Let $g_1 : X \rightarrow Y$ and $g_2 : Y \rightarrow Z$. For all $C \subseteq Z$, $\text{preimage } (g_2 \circ_{\text{map}} g_1) C = \text{preimage } g_1 (\text{preimage } g_2 C)$.*

Theorem 7.19 (pre distributes over (\circ_{map})). *Let $g_1 : X \rightarrow Y$ and $g_2 : Y \rightarrow Z$. Then $\text{pre } (g_2 \circ_{\text{map}} g_1) \equiv (\text{pre } g_2) \circ_{\text{pre}} (\text{pre } g_1)$.*

Proof. Let $\langle Z', p_2 \rangle := \text{pre } g_2$ and $C \subseteq Z$. Starting from the right-hand side of the equivalence,

$$\begin{aligned}
 & \text{ap}_{\text{pre}} ((\text{pre } g_2) \circ_{\text{pre}} (\text{pre } g_1)) C && (7.34) \\
 & \equiv \text{let } p := \lambda C. \text{ap}_{\text{pre}} (\text{pre } g_1) (p_2 C) && \text{Def of } \text{ap}_{\text{pre}}, (\circ_{\text{pre}}) \\
 & \quad \text{in } p (C \cap Z') \\
 & \equiv \text{ap}_{\text{pre}} (\text{pre } g_1) (p_2 (C \cap Z')) && \text{Def of } p \\
 & \equiv \text{ap}_{\text{pre}} (\text{pre } g_1) (\text{ap}_{\text{pre}} (\text{pre } g_2) C) && \text{Rewrite with } \text{ap}_{\text{pre}} \\
 & \equiv \text{preimage } g_1 (\text{preimage } g_2 C) && \text{Theorem 7.16} \\
 & \equiv \text{preimage } (g_2 \circ_{\text{map}} g_1) C && \text{Lemma 7.18} \\
 & \equiv \text{ap}_{\text{pre}} (\text{pre } (g_2 \circ_{\text{map}} g_1)) C && \text{Theorem 7.16} \quad \square
 \end{aligned}$$

7.5.2 Pairing

We have less luck with pairing than with composition, because **preimage** does not distribute over pairing. Fortunately, **preimage** distributes over unions, and over pairing and cartesian product together.

Lemma 7.20 (preimage distributes over $\langle \cdot, \cdot \rangle_{\text{map}}$ and (\times)). *Let $g_1 : X \rightarrow Y_1$ and $g_2 : X \rightarrow Y_2$. For all $B_1 \subseteq Y_1$ and $B_2 \subseteq Y_2$, $\text{preimage } \langle g_1, g_2 \rangle_{\text{map}} (B_1 \times B_2) = (\text{preimage } g_1 B_1) \cap (\text{preimage } g_2 B_2)$.*

Lemma 7.21 (preimage distributes over union). *Let $g : X \rightarrow Y$ and $B : J \Rightarrow \text{Set } Y$ be an indexed collection of subsets of Y . Then*

$$\bigcup_{j \in J} \text{preimage } g (B j) = \text{preimage } g \bigcup_{j \in J} B j \quad (7.35)$$

Theorem 7.22 (pre distributes over $\langle \cdot, \cdot \rangle_{\text{map}}$). *Let $g_1 : X \rightarrow Y_1$ and $g_2 : X \rightarrow Y_2$. Then $\text{pre } \langle g_1, g_2 \rangle_{\text{map}} \equiv \langle \text{pre } g_1, \text{pre } g_2 \rangle_{\text{pre}}$.*

Proof. Let $\langle Y'_1, p_1 \rangle := \text{pre } g_1$, $\langle Y'_2, p_2 \rangle := \text{pre } g_2$ and $B \subseteq Y_1 \times Y_2$. Starting from the right-hand side of the equivalence,

$$\begin{aligned}
& \text{ap}_{\text{pre}} \langle \text{pre } g_1, \text{pre } g_2 \rangle_{\text{pre}} B && (7.36) \\
& \equiv \text{let } p := \lambda B. \bigcup_{\langle y_1, y_2 \rangle \in B} p_1 \{y_1\} \cap p_2 \{y_2\} && \text{Def of } \text{ap}_{\text{pre}}, \langle \cdot, \cdot \rangle_{\text{pre}} \\
& \quad \text{in } p (B \cap (Y'_1 \times Y'_2)) \\
& \equiv \bigcup_{\langle y_1, y_2 \rangle \in B \cap (Y'_1 \times Y'_2)} p_1 \{y_1\} \cap p_2 \{y_2\} && \text{Def of } p \\
& \equiv \bigcup_{\langle y_1, y_2 \rangle \in B \cap (Y'_1 \times Y'_2)} \text{preimage } g_1 \{y_1\} \cap \text{preimage } g_2 \{y_2\} && \text{Theorem 7.16} \\
& \equiv \bigcup_{\langle y_1, y_2 \rangle \in B \cap (Y'_1 \times Y'_2)} \text{preimage } \langle g_1, g_2 \rangle_{\text{map}} (\{y_1\} \times \{y_2\}) && \text{Lemma 7.20} \\
& \equiv \bigcup_{\langle y_1, y_2 \rangle \in B \cap (Y'_1 \times Y'_2)} \text{preimage } \langle g_1, g_2 \rangle_{\text{map}} \{\langle y_1, y_2 \rangle\} && \text{Def of } (\times) \\
& \equiv \text{preimage } \langle g_1, g_2 \rangle_{\text{map}} (B \cap (Y'_1 \times Y'_2)) && \text{Lemma 7.21} \\
& \equiv \text{preimage } \langle g_1, g_2 \rangle_{\text{map}} B && \text{Lemma 7.15} \\
& \equiv \text{ap}_{\text{pre}} (\text{pre } \langle g_1, g_2 \rangle_{\text{map}}) B && \text{Theorem 7.16}
\end{aligned}$$

We have an unmet proof obligation from using Lemma 7.15: that $\text{range } \langle g_1, g_2 \rangle_{\text{map}} \subseteq Y'_1 \times Y'_2$.

Let $b \in \text{range } \langle g_1, g_2 \rangle_{\text{map}}$. By definition of $\langle \cdot, \cdot \rangle_{\text{map}}$, there exists $a \in \text{domain } g_1 \cap \text{domain } g_2$ such that $b = \langle g_1 a, g_2 a \rangle$. Thus, $b \in Y'_1 \times Y'_2$ if and only if $g_1 a \in Y'_1$ and $g_2 a \in Y'_2$.

By definition of pre , $Y'_1 = \text{range } g_1$ and $Y'_2 = \text{range } g_2$. Because $a \in \text{domain } g_1$, $g_1 a \in \text{range } g_1 = Y'_1$. Because $a \in \text{domain } g_2$, $g_2 a \in \text{range } g_2 = Y'_2$. \square

7.5.3 Disjoint Union

Like proving pre distributes over composition, the proof that it distributes over disjoint union simply lifts a lemma about preimage to lazy preimage mappings.

Lemma 7.23 (preimage distributes over (\uplus_{map})). *Let $g_1 : X \rightarrow Y$ and $g_2 : X \rightarrow Y$ have disjoint domains. For all $B \subseteq Y$, $\text{preimage } (g_1 \uplus_{\text{map}} g_2) B = (\text{preimage } g_1 B) \uplus (\text{preimage } g_2 B)$.*

Theorem 7.24 (pre distributes over (\uplus_{map})). *Let $g_1 : X \rightarrow Y$ and $g_2 : X \rightarrow Y$ have disjoint domains. Then $\text{pre } (g_1 \uplus_{\text{map}} g_2) \equiv (\text{pre } g_1) \uplus_{\text{pre}} (\text{pre } g_2)$.*

Proof. Let $Y'_1 := \text{range } g_1$, $Y'_2 := \text{range } g_2$ and $B \subseteq Y$. Starting from the right-hand side of the equivalence,

$$\begin{aligned}
& \text{ap}_{\text{pre}} ((\text{pre } g_1) \uplus_{\text{pre}} (\text{pre } g_2)) B && (7.37) \\
& \equiv \text{let } p := \lambda B. \text{ap}_{\text{pre}} (\text{pre } g_1) B \uplus \text{ap}_{\text{pre}} (\text{pre } g_2) B && \text{Def of } \text{ap}_{\text{pre}}, (\uplus_{\text{pre}}) \\
& \quad \text{in } p (B \cap (Y'_1 \cup Y'_2)) \\
& \equiv \text{ap}_{\text{pre}} (\text{pre } g_1) (B \cap (Y'_1 \cup Y'_2)) \uplus \text{ap}_{\text{pre}} (\text{pre } g_2) (B \cap (Y'_1 \cup Y'_2)) && \text{Def of } p \\
& \equiv \text{preimage } g_1 (B \cap (Y'_1 \cup Y'_2)) \uplus \text{preimage } g_2 (B \cap (Y'_1 \cup Y'_2)) && \text{Theorem 7.16} \\
& \equiv \text{preimage } (g_1 \uplus_{\text{map}} g_2) (B \cap (Y'_1 \cup Y'_2)) && \text{Lemma 7.23} \\
& \equiv \text{preimage } (g_1 \uplus_{\text{map}} g_2) B && \text{Lemma 7.15} \\
& \equiv \text{ap}_{\text{pre}} (\text{pre } (g_1 \uplus_{\text{map}} g_2)) B && \text{Theorem 7.16}
\end{aligned}$$

We have an unmet proof obligation from using Lemma 7.15: that $\text{range } (g_1 \uplus_{\text{map}} g_2) \subseteq Y'_1 \cup Y'_2$.

Let $b \in \text{range } (g_1 \uplus_{\text{map}} g_2)$. By definition of (\uplus_{map}) , there exists $a \in \text{domain } g_1 \uplus \text{domain } g_2$ such that if $a \in \text{domain } g_1$ then $b = g_1 a$ so $b \in \text{range } g_1 = Y'_1$, and if $a \in \text{domain } g_2$ then $b = g_2 a$ so $b \in \text{range } g_2 = Y'_2$. Thus $b \in Y'_1 \cup Y'_2$. \square

7.6 Deriving the Preimage Arrow

Now we can define an arrow that runs expressions backwards on sets of outputs. Its computations should produce preimage mappings or be preimage mappings.

As with the mapping arrow and mappings, we cannot have $X \xrightarrow[\text{pre}]{\sim} Y ::= X \xrightarrow[\text{pre}} Y$: we run into trouble trying to define arr_{pre} because a preimage mapping needs an observable range. To get one, it is easiest to parameterize preimage computations on a **Set** X ; therefore the

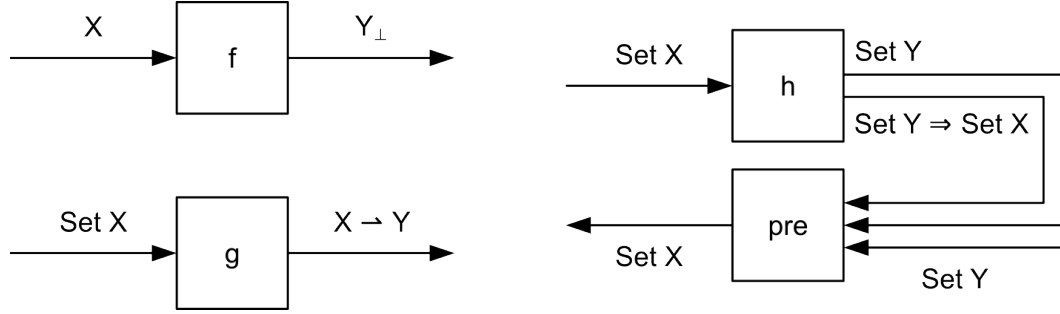


Figure 7.6: Comparison of arrows used as target categories. Computations $f : X \rightsquigarrow_{\perp} Y$ may return an error value \perp . Computations $g : X \rightsquigarrow_{\text{map}} Y$ produce partial mappings on a given $A \subseteq X$, leaving out inputs for which f would return an error. Computations $h : X \rightsquigarrow_{\text{pre}} Y$ produce lazy preimage mappings; i.e. $h \ A$ computes preimages under g .

preimage arrow type constructor is

$$X \rightsquigarrow_{\text{pre}} Y ::= \text{Set } X \Rightarrow (X \xrightarrow{\text{pre}} Y) \quad (7.38)$$

or $\text{Set } X \Rightarrow \langle \text{Set } Y, \text{Set } Y \Rightarrow \text{Set } X \rangle$. To deconstruct the type, a preimage arrow computation computes a range first, and returns the range and a lambda that computes preimages. Figure 7.6 illustrates this as a circuit diagram.

To use Theorem 7.6, we need to define correctness using a lift from the mapping arrow to the preimage arrow. A simple candidate with the right type is

$$\begin{aligned} \text{lift}_{\text{pre}} &: (X \rightsquigarrow_{\text{map}} Y) \Rightarrow (X \rightsquigarrow_{\text{pre}} Y) \\ \text{lift}_{\text{pre}} \ g \ A &:= \text{pre} \ (g \ A) \end{aligned} \quad (7.39)$$

By definition of lift_{pre} and Theorem 7.16, for all $g : X \rightsquigarrow_{\text{map}} Y$, and $A \subseteq X$ and $B \subseteq Y$,

$$\begin{aligned} \text{ap}_{\text{pre}} (\text{lift}_{\text{pre}} \ g \ A) \ B &\equiv \text{ap}_{\text{pre}} (\text{pre} \ (g \ A)) \ B \\ &\equiv \text{preimage} \ (g \ A) \ B \end{aligned} \quad (7.40)$$

Thus, lifted mapping arrow computations correctly compute preimages under restricted mappings, exactly as we should expect them to.

To derive the preimage arrow’s combinators in a way that makes lift_{pre} a homomorphism, we need preimage arrow equivalence to mean “computes the same preimages.”

$$\begin{array}{ll}
X_{\text{pre}} \rightsquigarrow Y ::= \text{Set } X \Rightarrow (X \rightrightarrows_{\text{pre}} Y) & \text{ifte}_{\text{pre}} : (X_{\text{pre}} \rightsquigarrow \text{Bool}) \Rightarrow (X_{\text{pre}} \rightsquigarrow Y) \Rightarrow (X_{\text{pre}} \rightsquigarrow Y) \Rightarrow (X_{\text{pre}} \rightsquigarrow Y) \\
\text{arr}_{\text{pre}} : (X \Rightarrow Y) \Rightarrow (X_{\text{pre}} \rightsquigarrow Y) & \text{ifte}_{\text{pre}} \ h_1 \ h_2 \ h_3 \ A := \text{let } h'_1 := h_1 \ A \\
\text{arr}_{\text{pre}} := \text{lift}_{\text{pre}} \circ \text{arr}_{\text{map}} & \quad h'_2 := h_2 \ (\text{ap}_{\text{pre}} \ h'_1 \ \{\text{true}\}) \\
& \quad h'_3 := h_3 \ (\text{ap}_{\text{pre}} \ h'_1 \ \{\text{false}\}) \\
& \quad \text{in } h'_2 \ \text{ap}_{\text{pre}} \ h'_3 \\
(\\>>>_{\text{pre}}) : (X_{\text{pre}} \rightsquigarrow Y) \Rightarrow (Y_{\text{pre}} \rightsquigarrow Z) \Rightarrow (X_{\text{pre}} \rightsquigarrow Z) & \text{lazy}_{\text{pre}} : (1 \Rightarrow (X_{\text{pre}} \rightsquigarrow Y)) \Rightarrow (X_{\text{pre}} \rightsquigarrow Y) \\
(h_1 \\>>>_{\text{pre}} h_2) \ A := \text{let } h'_1 := h_1 \ A & \text{lazy}_{\text{pre}} \ h \ A := \text{if } (A = \emptyset) \ (\text{pre } \emptyset) \ (h \ 0 \ A) \\
\quad h'_2 := h_2 \ (\text{range}_{\text{pre}} \ h'_1) & \\
\quad \text{in } h'_2 \circ_{\text{pre}} h'_1 & \\
(\\&&&_{\text{pre}}) : (X_{\text{pre}} \rightsquigarrow Y) \Rightarrow (X_{\text{pre}} \rightsquigarrow Z) \Rightarrow (X_{\text{pre}} \rightsquigarrow Y \times Z) & \text{lift}_{\text{pre}} : (X_{\text{map}} \rightsquigarrow Y) \Rightarrow (X_{\text{pre}} \rightsquigarrow Y) \\
(h_1 \\&&&_{\text{pre}} h_2) \ A := \langle h_1 \ A, h_2 \ A \rangle_{\text{pre}} & \text{lift}_{\text{pre}} \ g \ A := \text{pre } (g \ A)
\end{array}$$

Figure 7.7: Preimage arrow definitions.

Definition 7.25 (preimage arrow equivalence). *Two computations $h_1 : X_{\text{pre}} \rightsquigarrow Y$ and $h_2 : X_{\text{pre}} \rightsquigarrow Y$ are equivalent, or $h_1 \equiv h_2$, when $h_1 \ A \equiv h_2 \ A$ for all $A \subseteq X$.*

As with arr_{map} , defining arr_{pre} as a composition meets (7.8). The remainder of this section derives $(\\&&&_{\text{pre}})$, $(\\>>>_{\text{pre}})$, ifte_{pre} and lazy_{pre} from mapping arrow combinators, in a way that ensures lift_{pre} is an arrow homomorphism from the mapping arrow to the preimage arrow. Figure 7.7 contains the resulting definitions.

7.6.1 Composition

Starting with the left-hand side of (7.9),

$$\begin{aligned}
& \text{ap}_{\text{pre}} (\text{lift}_{\text{pre}} (g_1 \\>>>_{\text{map}} g_2) \ A) \ C & (7.41) \\
& \equiv \text{let } g'_1 := g_1 \ A & \text{Def of } \text{lift}_{\text{pre}}, (\\>>>_{\text{map}}) \\
& \quad g'_2 := g_2 \ (\text{range } g'_1) \\
& \quad \text{in } \text{ap}_{\text{pre}} (\text{pre } (g'_2 \circ_{\text{map}} g'_1)) \ C \\
& \equiv \text{let } g'_1 := g_1 \ A & \text{Theorem 7.19} \\
& \quad g'_2 := g_2 \ (\text{range } g'_1) \\
& \quad \text{in } \text{ap}_{\text{pre}} ((\text{pre } g'_1) \circ_{\text{pre}} (\text{pre } g'_2)) \ C \\
& \equiv \text{let } h_1 := \text{lift}_{\text{pre}} g_1 \ A & \text{Rewrite with } \text{lift}_{\text{pre}} \\
& \quad h_2 := \text{lift}_{\text{pre}} g_2 \ (\text{range}_{\text{pre}} h_1) \\
& \quad \text{in } \text{ap}_{\text{pre}} (h_2 \circ_{\text{pre}} h_1) \ C
\end{aligned}$$

Substituting h_1 for $\text{lift}_{\text{pre}} g_1$ and h_2 for $\text{lift}_{\text{pre}} g_2$, and removing the application of ap_{pre} from both sides of the equivalence gives a definition of (\ggg_{pre}) (Figure 7.7) for which (7.9) holds.

7.6.2 Pairing

Starting with the left-hand side of (7.10),

$$\begin{aligned}
& \text{ap}_{\text{pre}} (\text{lift}_{\text{pre}} (g_1 \&\&_{\text{map}} g_2) A) B && (7.42) \\
& \equiv \text{ap}_{\text{pre}} (\text{pre } \langle g_1 A, g_2 A \rangle_{\text{map}}) B && \text{Def of } \text{lift}_{\text{pre}}, (\&\&_{\text{map}}) \\
& \equiv \text{ap}_{\text{pre}} \langle \text{pre } (g_1 A), \text{pre } (g_2 A) \rangle_{\text{pre}} B && \text{Theorem 7.22} \\
& \equiv \text{ap}_{\text{pre}} \langle \text{lift}_{\text{pre}} g_1 A, \text{lift}_{\text{pre}} g_2 A \rangle_{\text{pre}} B && \text{Rewrite with } \text{lift}_{\text{pre}}
\end{aligned}$$

Substituting h_1 for $\text{lift}_{\text{pre}} g_1$ and h_2 for $\text{lift}_{\text{pre}} g_2$, and removing the application of ap_{pre} from both sides of the equivalence gives a definition of $(\&\&_{\text{pre}})$ (Figure 7.7) for which (7.10) holds.

7.6.3 Conditional

Starting with the left-hand side of (7.11),

$$\begin{aligned}
& \text{ap}_{\text{pre}} (\text{lift}_{\text{pre}} (\text{ifte}_{\text{map}} g_1 g_2 g_3) A) B && (7.43) \\
& \equiv \text{let } g'_1 := g_1 A && \text{Def of } \text{lift}_{\text{pre}}, \text{ifte}_{\text{map}} \\
& \quad g'_2 := g_2 (\text{preimage } g'_1 \{ \text{true} \}) \\
& \quad g'_3 := g_3 (\text{preimage } g'_1 \{ \text{false} \}) \\
& \quad \text{in } \text{ap}_{\text{pre}} (\text{pre } (g'_2 \uplus_{\text{map}} g'_3)) B \\
& \equiv \text{let } g'_1 := g_1 A && \text{Theorem 7.24} \\
& \quad g'_2 := g_2 (\text{preimage } g'_1 \{ \text{true} \}) \\
& \quad g'_3 := g_3 (\text{preimage } g'_1 \{ \text{false} \}) \\
& \quad \text{in } \text{ap}_{\text{pre}} ((\text{pre } g'_2) \uplus_{\text{pre}} (\text{pre } g'_3)) B \\
& \equiv \text{let } g'_1 := g_1 A && \text{Theorem 7.16} \\
& \quad g'_2 := g_2 (\text{ap}_{\text{pre}} (\text{pre } g'_1) \{ \text{true} \}) \\
& \quad g'_3 := g_3 (\text{ap}_{\text{pre}} (\text{pre } g'_1) \{ \text{false} \}) \\
& \quad \text{in } \text{ap}_{\text{pre}} ((\text{pre } g'_2) \uplus_{\text{pre}} (\text{pre } g'_3)) B
\end{aligned}$$

$$\begin{aligned}
&\equiv \text{let } h_1 := \text{lift}_{\text{pre}} g_1 A && \text{Rewrite with } \text{lift}_{\text{pre}} \\
&\quad h_2 := \text{lift}_{\text{pre}} g_2 (\text{ap}_{\text{pre}} h_1 \{\text{true}\}) \\
&\quad h_3 := \text{lift}_{\text{pre}} g_3 (\text{ap}_{\text{pre}} h_1 \{\text{false}\}) \\
&\text{in } \text{ap}_{\text{pre}} (h_2 \uplus_{\text{pre}} h_3) B
\end{aligned}$$

Substituting h_1 , h_2 and h_3 for $\text{lift}_{\text{pre}} g_1$, $\text{lift}_{\text{pre}} g_2$ and $\text{lift}_{\text{pre}} g_3$, and removing the application of ap_{pre} from both sides of the equivalence gives a definition of ifte_{pre} (Figure 7.7) for which (7.11) holds.

7.6.4 Laziness

Starting with the left-hand side of (7.12),

$$\begin{aligned}
&\text{ap}_{\text{pre}} (\text{lift}_{\text{pre}} (\text{lazy}_{\text{map}} g) A) B && (7.44) \\
&\equiv \text{let } g' := \text{if } (A = \emptyset) \emptyset (g \ 0 \ A) && \text{Def of } \text{lift}_{\text{pre}}, \text{lazy}_{\text{map}} \\
&\quad \text{in } \text{ap}_{\text{pre}} (\text{pre } g') B \\
&\equiv \text{let } h := \text{if } (A = \emptyset) (\text{pre } \emptyset) (\text{pre } (g \ 0 \ A)) && \text{Dist pre over if} \\
&\quad \text{in } \text{ap}_{\text{pre}} h B \\
&\equiv \text{let } h := \text{if } (A = \emptyset) (\text{pre } \emptyset) (\text{lift}_{\text{pre}} (g \ 0) A) && \text{Rewrite with } \text{lift}_{\text{pre}} \\
&\quad \text{in } \text{ap}_{\text{pre}} h B
\end{aligned}$$

Substituting $h \ 0$ for $\text{lift}_{\text{pre}} (g \ 0)$ and removing the application of ap_{pre} from both sides of the equivalence gives a definition for lazy_{pre} (Figure 7.7) for which (7.12) holds.

7.6.5 Correctness

Theorem 7.26 (preimage arrow correctness). *lift_{pre} is a homomorphism.*

Proof. By construction. □

Corollary 7.27 (semantic correctness). *For all e , $\llbracket e \rrbracket_{\text{pre}} \equiv \text{lift}_{\text{pre}} \llbracket e \rrbracket_{\text{map}}$.*

As with the mapping arrow, preimage arrow computations can be unruly. We would like to assume that each $h : X \xrightarrow[\text{pre}]{} Y$ acts as if it computes preimages under restricted mappings. The following equivalent property is easier to state, and makes proving the arrow laws simple.

Definition 7.28 (preimage arrow law). *Let $h : X \rightsquigarrow_{\text{pre}} Y$. If there exists a $g : X \rightsquigarrow_{\text{map}} Y$ such that $h \equiv \text{lift}_{\text{pre}} g$, then h obeys the **preimage arrow law**.*

By homomorphism of lift_{pre} , preimage arrow combinators preserve this law. It is therefore safe to assume that the preimage arrow law holds for all $h : X \rightsquigarrow_{\text{pre}} Y$.

Theorem 7.29. lift_{pre} is an arrow epimorphism.

Proof. Follows from Theorem 7.26 and restriction of $X \rightsquigarrow_{\text{pre}} Y$ to instances for which the preimage arrow law (Definition 7.28) holds. \square

Corollary 7.30. arr_{pre} , $(\&\&\&_{\text{pre}})$, $(\>\>\>_{\text{pre}})$, ifte_{pre} and lazy_{pre} define an arrow.

7.7 Preimages Under Partial, Probabilistic Functions

We have defined everything on the top of our roadmap:

$$\begin{array}{ccccc}
 X \rightsquigarrow_{\perp} Y & \xrightarrow{\text{lift}_{\text{map}}} & X \rightsquigarrow_{\text{map}} Y & \xrightarrow{\text{lift}_{\text{pre}}} & X \rightsquigarrow_{\text{pre}} Y \\
 \eta_{\perp*} \downarrow & & \downarrow \eta_{\text{map}*} & & \downarrow \eta_{\text{pre}*} \\
 X \rightsquigarrow_{\perp*} Y & \xrightarrow{\text{lift}_{\text{map}*}} & X \rightsquigarrow_{\text{map}*} Y & \xrightarrow{\text{lift}_{\text{pre}*}} & X \rightsquigarrow_{\text{pre}*} Y
 \end{array} \tag{7.45}$$

and proved that lift_{map} and lift_{pre} are homomorphisms. At this point, we can interpret an expression e in three ways using the same semantic function for first-order programs:

1. As $\llbracket e \rrbracket_{\perp} : X \rightsquigarrow_{\perp} Y$, an intensional function that may raise errors.
2. As $\llbracket e \rrbracket_{\text{map}} : X \rightsquigarrow_{\text{map}} Y$, which produces mappings, or extensional functions, on a restricted domain (correct by homomorphism of lift_{map}).
3. As $\llbracket e \rrbracket_{\text{pre}} : X \rightsquigarrow_{\text{pre}} Y$, which computes preimages under mappings produced by $\llbracket e \rrbracket_{\text{map}}$ (correct by homomorphism of lift_{pre}).

These interpretations have two shortcomings:

1. They do not pass an implicit random source through e 's subexpressions.
2. Using them requires knowing the set of inputs on which e terminates. If $\llbracket e \rrbracket_{\perp}$ does not terminate on just one input in $A \subseteq X$, neither $\llbracket e \rrbracket_{\text{map}} A$ nor $\llbracket e \rrbracket_{\text{pre}} A$ terminates.

In this section, we define the arrows on the bottom of the roadmap (7.45) by transforming the arrows on the top into arrows that pass an implicit random source and always terminate. Their correctness again comes down to proving that the combinators between them are homomorphisms, though guaranteed termination needs special treatment.

7.7.1 Motivation

Probabilistic functions that may not terminate, but do so with probability 1, are common. For example, suppose `random` retrieves numbers in $[0, 1]$ from an implicit random source. The following probabilistic function defines the well-known geometric distribution by counting the number of times `random` $< p$:

$$\text{geometric } p := \text{if } (\text{random} < p) \ 0 \ (1 + \text{geometric } p) \quad (7.46)$$

For any $p > 0$, `geometric p` may not terminate, but the probability of never taking the “else” branch is $(1 - p) \cdot (1 - p) \cdot (1 - p) \cdots = 0$. Thus, `geometric p` terminates with probability 1.

Suppose we interpret `geometric p` as $h : \Omega \rightsquigarrow_{\text{pre}} \mathbb{N}$, a preimage arrow computation from random sources $\omega \in \Omega$ to naturals, and we have a probability measure $P : \text{Set } \Omega \rightarrow [0, 1]$. The probability of $N \subseteq \mathbb{N}$ is then $P(\text{ap}_{\text{pre}}(h \ \Omega) \ N)$. To compute this, we must

- Ensure $\text{ap}_{\text{pre}}(h \ \Omega) \ N$ terminates.
- Ensure each $\omega \in \Omega$ contains enough random numbers.
- Determine how `random` indexes numbers in ω .

Ensuring $\text{ap}_{\text{pre}}(h \ \Omega) \ N$ terminates is the most difficult, but doing the other two will provide structure that makes it much easier.

7.7.2 Threading and Indexing

To ensure random sources contain enough numbers, they should be infinite.

Typically, to thread a random source $\omega \in \Omega$ through computations, ω is made an infinite stream. Each computation receives and returns an ω . The interpretation of `random`

as a computation takes ω 's head and returns its tail. Combinators pass ω unchanged to one subcomputation, and pass the resulting ω' unchanged to the next. This is typically done with a monad, and it imposes a total order on evaluation.

A little-used alternative that imposes only a partial order makes ω an infinite binary tree. Each computation receives an ω but does not return one. The interpretation of **random** as a computation simply returns ω 's root value. Combinators ignore the root, split ω into a left subtree ω_{left} and a right subtree ω_{right} , and pass each to their subcomputations.

Arrows can thread a stream or a tree in the same manner, but the resulting combinators have large definitions, and are conceptually difficult and hard to manipulate. Fortunately, it is relatively easy to assign each arrow computation a unique index into a tree-shaped random source and pass the random source unchanged. To do this, we need an indexing scheme.

Definition 7.31 (binary indexing scheme). *Let J be an index set, $j_0 \in J$ a distinguished element, and $\text{left} : J \Rightarrow J$ and $\text{right} : J \Rightarrow J$ be total, injective functions. If for all $j \in J$, $j = \text{next } j_0$ for some finite composition next of left and right , then J , j_0 , left and right define a **binary indexing scheme**.*

For example, let J be the set of lists of $\{0, 1\}$, $j_0 := \langle \rangle$, and $\text{left } j := \langle 0, j \rangle$ and $\text{right } j := \langle 1, j \rangle$. Alternatively, let J be the set of dyadic rationals in $(0, 1)$ (i.e. those with power-of-two denominators), $j_0 := \frac{1}{2}$ and

$$\begin{aligned} \text{left } (p/q) &:= (p - \tfrac{1}{2})/q \\ \text{right } (p/q) &:= (p + \tfrac{1}{2})/q \end{aligned} \tag{7.47}$$

With this alternative, left-to-right evaluation order can be made to correspond with the natural order ($<$) over J .

In any case, J is countable, and can be thought of as a set of indexes into an infinite binary tree. Values of type $J \rightarrow A$ encode such trees of values in A as total mappings (i.e. infinite vectors).

$x \rightsquigarrow_{a^*} y ::= \text{AStore } s \ (x \rightsquigarrow_a y) ::= J \Rightarrow (\langle s, x \rangle \rightsquigarrow_a y)$	$\text{ifte}_{a^*} : (x \rightsquigarrow_{a^*} \text{Bool}) \Rightarrow (x \rightsquigarrow_{a^*} y) \Rightarrow (x \rightsquigarrow_{a^*} y) \Rightarrow (x \rightsquigarrow_{a^*} y)$
$\text{arr}_{a^*} : (x \Rightarrow y) \Rightarrow (x \rightsquigarrow_{a^*} y)$	$\text{ifte}_{a^*} \ k_1 \ k_2 \ k_3 \ j := \text{ifte}_a \ (k_1 \ (\text{left } j))$
$\text{arr}_{a^*} := \eta_{a^*} \circ \text{arr}_a$	$\quad (k_2 \ (\text{left } (\text{right } j)))$
	$\quad (k_3 \ (\text{right } (\text{right } j)))$
$(\ggg_{a^*}) : (x \rightsquigarrow_{a^*} y) \Rightarrow (y \rightsquigarrow_{a^*} z) \Rightarrow (x \rightsquigarrow_{a^*} z)$	$\text{lazy}_{a^*} : (1 \Rightarrow (x \rightsquigarrow_{a^*} y)) \Rightarrow (x \rightsquigarrow_{a^*} y)$
$(k_1 \ggg_{a^*} k_2) j :=$	$\text{lazy}_{a^*} \ k \ j := \text{lazy}_a \ \lambda 0. k \ 0 \ j$
$\quad (\text{arr}_a \ \text{fst} \ \&\&\&_a \ k_1 \ (\text{left } j)) \ggg_a k_2 \ (\text{right } j)$	
$(\&\&\&_{a^*}) : (x \rightsquigarrow_{a^*} y_1) \Rightarrow (x \rightsquigarrow_{a^*} y_2) \Rightarrow (x \rightsquigarrow_{a^*} \langle y_1, y_2 \rangle)$	$\eta_{a^*} : (x \rightsquigarrow_a y) \Rightarrow (x \rightsquigarrow_{a^*} y)$
$(k_1 \ \&\&\&_{a^*} k_2) j := k_1 \ (\text{left } j) \ \&\&\&_a k_2 \ (\text{right } j)$	$\eta_{a^*} \ f \ j := \text{arr}_a \ \text{snd} \ \ggg_a \ f$

Figure 7.8: AStore (associative store) arrow transformer definitions.

7.7.3 Applicative, Associative Store Transformer

We thread infinite binary trees through bottom, mapping, and preimage arrow computations by defining an **arrow transformer**: a type constructor that receives and produces an arrow type, and combinators for arrows of the produced type. The applicative store arrow transformer's type constructor takes a store type s and an arrow type $x \rightsquigarrow_a y$:

$$\text{AStore } s \ (x \rightsquigarrow_a y) ::= J \Rightarrow (\langle s, x \rangle \rightsquigarrow_a y) \quad (7.48)$$

Reading the type, we see that computations receive an index $j \in J$ and produce a computation that receives a store as well as an x . Lifting extracts the x from the input pair and sends it on to the original computation, ignoring j :

$$\begin{aligned} \eta_{a^*} : (x \rightsquigarrow_a y) &\Rightarrow \text{AStore } s \ (x \rightsquigarrow_a y) \\ \eta_{a^*} \ f \ j &:= \text{arr}_a \ \text{snd} \ \ggg_a \ f \end{aligned} \quad (7.49)$$

Figure 7.8 defines the remaining combinators. Each subcomputation receives **left** j , **right** j , or some other unique binary index. We thus think of programs interpreted as AStore arrows as being completely unrolled into an infinite binary tree, with each subcomputation labeled with its tree index.

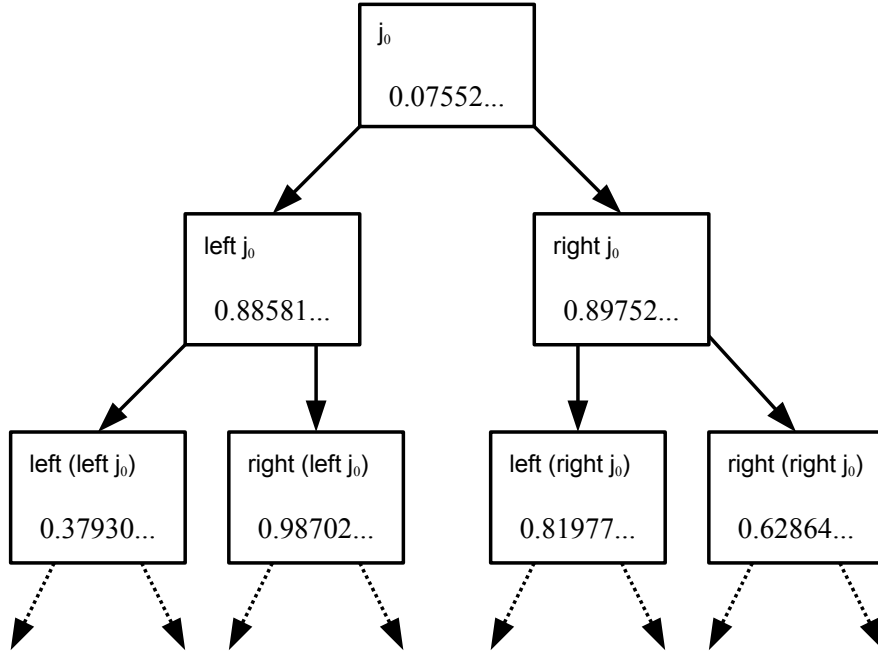


Figure 7.9: An $\omega \in \Omega$ is an infinite binary tree of random values encoded as a total mapping from tree indexes in J to real numbers in $[0, 1]$.

7.7.4 Partial, Probabilistic Programs

To interpret probabilistic programs, we put an infinite random tree in the store.

Definition 7.32 (random source). *Let $\Omega := J \rightarrow [0, 1]$. A **random source** is any infinite binary tree $\omega \in \Omega$.*

Figure 7.9 illustrates a single $\omega \in \Omega$.

To interpret partial programs, we need to ensure termination. One ultimately implementable way is to have the store dictate which branch of each conditional, if any, is taken.

Definition 7.33 (branch trace). *A **branch trace** is any $t : J \rightarrow \text{Bool}_\perp$ such that $t\ j = \text{true}$ or $t\ j = \text{false}$ for no more than finitely many $j \in J$.*

Let $T \subset J \rightarrow \text{Bool}_\perp$ be the largest set of branch traces.

Now $X \rightsquigarrow_{a^*} Y ::= \text{AStore } (\Omega \times T) (X \rightsquigarrow_a Y)$ is an **AStore** arrow type whose computations thread both random stores and branch traces.

For probabilistic programs, we define a combinator random_{a^*} that returns the number at its tree index in the random source, and extend $\llbracket \cdot \rrbracket_{a^*}$ for arrows a^* for which random_{a^*} is defined:

$$\begin{aligned} \text{random}_{a^*} &: X \rightsquigarrow_{a^*} [0, 1] \\ \text{random}_{a^*} j &:= \text{arr}_a (\text{fst} \ggg \text{fst} \ggg \pi j) \\ \llbracket \text{random} \rrbracket_{a^*} &:\equiv \text{random}_{a^*} \end{aligned} \tag{7.50}$$

Here, πj projects its argument onto the argument's j th coordinate:

$$\begin{aligned} \pi : J \Rightarrow (J \rightarrow X) \Rightarrow X \\ \pi j f &:= f j \end{aligned} \tag{7.51}$$

So πj is analogous to **fst** and **snd** for pairs, but for vectors at index j .

For partial programs, we define a combinator that reads branch traces, and an if-then-else combinator that ensures its test expression agrees with the trace:

$$\begin{aligned} \text{branch}_{a^*} &: X \rightsquigarrow_{a^*} \text{Bool} \\ \text{branch}_{a^*} j &:= \text{arr}_a (\text{fst} \ggg \text{snd} \ggg \pi j) \\ \text{ifte}_{a^*}^\Downarrow : (x \rightsquigarrow_{a^*} \text{Bool}) \Rightarrow (x \rightsquigarrow_{a^*} y) \Rightarrow (x \rightsquigarrow_{a^*} y) \Rightarrow (x \rightsquigarrow_{a^*} y) \\ \text{ifte}_{a^*}^\Downarrow k_1 k_2 k_3 j &:= \text{ifte}_a ((k_1 (\text{left } j) \&\& \text{branch}_{a^*} j) \ggg_a \text{arr}_a \text{agrees}) \\ &\quad (k_2 (\text{left } (\text{right } j))) \\ &\quad (k_3 (\text{right } (\text{right } j))) \end{aligned} \tag{7.52}$$

where $\text{agrees } \langle b_1, b_2 \rangle := \text{if } (b_1 = b_2) b_1 \perp$. Thus, if the branch trace does not agree with the test expression, it returns an error. We define a new semantic function $\llbracket \cdot \rrbracket_{a^*}^\Downarrow$ by replacing the if rule in $\llbracket \cdot \rrbracket_{a^*}$:

$$\llbracket \text{if } e_c \ e_t \ e_f \rrbracket_{a^*}^\Downarrow :\equiv \text{ifte}_{a^*}^\Downarrow \llbracket e_c \rrbracket_{a^*}^\Downarrow \llbracket \text{lazy } e_t \rrbracket_{a^*}^\Downarrow \llbracket \text{lazy } e_f \rrbracket_{a^*}^\Downarrow \tag{7.53}$$

For an **AStore** computation k , we obviously must run k on every branch trace in T and filter out \perp , or somehow find inputs $\langle \omega, t \rangle, a$ for which **agrees** never returns \perp . Preimage **AStore** arrows do the former by first computing an image, and the latter by computing preimages of sets that cannot contain \perp .

Definition 7.34 (terminating, probabilistic arrows). *Define*

$$\begin{aligned} X \rightsquigarrow_{\perp}^* Y &::= \text{AStore } (\Omega \times T) (X \rightsquigarrow_{\perp} Y) \\ X \rightsquigarrow_{\text{map}}^* Y &::= \text{AStore } (\Omega \times T) (X \rightsquigarrow_{\text{map}} Y) \\ X \rightsquigarrow_{\text{pre}}^* Y &::= \text{AStore } (\Omega \times T) (X \rightsquigarrow_{\text{pre}} Y) \end{aligned} \tag{7.54}$$

as the type constructors for the **bottom***, **mapping*** and **preimage*** arrows.

7.7.5 Correctness

We have two arrow lifts to prove homomorphic: one from pure computations to effectful (i.e. from those that do not access the store to those that do), and one from effectful computations to effectful. For both, we need **AStore** arrow equivalence to be more extensional.

Definition 7.35 (**AStore** arrow equivalence). *Two **AStore** arrow computations k_1 and k_2 are equivalent, or $k_1 \equiv k_2$, when $k_1 j \equiv k_2 j$ for all $j \in J$.*

Pure Expressions

Proving η_{a^*} is a homomorphism proves $\llbracket \cdot \rrbracket_{a^*}$ correctly interprets pure expressions. Because **AStore** accepts any arrow type $x \rightsquigarrow_a y$, we can do so using only the arrow laws. From here on, we assume every **AStore** arrow's base type's combinators obey the arrow laws listed in Section 7.2.1.

Theorem 7.36 (pure **AStore** arrow correctness). *η_{a^*} is a homomorphism.*

Proof. Defining arr_{a^*} as a composition clearly meets the first homomorphism law (7.8). For homomorphism laws (7.9)–(7.11), start from the right side, expand definitions, and use arrow laws (7.14)–(7.16) to factor out $\text{arr}_a \text{ snd}$.

For (7.12), additionally β -reduce within the outer thunk, then use the lazy distributive law (7.17) to extract $\text{arr}_a \text{snd}$. \square

Corollary 7.37 (pure semantic correctness). *For all pure e , $\llbracket e \rrbracket_{a^*} \equiv \eta_{a^*} \llbracket e \rrbracket_a$.*

Effectful Expressions

To prove all interpretations of effectful expressions correct, we need a lift between AStore arrows. Let $x \rightsquigarrow_{a^*} y ::= \text{AStore } s \ (x \rightsquigarrow_a y)$ and $x \rightsquigarrow_{b^*} y ::= \text{AStore } s \ (x \rightsquigarrow_b y)$. Define

$$\begin{aligned} \text{lift}_{b^*} : (x \rightsquigarrow_{a^*} y) &\Rightarrow (x \rightsquigarrow_{b^*} y) \\ \text{lift}_{b^*} f \ j &:= \text{lift}_b (f \ j) \end{aligned} \tag{7.55}$$

where $\text{lift}_b : (x \rightsquigarrow_a y) \Rightarrow (x \rightsquigarrow_b y)$. A commutative diagram shows the relationships more clearly:

$$\begin{array}{ccc} x \rightsquigarrow_a y & \xrightarrow{\text{lift}_b} & x \rightsquigarrow_b y \\ \eta_{a^*} \downarrow & & \downarrow \eta_{b^*} \\ x \rightsquigarrow_{a^*} y & \xrightarrow{\text{lift}_{b^*}} & x \rightsquigarrow_{b^*} y \end{array} \tag{7.56}$$

At minimum, we should expect to produce equivalent $x \rightsquigarrow_{b^*} y$ computations from $x \rightsquigarrow_a y$ computations whether a lift or an η is done first.

Theorem 7.38 (natural transformation). *If lift_b is an arrow homomorphism, then (7.56) commutes.*

Proof. Expand definitions and apply homomorphism laws (7.9) and (7.8) for lift_b :

$$\begin{aligned} \text{lift}_{b^*} (\eta_{a^*} f) &\equiv \lambda j. \text{lift}_b (\text{arr}_a \text{snd} \ggg_a f) \\ &\equiv \lambda j. \text{lift}_b (\text{arr}_a \text{snd}) \ggg_b \text{lift}_b f \\ &\equiv \lambda j. \text{arr}_b \text{snd} \ggg_b \text{lift}_b f \\ &\equiv \eta_{b^*} (\text{lift}_b f) \end{aligned} \tag{7.57} \quad \square$$

Theorem 7.39 (effectful AStore arrow correctness). *If lift_b is an arrow homomorphism from a to b , then lift_{b^*} is an arrow homomorphism from a^* to b^* .*

Proof. For each homomorphism property (7.8)–(7.12), expand the definitions of lift_{b^*} and the combinator, distribute lift_b , rewrite in terms of lift_{b^*} , and rewrite using the definition of the combinator. For example, for distribution over pairing:

$$\begin{aligned}
\text{lift}_{b^*} (k_1 \&\&\&_{a^*} k_2) j &\equiv \text{lift}_b ((k_1 \&\&\&_{a^*} k_2) j) \\
&\equiv \text{lift}_b (k_1 (\text{left } j) \&\&\&_a k_2 (\text{right } j)) \\
&\equiv \text{lift}_b (k_1 (\text{left } j)) \&\&\&_b \text{lift}_b (k_2 (\text{right } j)) \\
&\equiv (\text{lift}_{b^*} k_1) (\text{left } j) \&\&\&_b (\text{lift}_{b^*} k_2) (\text{right } j) \\
&\equiv (\text{lift}_{b^*} k_1 \&\&\&_{b^*} \text{lift}_{b^*} k_2) j
\end{aligned} \tag{7.58}$$

The remaining properties are similar, though distributing lift_{b^*} over lazy_{a^*} requires defining an extra thunk in the last step. \square

Corollary 7.40 (effectful semantic correctness). *If lift_b is an arrow homomorphism, then for all expressions e , $\llbracket e \rrbracket_{b^*} \equiv \text{lift}_{b^*} \llbracket e \rrbracket_{a^*}$ and $\llbracket e \rrbracket_{b^*}^\Downarrow \equiv \text{lift}_{b^*} \llbracket e \rrbracket_{a^*}^\Downarrow$.*

Corollary 7.41 (mapping* and preimage* arrow correctness). *The following diagram commutes:*

$$\begin{array}{ccccc}
X \rightsquigarrow_{\perp} Y & \xrightarrow{\text{lift}_{\text{map}}} & X \rightsquigarrow_{\text{map}} Y & \xrightarrow{\text{lift}_{\text{pre}}} & X \rightsquigarrow_{\text{pre}} Y \\
\eta_{\perp^*} \downarrow & & \downarrow \eta_{\text{map}^*} & & \downarrow \eta_{\text{pre}^*} \\
X \rightsquigarrow_{\perp^*} Y & \xrightarrow{\text{lift}_{\text{map}^*}} & X \rightsquigarrow_{\text{map}^*} Y & \xrightarrow{\text{lift}_{\text{pre}^*}} & X \rightsquigarrow_{\text{pre}^*} Y
\end{array} \tag{7.59}$$

Further, $\text{lift}_{\text{map}^*}$ and $\text{lift}_{\text{pre}^*}$ are arrow homomorphisms.

As with the correctness of interpretations using the mapping and preimage arrows, the correctness of interpretations using the mapping* and preimage* arrows follows from $\text{lift}_{\text{map}^*}$ and $\text{lift}_{\text{pre}^*}$ being arrow homomorphisms, and Theorem 7.6.

Corollary 7.42 (effectful semantic correctness). *For all expressions e ,*

$$\begin{aligned}
\llbracket e \rrbracket_{\text{pre}^*} &\equiv \text{lift}_{\text{pre}^*} (\text{lift}_{\text{map}^*} \llbracket e \rrbracket_{\perp^*}) \\
\llbracket e \rrbracket_{\text{pre}^*}^\Downarrow &\equiv \text{lift}_{\text{pre}^*} (\text{lift}_{\text{map}^*} \llbracket e \rrbracket_{\perp^*}^\Downarrow)
\end{aligned} \tag{7.60}$$

Unfortunately, because a statement such as “ $k_1 \equiv k_2$ ” implies k_1 terminates if and only if k_2 terminates, we cannot use the same tactics to prove an asymmetric statement such as “ k_2 terminates with the correct answer whenever k_1 terminates; otherwise returns \perp .” For these kinds of termination theorems, we need to reason about the interaction of programs with their supplied branch traces.

7.7.6 Termination

Here, we relate $\llbracket e \rrbracket_{a^*}^\downarrow$ computations, which are interpreted using $\text{ifte}_{a^*}^\downarrow$ and should always terminate, with $\llbracket e \rrbracket_{a^*}$ computations, which are interpreted using ifte_{a^*} and may not terminate. To do so, we need to find the largest domain on which $\llbracket e \rrbracket_{a^*}^\downarrow$ and $\llbracket e \rrbracket_{a^*}$ should agree.

Definition 7.43 (maximal domain). *A computation’s **maximal domain** is the largest A^* for which*

- For $f : X \rightsquigarrow_\perp Y$, $\text{domain}_\perp f A^* = A^*$.
- For $g : X \rightsquigarrow_{\text{map}} Y$, $\text{domain} (g A^*) = A^*$.
- For $h : X \rightsquigarrow_{\text{pre}} Y$, $\text{domain}_{\text{pre}} (h A^*) = A^*$.

The maximal domain of $k : X \rightsquigarrow_{a^} Y$ is that of $k j_0$.*

Because the above statements imply termination, A^* is a subset of the largest domain for which the computations terminate. It is not too hard to show (but is a bit tedious) that lifting computations preserves the maximal domain; e.g. the maximal domain of $\text{lift}_{\text{map}} f$ is the same as f ’s, and the maximal domain of $\text{lift}_{\text{pre}^*} g$ is the same as g ’s.

To ensure maximal domains exist, we need the domain operations above to have certain properties. For the mapping arrow, we must first make the intuition that computations “act as if they return restricted mappings” more precise. First, mapping restriction is defined by

$$\begin{aligned} \text{restrict} : (X \multimap Y) &\Rightarrow \text{Set } X \Rightarrow (X \multimap Y) \\ \text{restrict } g \ A &:= \lambda a \in (A \cap \text{domain } g). g \ a \end{aligned} \tag{7.61}$$

Theorem 7.44 (mapping arrow restriction). *Let $g : X \rightsquigarrow_{\text{map}} Y$, and $A^\Downarrow \subseteq X$ be the largest for which $g \ A^\Downarrow$ terminates. For all $A \subseteq A^\Downarrow$, $g \ A = \text{restrict } (g \ A^\Downarrow) \ A$.*

Proof. By the mapping arrow law (Definition 7.12) there is an $f : X \rightsquigarrow_{\perp} Y$ such that $g \equiv \text{lift}_{\text{map}} f$. Thus,

$$\begin{aligned}
\text{restrict } (g \ A^\Downarrow) \ A &\equiv \text{restrict } (\text{lift}_{\text{map}} f \ A^\Downarrow) \ A & (7.62) \\
&\equiv \text{restrict } (\{\langle a, b \rangle \in \text{mapping } f \ A^\Downarrow \mid b \neq \perp\}) \ A \\
&\equiv \{\langle a, b \rangle \in \text{mapping } f \ A \mid b \neq \perp\} \\
&\equiv \text{lift}_{\text{map}} f \ A \\
&\equiv g \ A & \square
\end{aligned}$$

Theorem 7.45 (domain closure operators). *If $f : X \rightsquigarrow_{\perp} Y$, $g : X \rightsquigarrow_{\text{map}} Y$ and $h : X \rightsquigarrow_{\text{pre}} Y$, then $\text{domain}_{\perp} f$, $\text{domain} \circ g$, and $\text{domain}_{\text{pre}} \circ h$ are monotone, decreasing, and idempotent in the subdomains on which they terminate.*

Proof. These properties follow from the same properties of selection, restriction, and of preimages of images. \square

Now we can relate $\llbracket e \rrbracket_{\perp}^{\Downarrow}$ computations to $\llbracket e \rrbracket_{\perp}$ computations. First, for any input for which $\llbracket e \rrbracket_{\perp}$ terminates, there should be a branch trace for which $\llbracket e \rrbracket_{\perp}^{\Downarrow}$ returns the correct output; it should otherwise return \perp .

Theorem 7.46. *Let $f := \llbracket e \rrbracket_{\perp} : X \rightsquigarrow_{\perp} Y$ with maximal domain A^* , and $f' := \llbracket e \rrbracket_{\perp}^{\Downarrow}$. For all $\langle \langle \omega, t \rangle, a \rangle \in A^*$, there exists a $T' \subseteq T$ such that*

- *If $t' \in T'$ then $f' \ j_0 \ \langle \langle \omega, t' \rangle, a \rangle = f \ j_0 \ \langle \langle \omega, t \rangle, a \rangle$.*
- *If $t' \in T \setminus T'$ then $f' \ j_0 \ \langle \langle \omega, t' \rangle, a \rangle = \perp$.*

Proof. Define T' as the set of all $t' \in J \rightarrow \text{Bool}_{\perp}$ such that $t' \ j = z$ if the subcomputation with index j is an if whose test returns z . Because $f \ j_0 \ \langle \langle \omega, t \rangle, a \rangle$ terminates, $t' \ j \neq \perp$ for at most finitely many j , so each $t' \in T$.

Let $\mathbf{t}' \in \mathbf{T}'$. Because the test of every if subcomputation at index j agrees with $\mathbf{t}' j$ and f ignores branch traces, $f' j_0 \langle \langle \omega, \mathbf{t}' \rangle, \mathbf{a} \rangle = f j_0 \langle \langle \omega, \mathbf{t} \rangle, \mathbf{a} \rangle$.

Let $\mathbf{t}' \in \mathbf{T} \setminus \mathbf{T}'$. There exists an if subexpression with a test that does not agree with \mathbf{t}' ; therefore $f' j_0 \langle \langle \omega, \mathbf{t}' \rangle, \mathbf{a} \rangle = \perp$. \square

Next, for any input for which $\llbracket e \rrbracket_{\perp^*}$ does not terminate or returns \perp , $\llbracket e \rrbracket_{\perp^*}^{\Downarrow}$ should return \perp . Proving this is a little easier if we first identify subsets of J that correspond with finite prefixes of an infinite binary tree.

Definition 7.47 (index prefix/suffix). *A finite $J' \subset J$ is an **index prefix** if $J' = \{j_0\}$ or, for some index prefix J'' and $j \in J''$, $J' = J'' \uplus \{\text{left } j\}$ or $J' = J'' \uplus \{\text{right } j\}$. The corresponding **index suffix** is $J \setminus J'$.*

It is not hard to show that every index suffix is closed under **left** and **right**.

For a given $\mathbf{t} \in \mathbf{T}$, an index prefix J' serves as a convenient bounding set for the finitely many indexes j for which $\mathbf{t} j \neq \perp$. Applying **left** and/or **right** repeatedly to any $j \in J'$ eventually yields a $j' \in J \setminus J'$, for which $\mathbf{t} j' = \perp$.

Theorem 7.48. *Let $f := \llbracket e \rrbracket_{\perp^*} : X \rightsquigarrow_{\perp^*} Y$ with maximal domain A^* , and $f' := \llbracket e \rrbracket_{\perp^*}^{\Downarrow}$. For all $\mathbf{a} \in ((\Omega \times \mathbf{T}) \times X) \setminus A^*$, $f' j_0 \mathbf{a} = \perp$.*

Proof. Let $\mathbf{t} := \text{snd}(\text{fst } \mathbf{a})$ be the branch trace element of \mathbf{a} .

Suppose $f j_0 \mathbf{a}$ terminates. If an if subcomputation's test does not agree with \mathbf{t} , then $f' j_0 \mathbf{a} = \perp$. If every if's test agrees, $f' j_0 \mathbf{a} = f j_0 \mathbf{a} = \perp$.

Suppose $f j_0 \mathbf{a}$ does not terminate. The set of all indexes j for which $\mathbf{t} j \neq \perp$ is contained within an index prefix J' . By hypothesis, there is an if subcomputation at some index j' such that $j' \in J \setminus J'$. Because $\mathbf{t} j' = \perp$, $f' j_0 \mathbf{a} = \perp$. \square

Corollary 7.49. *For all e , the maximal domain of $\llbracket e \rrbracket_{\perp^*}^{\Downarrow}$ is a subset of that of $\llbracket e \rrbracket_{\perp^*}$.*

Corollary 7.50. *Let $f' := \llbracket e \rrbracket_{\perp^*}^{\Downarrow} : X \rightsquigarrow_{\perp^*} Y$ with maximal domain A^* , and $f := \llbracket e \rrbracket_{\perp^*}$. For all $\mathbf{a} \in A^*$, $f' j_0 \mathbf{a} = f j_0 \mathbf{a}$.*

Corollary 7.51 (correct computation everywhere). *Let $\llbracket e \rrbracket_{\perp}^{\downarrow} : X \rightsquigarrow_{\perp}^* Y$ have maximal domain A^* , and $X' := (\Omega \times T) \times X$. For all $a \in X'$, $A \subseteq X'$ and $B \subseteq Y$,*

$$\begin{aligned} \llbracket e \rrbracket_{\perp}^{\downarrow} \text{ j}_0 a &= \text{if } (a \in A^*) (\llbracket e \rrbracket_{\perp}^{\downarrow} \text{ j}_0 a) \perp \\ \llbracket e \rrbracket_{\text{map}^*}^{\downarrow} \text{ j}_0 A &= \llbracket e \rrbracket_{\text{map}^*}^{\downarrow} \text{ j}_0 (A \cap A^*) \\ \text{ap}_{\text{pre}} (\llbracket e \rrbracket_{\text{pre}^*}^{\downarrow} \text{ j}_0 A) B &= \text{ap}_{\text{pre}} (\llbracket e \rrbracket_{\text{pre}^*}^{\downarrow} \text{ j}_0 (A \cap A^*)) B \end{aligned} \tag{7.63}$$

In other words, preimages computed using $\llbracket \cdot \rrbracket_{\text{pre}^*}^{\downarrow}$ always terminate, never include inputs that give rise to errors or nontermination, and are correct.

7.8 Output Probabilities and Measurability

Typically, for $g : \Omega \rightarrow Y$, the probability of $B \subseteq Y$ is $P(\text{preimage } g B)$, where $P : \text{Set } \Omega \rightarrow [0, 1]$ assigns probabilities to subsets of Ω .

A mapping* computation's domain is $(\Omega \times T) \times X$, not Ω . We assume each $\omega \in \Omega$ is randomly chosen, but not each $t \in T$ nor each $x \in X$; therefore, neither T nor X should affect the probabilities of output sets. We clearly must measure *projections* of preimage sets, or $P(\text{image } (\text{fst} \ggg \text{fst}) A)$ for preimage sets $A \subseteq (\Omega \times T) \times X$.

Not all preimage sets have sensible measures. Sets that do are called **measurable**. Computing preimages and projecting them onto Ω must preserve measurability.

Our main results are the best we could hope for. First, the interpretations of all expressions are measurable, regardless of nontermination.

Theorem 7.52. *For all expressions e , $\llbracket e \rrbracket_{\text{map}^*}^{\downarrow}$ is measurable.*

Second, projecting a program's preimages onto Ω results in a measurable set.

Theorem 7.53. *If $A \subseteq (\Omega \times T) \times \{\langle \rangle\}$ is measurable, then $\text{image } (\text{fst} \ggg \text{fst}) A$ is measurable.*

The proofs of these theorems are in Chapter 9.

$\text{id}_{\text{pre}} A := \langle A, \lambda B. B \rangle$	$\text{proj}_1 : \text{Set } \langle X_1, X_2 \rangle \Rightarrow \text{Set } X_1$
$\text{const}_{\text{pre}} b A := \langle \{b\}, \lambda B. \text{if } (B = \emptyset) \emptyset A \rangle$	$\text{proj}_1 := \text{image fst}$
$\text{fst}_{\text{pre}} A := \langle \text{proj}_1 A, \text{unproj}_1 A \rangle$	$\text{proj}_2 : \text{Set } \langle X_1, X_2 \rangle \Rightarrow \text{Set } X_2$
$\text{snd}_{\text{pre}} A := \langle \text{proj}_2 A, \text{unproj}_2 A \rangle$	$\text{proj}_2 := \text{image snd}$
$\pi_{\text{pre}} j A := \langle \text{proj } j A, \text{unproj } j A \rangle$	
<hr/>	
$\text{proj} : J \Rightarrow \text{Set } (J \rightarrow X) \Rightarrow \text{Set } X$	$\text{unproj}_1 : \text{Set } \langle X_1, X_2 \rangle \Rightarrow \text{Set } X_1 \Rightarrow \text{Set } \langle X_1, X_2 \rangle$
$\text{proj } j A := \text{image } (\pi j) A$	$\text{unproj}_1 A A_1 := \text{preimage } (\text{mapping fst } A) A_1$ $\equiv A \cap (A_1 \times \text{proj}_2 A)$
$\text{unproj} : J \Rightarrow \text{Set } (J \rightarrow X) \Rightarrow \text{Set } X \Rightarrow \text{Set } (J \rightarrow X)$	$\text{unproj}_2 : \text{Set } \langle X_1, X_2 \rangle \Rightarrow \text{Set } X_2 \Rightarrow \text{Set } \langle X_1, X_2 \rangle$
$\text{unproj } j A B := \text{preimage } (\text{mapping } (\pi j) A) B$ $\equiv A \cap \prod_{i \in J} \text{if } (j = i) B (\text{proj } j A)$	$\text{unproj}_2 A A_2 := \text{preimage } (\text{mapping snd } A) A_2$ $\equiv A \cap (\text{proj}_1 A \times A_2)$

Figure 7.10: Preimage arrow lifts needed to interpret probabilistic programs.

7.9 Approximating Semantics

If we were to confine preimage computation to finite sets, we could implement the preimage arrow directly. But we would like something that works efficiently on infinite sets, even if it means approximating. We focus on a specific method: approximating product sets with covering rectangles.

7.9.1 Implementable Lifts

We would like to be able to compute preimages of uncountable sets, such as real intervals—but $\text{preimage } (g A) B$ is uncomputable for most mappings g and uncountable sets A and B no matter how cleverly they are represented. Further, because pre , lift_{pre} and arr_{pre} are ultimately defined in terms of preimage , we cannot implement them.

Fortunately, we need to apply arr_{pre} only to certain functions. Figure 7.1 (which defines $\llbracket \cdot \rrbracket_a$) lifts id , $\text{const } b$, fst and snd . Section 7.7.4, which defines the combinators used to interpret partial, probabilistic programs, lifts πj and agrees . Measurable functions made available as language primitives, such as arithmetic, must be lifted to the preimage arrow—though to maintain generality, we put off lifting arithmetic functions until Chapter 8.

Figure 7.10 gives explicit definitions for $\text{arr}_{\text{pre}} \text{id}$, $\text{arr}_{\text{pre}} \text{fst}$, $\text{arr}_{\text{pre}} \text{snd}$, $\text{arr}_{\text{pre}} (\text{const } b)$ and $\text{arr}_{\text{pre}} (\pi j)$. (We will deal with **agrees** separately.) To implement them, we must model sets in a way that makes $A = \emptyset$ is decidable, and the following representable and finitely computable:

- $A \cap B$, \emptyset , $\{\text{true}\}$, $\{\text{false}\}$ and $\{b\}$ for every $\text{const } b$
- $A_1 \times A_2$, $\text{proj}_1 A$ and $\text{proj}_2 A$ (7.64)
- $J \rightarrow X$, $\text{proj } j A$ and $\text{unproj } j A B$

Before addressing representation and computability, we need to define families of sets under which these operations are closed.

Definition 7.54 (rectangular family). *Rect X denotes the **rectangular family** of subsets of X . Rect X must contain \emptyset and X , and be closed under finite intersections. Products must satisfy the following rules:*

$$\text{Rect } \langle X_1, X_2 \rangle = (\text{Rect } X_1) \boxtimes (\text{Rect } X_2) \quad (7.65)$$

$$\text{Rect } (J \rightarrow X) = (\text{Rect } X)^{\boxtimes J} \quad (7.66)$$

where the following operations lift cartesian products to sets of sets:

$$\mathcal{A}_1 \boxtimes \mathcal{A}_2 := \{A_1 \times A_2 \mid A_1 \in \mathcal{A}_1, A_2 \in \mathcal{A}_2\} \quad (7.67)$$

$$\mathcal{A}^{\boxtimes J} := \bigcup_{J' \subset J \text{ finite}} \left\{ \prod_{j \in J'} A_j \mid A_j \in \mathcal{A}, j \in J' \iff A_j \subset \bigcup \mathcal{A} \right\} \quad (7.68)$$

We additionally define $\text{Rect Bool} ::= \mathcal{P} \text{ Bool}$. It is easy to show the collection of all rectangular families is closed under products, projections, and **unproj**.

Further, all of the operations in (7.64) can be exactly implemented if finite sets are modeled directly, sets in ordered spaces (such as \mathbb{R}) are modeled by intervals, and sets in $\text{Rect } \langle X_1, X_2 \rangle$ are modeled by pairs of type $\langle \text{Rect } X_1, \text{Rect } X_2 \rangle$. By (7.68), sets in $\text{Rect } (J \rightarrow X)$ have no more than finitely many projections that are proper subsets of X . They can be modeled by *finite* binary trees, where unrepresented projections are implicitly X . Figure 7.11

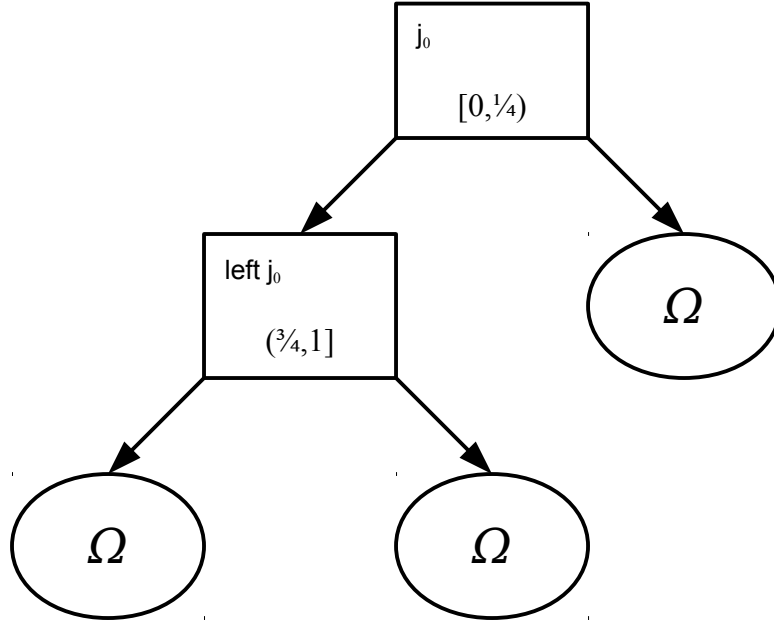


Figure 7.11: A finite binary tree model of $\text{unproj}(\text{left } j_0) (\text{unproj } j_0 \ \Omega \ [0, \frac{1}{4})) (\frac{3}{4}, 1]$. Because of Ω 's self-similarity, and because rectangles of $J \rightarrow [0, 1]$ are defined so that only finitely many projections are not $[0, 1]$, every rectangular subset of Ω has a finite binary tree model.

illustrates a model of a member of $\text{Rect}(J \rightarrow [0, 1])$; i.e. a rectangular subset of Ω .

The set of branch traces T is nonrectangular, but we can model T subsets by $J \rightarrow \text{Bool}_\perp$ rectangles, implicitly intersected with T .

Theorem 7.55 (T model). *If $T' \in \text{Rect}(J \rightarrow \text{Bool}_\perp)$ and $j \in J$, then $\text{proj } j \ T' = \text{proj } j \ (T' \cap T)$. If $B \subseteq \text{Bool}_\perp$, then $\text{unproj } j \ (T' \cap T) \ B = \text{unproj } j \ T' \ B \cap T$.*

Proof. Subset case is by projection monotonicity. For superset, let $b \in \text{proj } j \ T'$. Define t by $t \ j' = b$ if $j' = j$; $t \ j' = \perp$ if $\perp \in \text{proj } j' \ T'$; otherwise $t \ j' \in \text{proj } j' \ T'$.

By construction, $t \in T'$. For no more than finitely many $j' \in J$, $t \ j' \neq \perp$, so $t \in T$. Thus, there exists a $t \in T' \cap T$ such that $t \ j = b$, so $b \in \text{proj } j \ (T' \cap T)$.

The statement about unproj is an easy corollary. □

7.9.2 Approximate Preimage Mapping Operations

Implementing lazy_{pre} (defined in Figure 7.7) requires computing pre , but only for the empty mapping, which is trivial: $\text{pre } \emptyset \equiv \langle \emptyset, \lambda B. \emptyset \rangle$. Implementing the other combinators requires (\circ_{pre}) , $\langle \cdot, \cdot \rangle_{\text{pre}}$ and (\uplus_{pre}) .

From the preimage mapping definitions (Figure 7.5), we see that ap_{pre} is defined using (\cap) and that (\circ_{pre}) is defined using ap_{pre} , so (\circ_{pre}) is directly implementable. Unfortunately, we hit a snag with $\langle \cdot, \cdot \rangle_{\text{pre}}$: it loops over possibly uncountably many members of B in a big union. At this point, we need to approximate.

Theorem 7.56 (pair preimage approximation). *Let $g_1 : X \rightharpoonup Y_1$ and $g_2 : X \rightharpoonup Y_2$. For all $B \subseteq Y_1 \times Y_2$, $\text{preimage } \langle g_1, g_2 \rangle_{\text{map}} B \subseteq \text{preimage } g_1 (\text{proj}_1 B) \cap \text{preimage } g_2 (\text{proj}_2 B)$.*

Proof. By monotonicity of preimages and projections, and by Lemma 7.20. \square

It is not hard to use Theorem 7.56 to show that

$$\begin{aligned} \langle \cdot, \cdot \rangle'_{\text{pre}} : (X \rightharpoonup_{\text{pre}} Y_1) &\Rightarrow (X \rightharpoonup_{\text{pre}} Y_2) \Rightarrow (X \rightharpoonup_{\text{pre}} Y_1 \times Y_2) \\ \langle \langle Y'_1, p_1 \rangle, \langle Y'_2, p_2 \rangle \rangle'_{\text{pre}} &:= \langle Y'_1 \times Y'_2, \lambda B. p_1 (\text{proj}_1 B) \cap p_2 (\text{proj}_2 B) \rangle \end{aligned} \quad (7.69)$$

computes covering rectangles of preimages under pairing.

For (\uplus_{pre}) , we need an approximating replacement for (\cup) under which rectangular families are closed. In other words, we need a lattice join (\vee) with respect to (\subseteq) , with the following additional properties:

$$\begin{aligned} (A_1 \times A_2) \vee (B_1 \times B_2) &= (A_1 \vee B_1) \times (A_2 \vee B_2) \\ (\prod_{j \in J} A_j) \vee (\prod_{j \in J} B_j) &= \prod_{j \in J} A_j \vee B_j \end{aligned} \quad (7.70)$$

If for every nonproduct type X , $\text{Rect } X$ is closed under (\vee) , then rectangular families are clearly closed under (\vee) . Further, for any A and B , $A \cup B \subseteq A \vee B$.

Replacing each union in (\uplus_{pre}) with join yields the overapproximating (\uplus'_{pre}) :

$$\begin{aligned}
(\uplus'_{\text{pre}}) : (X \xrightarrow{\text{pre}} Y) &\Rightarrow (X \xrightarrow{\text{pre}} Y) \Rightarrow (X \xrightarrow{\text{pre}} Y) \\
h_1 \uplus'_{\text{pre}} h_2 &:= \text{let } Y' := \text{range}_{\text{pre}} h_1 \vee \text{range}_{\text{pre}} h_2 \\
&\quad p := \lambda B. \text{ap}_{\text{pre}} h_1 B \vee \text{ap}_{\text{pre}} h_2 B \\
&\quad \text{in } \langle Y', p \rangle
\end{aligned} \tag{7.71}$$

To interpret programs that may not terminate, or that terminate with probability 1, we need to approximate $\text{ifte}_{\text{pre}^*}^{\Downarrow}$ (7.52), which is defined in terms of **agrees**. Defining its approximation in terms of an approximation of **agrees** would not allow us to preserve the fact that expressions interpreted using $\text{ifte}_{\text{pre}^*}^{\Downarrow}$ always terminate. The best approximation of the preimage of **Bool** under **agrees** (as a mapping) is **Bool** \times **Bool**, which contains $\langle \text{true}, \text{false} \rangle$ and $\langle \text{false}, \text{true} \rangle$, and thus would not constrain the test to agree with the branch trace.

A lengthy (elided) sequence of substitutions to the defining expression for $\text{ifte}_{\text{pre}^*}^{\Downarrow}$ results in an **agrees**-free equivalence:

$$\begin{aligned}
\text{ifte}_{\text{pre}^*}^{\Downarrow} k_1 k_2 k_3 j A &\equiv \text{let } \langle C_k, p_k \rangle := k_1 j_1 A \\
&\quad \langle C_b, p_b \rangle := \text{branch}_{\text{pre}^*} j A \\
&\quad C_2 := C_k \cap C_b \cap \{\text{true}\} \\
&\quad C_3 := C_k \cap C_b \cap \{\text{false}\} \\
&\quad A_2 := p_k C_2 \cap p_b C_2 \\
&\quad A_3 := p_k C_3 \cap p_b C_3 \\
&\quad \text{in } k_2 j_2 A_2 \uplus_{\text{pre}} k_3 j_3 A_3
\end{aligned} \tag{7.72}$$

where $j_1 = \text{left } j$ and so on. Unfortunately, a straightforward approximation of this would still take unnecessary branches, when A_2 or A_3 overapproximates \emptyset .

C_b is the branch trace projection at j (with \perp removed). The set of indexes for which C_b is either $\{\text{true}\}$ or $\{\text{false}\}$ is finite, so it is bounded by an index prefix, outside of which branch trace projections are $\{\text{true}, \text{false}\}$. Therefore, if the approximating $\text{ifte}_{\text{pre}^*}^{\Downarrow}$ takes *no branches* when $C_b = \{\text{true}, \text{false}\}$, but approximates with a finite computation, expressions interpreted using $\text{ifte}_{\text{pre}^*}^{\Downarrow}$ will always terminate.

We need an overapproximation for the non-branching case. In the exact semantics, the returned preimage mapping's range is a subset of Y , and it returns subsets of $A_2 \uplus A_3$.

$$\begin{aligned}
X \xrightarrow{\text{pre}}' Y &::= \langle \text{Rect } Y, \text{Rect } Y \Rightarrow \text{Rect } X \rangle & \langle \cdot, \cdot \rangle'_{\text{pre}} : (X \xrightarrow{\text{pre}}' Y_1) \Rightarrow (X \xrightarrow{\text{pre}}' Y_2) \Rightarrow (X \xrightarrow{\text{pre}}' Y_1 \times Y_2) \\
\emptyset'_{\text{pre}} &::= \langle \emptyset, \lambda B. \emptyset \rangle & \langle \langle Y'_1, p_1 \rangle, \langle Y'_2, p_2 \rangle \rangle'_{\text{pre}} := \\
& & \langle Y'_1 \times Y'_2, \lambda B. p_1 (\text{proj}_1 B) \cap p_2 (\text{proj}_2 B) \rangle \\
\text{ap}'_{\text{pre}} : (X \xrightarrow{\text{pre}}' Y) \Rightarrow \text{Rect } Y \Rightarrow \text{Rect } X & & (\uplus'_{\text{pre}}) : (X \xrightarrow{\text{pre}}' Y) \Rightarrow (X \xrightarrow{\text{pre}}' Y) \Rightarrow (X \xrightarrow{\text{pre}}' Y) \\
\text{ap}'_{\text{pre}} \langle Y', p \rangle B &::= p (B \cap Y') & \langle Y'_1, p_1 \rangle \uplus'_{\text{pre}} \langle Y'_2, p_2 \rangle := \\
& & \langle Y'_1 \vee Y'_2, \lambda B. \text{ap}'_{\text{pre}} \langle Y'_1, p_1 \rangle B \vee \text{ap}'_{\text{pre}} \langle Y'_2, p_2 \rangle B \rangle \\
(\circ'_{\text{pre}}) : (Y \xrightarrow{\text{pre}}' Z) \Rightarrow (X \xrightarrow{\text{pre}}' Y) \Rightarrow (X \xrightarrow{\text{pre}}' Z) & & \\
\langle Z', p_2 \rangle \circ'_{\text{pre}} h_1 &::= \langle Z', \lambda C. \text{ap}'_{\text{pre}} h_1 (p_2 C) \rangle &
\end{aligned}$$

(a) Definitions for preimage mappings that compute rectangular covers.

$$\begin{aligned}
X \xrightarrow{\text{pre}}'' Y &::= \text{Rect } X \Rightarrow (X \xrightarrow{\text{pre}}' Y) & \text{ifte}'_{\text{pre}} : (X \xrightarrow{\text{pre}}'' \text{Bool}) \Rightarrow (X \xrightarrow{\text{pre}}'' Y) \Rightarrow (X \xrightarrow{\text{pre}}'' Y) \Rightarrow (X \xrightarrow{\text{pre}}'' Y) \\
(\ggg'_{\text{pre}}) : (X \xrightarrow{\text{pre}}'' Y) \Rightarrow (Y \xrightarrow{\text{pre}}'' Z) \Rightarrow (X \xrightarrow{\text{pre}}'' Z) & & \text{ifte}'_{\text{pre}} h_1 h_2 h_3 A := \text{let } h'_1 := h_1 A \\
(h_1 \ggg'_{\text{pre}} h_2) A &::= \text{let } h'_1 := h_1 A & h'_2 := h_2 (\text{ap}'_{\text{pre}} h'_1 \{\text{true}\}) \\
& & h'_3 := h_3 (\text{ap}'_{\text{pre}} h'_1 \{\text{false}\}) \\
& & \text{in } h'_2 \circ'_{\text{pre}} h'_1 \\
(\&\&\&'_{\text{pre}}) : (X \xrightarrow{\text{pre}}'' Y_1) \Rightarrow (X \xrightarrow{\text{pre}}'' Y_2) \Rightarrow (X \xrightarrow{\text{pre}}'' \langle Y_1, Y_2 \rangle) & & \text{lazy}'_{\text{pre}} : (1 \Rightarrow (X \xrightarrow{\text{pre}}'' Y)) \Rightarrow (X \xrightarrow{\text{pre}}'' Y) \\
(h_1 \&\&\&'_{\text{pre}} h_2) A &::= \langle h_1 A, h_2 A \rangle'_{\text{pre}} & \text{lazy}'_{\text{pre}} h A := \text{if } (A = \emptyset) \emptyset'_{\text{pre}} (h 0 A)
\end{aligned}$$

(b) Approximating preimage arrow, defined using approximating preimage mappings.

$$\begin{aligned}
X \xrightarrow{\text{pre}^*}'' Y &::= \text{AStore } (\Omega \times \mathbb{T}) (X \xrightarrow{\text{pre}}'' Y) & \text{ifte}^{\downarrow}_{\text{pre}^*} : (X \xrightarrow{\text{pre}^*}'' \text{Bool}) \Rightarrow (X \xrightarrow{\text{pre}^*}'' Y) \Rightarrow (X \xrightarrow{\text{pre}^*}'' Y) \Rightarrow (X \xrightarrow{\text{pre}^*}'' Y) \\
\text{random}'_{\text{pre}^*} : X \xrightarrow{\text{pre}^*}'' [0, 1] & & \text{ifte}^{\downarrow}_{\text{pre}^*} k_1 k_2 k_3 j := \\
\text{random}'_{\text{pre}^*} j &::= & \text{let } \langle C_k, p_k \rangle := k_1 (\text{left } j) A \\
\text{fst}_{\text{pre}} \ggg'_{\text{pre}} \text{fst}_{\text{pre}} \ggg'_{\text{pre}} \pi_{\text{pre}} j & & \langle C_b, p_b \rangle := \text{branch}_{\text{pre}^*} j A \\
& & C_2 := C_k \cap C_b \cap \{\text{true}\} \\
& & C_3 := C_k \cap C_b \cap \{\text{false}\} \\
& & A_2 := p_k C_2 \cap p_b C_2 \\
& & A_3 := p_k C_3 \cap p_b C_3 \\
\text{branch}'_{\text{pre}^*} : X \xrightarrow{\text{pre}^*}'' \text{Bool} & & \text{in if } (C_b = \{\text{true}, \text{false}\}) \\
\text{branch}'_{\text{pre}^*} j &::= & \langle \top, \lambda B. A_2 \vee A_3 \rangle \\
\text{fst}_{\text{pre}} \ggg'_{\text{pre}} \text{snd}_{\text{pre}} \ggg'_{\text{pre}} \pi_{\text{pre}} j & & (k_2 (\text{left } (\text{right } j)) A_2 \uplus'_{\text{pre}} k_3 (\text{right } (\text{right } j)) A_3) \\
\text{fst}'_{\text{pre}^*} &::= \eta'_{\text{pre}^*} \text{fst}_{\text{pre}}; \dots
\end{aligned}$$

(c) Preimage* arrow combinators for probabilistic choice and guaranteed termination. Figure 7.8 (AStore arrow transformer) defines η'_{pre^*} , (\ggg'_{pre^*}) , $(\&\&\&'_{\text{pre}^*})$, $\text{ifte}'_{\text{pre}^*}$ and $\text{lazy}'_{\text{pre}^*}$.

Figure 7.12: Implementable arrows that approximate preimage arrows. Specific lifts such as $\text{fst}_{\text{pre}} := \text{arr}_{\text{pre}} \text{fst}$ are computable (see Figure 7.10), but arr'_{pre} is not.

Therefore, $\text{ifte}^{\downarrow}_{\text{pre}^*}$ may return $\langle Y, \lambda B. A_2 \vee A_3 \rangle$ when $C_b = \{\text{true}, \text{false}\}$. We cannot refer to the type Y in the function definition, so we represent it using \top in the approximating semantics. Implementations can model it by a singleton “universe” instance for every $\text{Rect } Y$.

Figure 7.12b defines the final approximating preimage arrow. This arrow, the lifts in Figure 7.10, and the semantic function $\llbracket \cdot \rrbracket_a$ in Figure 7.1 define an approximating semantics for partial, probabilistic programs.

7.9.3 Correctness

From here on, $\llbracket \cdot \rrbracket_{\text{pre}^*}^{\Downarrow'}$ interprets programs as approximating preimage* arrow computations using $\text{ifte}_{\text{pre}^*}^{\Downarrow'}$. The following theorems assume $h := \llbracket e \rrbracket_{\text{pre}^*}^{\Downarrow} : X \rightsquigarrow_{\text{pre}^*} Y$ and $h' := \llbracket e \rrbracket_{\text{pre}^*}^{\Downarrow'} : X \rightsquigarrow_{\text{pre}^*}' Y$ for some expression e .

To use structural induction on the interpretation of e , we need a theorem that allows representing it as a finite expression (Definition 9.28). Because $\text{ifte}_{\text{pre}^*}^{\Downarrow'}$ does not branch when either branch could be taken, an equivalent finite expression exists for each rectangular domain subset A .

Theorem 7.57 (equivalent finite expression). *For all $A \in \text{Rect } \langle \langle \Omega, T \rangle, X \rangle$, there exists a finite expression e' for which, if $h'' := \llbracket e' \rrbracket_{\text{pre}^*}^{\Downarrow'}$, then $\text{ap}'_{\text{pre}}(h'' \text{ j}_0 A) B = \text{ap}'_{\text{pre}}(h' \text{ j}_0 A) B$ for all $B \in \text{Rect } Y$.*

Proof. Let $T' := \text{proj}_2(\text{proj}_1 A)$, and let the index prefix J' contain every j' for which $(\text{proj } j' T') \setminus \{\perp\}$ is either $\{\text{true}\}$ or $\{\text{false}\}$. To construct e' , exhaustively apply first-order functions in e , but replace any if $e_1 \ e_2 \ e_3$ whose index is not in J' with the equivalent expression if $e_1 \perp \perp$. Because e is well-defined, recurrences must be guarded by if, so this process terminates after finitely many applications. \square

Corollary 7.58 (terminating). *For all $A \in \text{Rect } \langle \langle \Omega, T \rangle, X \rangle$ and $B \in \text{Rect } Y$, $\text{ap}'_{\text{pre}}(h' \text{ j}_0 A) B$ terminates.*

Theorem 7.59 (sound). *For all $A \in \text{Rect } \langle \langle \Omega, T \rangle, X \rangle$ and $B \in \text{Rect } Y$, $\text{ap}_{\text{pre}}(h \text{ j}_0 A) B \subseteq \text{ap}'_{\text{pre}}(h' \text{ j}_0 A) B$.*

Proof. By construction and Corollary 7.58 (recall non-“ \equiv ” statements imply termination). \square

Theorem 7.60 (monotone). $\text{ap}'_{\text{pre}} (h' \text{ j}_0 A) B$ is monotone in both A and B .

Proof. Lattice operators (\cap) and (\vee) are monotone, as is (\times) . Therefore, id_{pre} and the other lifts in Figure 7.10 are monotone, and each approximating preimage arrow combinator preserves monotonicity. Approximating preimage* arrow combinators, which are defined in terms of approximating preimage arrow combinators (Figure 7.12b) likewise preserve monotonicity, as does η'_{pre^*} ; therefore id_{pre^*} and other lifts are monotone.

The definition of $\text{ifte}'_{\text{pre}^*}$ can be written in terms of lattice operators and approximating preimage arrow combinators for any A for which $C_b = \{\text{true}\}$ or $C_b = \{\text{false}\}$, and thus preserves monotonicity in those cases. If $C_b = \{\text{true}, \text{false}\}$, which is an upper bound for C_b , the returned value is an upper bound.

For monotonicity in A , suppose $A_1 \subseteq A_2$. Apply Theorem 7.57 with A_1 to yield e' ; clearly, it is also an equivalent finite expression for A_2 . Monotonicity follows from structural induction on the interpretation of e' .

For monotonicity in B , apply Theorem 7.57 with a fixed A . □

Theorem 7.61 (decreasing). If $A \in \text{Rect } \langle \langle \Omega, T \rangle, X \rangle$ and $B \in \text{Rect } Y$, $\text{ap}'_{\text{pre}} (h' \text{ j}_0 A) B \subseteq A$.

Proof. Because they compute exact preimages of rectangular sets under restriction to rectangular domains, id_{pre} and the other lifts in Figure 7.10 are decreasing.

By definition and applying basic lattice properties,

$$\begin{aligned}
\text{ap}'_{\text{pre}} ((h_1 \ggg'_{\text{pre}} h_2) A) B &\equiv \text{ap}'_{\text{pre}} (h_1 A) B' \text{ for some } B' & (7.73) \\
\text{ap}'_{\text{pre}} ((h_1 \&\&'_{\text{pre}} h_2) A) B &\equiv \text{ap}'_{\text{pre}} (h_1 A) (\text{proj}_1 B) \cap \text{ap}'_{\text{pre}} (h_2 A) (\text{proj}_2 B) \\
\text{ap}'_{\text{pre}} (\text{ifte}'_{\text{pre}} h_1 h_2 h_3 A) B &\equiv \text{let } A_2 := \text{ap}'_{\text{pre}} (h_1 A) \{\text{true}\} \\
&\quad A_3 := \text{ap}'_{\text{pre}} (h_1 A) \{\text{false}\} \\
&\quad \text{in } \text{ap}'_{\text{pre}} (h_2 A_2) B \vee \text{ap}'_{\text{pre}} (h_3 A_3) B \\
\text{ap}'_{\text{pre}} (\text{lazy}'_{\text{pre}} h A) B &\equiv \text{if } (A = \emptyset) \emptyset (\text{ap}'_{\text{pre}} (h \text{ 0 } A) B)
\end{aligned}$$

Thus, approximating preimage arrow combinators return decreasing computations when given decreasing computations. This property transfers trivially to approximating preimage* arrow

combinators. Use Theorem 7.57 with $A' := A$, and structural induction. \square

7.9.4 Preimage Refinement Algorithm

Given these properties, we might try to compute preimages of B by computing preimages with respect to increasingly fine discretizations of A .

Definition 7.62 (preimage refinement algorithm). *Let $B \in \text{Rect } Y$ and*

$$\begin{aligned} \text{refine} &: \text{Rect } \langle \langle \Omega, T \rangle, X \rangle \Rightarrow \text{Rect } \langle \langle \Omega, T \rangle, X \rangle \\ \text{refine } A &:= \text{ap}'_{\text{pre}} (h' \text{ j}_0 A) B \end{aligned} \tag{7.74}$$

Define $\text{partition} : \text{Rect } \langle \langle \Omega, T \rangle, X \rangle \Rightarrow \text{Set } (\text{Rect } \langle \langle \Omega, T \rangle, X \rangle)$ to produce positive-measure, disjoint rectangles, and define

$$\begin{aligned} \text{refine}^* &: \text{Set } (\text{Rect } \langle \langle \Omega, T \rangle, X \rangle) \Rightarrow \text{Set } (\text{Rect } \langle \langle \Omega, T \rangle, X \rangle) \\ \text{refine}^* \mathcal{A} &:= \text{image } \text{refine} \left(\bigcup_{A \in \mathcal{A}} \text{partition } A \right) \end{aligned} \tag{7.75}$$

For any $A \in \text{Rect } \langle \langle \Omega, T \rangle, X \rangle$, iterate refine^ on $\{A\}$.*

Theorem 7.61 (decreasing) guarantees $\text{refine } A$ is never larger than A . Theorem 7.60 (monotone) guarantees refining a *partition* of A never does worse than refining A itself. Theorem 7.59 (sound) guarantees the algorithm is **sound**: the preimage of B is always contained in the covering partition refine^* returns.

We would like it to be **complete** in the limit, up to null sets: covering partitions' measures should converge to the true measure. Unfortunately, preimage refinement appears to compute the **Jordan outer measure** of a preimage, which is not always its measure. A counterexample is the expression `rational? random`, where `rational?` returns `true` when its argument is rational and loops otherwise. (This is definable using a (\leq) primitive.) The preimage of $\{\text{true}\}$ (the rational numbers) has measure 0, but its Jordan outer measure is 1.

We conjecture that a minimal requirement for preimage refinement's measures to converge is that a program must converge with probability 1. There are certainly other

requirements. We leave these and proof of convergence of measures for future work.

For now, we use algorithms that depend only on soundness.

7.10 Implementations

We have four implementations: one of the exact semantics, two direct implementations of the approximating semantics, and a less direct but more efficient implementation of the approximating semantics, which we call *Dr. Bayes*.

If sets are restricted to be finite, the arrows used as translation targets in the exact semantics, defined in 7.2, 7.4, 7.5, 7.7 and 7.8, can be implemented directly in any practical λ -calculus. Computing exact preimages is very inefficient, even under the interpretations of very small programs. Still, we have found our Typed Racket [47] implementation useful for finding theorem candidates and counterexamples.

Given a rectangular set library, the approximating preimage arrows defined in Figures 7.10 and 7.12b can be implemented with few changes in any practical λ -calculus. We have done so in Typed Racket and Haskell [1]. Both implementations' arrow combinator definitions are almost line-for-line transliterations from the figures. They are at <https://github.com/ntoronto/writing/tree/master/2014esop-code>. [XXX: move]

Making the rectangular set type polymorphic seems to require the equivalent of a typeclass system. In Haskell, it also requires multi-parameter typeclasses or indexed type families [12] to associate set types with the types of their members. Using indexed type families, the only significant differences between the Haskell implementation and the approximating semantics are type contexts, `newtype` wrappers for arrow types, and using `Maybe` types as bottom arrow return types.

Typed Racket has no typeclass system on top of its type system, so the rectangular set type is monomorphic; thus, so are the arrow types. The lack of type variables in the combinator types is the only significant difference between the implementation and the approximating semantics.

Chapter 8 details the implementation of Dr. Bayes.

7.11 Conclusions

To allow recursion and arbitrary conditions in probabilistic programs, we combined the power of measure theory with the unifying elegance of arrows. We

1. Defined a transformation from first-order programs to arbitrary arrows.
2. Defined the bottom arrow as a standard translation target.
3. Derived the uncomputable preimage arrow as an alternative target.
4. Derived a sound, computable approximation of the preimage arrow, and enough computable lifts to transform programs.

Critically, the preimage arrow's lift from the bottom arrow distributes over bottom arrow computations. Our semantics thus generalizes this process to all programs: 1) encode a program as a bottom arrow computation; 2) lift this computation to get an uncomputable function that computes preimages; 3) distribute the lift; and 4) replace uncomputable expressions with sound approximations.

Using arrows drastically simplifies the correctness proofs. Almost every semantic correctness theorem proceeds from a proof that a lift distributes over five combinators. There are seven theorems in total corresponding to the morphisms in our roadmap (7.3), but the three center morphisms (pointing downward) are done in one proof, as are the two bottom morphisms (pointing rightward). In total, there are 20 cases, plus 11 for the original (and very simple) theorem that arrow homomorphisms distribute over program terms.

In contrast, the corresponding theorems with separate semantic functions would require seven proofs by structural induction over at least 11 rules (12 for programs that access the random store), for a total of at least 77 cases. This reduction in complexity by semantic abstraction not have been possible without targeting λ_{ZFC} , which allows such arrows to carry out uncountably infinite computations.

XXX: segue

Chapter 8

Preimage Computation Implementation

8.1 Set Representation

8.2 Preimages Under Real Functions

XXX: lifts are generally uncomputable, but many specific lifts are computable or approximately computable

XXX: something about computing functions backward by computing inverses forward

XXX: obvious how to do this with invertible functions (though infinite endpoints are a little tricky); not obvious how to do the same with two-argument functions

8.2.1 Invertible Primitives

We consider only strictly monotone functions on \mathbb{R} . Further on, we recover more generality by using language conditionals to implement *piecewise* monotone functions.

One reason we consider only strictly monotone functions is that they are easy to invert. Recall that a function is invertible (bijective) if and only if it is injective (one-to-one) and surjective (onto).

Lemma 8.1. *If $f : A \rightarrow B$ is strictly monotone, f is injective. If f is additionally surjective, f and its inverse are continuous.*

Preimages under invertible functions can be computed using their inverses. We are primarily interested in computing preimages under restricted functions.

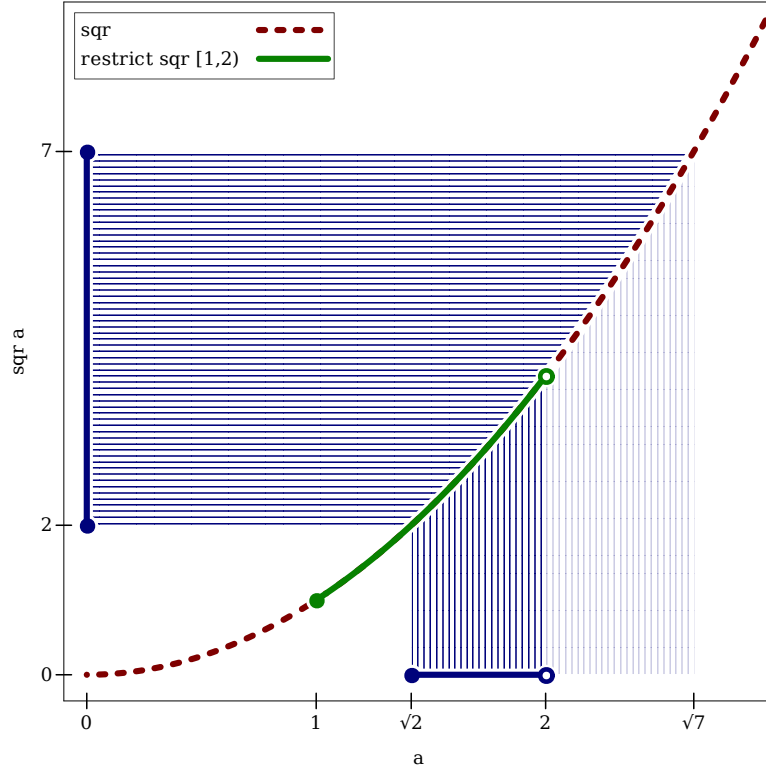


Figure 8.1: Computing the preimage of the interval $[2, 7]$ under sqr restricted to $[1, 2)$, by computing roots and intersecting with $[1, 2)$.

Lemma 8.2. *Let $A' \subseteq A$, $B' \subseteq B$, and $f : A \rightarrow B$ have inverse $f^{-1} : B \rightarrow A$. Let $f' := \text{restrict } f \ A'$. Then*

$$\text{preimage } f' \ B' = A' \cap \text{image } f^{-1} \ B' \quad (8.1)$$

These facts suggest that we can compute images (or preimages) of intervals under any strictly monotone, surjective f by applying f (or its inverse) to interval endpoints to yield an interval, as in Figure 8.1 This is evident for endpoints in A . Limit endpoints like $+\infty$ require a larger \bar{f} defined on a compact superset of A .

The next theorem is easier to state with interval notation in which the kind of interval is not baked into the syntax.

Definition 8.3 (interval). $[a_1, a_2, \alpha_1, \alpha_2]$ denotes an interval, where $a_1, a_2 \in \overline{\mathbb{R}}$ are extended real endpoints, and $\alpha_1, \alpha_2 \in \text{Bool}$ determine whether a_1 and a_2 are contained in the interval.

Example 8.4. Some intervals, using $[\cdot, \cdot, \cdot, \cdot]$ notation:

$$\begin{aligned}
[0, 1, \text{true}, \text{false}] &= [0, 1) \\
[-\infty, 0, \text{false}, \text{true}] &= (-\infty, 0] \\
[-\infty, +\infty, \text{false}, \text{false}] &= (-\infty, +\infty) = \mathbb{R} \\
[-\infty, +\infty, \text{true}, \text{true}] &= [-\infty, +\infty] = \overline{\mathbb{R}}
\end{aligned} \tag{8.2}$$

All but the last are subsets of \mathbb{R} . ◇

Theorem 8.5 (images of intervals by endpoints). *Let \bar{A} and \bar{B} be compact subsets of $\overline{\mathbb{R}}$, $\bar{f} : \bar{A} \rightarrow \bar{B}$ strictly monotone and surjective, and f the restriction of \bar{f} to some $A \subseteq \bar{A}$. For all nonempty $[a_1, a_2, \alpha_1, \alpha_2] \subseteq A$,*

- *If \bar{f} is increasing, image $f [a_1, a_2, \alpha_1, \alpha_2] = [\bar{f} a_1, \bar{f} a_2, \alpha_1, \alpha_2]$.*
- *If \bar{f} is decreasing, image $f [a_1, a_2, \alpha_1, \alpha_2] = [\bar{f} a_2, \bar{f} a_1, \alpha_2, \alpha_1]$.*

Proof. Because \bar{A} is compact and totally ordered, every subset of \bar{A} has a lower and an upper bound in \bar{A} . Therefore, the endpoints of every interval subset of A are in \bar{A} .

Let $(a_1, a_2] \subseteq A$. Suppose \bar{f} is strictly increasing; thus $a_1 < a \leq a_2$ if and only if $\bar{f} a_1 < \bar{f} a \leq \bar{f} a_2$, so $\text{image } f (a_1, a_2] = \text{image } \bar{f} (a_1, a_2] = (\bar{f} a_1, \bar{f} a_2]$. The remaining cases are similar. □

To use Theorem 8.5 to compute preimages under f by computing images under its inverse f^{-1} , we must know whether f^{-1} is increasing or decreasing. The following lemma can help.

Lemma 8.6. *If $f : A \rightarrow B$ is strictly monotone and surjective with inverse $f^{-1} : B \rightarrow A$, then f is increasing if and only if f^{-1} is increasing.*

Example 8.7. The extension of $\log : (0, +\infty) \rightarrow \mathbb{R}$ to the compact superdomain $[0, +\infty]$ is

defined by

$$\begin{aligned} \overline{\log} : [0, +\infty] &\rightarrow \overline{\mathbb{R}} \\ \overline{\log} a &:= \lim_{a' \rightarrow a} \log a' = \text{case } a \\ &\quad \begin{array}{ll} 0 &\longrightarrow -\infty \\ +\infty &\longrightarrow +\infty \\ \text{else} &\longrightarrow \log a \end{array} \end{aligned} \quad (8.3)$$

The extension of its inverse \exp is $\overline{\exp} : \overline{\mathbb{R}} \rightarrow [0, +\infty]$, defined similarly, which by Lemma 8.6 is also strictly increasing. Thus,

$$\begin{aligned} \text{image } \log (0, 1] &= (\overline{\log} 0, \overline{\log} 1] = (-\infty, 0] \\ \text{preimage } \log [0, +\infty) &= \text{image } \exp [0, +\infty) \\ &= [\overline{\exp} 0, \overline{\exp} +\infty) = [1, +\infty) \end{aligned} \quad (8.4)$$

by Theorem 8.5 and Lemma 8.2, ◇

8.2.2 Two-Argument Primitives

We do not expect to be able to compute preimages under $\mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ primitives by simply inverting them. Two-argument invertible real functions are difficult to define and are usually pathological.

Instead, we compute approximate preimages only, using inverses with respect to one argument (with the other held constant).

Definition 8.8 (axial inverse). *Let $f_c : A \times B \rightarrow C$. Functions $f_a : B \times C \rightarrow A$ and $f_b : C \times A \rightarrow B$ defined so that*

$$f_c \langle a, b \rangle = c \iff f_a \langle b, c \rangle = a \iff f_b \langle c, a \rangle = b \quad (8.5)$$

*are **axial inverses** with respect to f_c 's first and second arguments.*

We call f_c **axis-invertible** or **trijjective** when it has axial inverses f_a and f_b . We call f_a the **first axial inverse** of f_c because it is the inverse of f_c along the first axis: f_a with only

c varying (i.e. $\lambda c \in C$. $f_a \langle b, c \rangle$), is the inverse of f_c with only a varying (i.e. $\lambda a \in A$. $f_c \langle a, b \rangle$). Similarly, f_b is the *second axial inverse*.

Example 8.9. Let $\text{add}_c : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, $\text{add}_c \langle a, b \rangle := a + b$. Its axial inverses are $\text{add}_a \langle b, c \rangle := c - b$ and $\text{add}_b \langle c, a \rangle := c - a$. \diamond

We have chosen the axial inverse function types carefully: they are the only types for which f_c , f_a and f_b form a cyclic group.

Lemma 8.10. *The following statements are equivalent.*

- f_c has axial inverses f_a and f_b .
- f_a has axial inverses f_b and f_c .
- f_b has axial inverses f_c and f_a .

Equivalently, every axis-invertible function generates a cyclic group of order 3 by inversion in the first axis.

This fact is analogous to how mutual inverses f and f^{-1} also form a cyclic group (of order 2, generated by inversion). Similar to using mutual inversion to compute preimages under both \log and \exp , Lemma 8.10 allows computing preimages under two-argument functions related by axial inversion.

Example 8.11. Define $\text{sub}_c : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ by $\text{sub}_c \langle a, b \rangle := a - b$. Because $\text{sub}_c = \text{add}_b$, $\text{sub}_a = \text{add}_c$ and $\text{sub}_b = \text{add}_a$. \diamond

Unlike inverses, axial inverses do not provide a direct way to compute exact preimages. Instead, they provide a way to compute a preimage's smallest rectangular bounding set.

Theorem 8.12 (preimage bounds from axial inverse images). *Let $A' \subseteq A$, $B' \subseteq B$, $C' \subseteq C$, and $f_c : A \times B \rightarrow C$ with axial inverses f_a and f_b . If $f'_c = \text{restrict } f_c (A' \times B')$, then*

$$\text{preimage } f'_c C' \subseteq (A' \cap \text{image } f_a (B' \times C')) \times (B' \cap \text{image } f_b (C' \times A')) \quad (8.6)$$

Further, the right-hand side is the smallest rectangular superset.

Proof. The smallest rectangle containing $\text{preimage } f'_c C'$ is

$$\text{preimage } f'_c C' \subseteq (\text{image fst } (\text{preimage } f'_c C')) \times (\text{image snd } (\text{preimage } f'_c C')) \quad (8.7)$$

Starting with the first set in the product, expand definitions, distribute fst , replace $f_c \langle a, b \rangle = c$ by $f_a \langle b, c \rangle = a$, and simplify:

$$\begin{aligned} & \text{image fst } (\text{preimage } f'_c C') \\ &= \text{image fst } \{ \langle a, b \rangle \in A' \times B' \mid f_c \langle a, b \rangle \in C' \} \\ &= \{ a \in A' \mid \exists b \in B'. f_c \langle a, b \rangle \in C' \} \\ &= \{ a \in A' \mid \exists b \in B', c \in C'. f_c \langle a, b \rangle = c \} \\ &= \{ a \in A' \mid \exists b \in B', c \in C'. f_a \langle b, c \rangle = a \} \\ &= \{ f_a \langle b, c \rangle \mid b \in B', c \in C', f_a \langle b, c \rangle \in A' \} \\ &= A' \cap \{ f_a \langle b, c \rangle \mid b \in B', c \in C' \} \\ &= A' \cap \text{image } f_a (B' \times C') \end{aligned}$$

The second set in the product is similar. □

Example 8.13. See Figure 8.2. Let $\text{add}'_c := \text{restrict } \text{add}_c ([0, 1] \times [0, 2])$. By Theorem 8.12,

$$\begin{aligned} \text{preimage } \text{add}'_c [0, \tfrac{1}{2}] &\subseteq ([0, 1] \cap \text{image } \text{add}_a ([0, 2] \times [0, \tfrac{1}{2}])) \times \\ &\quad ([0, 2] \cap \text{image } \text{add}_b ([0, \tfrac{1}{2}] \times [0, 1])) \\ &= ([0, 1] \cap [-2, \tfrac{1}{2}]) \times ([0, 2] \cap [-1, \tfrac{1}{2}]) \\ &= [0, \tfrac{1}{2}] \times [0, \tfrac{1}{2}] \end{aligned}$$

is the smallest rectangular subset of $[0, 1] \times [0, 2]$ that contains the preimage of $[0, \frac{1}{2}]$ under add_c . ◇

At this point, we have an analogue of Lemma 8.2, in that we can compute (approximate) preimages by computing images under (axial) inverses. To compute images using interval

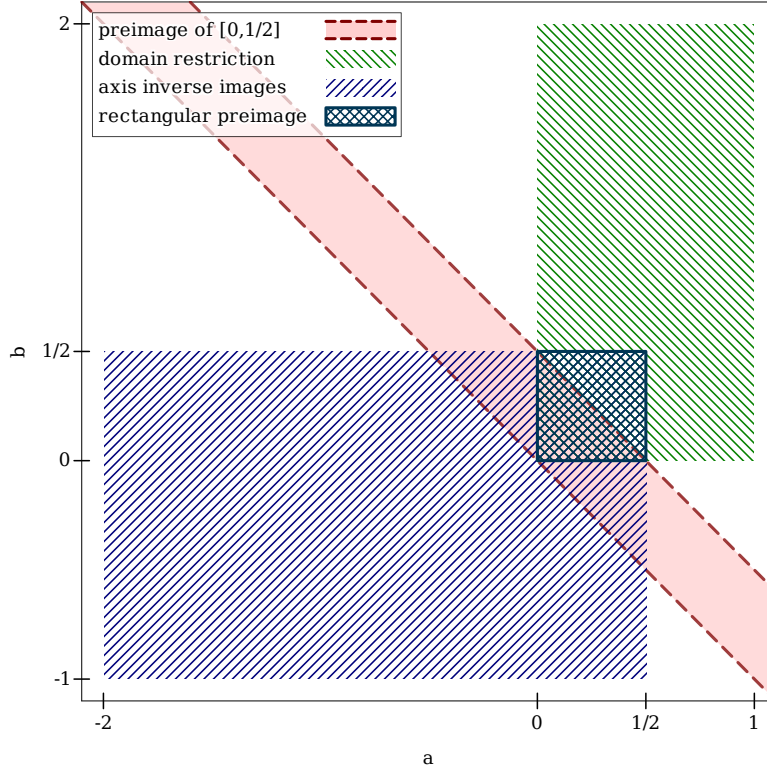


Figure 8.2: Computing an approximate preimage of $[0, \frac{1}{2}]$ under addition restricted to $[0, 1] \times [0, 2]$ (Example 8.13). The preimage is approximated by intersecting the domain with an overapproximation computed using axial inverses.

endpoints, we need analogues of Lemma 8.1 (strictly monotone, surjective functions are invertible and continuous), Theorem 8.5 (images of intervals by endpoints), and Lemma 8.6 (inverse direction).

We first need a notion of properties that hold along an axis for every fixed value of the other argument.

Definition 8.14 (uniform axis property). $f_c : A \times B \rightarrow C$ has property P **uniformly** in its first axis when $P(\text{flip}(\text{curry } f_c) \ b)$ for all $b \in B$, and uniformly in its second axis when $P(\text{curry } f_c \ a)$ for all $a \in A$. If the axis is not specified, P holds uniformly for both.

Now Lemma 8.1's analogue is an easy corollary.

Lemma 8.15. Let $f_c : A \times B \rightarrow C$ for totally ordered A , B and C . If f_c is uniformly surjective and either uniformly strictly increasing or uniformly strictly decreasing in each axis, then f_c

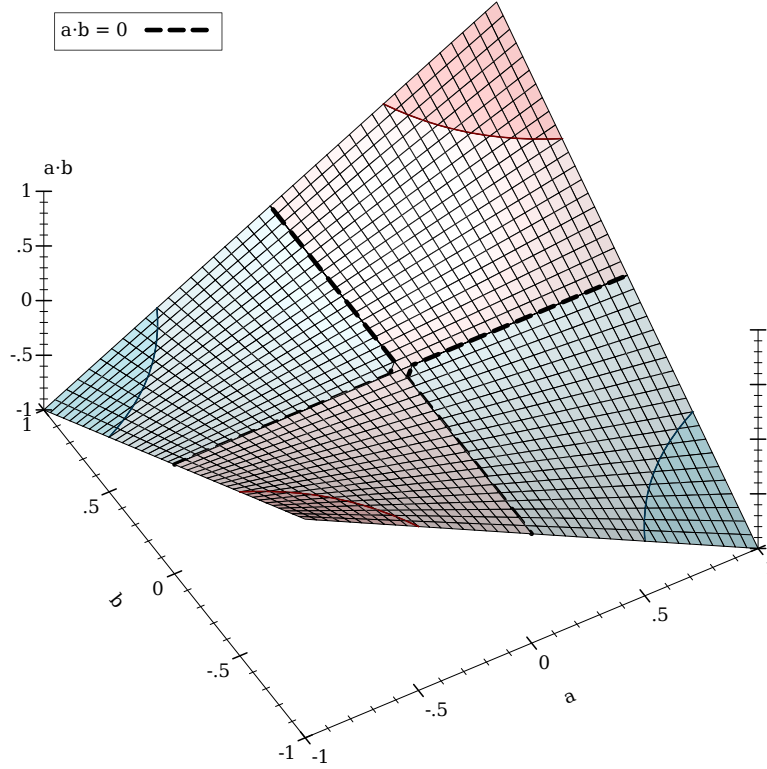


Figure 8.3: Multiplication on $\mathbb{R} \times \mathbb{R}$ is not uniformly surjective, nor uniformly strictly increasing or decreasing in each axis: $\mathbf{a} \cdot \mathbf{0} = 0$ and $\mathbf{0} \cdot \mathbf{b} = 0$ for all \mathbf{a} and \mathbf{b} (Example 8.18). Fortunately, it is uniformly surjective and strictly monotone in each open quadrant.

is axis-invertible; further, it and its axial inverses are continuous.

From here on, assume axis monotonicity properties are uniform unless otherwise stated.

Example 8.16. add_c is uniformly surjective and strictly increasing. sub_c is uniformly surjective and strictly increasing/decreasing in its first/second axis. Therefore, both are axis-invertible. \diamond

Restriction usually makes a function not uniformly surjective.

Example 8.17. Let $\text{add}'_c : [0, 1] \times [0, 1] \rightarrow [0, 2]$, defined by restricting add_c . It is strictly increasing, but not uniformly surjective: the range of $\text{curry } \text{add}'_c \mathbf{0}$ is $[0, 1]$, not $[0, 2]$. \diamond

Fortunately, restriction sometimes does the opposite.

Example 8.18. Define $\text{mul}_c : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ by $\text{mul}_c \langle a, b \rangle := a \cdot b$. It is not uniformly surjective nor strictly monotone because $\text{mul}_c \langle 0, b \rangle = 0$ for all $b \in \mathbb{R}$. (See Figure 8.3.) But $\text{mul}_c^{++} : (0, +\infty) \times (0, +\infty) \rightarrow (0, +\infty)$, and mul_c restricted to the other open quadrants, are uniformly surjective and strictly increasing or decreasing in each axis. \diamond

Theorem 8.5 justifies computing images of intervals with infinite endpoints by applying an extended function to the endpoints. Its two-argument analogue is more involved because extended, two-argument functions may not be defined at every point.

Example 8.19. add_c cannot be extended to $\overline{\text{add}_c} : \overline{\mathbb{R}} \times \overline{\mathbb{R}} \rightarrow \overline{\mathbb{R}}$ in the same way \log is extended to $\overline{\log}$ because

$$\lim_{\langle a', b' \rangle \rightarrow \langle a, b \rangle} \text{add}_c \langle a', b' \rangle \quad (8.8)$$

does not exist when $\langle a, b \rangle$ is $\langle -\infty, +\infty \rangle$ or $\langle +\infty, -\infty \rangle$. \diamond

The previous example suggests that extensions of increasing, two-argument functions are always well-defined except at off-diagonal corners. This is true, and similar statements hold for axes with other directions, and for more restricted domains.

Theorem 8.20 ($\overline{\mathbb{R}} \times \overline{\mathbb{R}}$ extension). *Let A, B, C be open subsets of \mathbb{R} , with $f_c : A \times B \rightarrow C$ uniformly surjective and strictly increasing or decreasing in each axis. Let $\overline{A}, \overline{B}$ and \overline{C} be the closures of A, B and C in $\overline{\mathbb{R}}$. The following extension is well-defined:*

$$\begin{aligned} \overline{f}_c : (\overline{A} \times \overline{B}) \setminus N &\rightarrow \overline{C} \\ \overline{f}_c \langle a, b \rangle &:= \lim_{\langle a', b' \rangle \rightarrow \langle a, b \rangle} f_c \langle a', b' \rangle \end{aligned} \quad (8.9)$$

where $N := \{\langle \min \overline{A}, \max \overline{B} \rangle, \langle \max \overline{A}, \min \overline{B} \rangle\}$ if f_c is increasing or decreasing, and $N := \{\langle \min \overline{A}, \min \overline{B} \rangle, \langle \max \overline{A}, \max \overline{B} \rangle\}$ if f_c is increasing/decreasing or decreasing/increasing.

Proof. Suppose f_c is increasing, and let $g : \mathbb{N} \rightarrow A \times B$ be a sequence of f_c 's domain values.

Any g that converges to $\langle \max \overline{A}, \max \overline{B} \rangle$ has a strictly increasing subsequence. By monotonicity, $\text{map } f_c \ g$ has a strictly increasing subsequence. It is bounded by $\max \overline{C}$, so $\overline{f}_c \langle \max \overline{A}, \max \overline{B} \rangle = \max \overline{C}$. A similar argument proves $\overline{f}_c \langle \min \overline{A}, \min \overline{B} \rangle = \min \overline{C}$.

For any g that converges to $\langle \max \bar{A}, b' \rangle$ for some $b' \in B$, define

$$g' := \text{map } (\lambda \langle a, b \rangle. \langle f_a \langle b', f_c \langle a, b \rangle \rangle', b' \rangle) g \quad (8.10)$$

where f_a is f_c 's first axial inverse, so that $\text{map } f_c g = \text{map } f_c g'$. Because g' has a subsequence that is strictly increasing in the first of each pair, and because the second of each pair is the constant b' , by monotonicity, $\text{map } f_c g'$ has a strictly increasing subsequence. It is bounded by $\max \bar{C}$, so $\bar{f}_c \langle \max \bar{A}, b' \rangle = \max \bar{C}$. By similar arguments, $\bar{f} \langle \min \bar{A}, b \rangle = \min \bar{C}$ and so on.

Arguments for f decreasing, etc., are similar. \square

Following the proof of Theorem 8.20, extensions of two-argument functions can be defined by two corner cases, four border cases, and an interior case.

Example 8.21. Define $\text{pow}_c : (0, 1) \times (0, +\infty) \rightarrow (0, 1)$ by $\text{pow}_c \langle a, b \rangle := \exp(b \cdot \log a)$, which is increasing/decreasing. Its extension to a subset of $\bar{\mathbb{R}} \times \bar{\mathbb{R}}$ is

$$\begin{aligned} \overline{\text{pow}_c} : ([0, 1] \times [0, +\infty]) \setminus N &\rightarrow [0, 1] \\ \overline{\text{pow}_c} \langle a, b \rangle &:= \text{case } \langle a, b \rangle \\ &\quad \langle 0, +\infty \rangle \longrightarrow 0 \\ &\quad \langle 1, 0 \rangle \longrightarrow 1 \\ &\quad \langle 0, b \rangle \longrightarrow 0 \\ &\quad \langle 1, b \rangle \longrightarrow 1 \\ &\quad \langle a, 0 \rangle \longrightarrow 1 \\ &\quad \langle a, +\infty \rangle \longrightarrow 0 \\ &\quad \text{else} \longrightarrow \text{pow}_c \langle a, b \rangle \end{aligned} \quad (8.11)$$

where $N := \{\langle 0, 0 \rangle, \langle 1, +\infty \rangle\}$. \diamond

The analogue of Theorem 8.5 is easiest to state if we have predicates that indicate a function's direction in each axis. Define $\text{inc}_1 : (A \times B \rightarrow C) \Rightarrow \text{Bool}$ so that $\text{inc}_1 f$ if and only if f is strictly increasing in its first axis, and similarly inc_2 so that $\text{inc}_2 f$ if and only if f is strictly increasing in its second axis.

Theorem 8.22 (images of rectangles by interval endpoints). *Let A, B, C be open subsets of \mathbb{R} , with $f_c : A \times B \rightarrow C$ uniformly surjective and strictly increasing or decreasing in*

each axis, with extension \bar{f}_c as defined in Theorem 8.20. If $A' := [a_1, a_2, \alpha_1, \alpha_2] \subseteq A$ and $B' := [b_1, b_2, \beta_1, \beta_2] \subseteq B$, then the image C' under f_c is

$$\begin{aligned}
& \text{image } f_c ([a_1, a_2, \alpha_1, \alpha_2] \times [b_1, b_2, \beta_1, \beta_2]) \\
&= \text{let } \langle a_1, a_2, \alpha_1, \alpha_2 \rangle := \text{cond } (\text{inc}_1 f_c) \longrightarrow \langle a_1, a_2, \alpha_1, \alpha_2 \rangle \\
&\quad \text{else} \longrightarrow \langle a_2, a_1, \alpha_2, \alpha_1 \rangle \\
&\quad \langle b_1, b_2, \beta_1, \beta_2 \rangle := \text{cond } (\text{inc}_2 f_c) \longrightarrow \langle b_1, b_2, \beta_1, \beta_2 \rangle \\
&\quad \text{else} \longrightarrow \langle b_2, b_1, \beta_2, \beta_1 \rangle \\
&\text{in } [\bar{f}_c \langle a_1, b_1 \rangle, \bar{f}_c \langle a_2, b_2 \rangle, \alpha_1 \text{ and } \beta_1, \alpha_2 \text{ and } \beta_2]
\end{aligned} \tag{8.12}$$

Proof. Because f_c is continuous and $A' \times B'$ is a connected set, C' is a connected set, which in \mathbb{R} is an interval. Thus, we need to determine only its bounds and whether it contains each endpoint.

Suppose f_c is increasing. By monotonicity, C' is contained in $[\bar{f}_c \langle a_1, b_1 \rangle, \bar{f}_c \langle a_2, b_2 \rangle]$. If either α_1 or β_1 is false, C' cannot contain $\bar{f}_c \langle a_1, b_1 \rangle$. If either α_2 or β_2 is false, C' cannot contain $\bar{f}_c \langle a_2, b_2 \rangle$. Therefore $C' = [\bar{f}_c \langle a_1, b_1 \rangle, \bar{f}_c \langle a_2, b_2 \rangle, \alpha_1 \text{ and } \beta_1, \alpha_2 \text{ and } \beta_2]$.

We still must prove $\langle a_1, b_1 \rangle$ and $\langle a_2, b_2 \rangle$ are in \bar{f}_c 's domain. First, recall $\bar{f}_c : (\bar{A} \times \bar{B}) \setminus N \rightarrow \bar{C}$, where \bar{A} , \bar{B} and \bar{C} are the closures of A , B and C in \mathbb{R} , and $N = \{\langle \min \bar{A}, \max \bar{B} \rangle, \langle \max \bar{A}, \min \bar{B} \rangle\}$.

Because $A' \subseteq A$ and $B' \subseteq B$, and A and B are open sets, $a_1 \neq \max \bar{A}$, $a_2 \neq \min \bar{A}$, $b_1 \neq \max \bar{B}$, and $b_2 \neq \min \bar{B}$, so

$$\begin{aligned}
\langle a_1, b_1 \rangle &\neq \langle \max \bar{A}, b \rangle \quad \text{for all } b \in \bar{B} \\
\langle a_1, b_1 \rangle &\neq \langle a, \max \bar{B} \rangle \quad \text{for all } a \in \bar{A} \\
\langle a_2, b_2 \rangle &\neq \langle \min \bar{A}, b \rangle \quad \text{for all } b \in \bar{B} \\
\langle a_2, b_2 \rangle &\neq \langle a, \min \bar{B} \rangle \quad \text{for all } a \in \bar{A}
\end{aligned} \tag{8.13}$$

Therefore, $\langle a_1, b_1 \rangle \notin N$ and $\langle a_2, b_2 \rangle \notin N$, as desired.

The remaining cases for f_c are similar. □

Example 8.23. Because $\text{inc}_1 \text{ pow}_c$ and not $(\text{inc}_2 \text{ pow}_c)$,

$$\begin{aligned}
& \text{image } \text{pow}_c ((0, \tfrac{1}{2}] \times [2, +\infty)) \\
&= \text{let } \langle a_1, a_2, \alpha_1, \alpha_2 \rangle := \langle 0, \tfrac{1}{2}, \text{false}, \text{true} \rangle \\
&\quad \langle b_1, b_2, \beta_1, \beta_2 \rangle := \langle +\infty, 2, \text{false}, \text{true} \rangle \\
&\quad \text{in } [\overline{\text{pow}_c} \langle a_1, b_1 \rangle, \overline{\text{pow}_c} \langle a_2, b_2 \rangle, \alpha_1 \text{ and } \beta_1, \alpha_2 \text{ and } \beta_2] \\
&= [\overline{\text{pow}_c} \langle 0, +\infty \rangle, \overline{\text{pow}_c} \langle \tfrac{1}{2}, 2 \rangle, \text{false and false, true and true}] \\
&= [0, \tfrac{1}{4}, \text{false}, \text{true}] \\
&= (0, \tfrac{1}{4}] \quad \diamond
\end{aligned}$$

To use Theorem 8.22 to compute approximate preimages under f_c by computing images under its axial inverses, we must know whether each axis of f_a and f_b is increasing or decreasing. It helps to have an analogue of Lemma 8.6 (inverse direction).

Theorem 8.24. *Let $f_c : A \times B \rightarrow C$ be uniformly surjective and strictly increasing or decreasing in each axis, with axial inverses f_a and f_b . The following statements hold:*

1. $\text{inc}_1 f_a$ if and only if $(\text{inc}_1 f_c) \text{ xor } (\text{inc}_2 f_c)$.
2. $\text{inc}_2 f_a$ if and only if $\text{inc}_1 f_c$.

Proof. For statement 1, let $c \in C$, $b_1, b_2 \in B$, $a_1 := f_a \langle b_1, c \rangle$ and $a_2 := f_a \langle b_2, c \rangle$. Let $c' := f_c \langle a_1, b_2 \rangle$; note $c = f_c \langle a_1, b_1 \rangle = f_c \langle a_2, b_2 \rangle$. Suppose $\text{inc}_1 f_c$ and $\text{inc}_2 f_c$; then $a_1 > a_2 \iff c < c'$ and $b_1 < b_2 \iff c < c'$, so $b_1 < b_2 \iff a_1 > a_2$. The remaining cases are similar.

For statement 2, fix $b \in B$ and apply Lemma 8.6. \square

By Lemma 8.10, $\text{inc}_1 f_b \iff \text{inc}_2 f_c$, and $\text{inc}_2 f_b \iff \text{inc}_1 f_a$. We can therefore easily determine the uniform directions of f_a 's and f_b 's axes from the uniform directions of f_c 's axes.

We are certain the preceeding definitions and theorems extend naturally to functions with any number of arguments, but have not needed to extend them yet.

8.2.3 Discontinuous Primitives

8.3 Primitive Implementation

Because floating-point functions are defined on subsets of $\overline{\mathbb{R}}$, it would seem we could compute preimages under strictly monotone, real functions by applying their floating-point counterparts to interval endpoints. This is mostly true, but we must take care to round in the right directions and account for floating-point negative zero.

$$\begin{aligned} \overline{\text{pos-recip}} : [0, +\infty] &\rightarrow [0, +\infty] \\ \overline{\text{pos-recip}} \ a &= \text{cond } \begin{array}{ll} 0 < a < +\infty & \longrightarrow 1/a \\ a = 0 & \longrightarrow +\infty \\ a = +\infty & \longrightarrow 0 \end{array} \end{aligned} \quad (8.14)$$

8.4 Partitioned Sampling

Suppose that we want to sample values in a probability space X, P by first sampling a set from a *partition* of X and then sampling from that set.¹ [XXX: this transition should mention why we want to do partition sampling instead of just supposing that we do]

First, we define

$$\begin{aligned} \text{condition} : \text{Set } X \rightarrow [0, 1] &\Rightarrow \text{Set } X \Rightarrow \text{Set } X \rightarrow [0, 1] \\ \text{condition } P \ A &:= \lambda A' \in \text{domain } P. P \ (A' \cap A) / P \ A \end{aligned} \quad (8.15)$$

to restrict a probability measure P to a measurable, positive-probability set A and renormalize it.

Definition 8.25 (partitioned sampling). *Let X, P be an arbitrary probability space, N be an at-most-countable index set, and $s : N \rightarrow \text{Set } X$ be a partition of X into $|N|$ measurable parts. The following procedure samples from X :*

1. Choose $n \in N$ with probability $P \ (s \ n)$.

¹This is not *stratified* sampling, which samples a fixed number of times from each partition.

2. Choose $a \in s_n$ according to condition $P(s_n)$.

It is not hard to show that partitioned sampling chooses an $a \in X$ according to P .

Example 8.26 (partitioned sampling from a standard normal). Let P be the standard normal distribution's probability measure. To sample according to P , let $N := \{\text{neg}, \text{pos}\}$ and $s = [\text{neg} \mapsto (-\infty, 0], \text{pos} \mapsto (0, +\infty)]$, and define $Q : N \rightarrow \text{Set } \mathbb{R} \rightarrow [0, 1]$ by

$$\begin{aligned} Q \text{ neg } A &= P((-\infty, 0] \cap A) / \frac{1}{2} \\ Q \text{ pos } A &= P((0, +\infty) \cap A) / \frac{1}{2} \end{aligned} \tag{8.16}$$

Then

1. Choose $n = \text{neg}$ or $n = \text{pos}$, each with probability $\frac{1}{2}$.
2. Choose $a \in s_n$ according to Q_n . ◇

Partitioned sampling has two weaknesses. First, it requires $P(s_n)$ to be easy to compute for all $n \in N$. If this were true, we would not need to sample in the first place—i.e. it assumes a solution to the overall problem we are trying to solve. Second, it assumes sampling according to condition $P(s_n)$ is easy, which is also not reasonable, as sampling according to a conditioned distribution is a subproblem we are trying to solve.

But suppose we could easily sample a partition index according to a different distribution over N , and according to a different distribution over part s_n for each $n \in N$. Doing so and returning weighted samples to adjust for the differences in distribution comprises *partitioned importance sampling*.

First, we define

$$\begin{aligned} \text{subcond} : \text{Set } X \rightarrow [0, 1] &\Rightarrow \text{Set } X \Rightarrow \text{Set } X \rightarrow [0, 1] \\ \text{subcond } P \ A &:= \lambda A' \in \text{domain } P. P(A' \cap A) \end{aligned} \tag{8.17}$$

to restrict a probability measure P to a measurable set A *without* renormalizing it; i.e. it returns a subprobability measure.

Definition 8.27 (partitioned importance sampling). *Suppose we have*

- An arbitrary probability space X, P .
- An at-most-countable index set N .
- A probability mass function $p : N \rightarrow [0, 1]$ such that $p\ n > 0$ for all $n \in N$.
- A partition $s : N \rightarrow \text{Set } X$ of X into $|N|$ measurable parts.
- Candidate probability measures $Q : N \rightarrow \text{Set } X \rightarrow [0, 1]$, one for each partition.

To sample from X according to P ,

1. Choose $n \in N$ with probability $p\ n$.
2. Choose $a \in X$ according to $Q\ n$.
3. Compute $w := \frac{1}{p\ n} \cdot \text{diff}^+ (\text{subcond } P\ (s\ n))\ (Q\ n)\ a$.
4. Return the weighted sample $\langle a, w \rangle$.

The function $\text{diff}^+ (\text{subcond } P\ (s\ n))\ (Q\ n)$, with type $X \rightarrow [0, +\infty)$, is a **Radon-Nikodým² derivative**. If P has density f , $Q\ n$ has density g with respect to the same base measure (e.g. same-dimensional length, area or volume), and $a \in s\ n$ implies $g\ a > 0$, then³

$$\text{diff}^+ (\text{subcond } P\ (s\ n))\ (Q\ n)\ a = \text{if } (a \in s\ n) (f\ a / g\ a) 0 \quad (8.18)$$

Chapter 10 has definitions and more details. We use diff^+ in a more general sense, but in this section, it is usually fine to think of its return values as quotients of densities.

An importance sampling algorithm is correct when all expected values computed using its weighted samples are the equal to the true expected values. The next theorem states that this is true of partitioned importance sampling under reasonable conditions, which are analogous to the support of $\text{subcond } P\ (s\ n)$ being no larger than that of $Q\ n$.

Theorem 8.28 (partitioned importance sampling correctness). *Let X, P, N, s, p , and Q as in Definition 8.27 (partitioned importance sampling) such that $\text{subcond } P\ (s\ n) \ll Q\ n$ for all $n \in N$. Define $P_N : \text{Set } N \rightarrow [0, 1]$ by extending p to a measure.*

²Pronounced “RADon neekohDIM,” and named after Austrian mathematician Johann Radon and Polish mathematician Otto Nikodým.

³The equality in (8.18) holds $(Q\ n)$ -almost everywhere.

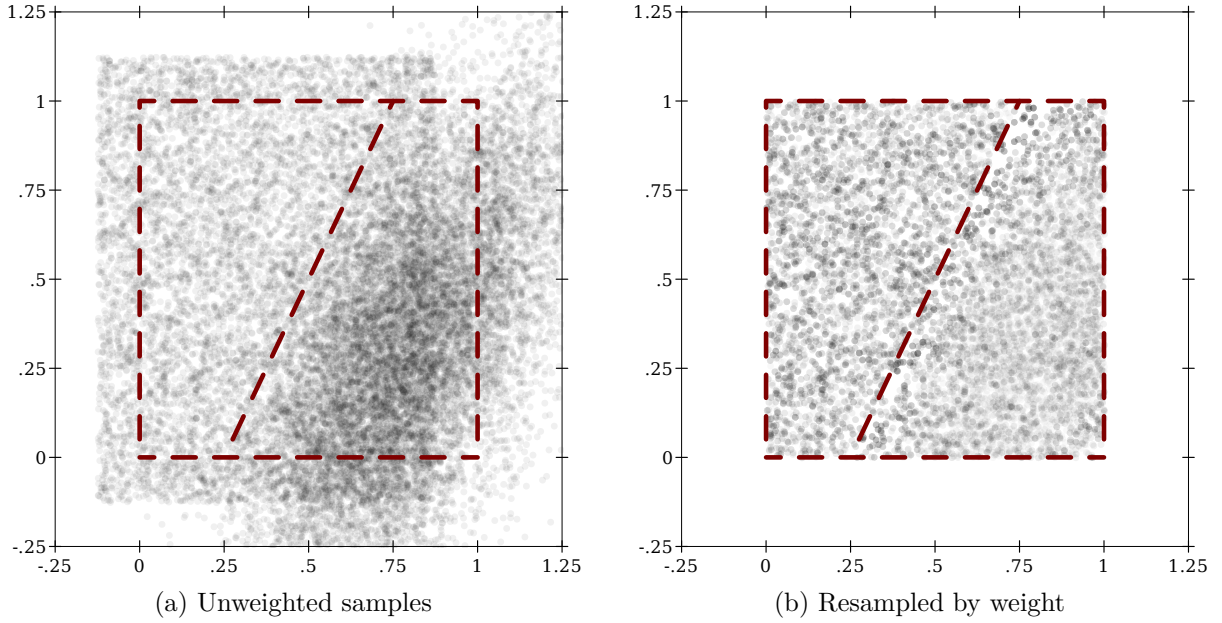


Figure 8.4: Partitioned importance sampling used to sample uniformly in a partition of the unit square, using two overlapping candidate distributions.

If $g : X \rightarrow \mathbb{R}$ is a P -integrable mapping, and

$$\begin{aligned}
 g' : N \times X &\rightarrow \mathbb{R} \\
 g' \langle n, a \rangle &:= g(a) \cdot \frac{1}{p_n} \cdot \text{diff}^+ (\text{subcond } P(s_n)) (Q_n)(a)
 \end{aligned} \tag{8.19}$$

then $\int g' (P_N \times Q) (N \times X) = \int g P X$.

Proof. See Chapter 10. □

Partitioned importance sampling allows quite a lot of freedom: parts can be chosen with arbitrary nonzero probability, and each part can have its own candidate distribution, which may be defined on a superset of the part.

Example 8.29 (2D partitioned importance sampling). Figure 8.4 shows the result of partitioned importance sampling in a partition of the unit square. In this instance,

- $X := [0, 1] \times [0, 1]$ and P is the uniform measure on X (i.e. area).
- $N := \{\text{left}, \text{right}\}$.
- $p := [\text{left} \mapsto 0.4, \text{right} \mapsto 0.6]$.

- $s \text{ left} = \{\langle x, y \rangle \in X \mid y > 2 \cdot x - \frac{1}{2}\}$, similarly for $s \text{ right}$.
- $Q \text{ left}$ is the uniform measure on a superset of $s \text{ left}$, and $Q \text{ right}$ is a multivariate Gaussian measure centered at $\langle 0.8, 0.3 \rangle$.

The implementation does not actually construct most of these objects. It constructs

- A density function $f : \mathbb{R} \times \mathbb{R} \Rightarrow [0, +\infty)$ to represent P .
- A family of predicates $s? : N \Rightarrow X \Rightarrow \text{Bool}$ to decide $a \in s \text{ n}$.
- Candidate densities $g : N \Rightarrow X \Rightarrow [0, +\infty)$ to represent Q .

It computes weights using $\text{diff}^+ (\text{subcond } P (s \text{ n})) (Q \text{ n}) a = f a / g n a$ when $s? n a = \text{true}$.

It directly represents only N and p , but we will find even this to be infeasible shortly. \diamond

Two properties make the preceding example simple. First, the partition has finitely many parts. Second, the measures $\text{subcond } P (s \text{ n})$ and $Q \text{ n}$ have densities with respect to the same base measure, which ensures $\text{diff}^+ (\text{subcond } P (s \text{ n})) (Q \text{ n})$ exists and is easy to compute.

When sampling in the domain of programs, neither property holds in general.

8.4.1 Partitioning Probabilistic Program Domains

For the random source part $R := J \rightarrow [0, 1]$ of probabilistic language domains, which consists of infinite binary trees of reals, it is not clear that Theorem 8.28 is applicable. The main problem is that infinite-dimensional Radon-Nikodým derivatives do not generally exist.

Fortunately, they *can* exist if the two measures differ in only finitely many axes. More precisely, let $P_1 : \text{Set } R \rightarrow [0, 1]$ and $P_2 : \text{Set } R \rightarrow [0, 1]$ be probability measures, and $J' \subseteq J$ be a finite set of tree indexes. Suppose P_1 can be factored into a distribution P'_1 over $J' \rightarrow [0, 1]$ and a distribution over $(J \setminus J') \rightarrow [0, 1]$, and P'_2 can be similarly factored into P'_2 and *the same* distribution over $(J \setminus J') \rightarrow [0, 1]$. Then, under reasonable conditions (which are analogous to the support of P'_1 being no larger than that of P'_2), $\text{diff}^+ P_1 P_2$ exists and can be computed using $\text{diff}^+ P'_1 P'_2$. Chapter 10 contains a formal statement of this fact and a proof.

To ensure $\text{subcond } P (s \text{ n})$ and $Q \text{ n}$ differ in only finitely many axes, we partition R according to branch traces. Each branch trace corresponds with a program that reads any

$r \in R$ at only finitely many indexes $J' \subseteq J$. [XXX: connect idea better]

In the remainder of this subsection, assume a fixed program p . Let $f := \llbracket p \rrbracket_{\perp}^{\downarrow} : \langle\langle R, T \rangle, \langle \rangle \rangle \rightsquigarrow_{\perp}^* Y$ be its interpretation as a bottom* arrow computation, with maximal domain A^* . Define $T^* := \text{image}(\text{fst} \ggg \text{snd}) A^*$ as its *maximal branch trace set* and $R^* := \text{image}(\text{fst} \ggg \text{fst}) A^*$ as its *maximal random source set*.

We need a notion of the random sources that agree with a given branch trace $t \in T$; i.e. those $r \in R$ for which $f \langle\langle r, t \rangle, \langle \rangle \rangle \neq \perp$.

Definition 8.30 (induced random sources). *Let $t \in T$ be a branch trace. The **random sources induced by t** are a subset of R defined by $R' := \{r \in R \mid \langle\langle r, t \rangle, \langle \rangle \rangle \in A^*\}$.*

[XXX: graph of induced partition from program if (random < random) true false?]

Using T as the partition index set and defining the partition's parts as induced random sources *almost* works, in the sense that the required Radon-Nikodým derivatives exist. Unfortunately, we cannot use T or T^* as the partition index set because many branch traces can induce the same random sources.

Example 8.31. For the program

$$\text{if } (\text{random} < p) \ 0 \ \text{random} \tag{8.20}$$

there are at least two branch traces in the program's maximal domain: $t_0 := [j_0 \mapsto \text{true}, * \mapsto \perp]$ and $t_1 := [j_0 \mapsto \text{false}, * \mapsto \perp]$. There is also $[j_0 \mapsto \text{false}, \text{left } j_0 \mapsto \text{true}, * \mapsto \perp]$, because it agrees with every execution that $[j_0 \mapsto \text{false}, * \mapsto \perp]$ agrees with. In fact, there are infinitely many branch traces in T^* that induce the same random sources as either t_0 or t_1 . [XXX: tree diagrams?] ◇

We need to find a subset of branch traces whose induced random sources are disjoint. The main idea is to define equivalence classes of branch traces that induce the same random sources, and use the “smallest” branch trace in each class as a part index.

First, we need to be sure that such equivalence classes can induce a partition.

Theorem 8.32. *Let $t_1, t_2 \in T^*$. If t_1 induces R'_1 and t_2 induces R'_2 , then $R'_1 = R'_2$ or $R'_1 \cap R'_2 = \emptyset$.*

Proof. XXX: do this □

To identify the “smallest” trace in each class, we must define an ordering over them. One fairly natural way is to say a branch trace is smaller than another when it describes fewer branch decisions; i.e. its tree has fewer non- \perp elements. Two branch traces that differ by returning respectively **true** and **false** for the same j may represent different execution paths, so they must be incomparable.

Definition 8.33 (branch trace partial order). $t_1 \leq t_2$ when for all $j \in J$, if $t_1 j \neq t_2 j$, then $t_1 j = \perp$.

If T^* is partitioned into equivalence classes of traces that induce the same random sources, each part in the partition contains a least member with respect to (\leq) .

Theorem 8.34. *Let $t \in T^*$ induce R' , and let T' be the largest subset of T^* that induces R' . T' has a least member t_* .*

Proof. Define $t_* \in T$ by

$$\begin{aligned} t_* &:= \text{cond } \forall t' \in T'. t' j = \text{true} \longrightarrow \text{true} \\ &\quad \forall t' \in T'. t' j = \text{false} \longrightarrow \text{false} \\ &\quad \text{else} \longrightarrow \perp \end{aligned} \tag{8.21}$$

Thus, $t_* \leq t'$ for all $t' \in T'$, or t_* is a lower bound for T' .

Let $r \in R'$. By construction, every conditional subcomputation at index $j \in J$ agrees with $t_* j$. Therefore $f \langle \langle r, t_* \rangle, \langle \rangle \rangle \neq \perp$, so $\langle \langle r, t_* \rangle, \langle \rangle \rangle \in A^*$, and thus $t_* \in T'$. □

An easy consequence is $T' = \{t' \in T^* \mid t_* \leq t'\}$. [XXX: is it easy?] Thus, each least member represents a set of larger branch traces that induce the same set of random sources.

We can get our sought-after index set by defining the set of smallest branch traces.

Definition 8.35 (minimal branch traces). *The set of **minimal branch traces** T_* is the set of minimal elements in T^* , or*

$$T_* := \{t_1 \in T^* \mid \forall t_2 \in T^*. t_2 \leq t_1 \implies t_2 = t_1\} \quad (8.22)$$

Theorem 8.36. T_* induces a partition of R^* .

Proof. XXX: todo □

We can infer from Theorem 8.36 that a program’s minimal branch trace set T_* contains only the actual branches taken when running the program on every $r \in R^*$. Therefore, one way to sample from T_* with the correct probability—at least, for programs that halt with probability 1—would be to choose an $r \in R$ uniformly, and run the program on r while recording each branch decision.

But this sampling scheme has problems similar to those of partitioned sampling (Definition 8.25). First, it assumes the probabilities of branch traces, which are the probabilities of the R^* subsets they induce, are easy to compute. Second, we are interested in sampling from an *arbitrarily low-probability subset* of R^* , which may be covered by the partition induced by an *arbitrarily low-probability subset* of T_* .

It appears we have a chicken-and-egg problem, in that

1. Sampling in a small subset of R^* requires sampling in a small subset of T_* .
2. Sampling in a small subset of T_* requires sampling in a small subset of R^* .

Fortunately, if we allow ourselves subsets of a slightly larger set than T_* , and allow ourselves to sample within overapproximating *covers* of R^* subsets, we can use approximate preimage computation to sample from T_* and R^* subsets simultaneously.

8.4.2 Approximate Partitions of Probabilistic Program Domains

We first define the set of branch traces that is slightly larger than T_* .

The idea is to define a set of *feasible* branch traces T_+ that is derived only from the

shape of a program, not its actual executions. We must then ensure that every additional branch trace (i.e. every $t \in T_+ \setminus T_*$) induces \emptyset , so that T_+ induces a partition. After that, we define an algorithm for sampling from T_+ , which does not require running a probabilistic program on any $r \in R$, and prove the algorithm correct. We can then extend this algorithm to use preimage computation to sample in arbitrarily good approximations of small subsets of T_* and R^* .

Defining T_+ in terms of a program's abstract branching shape requires an additional nonstandard interpretation. Figure 8.5a defines the **branch index** arrow. Its type is $J \Rightarrow \text{Idxs}$, which does not refer to a domain or codomain type of program values because its computations do not receive or compute program values. Instead, they build lazy trees of possible branching decisions, ignoring the actual values of if conditions. For example, lifted, pure functions are interpreted as $\lambda j. \langle \rangle$, which takes the function's computation index and returns no decisions. Composition and pairing of subcomputations k_1 and k_2 both return $\langle k_1 \text{ (left } j), k_2 \text{ (right } j) \rangle$: a node with two children that contain the feasible branch decisions in their subcomputations.

Only $\text{ifte}_{\text{idxs}^*}$ does more than simple structural recursion: it returns $\text{if-idxs } j \text{ idxs}_2 \text{ idxs}_3$ to represent a decision at computation index j . The children idxs_2 and idxs_3 are lazy, abstract representations of the if's branches. Like a concrete execution, a branch trace sampler is expected to compute and recur through only one of them.

Figure 8.5b defines the function **traces**, which transforms these lazy trees of type Idxs into sets of branch traces. To ensure **trace** always terminates, it is defined using a depth-limited helper function, which it calls with every depth $n \in \mathbb{N}$, and collects the results in a countable union. It is thus unimplementable, but we will use it only to precisely define the range of a random value.

We define a program's feasible branch traces using the standard first-order semantic function $\llbracket \cdot \rrbracket_a^\downarrow$ with arrow $a = \text{idxs}^*$.

Definition 8.37 (feasible branch traces). *A program p 's **feasible branch traces** are those*

$$\begin{aligned}
\text{Idxs} &::= \langle \rangle \mid \langle \text{Idxs}, \text{Idxs} \rangle \\
&\mid \text{if-idxs } J \ (1 \Rightarrow \text{Idxs}) \ (1 \Rightarrow \text{Idxs}) \\
\text{arr}_{\text{idxs}^*} &: (x \Rightarrow y) \Rightarrow (J \Rightarrow \text{Idxs}) \\
\text{arr}_{\text{idxs}^*} \ f \ j &:= \langle \rangle \\
(\ggg_{\text{idxs}^*}) &: (J \Rightarrow \text{Idxs}) \Rightarrow (J \Rightarrow \text{Idxs}) \Rightarrow (J \Rightarrow \text{Idxs}) \\
(k_1 \ggg_{\text{idxs}^*} k_2) \ j &:= \langle k_1 \ (\text{left } j), k_2 \ (\text{right } j) \rangle \\
(\&\&\&_{\text{idxs}^*}) &: (J \Rightarrow \text{Idxs}) \Rightarrow (J \Rightarrow \text{Idxs}) \Rightarrow (J \Rightarrow \text{Idxs}) \\
(k_1 \&\&\&_{\text{idxs}^*} k_2) \ j &:= \langle k_1 \ (\text{left } j), k_2 \ (\text{right } j) \rangle \\
\text{ifte}_{\text{idxs}^*} &: (J \Rightarrow \text{Idxs}) \Rightarrow (J \Rightarrow \text{Idxs}) \Rightarrow (J \Rightarrow \text{Idxs}) \Rightarrow (J \Rightarrow \text{Idxs}) \\
\text{ifte}_{\text{idxs}^*} \ k_1 \ k_2 \ k_3 \ j &:= \text{let } \text{idxs}_2 := \lambda 0. k_2 \ (\text{left } (\text{right } j)) \\
&\quad \text{idxs}_3 := \lambda 0. k_3 \ (\text{right } (\text{right } j)) \\
&\quad \text{in } \langle k_1 \ j, \text{if-idxs } j \ \text{idxs}_2 \ \text{idxs}_3 \rangle \\
\text{lazy}_{\text{idxs}^*} &: (1 \Rightarrow (J \Rightarrow \text{Idxs})) \Rightarrow (J \Rightarrow \text{Idxs}) \\
\text{lazy}_{\text{idxs}^*} \ k \ j &:= k \ 0 \ j \\
\text{random}_{\text{idxs}^*} &: J \Rightarrow \text{Idxs} \\
\text{random}_{\text{idxs}^*} \ j &:= \langle \rangle
\end{aligned}$$

(a) Branch index arrow. Computations return a lazy tree of type Idxs , of feasible branch decisions, ignoring the actual values of if conditions. The arrow is directly implementable in any λ -calculus.

$$\begin{aligned}
\text{traces} &: \text{Idxs} \Rightarrow \text{Set } T \\
\text{traces } \text{idxs} &:= \bigcup_{n \in \mathbb{N}} \text{traces}^* \ n \ \text{idxs} \ [* \mapsto \perp] \\
\text{traces}^* &: \mathbb{N} \Rightarrow \text{Idxs} \Rightarrow T \Rightarrow \text{Set } T \\
\text{traces}^* \ n \ \langle \rangle \ t &:= \{t\} \\
\text{traces}^* \ n \ \langle \text{idxs}_1, \text{idxs}_2 \rangle \ t &:= \text{let } T' := \text{traces}^* \ n \ \text{idxs}_1 \ t \\
&\quad \text{in } \bigcup_{t' \in T'} \text{traces}^* \ n \ \text{idxs}_2 \ t' \\
\text{traces}^* \ 0 \ (\text{if-idxs } j \ \text{idxs}_2 \ \text{idxs}_3) \ t &:= \{t\} \\
\text{traces}^* \ n \ (\text{if-idxs } j \ \text{idxs}_2 \ \text{idxs}_3) \ t &:= \text{traces}^* \ (n-1) \ (\text{idxs}_2 \ 0) \ t[j \mapsto \text{true}] \cup \\
&\quad \text{traces}^* \ (n-1) \ (\text{idxs}_3 \ 0) \ t[j \mapsto \text{false}] \cup \\
&\quad \{t\}
\end{aligned} \tag{8.23}$$

(b) The λ_{ZFC} function `traces` turns a lazy tree of feasible branch decisions into a set of feasible branch traces.

$$\begin{aligned}
\text{sample-trace} &: \text{Idxs} \rightarrow \langle \mathbb{R}, T \rangle \\
\text{sample-trace } \text{idxs} &:= \text{sample-trace}^* \ \text{idxs} \ \langle 1, [* \mapsto \perp] \rangle \\
\text{sample-trace}^* &: \text{Idxs} \Rightarrow \langle \mathbb{R}, T \rangle \Rightarrow \langle \mathbb{R}, T \rangle \\
\text{sample-trace}^* \ \langle \rangle \ pt &:= pt \\
\text{sample-trace}^* \ \langle \text{idxs}_1, \text{idxs}_2 \rangle \ pt &:= \text{let } pt' := \text{sample-trace}^* \ \text{idxs}_1 \ pt \\
&\quad \text{in } \text{sample-trace}^* \ \text{idxs}_2 \ pt' \\
\text{sample-trace}^* \ (\text{if-idxs } j \ \text{idxs}_2 \ \text{idxs}_3) \ \langle p_t, t \rangle &:= \text{let } \langle p_b, b \rangle := \text{sample-branch } \{\text{true}, \text{false}, \perp\} \\
&\quad pt' := \langle p_t \cdot p_b, t[j \mapsto b] \rangle \\
&\quad \text{in case } b \\
&\quad \text{true} \longrightarrow \text{sample-trace}^* \ (\text{idxs}_2 \ 0) \ pt' \\
&\quad \text{false} \longrightarrow \text{sample-trace}^* \ (\text{idxs}_3 \ 0) \ pt' \\
&\quad \perp \longrightarrow pt'
\end{aligned} \tag{8.24}$$

(c) The stochastic function `sample-trace` samples t from the set returned by `traces`, and returns t and its probability. It is directly implementable in any λ -calculus with probabilistic choice.

Figure 8.5: Branch index collecting semantics.

in the set $T_+ := \text{traces } (\llbracket p \rrbracket_{\text{idxs}^*}^\downarrow j_0)$.

If we are to sample $t \in T_+$ and then sample within the R' induced by t , we must ensure that T_+ includes T_* and induces a partition.

Theorem 8.38. *Let $t \in T$ induce R' . $R' \neq \emptyset$ if and only if $t \in T^*$. [XXX: needed?]*

Proof. $R' \neq \emptyset$ if and only if there is an $r \in R'$ such that $\langle \langle r, t \rangle, \langle \rangle \rangle \in A^*$. □

Theorem 8.39. *For all $t \in T_+$, either $t \in T_*$ or t induces \emptyset .*

Proof. XXX: todo □

Corollary 8.40. T_+ induces a partition of R^* .

XXX: now we define a stochastic procedure for sampling in T_+

XXX: apologize for being less precise (e.g. stochastic procedures look like λ_{ZFC} but aren't; we assume every operation is lifted to receive a random source as input, etc.)

XXX: suppose we have a stochastic procedure $\text{sample-branch} : \text{Set Bool}_\perp \Rightarrow \langle \mathbb{R}, \text{Bool}_\perp \rangle$, where $\text{sample-branch } B$ returns any member of B with some nonzero, constant probability

Figure 8.5c defines sample-trace , a stochastic procedure defined in terms of sample-branch . Given an $\text{idxs} : \text{Idxs}$, it returns a member of traces idxs and its probability. In other words, for a program p , the probability that $\langle p_t, t \rangle = \text{sample-trace } (\llbracket p \rrbracket_{\text{idxs}^*}^\downarrow j_0)$ is p_t , and $t \in T_+$. Further, every $t \in T_+$ has a constant, nonzero probability of being chosen.

Of course, we must prove these facts.

XXX: move the following soundness and completeness proofs to Chapter 10, or define T_+ in terms of sample-trace to make the proofs unnecessary; whether the latter will work (or whether it will work better) will probably become clear when I do the missing proof above

Theorem 8.41 (sample-trace* soundness). *For all $t \in T$ and $\text{idxs} : \text{Idxs}$, there exists an $n \in \mathbb{N}$ such that $\text{snd } (\text{sample-trace}^* \text{ idxs } \langle p_t, t \rangle) \in \text{traces}^* n \text{ idxs } t$.*

Proof. By structural induction on Idxs .

Case $\text{idxs} = \langle \rangle$. For all $n \in \mathbb{N}$, these two statements are equivalent by substituting definitions of sample-trace^* and traces^* :

$$\begin{aligned} \text{snd} (\text{sample-trace}^* \langle \rangle \langle p_t, t \rangle) &\in \text{traces}^* n \langle \rangle t \\ t &\in \{t\} \end{aligned} \tag{8.25}$$

Case $\text{idxs} = \langle \text{idxs}_1, \text{idxs}_2 \rangle$. Let

$$\begin{aligned} \langle p'_t, t' \rangle &:= \text{sample-trace}^* \text{idxs}_1 \langle p_t, t \rangle \\ \langle p''_t, t'' \rangle &:= \text{sample-trace}^* \text{idxs}_2 \langle p'_t, t' \rangle \end{aligned} \tag{8.26}$$

as in the definition of sample-trace^* . By hypothesis, there exists an $n' \in \mathbb{N}$ such that $t' \in \text{traces}^* n' \text{idxs}_1 t$. Again by hypothesis, there exists an $n'' \in \mathbb{N}$ such that $t'' \in \text{traces}^* n'' \text{idxs}_2 t'$. Let $n := \max n' n''$.

Case $\text{idxs} = \text{if-idxs } j \text{ idxs}_2 \text{ idxs}_3$. Let

$$\begin{aligned} \langle p_b, b \rangle &:= \text{sample-branch} \{ \text{true}, \text{false}, \perp \} \\ p_{t'} &:= \langle p_t \cdot p_b, t[j \mapsto b] \rangle \end{aligned} \tag{8.27}$$

as in the definition of sample-trace^* , and let $\langle p'_t, t' \rangle = p_{t'}$. If $b = \text{true}$, by hypothesis, there exists an $n' \in \mathbb{N}$ such that $\text{snd} (\text{sample-trace}^* (\text{idxs}_2 0) p_{t'}) \in \text{traces}^* n' (\text{idxs}_2 0) t'$; similarly for $b = \text{false}$. Let $n := n' + 1$, so $n' = n - 1$. If $b = \perp$, then $t' = t$, and $t \in \{t\}$ for all $n \in \mathbb{N}$. \square

Corollary 8.42 (*sample-trace soundness*). *For all $\text{idxs} : \text{Idxs}$, $\text{snd} (\text{sample-trace } \text{idxs}) \in \text{traces } \text{idxs}$.*

Theorem 8.43 (*sample-trace* completeness*). *For all $\text{idxs} : \text{Idxs}$, $t \in \mathbb{T}$ and $n \in \mathbb{N}$, $t' \in \text{traces}^* n \text{idxs } t$ implies $\Pr[\text{snd} (\text{sample-trace}^* \text{idxs } \langle p_t, t \rangle) = t'] > 0$ (or just $\Pr[t'] > 0$).*

Proof. By structural induction on Idxs .

Case $\text{idxs} = \langle \rangle$. By definition of traces^* , $t' \in \text{traces}^* n \langle \rangle t$ if and only if $t' \in \{t\}$, or $t' = t$. By definition of sample-trace^* , $\Pr[\text{snd} (\text{sample-trace}^* \text{idxs } \langle p_t, t \rangle) = t] = 1$.

Case $\text{idxs} = \langle \text{idxs}_1, \text{idxs}_2 \rangle$. Let

$$\begin{aligned} T' &:= \text{traces}^* \text{ n idxs}_1 t \\ T'' &:= \bigcup_{t' \in T'} \text{traces}^* \text{ n idxs}_2 t' \end{aligned} \tag{8.28}$$

as in the definition of traces^* . Let

$$\begin{aligned} \langle p'_t, t' \rangle &:= \text{sample-trace}^* \text{ idxs}_1 \langle p_t, t \rangle \\ \langle p''_t, t'' \rangle &:= \text{sample-trace}^* \text{ idxs}_2 \langle p'_t, t' \rangle \end{aligned} \tag{8.29}$$

as in the definition of sample-trace^* . Suppose $t'' \in T''$. Then there exists a $t' \in T'$ such that $t'' \in \text{traces}^* \text{ n idxs}_2 t'$. By hypothesis, $\Pr[t'] > 0$, therefore by hypothesis, $\Pr[t''] > 0$.

Case $\text{idxs} = \text{if-idxs } j \text{ idxs}_2 \text{ idxs}_3$. Let

$$\begin{aligned} \langle p_b, b \rangle &:= \text{sample-branch} \{ \text{true}, \text{false}, \perp \} \\ p_{t'} &:= \langle p_t \cdot p_b, t[j \mapsto b] \rangle \end{aligned} \tag{8.30}$$

as in the definition of sample-trace^* , and let $\langle p'_t, t' \rangle = p_{t'}$.

If $n = 0$, then let

$$T' := \text{traces}^* 0 \text{ idxs } t = \{t\} \tag{8.31}$$

as in the definition of traces^* . Suppose $t' \in T'$, equiv. $t' = t$, so $\Pr[t'] > 0$ iff $\Pr[t] > 0$. Because $\Pr[b = \perp] > 0$ and $t[j \mapsto \perp] = t$, $\Pr[t] > 0$.

If $n > 0$, then expand $\text{traces}^* \text{ n idxs } t$ to

$$\begin{aligned} T'_{\text{true}} &:= \text{traces}^* (n-1) (\text{idxs}_2 0) t[j \mapsto \text{true}] \\ T'_{\text{false}} &:= \text{traces}^* (n-1) (\text{idxs}_3 0) t[j \mapsto \text{false}] \\ T'_{\perp} &:= \{t\} \\ T' &:= T'_{\text{true}} \cup T'_{\text{false}} \cup T'_{\perp} \end{aligned} \tag{8.32}$$

as in the definition of traces^* . Suppose $t' \in T'$. If $t' \in T'_{\text{true}}$, then because $\Pr[b = \text{true}] > 0$, by hypothesis, $\Pr[t'] > 0$; similarly for $t' \in T'_{\text{false}}$ and $t' \in T'_{\perp}$. \square

Corollary 8.44 (sample-trace completeness). *For all $\text{idxs} : \text{Idx}$, $t \in \text{traces idxs}$ implies $\Pr[\text{snd}(\text{sample-trace idxs}) = t] > 0$.*

Theorem 8.45 (sample-trace* correctness). *Let $\text{idxs} : \text{Idx}$, $t \in T_+$, $p_t \in [0, 1]$, and $\langle p'_t, t' \rangle := \text{sample-trace}^* \text{ idxs } \langle p_t, t \rangle$. If $\Pr[t] = p_t$, then $\Pr[t'] = p'_t$.*

Proof. By structural induction on Idxs .

Case $\text{idxs} = \langle \rangle$. By definition of sample-trace^* , $\langle p'_t, t' \rangle = \langle p_t, t \rangle$. Clearly $\Pr[t] = p_t$ implies $\Pr[t'] = p'_t$.

Case $\text{idxs} = \langle \text{idxs}_1, \text{idxs}_2 \rangle$. Let

$$\begin{aligned} \langle p'_t, t' \rangle &:= \text{sample-trace}^* \text{ idxs}_1 \langle p_t, t \rangle \\ \langle p''_t, t'' \rangle &:= \text{sample-trace}^* \text{ idxs}_2 \langle p'_t, t' \rangle \end{aligned} \tag{8.33}$$

as in the definition of sample-trace^* . By hypothesis, $\Pr[t] = p_t$ implies $\Pr[t'] = p'_t$, which (by hypothesis) implies $\Pr[t''] = p''_t$.

Case $\text{idxs} = \text{if-idxs } j \text{ idxs}_2 \text{ idxs}_3$. Let

$$\begin{aligned} \langle p_b, b \rangle &:= \text{sample-branch } \{\text{true}, \text{false}, \perp\} \\ p_{t'} &:= \langle p_t \cdot p_b, t[j \mapsto b] \rangle \end{aligned} \tag{8.34}$$

as in the definition of sample-trace^* , and let $\langle p'_t, t' \rangle = p_{t'}$. By definition of sample-branch , $\Pr[b] = p_b$. Because b and t are independently chosen and t' can be generated by no other execution path, $\Pr[t'] = \Pr[t] \cdot \Pr[b] = p_t \cdot p_b$. If $b = \text{true}$, then in $\langle p'_t, t' \rangle := \text{sample-trace}^* (\text{idxs}_2 \ 0) \text{ pr}'$, by hypothesis, $\Pr[t''] = p''_t$; similarly for $b = \text{false}$. \square

Corollary 8.46 (sample-trace correctness). *If $\langle p_t, t \rangle := \text{sample-trace idxs}$, then $\Pr[t] = p_t$.*

XXX: sampling from a *subset* of T_+ that induces a part of the partition of R^* :

$\text{sample-part} : \text{Idxs} \Rightarrow \langle \mathbb{R}, \text{Rect } \langle R, T \rangle \rangle \Rightarrow \langle \mathbb{R}, \text{Rect } \langle R, T \rangle \rangle$

$$\begin{aligned}
\text{sample-part idxs } \langle p, R \times T \rangle &:= \\
\text{case } \langle \text{idxs}, \langle p, \text{refine } (R \times T) \rangle \rangle & \\
\langle \langle \rangle, pr \rangle &\longrightarrow pr \\
\langle \langle \text{idxs}_1, \text{idxs}_2 \rangle, pr \rangle &\longrightarrow \text{let } pr' := \text{sample-part idxs}_1 pr \\
&\quad \text{in sample-part idxs}_2 pr' \\
\langle \text{if-idxs } j \text{ idxs}_2 \text{ idxs}_3, \langle p, R \times T \rangle \rangle &\longrightarrow \text{let } \langle p_b, b \rangle := \text{sample-branch } (\text{proj } j \text{ } T) \\
&\quad pr' := \langle p \cdot p_b, R \times \text{unproj } j \text{ } T \{b\} \rangle \\
&\quad \text{in case } b \\
&\quad \text{true} \longrightarrow \text{sample-part } (\text{idxs}_2 \text{ } 0) pr' \\
&\quad \text{false} \longrightarrow \text{sample-part } (\text{idxs}_3 \text{ } 0) pr' \\
&\quad \perp \longrightarrow \langle p \cdot p_b, \emptyset \rangle
\end{aligned} \tag{8.35}$$

XXX: sampling a weighted point from $\text{Rect } \langle R, T \rangle$:

$$\begin{aligned}
&\text{let } \langle r, q \rangle := \text{sample-point } R \\
&\quad t := \iota t \in T. \forall t' \in T. t \leq t' \\
&\quad b := f \langle \langle r, t \rangle, \langle \rangle \rangle \\
&\quad w := \text{if } (b \in B) (1 / p \cdot 1 / q) 0 \\
&\quad \text{in } \langle r, w \rangle
\end{aligned} \tag{8.36}$$

XXX: the least $t \in T$, or $\iota t \in T. \forall t' \in T. t \leq t'$, is easy to compute: set each unrestricted axis to \perp

XXX: decide what theorems are needed and prove them

XXX: extension: interval splitting

XXX: extension: sampling a weighted point using refinement

XXX: extension: not sampling any decision as \perp when sampling in T_+ ; characterize the programs for which this extension doesn't terminate; explain options to detect or fix

Chapter 9

Measurability Theorems

Proving measurability is critical in proving correctness, in that it establishes that the outputs of all programs have sensible distributions. While critical, it is somewhat distracting to the main narrative. Instead of ignoring measurability, however, as is so often done, we have moved it near the end, where readers that are somehow *still starving for even more mathematics* can devour it and—possibly—finally be satiated.

9.1 Basic Definitions

For readers familiar with topology, we review the necessary fundamentals by analogy to topology. However, we have tried to include enough of the fundamentals that readers not familiar with basic topology can verify the proofs. [XXX: cite book for imported lemmas]

The analogue of a topology of open sets is a σ -algebra of measurable sets.

Definition 9.1 (σ -algebra, measurable set). *A collection of sets $\mathcal{A} \subseteq \mathcal{P} X$ is called a σ -algebra on X if it contains X and is closed under complements and countable unions. The sets in \mathcal{A} are called **measurable sets**.*

$X \setminus X = \emptyset$, so $\emptyset \in \mathcal{A}$. Additionally, it follows from De Morgan's law that \mathcal{A} is closed under countable intersections.

The analogue of continuity is measurability.

Definition 9.2 (measurable mapping). *Let \mathcal{A} and \mathcal{B} be σ -algebras on X and Y . A mapping $g : X \rightarrow Y$ is **\mathcal{A} - \mathcal{B} -measurable** if for all $B \in \mathcal{B}$, preimage $g^{-1} B \in \mathcal{A}$.*

When the domain and codomain σ -algebras \mathcal{A} and \mathcal{B} are clear from context, we will simply say g is **measurable**.

Measurability is usually a weaker condition than continuity. For example, with respect to the σ -algebra generated from \mathbb{R} 's standard topology (i.e. using the standard topology as a sort of “base”), measurable $\mathbb{R} \rightarrow \mathbb{R}$ functions may have countably many discontinuities. Likewise, real comparison functions such as $(=)$, $(<)$, $(>)$ and their negations are measurable, but not continuous.

Product σ -algebras are defined analogously to product topologies.

Definition 9.3 (finite product σ -algebra). *Let \mathcal{A}_1 and \mathcal{A}_2 be σ -algebras on X_1 and X_2 , and define $X := X_1 \times X_2$. The **product σ -algebra** $\mathcal{A}_1 \otimes \mathcal{A}_2$ is the smallest (i.e. coarsest) σ -algebra for which mapping $\text{fst } X$ and mapping $\text{snd } X$ are measurable.*

Definition 9.4 (arbitrary product σ -algebra). *Let \mathcal{A} be a σ -algebra on X . The **product σ -algebra** $\mathcal{A}^{\otimes J}$ is the smallest σ -algebra for which, for all $j \in J$, mapping $(\pi_j) (J \rightarrow X)$ is measurable.*

9.2 Measurable Pure Computations

It is easier to prove measurability of pure computations than to prove measurability of partial, probabilistic ones. Further, we can use the results to prove that the interpretations of all partial, probabilistic expressions are measurable.

We must first define what it means for a *computation* to be measurable.

Definition 9.5 (measurable mapping arrow computation). *Let \mathcal{A} and \mathcal{B} be σ -algebras on X and Y . A computation $g : X \xrightarrow[\text{map}]{} Y$ is **\mathcal{A} - \mathcal{B} -measurable** if $g \ A^*$ is an \mathcal{A} - \mathcal{B} -measurable mapping, where A^* is g 's maximal domain.*

Theorem 9.6 (maximal domain measurability). *Let $g : X \xrightarrow[\text{map}]{} Y$ be an \mathcal{A} - \mathcal{B} -measurable mapping arrow computation. Its maximal domain A^* is in \mathcal{A} .*

Proof. Because $g : A^* \rightarrow Y$ is measurable, $\text{preimage } (g : A^* \rightarrow Y) A = A^*$ is in \mathcal{A} . □

Mapping arrow computations can be applied to sets other than their maximal domains. We need to ensure doing so yields a measurable mapping, at least for measurable subsets of A^* . Fortunately, that is true without any extra conditions.

Lemma 9.7. *Let $g : X \rightarrow Y$ be an \mathcal{A} - \mathcal{B} -measurable mapping. For any $A \in \mathcal{A}$, $\text{restrict } g : A$ is \mathcal{A} - \mathcal{B} -measurable.*

Theorem 9.8. *Let $g : X \xrightarrow[\text{map}]{} Y$ be an \mathcal{A} - \mathcal{B} -measurable mapping arrow computation with maximal domain A^* . For all $A \subseteq A^*$ with $A \in \mathcal{A}$, $g : A$ is \mathcal{A} - \mathcal{B} -measurable.*

Proof. By Theorem 7.44 (mapping arrow restriction) and Lemma 9.7. □

We do not need to prove all interpretations using $\llbracket \cdot \rrbracket_a$ are measurable. However, we do need to prove mapping arrow combinators preserve measurability.

9.2.1 Composition

Proving compositions are measurable takes the most work. The main complication is that, under measurable mappings, while *preimages* of measurable sets are measurable, *images* of measurable sets may not be. We need the following four extra theorems to get around this.

Lemma 9.9 (images of preimages). *If $g : X \rightarrow Y$ and $B \subseteq Y$, $\text{image } g (\text{preimage } g B) \subseteq B$.*

Lemma 9.10 (expanded post-composition). *Let $g_1 : X \rightarrow Y$ and $g_2 : Y \rightarrow Z$ such that $\text{range } g_1 \subseteq \text{domain } g_2$, and let $g'_2 : Y \rightarrow Z$ such that $g_2 \subseteq g'_2$. Then $g_2 \circ_{\text{map}} g_1 = g'_2 \circ_{\text{map}} g_1$.*

Theorem 9.11 (mapping arrow monotonicity). *Let $g : X \xrightarrow[\text{map}]{} Y$. For any $A' \subseteq A \subseteq A^*$, $g : A' \subseteq g : A$.*

Proof. By Theorem 7.44 (mapping arrow restriction). □

Theorem 9.12 (maximal domain subsets). *Let $g : X \xrightarrow[\text{map}]{} Y$. For all $A \subseteq A^*$, $\text{domain } (g : A) = A$.*

Proof. Follows from Theorem 7.45. □

Now we can prove measurability.

Lemma 9.13 ((\circ_{map}) measurability). *If $g_1 : X \rightarrow Y$ is \mathcal{A} - \mathcal{B} -measurable and $g_2 : Y \rightarrow Z$ is \mathcal{B} - \mathcal{C} -measurable, then $g_2 \circ_{\text{map}} g_1$ is \mathcal{A} - \mathcal{C} -measurable.*

Theorem 9.14 ((\ggg_{map}) measurability). *If $g_1 : X \xrightarrow{\sim}_{\text{map}} Y$ is \mathcal{A} - \mathcal{B} -measurable and $g_2 : Y \xrightarrow{\sim}_{\text{map}} Z$ is \mathcal{B} - \mathcal{C} -measurable, then $g_1 \ggg_{\text{map}} g_2$ is \mathcal{A} - \mathcal{C} -measurable.*

Proof. Let $A^* \in \mathcal{A}$ and $B^* \in \mathcal{B}$ be respectively g_1 's and g_2 's maximal domains. The maximal domain of $g_1 \ggg_{\text{map}} g_2$ is $A^{**} := \text{preimage}(g_1 A^*) B^*$, which is in \mathcal{A} . By definition,

$$\begin{aligned} (g_1 \ggg_{\text{map}} g_2) A^{**} &= \text{let } g'_1 := g_1 A^{**} \\ &\quad g'_2 := g_2 (\text{range } g'_1) \\ &\quad \text{in } g'_2 \circ_{\text{map}} g'_1 \end{aligned} \tag{9.1}$$

By Theorem 9.8, g'_1 is an \mathcal{A} - \mathcal{B} -measurable mapping. Unfortunately, g'_2 may not be \mathcal{B} - \mathcal{C} -measurable when $\text{range } g'_1 \notin \mathcal{B}$.

Let $g''_2 := g_2 B^*$, which is a \mathcal{B} - \mathcal{C} -measurable mapping. By Lemma 9.13, $g''_2 \circ_{\text{map}} g'_1$ is \mathcal{A} - \mathcal{C} -measurable. We need only show that $g'_2 \circ_{\text{map}} g'_1 = g''_2 \circ_{\text{map}} g'_1$, which by Lemma 9.10 is true if $\text{range } g'_1 \subseteq \text{domain } g'_2$ and $g'_2 \subseteq g''_2$.

By Theorem 9.12, $A^{**} \subseteq A^*$ implies $\text{domain } g'_1 = A^{**}$. By Theorem 9.11 and Lemma 9.9,

$$\begin{aligned} \text{range } g'_1 &= \text{image}(g_1 A^{**}) (\text{preimage}(g_1 A^*) B^*) \\ &= \text{image}(g_1 A^*) (\text{preimage}(g_1 A^*) B^*) \\ &\subseteq B^* \end{aligned} \tag{9.2}$$

$\text{range } g'_1 \subseteq B^*$ implies (by Theorem 9.12) that $\text{domain } g'_2 = \text{range } g'_1$, and (by Theorem 9.11) that $g'_2 \subseteq g''_2$. □

9.2.2 Pairing

Proving pairing preserves measurability is straightforward given a corresponding theorem about mappings.

Lemma 9.15 ($\langle \cdot, \cdot \rangle_{\text{map}}$ measurability). *If $g_1 : X \rightarrow Y_1$ is \mathcal{A} - \mathcal{B}_1 -measurable and $g_2 : X \rightarrow Y_2$ is \mathcal{A} - \mathcal{B}_2 -measurable, then $\langle g_1, g_2 \rangle_{\text{map}}$ is \mathcal{A} - $(\mathcal{B}_1 \otimes \mathcal{B}_2)$ -measurable.*

Theorem 9.16 (\mathbb{E}_{map} measurability). *If $g_1 : X \rightsquigarrow_{\text{map}} Y_1$ is \mathcal{A} - \mathcal{B}_1 -measurable and $g_2 : X \rightsquigarrow_{\text{map}} Y_2$ is \mathcal{A} - \mathcal{B}_2 -measurable, then $g_1 \mathbb{E}_{\text{map}} g_2$ is \mathcal{A} - $(\mathcal{B}_1 \otimes \mathcal{B}_2)$ -measurable.*

Proof. Let A_1^* and A_2^* be respectively g_1 's and g_2 's maximal domains. The maximal domain of $g_1 \mathbb{E}_{\text{map}} g_2$ is $A^{**} := A_1^* \cap A_2^*$, which is in \mathcal{A} . By definition, $(g_1 \mathbb{E}_{\text{map}} g_2) A^{**} = \langle g_1 A^{**}, g_2 A^{**} \rangle_{\text{map}}$, which by Lemma 9.15 is \mathcal{A} - $(\mathcal{B}_1 \otimes \mathcal{B}_2)$ -measurable. \square

9.2.3 Conditional

Conditionals can be proved measurable given a theorem that ensures the measurability of *finite* unions of disjoint, measurable mappings. We will need the corresponding theorem for *countable* unions further on, however.

Lemma 9.17 (union of measurable mappings). *The union of a countable set of \mathcal{A} - \mathcal{B} -measurable mappings with disjoint domains is \mathcal{A} - \mathcal{B} -measurable.*

Theorem 9.18 (ifte_{map} measurability). *If $g_1 : X \rightsquigarrow_{\text{map}} \text{Bool}$, and $g_2 : X \rightsquigarrow_{\text{map}} Y$ and $g_3 : X \rightsquigarrow_{\text{map}} Y$ are respectively \mathcal{A} - $(\mathcal{P} \text{ Bool})$ -measurable and \mathcal{A} - \mathcal{B} -measurable, then $\text{ifte}_{\text{map}} g_1 g_2 g_3$ is \mathcal{A} - \mathcal{B} -measurable.*

Proof. Let \mathcal{A}_1^* , \mathcal{A}_2^* and \mathcal{A}_3^* be g_1 's, g_2 's and g_3 's maximal domains. The maximal domain of $\text{ifte}_{\text{map}} g_1 g_2 g_3$ is A^{**} , defined by

$$\begin{aligned} A_2^{**} &:= A_2^* \cap \text{preimage } (g_1 \mathcal{A}_1^*) \{\text{true}\} \\ A_3^{**} &:= A_3^* \cap \text{preimage } (g_1 \mathcal{A}_1^*) \{\text{false}\} \\ A^{**} &:= A_2^{**} \uplus A_3^{**} \end{aligned} \tag{9.3}$$

Because $\text{preimage } (g_1 \mathcal{A}_1^*) B \in \mathcal{A}$ for any $B \subseteq \text{Bool}$, $A^{**} \in \mathcal{A}$. By definition,

$$\begin{aligned} \text{ifte}_{\text{map}} g_1 g_2 g_3 A^{**} = & \text{let } g'_1 := g_1 A^{**} \\ & g'_2 := g_2 (\text{preimage } g'_1 \{\text{true}\}) \\ & g'_3 := g_3 (\text{preimage } g'_1 \{\text{false}\}) \\ & \text{in } g'_2 \uplus_{\text{map}} g'_3 \end{aligned} \tag{9.4}$$

By hypothesis, g'_1 , g'_2 and g'_3 are measurable mappings. By Theorem 7.44 (mapping arrow restriction), g'_2 and g'_3 have disjoint domains. Apply Lemma 9.17. \square

9.2.4 Laziness

We must first prove measurability of an often-ignored corner case.

Theorem 9.19 (measurability of \emptyset). *For any σ -algebras \mathcal{A} and \mathcal{B} , the empty mapping \emptyset is \mathcal{A} - \mathcal{B} -measurable.*

Proof. For any $B \in \mathcal{B}$, $\text{preimage } \emptyset B = \emptyset$, and $\emptyset \in \mathcal{A}$. \square

Theorem 9.20 (measurability under lazy_{map}). *Let $g : 1 \Rightarrow (X \xrightarrow{\text{map}} Y)$. If $g \ 0$ is \mathcal{A} - \mathcal{B} -measurable, then $\text{lazy}_{\text{map}} g$ is \mathcal{A} - \mathcal{B} -measurable.*

Proof. The maximal domain A^{**} of $\text{lazy}_{\text{map}} g$ is that of $g \ 0$. By definition,

$$\text{lazy}_{\text{map}} g A^{**} = \text{if } (A^{**} = \emptyset) \ \emptyset \ (g \ 0 A^{**}) \tag{9.5}$$

If $A^{**} = \emptyset$, then $\text{lazy}_{\text{map}} g A^{**} = \emptyset$; apply Theorem 9.19. If $A^{**} \neq \emptyset$, then $\text{lazy}_{\text{map}} g = g \ 0$, which is \mathcal{A} - \mathcal{B} -measurable. \square

9.3 Measurable Probabilistic Computations

As with pure computations, we must first define what it means for an effectful computation to be measurable.

Definition 9.21 (measurable mapping* arrow computation). Let \mathcal{A} and \mathcal{B} be σ -algebras on $(\Omega \times \mathsf{T}) \times \mathsf{X}$ and Y . A computation $g : \mathsf{X} \rightsquigarrow_{\text{map}^*} \mathsf{Y}$ is **\mathcal{A} - \mathcal{B} -measurable** if $g \text{ j}_0$ is an \mathcal{A} - \mathcal{B} -measurable mapping arrow computation.

Theorem 9.22. If $g : \mathsf{X} \rightsquigarrow_{\text{map}^*} \mathsf{Y}$ is \mathcal{A} - \mathcal{B} -measurable, then for all $j \in \mathsf{J}$, $g \text{ j}$ is an \mathcal{A} - \mathcal{B} -measurable mapping arrow computation.

Proof. By induction on J : if $g \text{ j}$ is measurable, so are $g \text{ (left j)}$ and $g \text{ (right j)}$. \square

To make general measurability statements about computations, whether they have flat or product types, it helps to have a notion of a standard σ -algebra.

Definition 9.23 (standard σ -algebra). For a set X used as a type, $\Sigma \mathsf{X}$ denotes its **standard σ -algebra**, which must be defined under the following constraints:

$$\Sigma \langle \mathsf{X}_1, \mathsf{X}_2 \rangle = \Sigma \mathsf{X}_1 \otimes \Sigma \mathsf{X}_2 \quad (9.6)$$

$$\Sigma (\mathsf{J} \rightarrow \mathsf{X}) = (\Sigma \mathsf{X})^{\otimes \mathsf{J}} \quad (9.7)$$

From here on, when no σ -algebras are given, “measurable” means “measurable with respect to standard σ -algebras.”

The following definitions allow distinguishing the results of conditional expressions and any two branch traces:

$$\Sigma \text{ Bool} ::= \mathcal{P} \text{ Bool} \quad (9.8)$$

$$\Sigma \mathsf{T} ::= \mathcal{P} \mathsf{T} \quad (9.9)$$

Lemma 9.24 (measurable mapping arrow lifts). $\text{arr}_{\text{map}} \text{ id}$, $\text{arr}_{\text{map}} \text{ fst}$ and $\text{arr}_{\text{map}} \text{ snd}$ are measurable. $\text{arr}_{\text{map}} (\text{const } b)$ is measurable if $\{b\}$ is a measurable set. For all $j \in \mathsf{J}$, $\text{arr}_{\text{map}} (\pi \text{ j})$ is measurable.

Corollary 9.25 (measurable mapping* arrow lifts). $\text{arr}_{\text{map}^*} \text{ id}$, $\text{arr}_{\text{map}^*} \text{ fst}$ and $\text{arr}_{\text{map}^*} \text{ snd}$ are measurable. $\text{arr}_{\text{map}^*} (\text{const } b)$ is measurable if $\{b\}$ is a measurable set. $\text{random}_{\text{map}^*}$ and $\text{branch}_{\text{map}^*}$ are measurable.

Theorem 9.26 (AStore combinators preserve measurability). *Every AStore arrow combinator produces measurable mapping* computations from measurable mapping* computations.*

Proof. AStore's combinators are defined in terms of the base arrow's combinators and arr_{map} fst and arr_{map} snd . □

Theorem 9.27 ($\text{ifte}_{\text{map}^*}^{\Downarrow}$ measurability). *$\text{ifte}_{\text{map}^*}^{\Downarrow}$ is measurable.*

Proof. $\text{branch}_{\text{map}^*}$ is measurable, and arr_{map} agrees is measurable by (9.8). □

We can now prove all nonrecursive programs measurable by induction.

Definition 9.28 (finite expression). *A **finite expression** is any expression for which no subexpression is a first-order application.*

Theorem 9.29 (all finite expressions are measurable). *For all finite expressions e , $\llbracket e \rrbracket_{\text{map}^*}$ is measurable.*

Proof. By structural induction and the above theorems. □

Now all we need to do is represent recursive programs as a net of finite expressions, and take a sort of limit.

Theorem 9.30 (approximation with finite expressions). *Let $g := \llbracket e \rrbracket_{\text{map}^*}^{\Downarrow} : X \rightsquigarrow_{\text{map}^*} Y$ and $t \in T$. Define $A := (\Omega \times \{t\}) \times X$. There is a finite expression e' for which $\llbracket e' \rrbracket_{\text{map}^*} j_0 A = g j_0 A$.*

Proof. Let the index prefix J' contain every j for which $t j \neq \perp$. To construct e' , exhaustively apply first-order functions in e , but replace any $\text{ifte}_{\text{map}^*}^{\Downarrow}$ whose index j is not in J' with the equivalent expression \perp . Because e is well-defined, recurrences must be guarded by if , so this process terminates after finitely many first-order applications. □

Theorem 9.31 (all probabilistic expressions are measurable). *For all expressions e , $\llbracket e \rrbracket_{\text{map}^*}^{\Downarrow}$ is measurable.*

Proof. Let $g := \llbracket e \rrbracket_{\text{map}^*}^\downarrow$ and $g' := g \text{ j}_0 ((\Omega \times T) \times X)$. By Corollary 7.51 (correct computation everywhere), $g' = g \text{ j}_0 A^*$ where A^* is g 's maximal domain; thus we need only show g' is a measurable mapping.

By Theorem 7.44 (mapping arrow restriction),

$$g' = \bigcup_{t \in T} g \text{ j}_0 ((\Omega \times \{t\}) \times X) \quad (9.10)$$

By Theorem 9.30 (approximation with finite expressions), for every $t \in T$, there is a finite expression whose interpretation agrees with g on $(\Omega \times \{t\}) \times X$. Therefore, by Theorem 9.29 (all finite expressions are measurable), $g \text{ j}_0 ((\Omega \times \{t\}) \times X)$ is a measurable mapping. By Theorem 7.44 (mapping arrow restriction), they have disjoint domains. By Lemma 9.17 (union of measurable mappings), their union is measurable. \square

Theorem 9.31 remains true when $\llbracket \cdot \rrbracket_a$ is extended with any rule whose right side is measurable, including rules for real arithmetic, equality, inequality and limits. More generally, any continuous or (countably) piecewise continuous function can be made available as a language primitive, as long as its domain's and codomain's standard σ -algebras are generated from their topologies.

It is not difficult to compose $\llbracket \cdot \rrbracket_a$ with another semantic function that defunctionalizes lambda expressions. Thus, the interpretations of all expressions in higher-order languages are measurable.

9.4 Measurable Projections

If $g := \llbracket e \rrbracket_{\text{map}^*}^\downarrow : X \rightsquigarrow_{\text{map}^*} Y$, the probability of a measurable output set $B \in \Sigma Y$ is

$$P (\text{image } (fst \ggg fst) (\text{preimage } (g \text{ j}_0 A^*) B)) \quad (9.11)$$

Unfortunately, projections are generally not measurable. Fortunately, for interpretations of programs $\llbracket p \rrbracket_{\text{map}^*}^\downarrow$, for which $X = \{\langle \rangle\}$, we have a special case.

Theorem 9.32 (measurable finite projections). *Let $A \in \Sigma \langle X_1, X_2 \rangle$. If X_2 is at most countable and $\Sigma X_2 = \mathcal{P} X_2$, then $\text{image fst } A \in \mathcal{A}_1$.*

Proof. Because $\Sigma X_2 = \mathcal{P} X_2$, A is a countable union of rectangles of the form $A_1 \times \{a_2\}$, where $A_1 \in \Sigma X_1$ and $a_2 \in X_2$. Because image fst distributes over unions, $\text{image fst } A$ is a countable union of sets in ΣX_1 . \square

Theorem 9.33. *Let $g : X \rightsquigarrow_{\text{map}^*} Y$ be measurable. If X is at most countable and $\Sigma X = \mathcal{P} X$, for all $B \in \Sigma Y$, $\text{image (fst } \ggg \text{ fst) (preimage (g } j_0 \text{ } A^*) B) \in \Sigma \Omega$.*

Proof. T is countable and $\Sigma T = \mathcal{P} T$ by (9.9); apply Theorem 9.32 twice. \square

In particular, for $\llbracket p \rrbracket_{\text{map}^*}^{\downarrow} : \{\langle \rangle\} \rightsquigarrow_{\text{map}^*} Y$, the probabilities of ΣY are well-defined.

Chapter 10

Sampling Theorems

This chapter contains proofs of measure-theoretic theorems stated in Chapter 8.

10.1 Basic Definitions

While the following review is necessarily incomplete, we have tried to include enough discussion for readers unfamiliar with measure theory, and enough formalism that the proofs can be verified without consulting an outside text. For example, we do not define measure-theoretic integration, but we contrast it with integration typically learned in differential calculus, and import theorems about its properties and interactions with other operations we use.

10.1.1 Measures

Measure theory is named for its primary abstraction of length, area, volume and probability—and anything else for which assigning reals to sets in an additive way makes sense.

Definition 10.1 (measure). *A partial function $m : \text{Set } X \rightarrow [0, +\infty]$ with domain $\mathcal{A} := \text{domain } m$ is a **measure** if*

- \mathcal{A} is a σ -algebra
- $m \emptyset = 0$
- It is σ -**additive**: for any disjoint collection $A : \mathbb{N} \Rightarrow \text{Set } X$ of sets in \mathcal{A} ,

$$m \left(\bigcup_{n \in \mathbb{N}} A_n \right) = \sum_{n \in \mathbb{N}} m(A_n) \quad (10.1)$$

From here on, we rely again on the notion of a set X 's standard σ -algebra ΣX , and assume the domain of a measure $m : \text{Set } X \rightarrow [0, 1]$ is ΣX .

We will need to distinguish three kinds of measures.

Definition 10.2 (probability, finite, and σ -finite measures). *A measure $m : \text{Set } X \rightarrow [0, +\infty]$ may be*

- *A **probability measure** if $m X = 1$.*
- *A **finite measure** if $m X < +\infty$.*
- *A **σ -finite measure** if there is a collection $A : \mathbb{N} \Rightarrow \Sigma X$ such that $m (A n) < +\infty$ for all $n \in \mathbb{N}$, and $\bigcup_{n \in \mathbb{N}} A n = X$.*

Trivially, probability measures are also finite measures, which in turn are also σ -finite.

A ubiquitous example of a σ -finite measure is **Lebesgue measure**,¹ which maps sets of \mathbb{R}^n (for $n \geq 1$) to their lengths, areas and volumes. Indeed, the Lebesgue measure of \mathbb{R} is $+\infty$, but \mathbb{R} is the union of countably many sets with finite measure; e.g. $\mathbb{R} = \bigcup_{n \in \mathbb{N}} [-n, n]$.

Counting measure simply returns the cardinality of a set. If X is countable and $m : \text{Set } X \rightarrow [0, +\infty]$ is counting measure, then m is σ -finite. If X is finite, m is finite.

Image measure defines measures over the outputs of functions in terms of measures over their inputs.

Definition 10.3 (image measure). *Let $m : \text{Set } X \rightarrow [0, +\infty]$ be a measure and $g : X \rightarrow Y$ be measurable. Then g 's **image measure** with respect to m is $m' : \text{Set } Y \rightarrow [0, +\infty]$, defined by $m' B = m (\text{preimage } g B)$.*

Measures provide a way to differentiate between propositions that are always true, and propositions that may be false, but are true for certain practical purposes. To determine the latter, we need the concept of a **null set**: a set of measure zero. For example, with Lebesgue measure, $\{4\}$, or any other singleton, is a null set. In general, so is any countable union of null sets.

¹Pronounced “lehBEG,” and named after French mathematician Henri Lebesgue.

Definition 10.4 (almost everywhere). *A measurable predicate $p?$ holds **almost everywhere** with respect to measure $m : \text{Set } X \rightarrow [0, +\infty]$ when it holds on the complement of a null set:*

$$\begin{aligned} \text{ae?} : (\text{Set } X \rightarrow [0, +\infty]) &\Rightarrow (X \rightarrow \text{Bool}) \Rightarrow \text{Bool} \\ \text{ae? } m \text{ } p? &:= m (\text{preimage } p? \text{ } \{\text{false}\}) = 0 \end{aligned} \tag{10.2}$$

If m is a probability measure, $\text{ae? } m \text{ } p?$ is equivalent to $m (\text{preimage } p? \text{ } \{\text{true}\}) = 1$, or to $p?$ holding on a set of measure 1. If m is a finite measure, it is equivalent to $p?$ holding on a set of measure $m X$. If m is any other kind of measure, $\text{ae? } m \text{ } p?$ must be determined using null sets. If we were to say $p?$ holds almost everywhere when $m (\text{preimage } p? \text{ } \{\text{true}\}) = m X = +\infty$, we would have to say $\lambda x \in \mathbb{R}. x > 0$ holds almost everywhere with respect to Lebesgue measure.

In this chapter, we are most interested in almost-everywhere equality of mappings.

Definition 10.5 (almost-everywhere equality). *Two total mappings are **equal almost everywhere** with respect to a measure $m : \text{Set } X \rightarrow [0, +\infty]$ when they are not equal only on a null set, or $\text{ae-equal? } m \text{ } g_1 \text{ } g_2$ where*

$$\begin{aligned} \text{ae-equal?} : (\text{Set } X \rightarrow [0, +\infty]) &\Rightarrow (X \rightarrow Y) \Rightarrow (X \rightarrow Y) \Rightarrow \text{Bool} \\ \text{ae-equal? } m \text{ } g_1 \text{ } g_2 &:= \text{ae? } m \text{ } \lambda a \in \text{domain } g_1. g_1 \text{ } a = g_2 \text{ } a \end{aligned} \tag{10.3}$$

From here on, we use the more common “ $g_1 = g_2$ (m-a.e.)” instead of $\text{ae-equal? } m \text{ } g_1 \text{ } g_2$.

10.1.2 Integration

While measure-theoretic integration—called **Lebesgue integration**—is λ_{ZFC} -definable, defining it will not illuminate the proofs further on. The main things to know are:

- Lebesgue integration is done with respect to any base measure for the domain of integration. In contrast, Riemann integration² (as taught in differential calculus) is done only with respect to length, area or volume in \mathbb{R} , \mathbb{R}^2 and other finite products \mathbb{R}^n .
- Lebesgue integration with respect to Lebesgue measure is strictly more general than

²Pronounced “REEmahn,” and named after the German mathematician Bernhard Riemann.

Riemann integration on \mathbb{R}^n , in that it can integrate more functions. Further, a Lebesgue integral is equivalent to a corresponding Riemann integral when the latter exists.

- Lebesgue integration with respect to *counting* measure is summation.

Rather than define Lebesgue integration from first principles, we only *functionalize* it: we assume it has been defined, and turn it from special notation into a lambda. Using the notation for Lebesgue integration, we define

$$\begin{aligned} \text{int} : (X \rightarrow \mathbb{R}) &\Rightarrow (\text{Set } X \rightarrow [0, +\infty]) \Rightarrow (\text{Set } X \rightarrow [-\infty, +\infty]) \\ \text{int } g \text{ } m &:= \lambda A \in \text{domain } m. \int_A g \, dm \end{aligned} \quad (10.4)$$

Now $\text{int } g \text{ } m$ is an *indefinite* integral of g : another partial function, defined on the domain of m , that returns the definite integral on a given set A . For example, if $g \, x := x^2$ and $m : \text{Set } \mathbb{R} \rightarrow [0, +\infty]$ is Lebesgue measure, $\text{int } g \text{ } m$ measures areas under the curve $y = x^2$.

We can compute areas under the curve $y = x^2$ using Riemann integration:

$$\text{int } g \text{ } m \, [0, 1) = \int_{[0,1)} g \, dm = \int_0^1 x^2 \, dx = \frac{1^3}{3} - \frac{0^3}{3} = \frac{1}{3} \quad (10.5)$$

Of course, $\text{int } g \text{ } m$ accepts any $A \in \text{domain } m$. Because $\text{domain } m$ is a σ -algebra, this includes countable unions, countable intersections, and complements of intervals.

For real-valued functions, Lebesgue integration gives another, sometimes more convenient way to characterize almost-everywhere equality: two functions are equal almost everywhere if and only if their indefinite integrals are equal.

Lemma 10.6 (real function a.e. equality). *If $m : \text{Set } X \rightarrow [0, +\infty]$ is a σ -finite measure and $g_1, g_2 : X \rightarrow \mathbb{R}$ are measurable, then $g_1 = g_2$ (m -a.e.) if and only if $\text{int } g_1 \text{ } m = \text{int } g_2 \text{ } m$.*

The type of int might suggest its intended use; in particular, $\text{Set } X \rightarrow [-\infty, +\infty]$ is similar to $\text{Set } X \rightarrow [0, +\infty]$, which we use as the type of measures.³ We have functionalized indefinite integration to emphasize that, in this chapter and much of measure-theoretic

³In fact, $\text{Set } X \rightarrow [-\infty, +\infty]$ is the type we would use for *signed* measures if we needed them.

practice, integration's primary purpose is not to compute concrete areas and volumes, but to *transform measures*. This is supported by the following imported theorem.

Lemma 10.7 (indefinite integration yields measures). *If $g : X \rightarrow [0, +\infty)$ is measurable and $m : \text{Set } X \rightarrow [0, +\infty]$ is a measure, then $\text{int } g \, m$ is a measure.*

For example, if $g : \mathbb{R} \rightarrow [0, +\infty)$ is a probability density function and m is Lebesgue measure on \mathbb{R} , then $\text{int } g \, m$ is a probability measure.

Lemma 10.7 implies there is a function

$$\text{int}^+ : (X \rightarrow [0, +\infty)) \Rightarrow (\text{Set } X \rightarrow [0, +\infty]) \Rightarrow (\text{Set } X \rightarrow [0, +\infty]) \quad (10.6)$$

which agrees with int for all arguments $g : X \rightarrow [0, +\infty)$. We thus begin defining an algebra of measures and operations on them with int^+ .

We should expect integration to be positive linear, and it is. In the following, assume that arithmetic is lifted to operate pointwise on mappings.

Lemma 10.8. *Let $g_1, g_2 : X \rightarrow [0, +\infty)$ be measurable and $m : \text{Set } X \rightarrow [0, +\infty]$ be a measure. Then $\text{int } (g_1 + g_2) \, m = \text{int } g_1 \, m + \text{int } g_2 \, m$.*

Lemma 10.9. *Let $g : X \rightarrow [0, +\infty)$ be measurable, $\alpha \geq 0$ and $m : \text{Set } X \rightarrow [0, +\infty]$ be a measure. Then $\text{int } (\alpha \cdot g) \, m = \alpha \cdot \text{int } g \, m$.*

Lastly, compositions within integrals can be moved into the base measure.

Lemma 10.10 (image measure integration). *Let $m : \text{Set } X \rightarrow [0, +\infty]$ be a measure, $g_1 : X \rightarrow Y$ and $g_2 : Y \rightarrow \mathbb{R}$ be measurable, and m' be g_1 's image measure with respect to m . For all $B \in \Sigma Y$, $\text{int } (g_2 \circ_{\text{map}} g_1) \, m (\text{preimage } g_1 \, B) = \text{int } g_2 \, m' \, B$. [XXX: need more conditions]*

10.1.3 Differentiation

In differential calculus, indefinite integration has an inverse: differentiation. In measure theory, indefinite Lebesgue integration also has an inverse, which is also called differentiation.

One significant difference is that, because indefinite Lebesgue integration returns measures, differentiation operates on measures.

In differential calculus, differentiation is defined only for differentiable functions. In measure theory, the analogous property is absolute continuity.

Definition 10.11 (absolute continuity). *Given measures $m_1, m_2 : \text{Set } X \rightarrow [0, +\infty]$, we say m_1 is **absolutely continuous** with respect to m_2 if $m_1 \ll m_2$, where*

$$\begin{aligned} (\ll) : (\text{Set } X \rightarrow [0, +\infty]) &\Rightarrow (\text{Set } X \rightarrow [0, +\infty]) \Rightarrow \text{Bool} \\ m_1 \ll m_2 &:= \forall A \in \text{domain } m_2. m_2 A = 0 \implies m_1 A = 0 \end{aligned} \tag{10.7}$$

By Definition 10.11, $m_1 \ll m_2$ means m_1 has at least as many measure-zero sets as m_2 , and is therefore, in a sense, smaller. If P and Q are probability measures, $P \ll Q$ essentially means P 's support is no larger than Q 's support.

As for integration, for differentiation, we functionalize special notation:

$$\begin{aligned} \text{diff}^+ : (\text{Set } X \rightarrow [0, +\infty]) &\Rightarrow (\text{Set } X \rightarrow [0, +\infty]) \Rightarrow (X \rightarrow [0, +\infty]) \\ \text{diff}^+ m_1 m_2 &:= \frac{dm_1}{dm_2} \end{aligned} \tag{10.8}$$

This returns a **Radon-Nikodým derivative**. Such derivatives are named after the following theorem, which gives circumstances under which $\text{diff}^+ m_1 m_2$ exists, and states that int^+ is the left inverse of diff^+ (with second arguments held constant).

Lemma 10.12 (Radon-Nikodým). *If $m_1, m_2 : \text{Set } X \rightarrow [0, +\infty]$ are σ -finite measures and $m_1 \ll m_2$, then $\text{diff}^+ m_1 m_2$ exists and $m_1 = \text{int}^+ (\text{diff}^+ m_1 m_2) m_2$.*

The function $\text{diff}^+ m_1 m_2 : X \rightarrow [0, +\infty]$ is often called the *density* of m_1 with respect to m_2 , but we call them *derivatives*, reserving *density* for derivatives with respect to Lebesgue measure. By Lemma 10.6, any $g : X \rightarrow [0, +\infty]$ for which $g = \text{diff}^+ m_1 m_2$ (m_2 -a.e.) meets the Radon-Nikodým theorem's conclusion $m_1 = \text{int}^+ g m_2$. We therefore say that Radon-Nikodým derivatives are unique up to equality m_2 -a.e.

By analogy to differential calculus, we should expect diff^+ to be the left inverse of int^+ (with second arguments held constant). It is, up to equality m_2 -a.e.

Lemma 10.13. *If $g_1 : X \rightarrow [0, +\infty)$ is measurable and $m_2 : \text{Set } X \rightarrow [0, +\infty]$ is a σ -finite measure, then $\text{int}^+ g_1 m_2 \ll m_2$ and $g_1 = \text{diff}^+ (\text{int}^+ g_1 m_2) m_2$ (m_2 -a.e.).*

The preceding two theorems are analogous to the fundamental theorem of calculus.

We should expect differentiation to be positive linear, and it is.

Lemma 10.14. *Let $m_1, m_2, m : \text{Set } X \rightarrow [0, +\infty]$ be σ -finite measures with $m_1 \ll m$ and $m_2 \ll m$. Then $m_1 + m_2 \ll m$ and $\text{diff}^+ (m_1 + m_2) m = \text{diff}^+ m_1 m + \text{diff}^+ m_2 m$ (m -a.e.).*

Lemma 10.15. *Let $m_1, m_2 : \text{Set } X \rightarrow [0, +\infty]$ be σ -finite measures with $m_1 \ll m_2$. For all $\alpha \geq 0$ and $\beta > 0$, $\alpha \cdot m_1 \ll \beta \cdot m_2$ and $\text{diff}^+ (\alpha \cdot m_1) (\beta \cdot m_2) = \frac{\alpha}{\beta} \cdot \text{diff}^+ m_1 m_2$ (m -a.e.).*

As in differential calculus, there is a chain rule.

Lemma 10.16 (chain rule). *Let $m_1, m_2, m_3 : \text{Set } X \rightarrow [0, +\infty]$ be σ -finite measures with $m_1 \ll m_2$ and $m_2 \ll m_3$. Then $m_1 \ll m_3$ and $\text{diff}^+ m_1 m_2 \cdot \text{diff}^+ m_2 m_3 = \text{diff}^+ m_1 m_3$ (m_3 -a.e.).*

We need two more differentiation rules, which have no direct analogues in differential calculus. Importing them makes our algebra of measures complete enough to prove importance sampling correct. The first is a rule for reciprocals.

Lemma 10.17 (reciprocal rule). *Let $m_1, m_2 : \text{Set } X \rightarrow [0, +\infty]$ be σ -finite measures with $m_2 \ll m_1$ and $m_1 \ll m_2$. Then $\text{diff}^+ m_1 m_2 = 1 / \text{diff}^+ m_2 m_1$ (m_1 -a.e. and m_2 -a.e.).*

The second provides a way to integrate out derivatives, or to use differentiation to change the base measure in Lebesgue integration.

Lemma 10.18 (change of measure). *Let $m_1, m_2 : \text{Set } X \rightarrow [0, +\infty]$ be σ -finite measures with $m_1 \ll m_2$, and $g : X \rightarrow \mathbb{R}$ be measurable. Then $\text{int } g m_1 = \text{int } (g \cdot \text{diff}^+ m_1 m_2) m_2$.*

Suppose we have a joint and candidate probability densities $\mathbf{p}, \mathbf{q} : \mathbb{R}^n \rightarrow [0, +\infty)$, and we sample according to \mathbf{q} and weight the samples by \mathbf{p} / \mathbf{q} . The weighted samples represent \mathbf{p} if expected values estimated using them are correct; i.e. for all measurable $g : \mathbb{R}^n \rightarrow \mathbb{R}$,

$$\text{int } g \, P = \text{int } (g \cdot \mathbf{p} / \mathbf{q}) \, Q \quad (10.9)$$

where $P := \text{int}^+ \mathbf{p} \, m$ and $Q := \text{int}^+ \mathbf{q} \, m$, and m is Lebesgue measure on \mathbb{R}^n .

The density route to a proof is simple, and requires that \mathbf{q} be nonzero everywhere. By definition and Lemma 10.13 (diff^+ is a left inverse of int^+), $\mathbf{q} = \text{diff}^+ Q \, m$ (m -a.e.), so by Lemma 10.18 (change of measure) with $m_1 = Q$,

$$\begin{aligned} \text{int } (g \cdot \mathbf{p} / \mathbf{q}) \, Q &= \text{int } (g \cdot \mathbf{p} / \mathbf{q} \cdot \mathbf{q}) \, m \\ &= \text{int } (g \cdot \mathbf{p}) \, m \end{aligned} \quad (10.10)$$

which again by Lemmas 10.13 and 10.18 is $\text{int } g \, P$.

Taking the measure route demonstrates how to prove more general importance sampling theorems. We again require \mathbf{q} to be nonzero everywhere; then

$$\begin{aligned} \mathbf{p} / \mathbf{q} &= \mathbf{p} \cdot 1 / \mathbf{q} = \text{diff}^+ P \, m \cdot 1 / \text{diff}^+ Q \, m \quad (m\text{-a.e.}) && \text{Lemma 10.13} \\ &= \text{diff}^+ P \, m \cdot \text{diff}^+ m \, Q \quad (m, Q\text{-a.e.}) && \text{Lemma 10.17} \\ &= \text{diff}^+ P \, Q \quad (Q\text{-a.e.}) && \text{Lemma 10.16} \end{aligned} \quad (10.11)$$

Because $g \cdot \mathbf{p} / \mathbf{q} = g \cdot \text{diff}^+ P \, Q$ (Q -a.e.),

$$\begin{aligned} \text{int } (g \cdot \mathbf{p} / \mathbf{q}) \, Q &= \text{int } (g \cdot \text{diff}^+ P \, Q) \, Q && \text{Lemma 10.6} \\ &= \text{int } g \, P && \text{Lemma 10.18} \end{aligned} \quad (10.12)$$

The more general method is this: instead of densities \mathbf{p} and \mathbf{q} , define measures P and Q , derive $\text{diff}^+ P \, Q$, and apply Lemma 10.18.

The proof of correctness of *partitioned* importance sampling proceeds this way, but requires more machinery to construct a measure-theoretic model of the sampling process.

10.1.4 Transition Kernels

In naïve probability theory, conditional density functions model probabilistic processes that depend on the outcome of another. In measure-theoretic probability, this is accomplished using transition kernels, which are not much more than functions that return measures.

Definition 10.19 (transition kernel). *A function $k : X \rightarrow \text{Set } Y \rightarrow [0, +\infty]$ is a **transition kernel** when both of the following hold.*

- *For all $a \in X$, $k\ a$ is a measure.*
- *For all $B \in \Sigma\ Y$, $\lambda a \in X. k\ a\ B$ is measurable.*

For any measure property, we say k has that property when it holds for all $k\ a$. Therefore, k is a **probability kernel**, **finite kernel**, or **σ -finite kernel** when for all $a \in X$, $k\ a$ is respectively a probability measure, finite measure, or σ -finite measure.

Product models of dependent processes can be built by starting with a probability measure and iteratively extending it using probability kernels.

Lemma 10.20 (finite kernel products). *Let $m : \text{Set } X \rightarrow [0, +\infty]$ be a finite measure and $k : X \rightarrow \text{Set } Y \rightarrow [0, +\infty]$ be a finite kernel. There exists a unique σ -finite measure $m \times k : \text{Set } \langle X, Y \rangle \rightarrow [0, +\infty]$ that is determined by its output on rectangles; i.e. defined by extending the following to a product measure: for all $A \in \Sigma\ X$ and $B \in \Sigma\ Y$,*

$$(m \times k)(A \times B) = \int^+ (\lambda a \in X. k\ a\ B) m\ A \quad (10.13)$$

If m is a probability measure and k a probability kernel, $m \times k$ is a probability measure.

For example, if $N : \mathbb{R} \rightarrow \text{Set } \mathbb{R} \rightarrow [0, 1]$ takes a mean μ and returns a normal probability measure centered on μ with standard deviation 1, then the interpretation of

$$\begin{aligned} X &\sim \text{Normal}(0, 1) \\ Y &\sim \text{Normal}(X, 1) \end{aligned} \quad (10.14)$$

as a measure-theoretic joint distribution is $(N\ 0) \times N$.

For the proofs in the next section, we need a way to turn integrals with respect to $m \times k$ measures into nested integrals.

Lemma 10.21 (Fubini's for transition kernels). *Let $m : \text{Set } X \rightarrow [0, +\infty]$ be a finite measure and $k : X \rightarrow \text{Set } Y \rightarrow [0, +\infty]$ be a finite kernel. If $g : X \times Y \rightarrow \overline{\mathbb{R}}$ is measurable, and nonnegative or $(m \times k)$ -integrable, then*

$$\begin{aligned} \text{int } g \ (m \times k) \ (X \times Y) \\ = \text{int } (\lambda a \in X. \text{int } (\lambda b \in Y. g \ \langle a, b \rangle)) \ (k \ a) \ Y) \ m \ X \end{aligned} \tag{10.15}$$

10.2 Sampling Proofs

Recall that $\text{subcond } P \ A' := \lambda A \in \text{domain } P. P \ (A' \cap A)$. [XXX: define in Chapter 8]

Theorem 10.22 (partitioned importance sampling correctness). *Let X , P , N , s , p , and Q as in Definition 8.27 (partitioned importance sampling) such that $\text{subcond } P \ (s \ n) \ll Q \ n$ for all $n \in N$. Define $P_N : \text{Set } N \rightarrow [0, 1]$ by integrating p with respect to counting measure.*

If $g : X \rightarrow \mathbb{R}$ is a P -integrable mapping, and

$$\begin{aligned} g' : N \times X &\rightarrow \mathbb{R} \\ g' \ \langle n, a \rangle &:= g \ a \cdot \frac{1}{p \ n} \cdot \text{diff}^+ \ (\text{subcond } P \ (s \ n)) \ (Q \ n) \ a \end{aligned} \tag{10.16}$$

then $\text{int } g' \ (P_N \times Q) \ (N \times X) = \text{int } g \ P \ X$.

Proof. Let $w_1 \ n := \frac{1}{p \ n}$ and $w_2 \ n := \text{diff}^+ \ (\text{subcond } P \ (s \ n)) \ (Q \ n)$. Starting from the left side,

$$\begin{aligned} \text{int } g' \ (P_N \times Q) \ (N \times X) & \tag{10.17} \\ = \text{int } (\lambda \langle n, a \rangle \in N \times X. g \ a \cdot w_1 \ n \cdot w_2 \ n \ a) \ (P_N \times Q) \ (N \times X) & \quad \text{Def of } g' \\ = \text{int } (\lambda n \in N. \text{int } (\lambda a \in X. g \ a \cdot w_1 \ n \cdot w_2 \ n \ a) \ (Q \ n) \ X) \ P_N \ N & \quad \text{Lemma 10.21} \\ = \text{int } (\lambda n \in N. \text{int } (g \cdot w_1 \ n \cdot w_2 \ n) \ (Q \ n) \ X) \ P_N \ N & \quad \text{Lift } (\cdot) \\ = \text{int } (\lambda n \in N. w_1 \ n \cdot \text{int } (g \cdot w_2 \ n) \ (Q \ n) \ X) \ P_N \ N & \quad \text{Lemma 10.9} \\ = \text{int } (\lambda n \in N. w_1 \ n \cdot \text{int } g \ (\text{subcond } P \ (s \ n)) \ X) \ P_N \ N & \quad \text{Def } w_2, \text{ Lemma 10.18} \end{aligned}$$

Because P_N is defined with respect to counting measure, turn integration into summation:

$$\begin{aligned}
&= \sum_{n \in N} p_n \cdot \frac{1}{p_n} \cdot \int g(\text{subcond } P(s_n)) X && \text{Def of } w_1 \\
&= \sum_{n \in N} \int g(\text{subcond } P(s_n)) X && p_n > 0 \\
&= \sum_{n \in N} \int g P(s_n) && \text{XXX: justify} \\
&= \int g P \left(\bigcup_{n \in N} (s_n) \right) && \sigma\text{-additivity} \\
&= \int g P X && \text{Def of } s \quad \square
\end{aligned}$$

When m_1 and m_2 are measures on infinite spaces, it is not clear that $\text{diff}^+ m_1 m_2$ exists or how to compute it. It seems it should exist when m_1 and m_2 differ only on a finite projection of their domains, and that it should be easy to compute when the distributions of those finite projections can be defined by densities.

We start with a theorem for $\text{diff}^+ (m_1 \times k) (m_2 \times k)$ that says we may ignore k .

Theorem 10.23. *Let $m_1, m_2 : \text{Set } X \rightarrow [0, +\infty]$ be finite measures such that $m_1 \ll m_2$, and $k : X \rightarrow \text{Set } Y \rightarrow [0, +\infty]$ be a finite kernel. Then $m_1 \times k \ll m_2 \times k$ and $\text{diff}^+ (m_1 \times k) (m_2 \times k) = (\lambda \langle a, b \rangle \in X \times Y. \text{diff}^+ m_1 m_2 a) (m_2 \times k\text{-a.e.})$.*

Proof. XXX: prove absolute continuity

Let $A \in \Sigma X$ and $B \in \Sigma Y$. Integrating the left-hand side,

$$\begin{aligned}
&\int (\text{diff}^+ (m_1 \times k) (m_2 \times k)) (m_2 \times k) (A \times B) && (10.18) \\
&= (m_1 \times k) (A \times B) && \text{Lemma 10.12}
\end{aligned}$$

Integrating the right-hand side,

$$\begin{aligned}
& \text{int } (\lambda \langle a, b \rangle \in X \times Y. \text{diff}^+ m_1 m_2 a) (m_2 \times k) (A \times B) & (10.19) \\
&= \text{int } (\lambda a \in X. \text{int } (\lambda b \in Y. \text{diff}^+ m_1 m_2 a) (k a) B) m_2 A & \text{Lemma 10.21} \\
&= \text{int } (\text{diff}^+ m_1 m_2 \cdot \lambda a \in X. \text{int } (\lambda b \in Y. 1) (k a) B) m_2 A & \text{Lemma 10.9, Lift } (\cdot) \\
&= \text{int } (\lambda a \in X. k a B) m_1 A & \text{Lemma 10.18} \\
&= (m_1 \times k) (A \times B) & \text{Lemma 10.20}
\end{aligned}$$

Therefore, because $m_1 \times k$ is uniquely defined by its output on all such $A \times B$,

$$\begin{aligned}
& \text{int } (\text{diff}^+ (m_1 \times k) (m_2 \times k)) (m_2 \times k) \\
&= \text{int } (\lambda \langle a, b \rangle \in X \times Y. \text{diff}^+ m_1 m_2 a) (m_2 \times k)
\end{aligned} \tag{10.20}$$

Apply Lemma 10.6 (real function a.e. equality). \square

It is not hard to extend the preceeding theorem to arbitrary sublists of finite lists, or to arbitrary finite substructures of any algebraic data type, by induction. But we need a version of it for arbitrary finite substructures of infinite binary trees, which we have defined non-inductively as mappings $\Omega := J \rightarrow [0, 1]$ from tree indexes to reals.

One solution is to define an injective transformation g from any $\omega \in \Omega$ to a pair $\langle \omega_{\text{fin}}, \omega_{\text{inf}} \rangle$, where ω_{fin} is a finite substructure of ω and ω_{inf} is the rest of it, and apply Theorem 10.23. The proof is easier to do first in generality, without specifying the structure of ω , requiring the substructure to be finite, or requiring the pairs to contain projections.

Theorem 10.24. *Let $\mu_1, \mu_2 : \text{Set } Z \rightarrow [0, +\infty]$ be σ -finite measures. If there exist finite measures $m_1, m_2 : \text{Set } X \rightarrow [0, +\infty]$ such that $m_1 \ll m_2$, a finite kernel $k : X \rightarrow \text{Set } Y \rightarrow [0, +\infty]$, and an injective, measurable function $g : Z \rightarrow X \times Y$ such that for all $D \in \Sigma \langle X, Y \rangle$,*

$$\begin{aligned}
(m_1 \times k) D &= \mu_1 (\text{preimage } g D) \\
(m_2 \times k) D &= \mu_2 (\text{preimage } g D)
\end{aligned} \tag{10.21}$$

then $\mu_1 \ll \mu_2$ and $\text{diff}^+ \mu_1 \mu_2 = \lambda z \in Z. \text{diff}^+ m_1 m_2 (\text{fst } (g z))$ (μ_2 -a.e.).

Proof. Suppose m_1, m_2, k and g as stated. XXX: prove absolute continuity

Let $C \in \Sigma Z$ and $D := \text{image } g C$; then

$$\begin{aligned}
& \int (\text{diff}^+ \mu_1 \mu_2) \mu_2 C && (10.22) \\
&= \int \mu_1 C && \text{Lemma 10.12} \\
&= (m_1 \times k) D && \text{Def } (m_1 \times k), \text{ injectivity} \\
&= \int (\text{diff}^+ (m_1 \times k) (m_2 \times k)) (m_2 \times k) D && \text{Lemma 10.12} \\
&= \int (\lambda \langle a, b \rangle \in X \times Y. \text{diff}^+ m_1 m_2 a) (m_2 \times k) D && \text{Theorem 10.23} \\
&= \int ((\lambda \langle a, b \rangle \in X \times Y. \text{diff}^+ m_1 m_2 a) \circ_{\text{map}} g) \mu_2 C && \text{Lemma 10.10} \\
&= \int (\lambda z \in Z. \text{diff}^+ m_1 m_2 (\text{fst } (g z))) \mu_2 C && \text{Def of } \circ_{\text{map}}
\end{aligned}$$

Apply Lemma 10.6 (real function a.e. equality). □

Thus, two measures μ_1 and μ_2 on infinite structures that can be decomposed into products $m_1 \times k$ and $m_2 \times k$ such that $\text{diff}^+ m_1 m_2$ exists—using any injective transformation—have a Radon-Nikodým derivative that can be uniquely defined in terms of $\text{diff}^+ m_1 m_2$.

Application to infinite binary trees mostly requires defining the transformation.

Theorem 10.25. *Let $J' \subseteq J$ be finite, and define $X := J' \rightarrow [0, 1]$ and $Y := (J \setminus J') \rightarrow [0, 1]$. Let $P', Q' : \text{Set } X \rightarrow [0, 1]$ be finite measures such that $P' \ll Q'$, and let $k : X \rightarrow \text{Set } Y \rightarrow [0, 1]$ be a finite kernel. Let $g : \Omega \rightarrow X \times Y$ be defined by $g \omega = \langle \text{restrict } \omega J', \text{restrict } \omega (J \setminus J') \rangle$. If $P, Q : \text{Set } \Omega \rightarrow [0, 1]$ are defined so that for all $\Omega' \in \Sigma \Omega$,*

$$\begin{aligned}
P \Omega' &= (P' \times k) (\text{image } g \Omega') \\
Q \Omega' &= (Q' \times k) (\text{image } g \Omega')
\end{aligned} \tag{10.23}$$

then $P \ll Q$ and $\text{diff}^+ P Q = \lambda \omega \in \Omega. \text{diff}^+ P' Q' (\text{restrict } \omega J')$ (Q -a.e.).

Proof. The inverse of g is $g^{-1} : X \times Y \rightarrow \Omega$, defined by

$$g^{-1} \langle \omega_{\text{fin}}, \omega_{\text{inf}} \rangle = \lambda j \in J. \text{ if } (j \in J') (\omega_{\text{fin}} j) (\omega_{\text{inf}} j) \quad (10.24)$$

Thus $(P' \times k) D = P (\text{preimage } g D)$ for all $D \in \Sigma \langle X, Y \rangle$; similarly for $(Q' \times k) D$. Apply Theorem 10.24. \square

In particular, if additionally P' and Q' can be defined by densities $p : (J' \rightarrow [0, 1]) \rightarrow [0, +\infty)$ and $q : (J' \rightarrow [0, 1]) \rightarrow [0, +\infty)$, then

$$\text{diff}^+ P Q = \lambda \omega \in \Omega. \frac{p (\text{restrict } \omega J')}{q (\text{restrict } \omega J')} \quad (Q\text{-a.e.}) \quad (10.25)$$

Chapter 11

Related Work

Probabilistic languages can be approximately placed into two groups: those defined by an implementation, and those defined by a semantics.

11.1 Implementations

Almost all of the languages defined by their implementations support conditional queries and compute converging approximations. The reports on these languages generally describe interpreters, compilers, and algorithms for sampling with probabilistic conditions. When they work correctly, they are useful.

Koller and Pfeffer [28] efficiently compute exact, discrete distributions for the outputs of programs in a Scheme-like language. BUGS [33] focuses on efficient approximate computation for probabilistic theories with a finitely many statements, and uses approximation methods that Bayesians typically use. BLOG [36] exists specifically to allow stating distributions over countably infinite vectors. BLAISE [10] allows stating both distributions and approximation methods for random variables. Church [17] is a Scheme-like probabilistic language with approximate inference, and focuses on expressiveness.

Kiselyov [26] embeds a probabilistic language in O’Caml for efficient computation. It uses continuations to enumerate or sample random variable values, and has a **fail** construct for the *complement* of conditioning. The sampler looks ahead for **fail** and can handle it efficiently. This may be justified by commutativity (Thm. 5.10), depending on interaction with other language features.

Recently, Wingate et al [56, 57] have defined the semantics of *nonstandard interpretations* that enable efficient inference, but do not define the languages. [XXX: recheck, relate to computation indexing]

11.2 Semantics

Early work in probabilistic language semantics is not motivated by Bayesian concerns, and thus does not address conditioning. Kozen [29] defines the meaning of bounded-space, imperative “while” programs as functions from probability measures to probability measures. Hurd [22] proves properties about programs with binary random choice by encoding programs and portions of measure theory in HOL. Jones [23] develops a domain-theoretic account of probability, and with it defines the probability monad, whose discrete version is a distribution-valued variation of the set or list monad.

Ramsey and Pfeffer [45] define the probability monad measure-theoretically and implement a language for finite probability. We do not build on this work because the probability monad does not build a global probability space, making it difficult to reason about conditioning.

Using inverse transform sampling, Park [41] extends a λ -calculus with probabilistic choice according to any of a general class of probability measures. This is the same technique we use in Chapter 8 to turn uniform probabilistic choice into choice according to other distributions, though preimage computation makes doing so more difficult.

Some recent work in probabilistic language semantics tackles conditioning. Pfeffer’s IBAL [44] is the earliest λ -calculus with finite probabilistic choice that also defines conditional queries. Borgström et al [11] develop Fun, a first-order functional language without recursion, extended with probabilistic choice and conditioning. Its semantics interprets programs as *measure transformers* by transforming expressions into arrow-like combinators. The implementation generates a decomposition of the probability density represented by the program, if it exists. Bhat et al [9] replaces Fun’s `if` with `match`, and interprets programs

more directly as probability density functions by compositionally transforming expressions into an extension of the probability monad.

11.3 Techniques

The forward phase in computing preimages takes a subdomain and returns an overapproximation of the function’s range for that subdomain. This clearly generalizes interval arithmetic [24] to all first-order algebraic types. We further generalize interval arithmetic by doing it backwards as well as forwards.

Our general approach—creating an exact semantics and deriving an implementable, approximating semantics—is similar to abstract interpretation [13]. Most functional programming researchers would hesitate to call Chapter 5’s approximating semantics abstract, however. While the approximations it computes converge, they are not conservative.

The work in Chapter 7 is much more in line with typical abstract interpretation: $\llbracket \cdot \rrbracket_{\text{pre}^*}^\downarrow$ is the concrete semantics, $\llbracket \cdot \rrbracket_{\text{pre}^*}^{\downarrow'}$ is an abstract semantics, the concrete values are arbitrary sets, and the abstract values are rectangles. It has many properties common to abstract semantics: it is sound, it labels expressions, the abstract domain is a lattice, and the exact semantics performs infinite computations.

However, it is far from typical in other ways. It is used to run programs, not for static analysis. The abstraction boundaries are the conditional branches of completely unrolled, infinite programs, and are not fixed. There is no Kleene iteration to find a fixed point. In our abstract semantics, infinite computations are done in a library of λ_{ZFC} -computable combinators, not by a semantic function. This cleanly separates the syntax from the semantics, and allows us to prove the exact semantics correct mostly by proving simple categorical properties.

11.4 Somewhat Related Work

Any programming language research described by the words “bijective” or “reversible” might seem to have much in common with ours. Unfortunately, when we look more closely, we can usually draw only loose analogies and perhaps inspiration. An example is lenses [19], which are transformations from X to Y that can be run forward and backward, in a way that maintains some relationship between X and Y . Usually, a destructive, external process is assumed, so that, for example, a change from $y \in Y$ to $y' \in Y$ induces a corresponding change from $x \in X$ to some $x' \in X$. When transformations lose information, lenses must satisfy certain behavioral laws. In our work, no input or output is updated, and preimages are always definable regardless of non-injectivity.

Many multi-paradigm languages [18], especially constraint functional languages, bear a strong resemblance to our work. In fact, it is easy to add a `fail` expression to our semantics, or to transform constraints into boolean program outputs. The most obvious difference is evaluation strategy. The most important difference is that our interpretation of programs returns *distributions* of constrained outputs, rather than arbitrary single values that meet constraints.

Chapter 12

Conclusions and Future Work

12.1 Conclusions

We started by defining λ_{ZFC} , a call-by-value λ -calculus with infinite sets and set operations, so that we could interpret Bayesian notation categorically.

We then investigated a general approach to trustworthy Bayesian languages: defining an exact semantics that interprets notation as measure-theoretic models, and then deriving a directly implementable approximating semantics. We restricted our investigation to countable distributions and theories with finitely many statements, as it is the first point in the design space where approximation is necessary, and requires no deep measure theory.

In a slight change of tactics, we fixed a canonical probability space of uniformly random, infinite binary trees, and interpreted programs as measure-theoretic random variables, and then as computations that compute exact preimages. The approximating semantics interprets programs as computations that compute conservative approximations of preimages, using rectangles instead of sets. We demonstrated that the language is useful by implementing the approximating semantics and encoding typical Bayesian theories. We also encoded theories without density models, which can only be interpreted using measure theory.

In short, we developed trustworthy, useful languages for Bayesian modeling and inference by founding them solidly on functional programming theory and measure-theoretic probability.

12.2 Future Work

Future work falls into two categories: expressiveness and optimization.

12.2.1 Expressiveness

Adding a new feature and its semantics to a Turing-equivalent language makes the language more **expressive** if the only way to encode the new feature in the original language with the original semantics is by a global transformation [XXX: cite Felleisen].

An example is adding lambdas. For a first-order language like we have defined, adding lambdas requires closure conversion and defunctionalization: turning every lambda value into a data structure containing bound variable values and a function pointer, and changing every application site to apply a global dispatching function that decodes such closures [XXX: cite]. It may be simpler, more efficient, or more elegant to add lambda terms to the language itself, despite the fact that ensuring higher-order application is measurable is difficult.

Other examples of expressive new features are mutation, exceptions and parameters [XXX: cite] (or more generally continuations [XXX: cite] and continuation marks [XXX: cite]). Once lambdas are available, these can be encoded by globally transforming programs into monadic computations [XXX: cite; include Kimball]. We want to know whether such global transformations are the simplest, most efficient, or most elegant ways to extend our probabilistic language's expressiveness.

Non-examples are loops and other constructs that are merely syntactic sugar, which can be implemented using local transformations. We will be able to provide all such features by making Racket syntax transformers available to programs.

It is not clear that lambdas are the best recursive abstraction for encoding Bayesian theories. Lambdas are good for creating abstractions because their bodies and the computations they carry out are generally unobservable. However, these facts may cause them to not meet Bayesian needs very well: Bayesians define theories in order to study them, and often to discover the behavior of a process cannot be directly observed. Abstracting such processes

using lambdas would require returning every intermediate value. We do not know what the right recursive abstraction is for such cases, let alone whether it can be implemented using a local transformation. This dimension of probabilistic language design clearly needs study.

In our experience, probabilistic programs are more difficult to debug than other kinds of programs. Therefore, one way to increase expressiveness while reducing errors is by adding a type system. A type system similar to Typed Racket [XXX: cite] would be a good choice. Because Typed Racket was originally meant for converting untyped programs into typed programs by adding a few annotations, it has true union types, and **occurrence typing**, which allows identifiers to have different types in each branch of a conditional based on the result of its test expression. Occurrence typing is similar to how the forward phase in preimage computation applies the interpretation of each *true* branch to the preimage of `{true}` under the test, and the *false* branch to the preimage of `{false}`. Union types are similar to our representation of disjoint unions of sets of tagged data structures.

In fact, there are many similarities between preimage computation and type checking and inference, which we plan to investigate.

12.2.2 Optimization

A type system would not only help reduce programmer error, but would provide information about program terms that an optimizer could use. For example, because our implementation's sets are monomorphic, every set operation must dispatch to a more specific operation based on the runtime data types of its arguments. If a type system determined that a certain computation consumed and produced only pairs, such dispatch would be unnecessary.

Importance sampling is known to suffer from high variance in its estimates when used on high-dimensional models [XXX: cite]. The higher the variance, the more samples are necessary to get good estimates. A major piece of future optimization work is to use different sampling algorithms, such as Markov Chain Monte Carlo (MCMC) [XXX: cite]. We expect the split-and-refine part of sampling to be particularly amenable, because it divides the

program domain into an at-most-countable partition. Markov-chain sampling of parts from this partition may easily sidestep a common problem with MCMC methods: that they often mix poorly when sampling within narrow, non-convex shapes. In any case, to use MCMC, we would need to solve this problem, because the preimages of conditions in typical Bayesian queries are narrow and non-convex.

Our semantics trades efficiency for simplicity by threading a constant, tree-shaped random source. This makes each `random` computation linear-time in the depth it appears in the completely inlined program, which can turn functions that should be linear-time into quadratic-time. Passing subtrees instead would make `random` constant-time, restoring these functions' apparent time complexity. Passing subtrees would also allow combinators to detect lack of change in their received subtrees and return cached values.

Besides rectangles, we intend to try other simple, but more expressive set representations, such as parallelotopes [4]. For representing sets of strings, regular expressions look promising [XXX: cite paper from CS 686], as does the idea of using *abstract* regular expressions to represent sets of lengths and other integral values.

Model equivalence in distribution extends readily to uncountable spaces. It defines a standard for measure-theoretic optimizations, which can only be done in the exact semantics. Examples of optimizations we could prove correct and implement using this equivalence are variable collapse, a probabilistic analogue of constant folding that can increase efficiency by an order of magnitude, and a probabilistic analogue of constraint propagation to speed up conditional queries.

More broadly, we hope to advance Bayesian practice by providing rich modeling languages with an efficient, correct implementations, which allow general recursion and every computable, probabilistic condition.

References

- [1] Haskell 98 language and libraries, the revised report, December 2002. URL <http://www.haskell.org/onlinereport/>.
- [2] Stephen Abbott. *Understanding Analysis*. Springer, 2001.
- [3] Peter Aczel. An introduction to inductive definitions. *Studies in Logic and the Foundations of Mathematics*, 90:739–782, 1977.
- [4] G. Amato and F. Scozzari. The abstract domain of parallelotopes. *Electronic Notes in Theoretical Computer Science*, 287:17–28, November 2012.
- [5] Robert J. Aumann. Borel structures for function spaces. *Illinois Journal of Mathematics*, 5:614–630, 1961.
- [6] Bruno Barras. Sets in Coq, Coq in sets. *Journal of Formalized Reasoning*, 3(1), 2010.
- [7] C. Berline and K. Grue. A κ -denotational semantics for Map Theory in ZFC+SI. *Theoretical Computer Science*, 179(1–2):137–202, 1997.
- [8] Yves Bertot and Pierre Castéran. *Interactive Theorem Proving and Program Development. Coq’Art: The Calculus of Inductive Constructions*. Texts in Theoretical Computer Science. Springer Verlag, 2004. URL <http://www.labri.fr/publications/13a/2004/BC04>.
- [9] Sooraj Bhat, Johannes Borgström, Andrew D. Gordon, and Claudio Russo. Deriving probability density functions from probabilistic functional programs. In *Tools and Algorithms for the Construction and Analysis of Systems*, 2013.
- [10] Keith A Bonawitz. *Composable Probabilistic Inference with Blaise*. PhD thesis, Massachusetts Institute of Technology, 2008.
- [11] Johannes Borgström, Andrew D. Gordon, Michael Greenberg, James Margetson, and Jurgen Van Gael. Measure transformer semantics for Bayesian machine learning. In *European Symposium on Programming*, pages 77–96, 2011.
- [12] Manuel M. T. Chakravarty, Gabriele Keller, Simon Peyton Jones, and Simon Marlow. Associated types with class. In *Principles of Programming Languages*, pages 1–13, 2005.

- [13] Patrick Cousot and Radhia Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Principles of Programming Languages*, pages 238–252, 1977.
- [14] Ryan Culpepper. *Refining Syntactic Sugar: Tools for Supporting Macro Development*. PhD thesis, Northeastern University, 2010. To Appear.
- [15] R. C. Flagg and J. Myhill. A type-free system extending ZFC. *Annals of Pure and Applied Logic*, 43:79–97, 1989.
- [16] Matthew Flatt and PLT. Reference: Racket. Technical Report PLT-TR-2010-1, PLT Inc., 2010. <http://racket-lang.org/tr1/>.
- [17] Noah Goodman, Vikash Mansinghka, Daniel Roy, Keith Bonawitz, and Joshua Tenenbaum. Church: a language for generative models. In *Uncertainty in Artificial Intelligence*, 2008.
- [18] Michael Hanus. Multi-paradigm declarative languages. In *Logic Programming*, pages 45–75. 2007.
- [19] Martin Hofmann, Benjamin C. Pierce, , and Daniel Wagner. Edit lenses. In *Principles of Programming Languages*, 2012.
- [20] K. Hrbacek and T.J. Jech. *Introduction to set theory*. Pure and Applied Mathematics. M. Dekker, 1999.
- [21] John Hughes. Generalizing monads to arrows. In *Science of Computer Programming*, volume 37, pages 67–111, 2000.
- [22] Joe Hurd. *Formal Verification of Probabilistic Algorithms*. PhD thesis, University of Cambridge, 2002.
- [23] Claire Jones. *Probabilistic Non-Determinism*. PhD thesis, Univ. of Edinburgh, 1990.
- [24] R. B. Kearfott. Interval computations: Introduction, uses, and resources. *Euromath Bulletin*, 2:95–112, 1996.
- [25] Richard Kennaway, Jan Willem Klop, M. Ronan Sleep, and Fer jan De Vries. Infinitary lambda calculus. *Theoretical Computer Science*, 175:93–125, 1997.
- [26] Oleg Kiselyov and Chung-chieh Shan. Monolingual probabilistic programming using generalized coroutines. In *Uncertainty in Artificial Intelligence*, 2008.

- [27] Achim Klenke. *Probability Theory: A Comprehensive Course*. Springer, 2006. ISBN 978-1-84800-047-6.
- [28] Daphne Koller, David McAllester, and Avi Pfeffer. Effective Bayesian inference for stochastic programs. In *14th National Conference on Artificial Intelligence*, August 1997.
- [29] Dexter Kozen. Semantics of probabilistic programs. In *Foundations of Computer Science*, 1979.
- [30] Daniel Leivant. Higher order logic. In *In Handbook of Logic in Artificial Intelligence and Logic Programming*, pages 229–321. Clarendon Press, 1994.
- [31] Sam Lindley, Philip Wadler, and Jeremy Yallop. Idioms are oblivious, arrows are meticulous, monads are promiscuous. *Electronic Notes in Theoretical Computer Science*, 2008.
- [32] Sam Lindley, Philip Wadler, and Jeremy Yallop. The arrow calculus. *Journal of Functional Programming*, 20:51–69, 2010.
- [33] David J. Lunn, Andrew Thomas, Nicky Best, and David Spiegelhalter. WinBUGS – a Bayesian modelling framework. *Statistics and Computing*, 10(4), 2000.
- [34] Robert Mateescu and Rina Dechter. Mixed deterministic and probabilistic networks. *Annals of Mathematics and Artificial Intelligence*, 2008.
- [35] Conor McBride and Ross Paterson. Applicative programming with effects. *Journal of Functional Programming*, 18(1), 2008.
- [36] Brian Milch, Bhaskara Marthi, Stuart Russell, David Sontag, Daniel Ong, and Andrey Kolobov. BLOG: Probabilistic models with unknown objects. In *International Joint Conference on Artificial Intelligence*, 2005.
- [37] James R. Munkres. *Topology*. Prentice Hall, second edition, 2000.
- [38] Paul J. Nahin. *Duelling Idiots and Other Probability Puzzlers*. Princeton University Press, 2000.
- [39] Russell O’Connor. Certified exact transcendental real number computation in Coq. In *TPHOLs’08*, pages 246–261, 2008.
- [40] Toby Ord. The many forms of hypercomputation. *Applied Mathematics and Computation*, 178:143–153, 2006.

- [41] Sungwoo Park, Frank Pfenning, and Sebastian Thrun. A probabilistic language based upon sampling functions. *Transactions on Programming Languages and Systems*, 31(1), 2008.
- [42] L. C. Paulson. Set theory for verification: II. Induction and recursion. *Journal of Automated Reasoning*, 15:167–215, 1995.
- [43] Lawrence C. Paulson. Set theory for verification: I. From foundations to functions. *Journal of Automated Reasoning*, 11:353–389, 1993.
- [44] Avi Pfeffer. The design and implementation of IBAL: A general-purpose probabilistic language. In *Statistical Relational Learning*. MIT Press, 2007.
- [45] Norman Ramsey and Avi Pfeffer. Stochastic lambda calculus and monads of probability distributions. In *Principles of Programming Languages*, 2002.
- [46] S M Samuels. The Radon-Nikodym theorem as a theorem in probability. *The American Mathematical Monthly*, 85(3):155–165, March 1978.
- [47] Sam Tobin-Hochstadt and Matthias Felleisen. The design and implementation of typed Scheme. In *Principles of Programming Languages*, pages 395–406, 2008.
- [48] Neil Toronto and Jay McCarthy. From Bayesian notation to pure Racket, via measure-theoretic probability in λ_{ZFC} . In *Impl. and Appl. of Functional Languages*, 2010.
- [49] Neil Toronto and Jay McCarthy. Computing in Cantor’s paradise with λ_{ZFC} . In *Functional and Logic Programming Symposium*, pages 290–306, 2012.
- [50] Neil Toronto, Bryan S. Morse, Kevin Seppi, and Dan Ventura. Super-resolution via recapture and Bayesian effect modeling. In *Computer Vision and Pattern Recognition*, 2009.
- [51] Alan M. Turing. On computable numbers, with an application to the Entscheidungsproblem. In *Proceedings of the London Mathematical Society*, volume 42, pages 230–265, 1936.
- [52] Athanassios Tzouvaras. Cardinality without enumeration. *Studia Logica: An International Journal for Symbolic Logic*, 80(1):121–141, June 2005.
- [53] Gabriel Uzquiano. Models of second-order Zermelo set theory. *The Bulletin of Symbolic Logic*, 5(3):289–302, 1999.

- [54] Philip Wadler. Monads for functional programming. In J. Jeuring and E. Meijer, editors, *Advanced Functional Programming*. 2001.
- [55] Benjamin Werner. Sets in types, types in sets. In *TACS'97*, pages 530–546, 1997.
- [56] David Wingate, Noah D. Goodman, Andreas Stuhlmüller, and Jeffrey M. Siskind. Nonstandard interpretations of probabilistic programs for efficient inference. In *Neural Information Processing Systems*, pages 1152–1160, 2011.
- [57] David Wingate, Andreas Stuhlmüller, and Noah D. Goodman. Lightweight implementations of probabilistic programming languages via transformational compilation. In *Artificial Intelligence and Statistics*, 2011.