# From Bayesian Notation to Pure Racket[*]

## via Discrete Measure-Theoretic Probability in $\lambda_{ZFC}$

**Implementation and Application of Functional Languages**

**September 1-3, 2010**

Neil Toronto and Jay McCarthy

PLT @ Brigham Young University, Utah, USA

[*] Formerly PLT-Scheme

# Bayesian Practice and Philosophy

- State a probabilistic **model** of a process; then pose a **query** that runs the process backwards

Doing Bayesian statistics is like doing physics, but for fuzzy, uncertain, or random things. In physics, you might use a Newtonian model that gives distance in terms of time (i.e. d=rt) and ask how much time it will take to travel given a certain distance.

2

# Bayesian Practice and Philosophy

- State a probabilistic **model** of a process; then pose a **query** that runs the process backwards

  - Document generation / Is this email spam?

2

# Bayesian Practice and Philosophy

- State a probabilistic **model** of a process; then pose a **query** that runs the process backwards

  ○ Document generation / Is this email spam?

  ○ Real-world scenes and image capture / Likely scene given a photograph

2

# Bayesian Practice and Philosophy

- State a probabilistic **model** of a process; then pose a **query** that runs the process backwards

  - Document generation / Is this email spam?

  - Real-world scenes and image capture / Likely scene given a photograph

---

"An approximate answer to the right question is worth a great deal more than a precise answer to the wrong question."

**John Wilder Tukey**

# Philosophy Into Practice (1)

- Approximations must be put off as long as possible

3

# Philosophy Into Practice (1)

- Approximations must be put off as long as possible

    ○ Models and queries are exact, and generally not closed-form nor finitely computable

3

# Philosophy Into Practice (1)

- Approximations must be put off as long as possible

  ○ Models and queries are exact, and generally not closed-form nor finitely computable

  ○ Compute answers as converging approximations

# Philosophy Into Practice (1)

- Approximations must be put off as long as possible

  - Models and queries are exact, and generally not closed-form nor finitely computable

  - Compute answers as converging approximations

- Example: enlarging images

Model and query

$$S_{i,j}^{\theta} \sim \text{Uniform}(-\pi, \pi) \quad S_{i,j}^{v^+} \sim \text{Uniform}(0, 1)$$
$$S_{i,j}^{d} \sim \text{Uniform}(-3, 3) \quad S_{i,j}^{v^-} \sim \text{Uniform}(0, 1)$$
$$S_{i,j}^{\sigma} \sim \text{Beta}(1.6, 1)$$

$$I_{i,j} | S_{N9(i,j)} \sim \text{Normal}(E[S_{i,j}], \omega)$$

$$\Phi_{i,j}(S_{N9(i,j)}) \equiv \exp\left(-\frac{\text{Var}[S_{i,j}]}{2\gamma^2}\right)$$

What is the distribution of $I'|I$?

3

# Philosophy Into Practice (1)

- Approximations must be put off as long as possible

  - Models and queries are exact, and generally not closed-form nor finitely computable

  - Compute answers as converging approximations

- Example: enlarging images

| Model and query | The answer's approximation |
|---|---|

$$S_{i,j}^{\theta} \sim \text{Uniform}(-\pi, \pi) \qquad S_{i,j}^{v^+} \sim \text{Uniform}(0,1)$$
$$S_{i,j}^{d} \sim \text{Uniform}(-3, 3) \qquad S_{i,j}^{v^-} \sim \text{Uniform}(0,1)$$
$$S_{i,j}^{\sigma} \sim \text{Beta}(1.6, 1)$$

$$\mathbf{I}_{i,j} | \mathbf{S}_{N9(i,j)} \sim \text{Normal}(\mathbf{E}[\mathbf{S}_{i,j}], \omega)$$

$$\Phi_{i,j}(\mathbf{S}_{N9(i,j)}) \equiv \exp\left(-\frac{\text{Var}[\mathbf{S}_{i,j}]}{2\gamma^2}\right)$$

What is the distribution of $\mathbf{I}'|\mathbf{I}$?



3

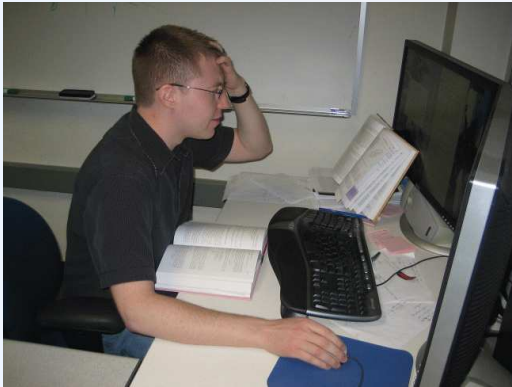# Philosophy Into Practice (2)

- Grads become compilers (and are just as grumpy)

4

# Philosophy Into Practice (2)

- Grads become compilers (and are just as grumpy)



- Our motivation: free (fire?) the grad students

4

# Philosophy Into Practice (2)

- Grads become compilers (and are just as grumpy)

- Our motivation: free (fire?) the grad students

- Our primary constraints

  - Do not approximate earlier than users would

  - Do not force users to approximate early

# Philosophical Constraints

- Compatible approach

  1. Informally determine meaning of notation

5

# Philosophical Constraints

- Compatible approach

  1. Informally determine meaning of notation

  2. Develop exact $[\![ \cdot ]\!]$ : "notation" $\rightarrow$ "calculations"

5

# Philosophical Constraints

- Compatible approach

  1. Informally determine meaning of notation

  2. Develop exact $[\![ \cdot ]\!]$ : "notation" $\rightarrow$ "calculations"

  3. Approximate $[\![ \cdot ]\!]$, prove convergence

5

# Philosophical Constraints

- Compatible approach

  1. Informally determine meaning of notation

  2. Develop exact $[\![ \cdot ]\!]$ : "notation" $\rightarrow$ "calculations"

  3. Approximate $[\![ \cdot ]\!]$, prove convergence

  4. Implement approximating $[\![ \cdot ]\!]$

5

# Philosophical Constraints

- Compatible approach

    1. Informally determine meaning of notation

    2. Develop exact $[\![\cdot]\!]$ : "notation" $\rightarrow$ "calculations"

    3. Approximate $[\![\cdot]\!]$, prove convergence

    4. Implement approximating $[\![\cdot]\!]$

- Analogous to abstract interpretation

    concrete/exact, abstract/approximating

# Which Probability Theory?

- Naive/undergraduate/informal probability theory

  - How Bayesians tend to think about probability

To turn notation into exact calculations, we need a theory of probability that tells us what those calculations should be. Bayesians tend to think and calculate using naive probability, which you probably learned if you had to take an undergraduate statistics course. But we can't use naive probability.

6

# Which Probability Theory?

- Naive/undergraduate/informal probability theory

  - How Bayesians tend to think about probability

  - But can't properly express infinities

# Which Probability Theory?

- Naive/undergraduate/informal probability theory

  ○ How Bayesians tend to think about probability

  ○ But can't properly express infinities

  ○ "Spooky interaction at a distance"

# Which Probability Theory?

- Naive/undergraduate/informal probability theory

    ○ How Bayesians tend to think about probability

    ○ But can't properly express infinities

    ○ "Spooky interaction at a distance"

- Measure-theoretic probability: global $(\Omega, \Sigma, \mathbb{P})$

Measure-theoretic probability explains non-local interaction by having all random variables interact through a single, global object. You might call it a `store,' but its actual name is `probability space.'

6

# Which Probability Theory?

- Naive/undergraduate/informal probability theory

  ○ How Bayesians tend to think about probability

  ○ But can't properly express infinities

  ○ "Spooky interaction at a distance"

- Measure-theoretic probability: global $(\Omega, \Sigma, \mathbb{P})$

  ○ $\Omega$ : set: all possible "worlds"

6

# Which Probability Theory?

- Naive/undergraduate/informal probability theory

  ○ How Bayesians tend to think about probability

  ○ But can't properly express infinities

  ○ "Spooky interaction at a distance"

- Measure-theoretic probability: global $(\Omega, \Sigma, \mathbb{P})$

  ○ $\Omega$ : set: all possible "worlds"

  ○ $X : \Omega \to S$: "getters" or "readers" for worlds
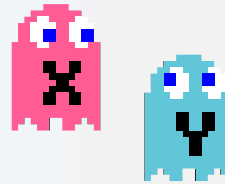
6

# Which Probability Theory?

- Naive/undergraduate/informal probability theory

  - How Bayesians tend to think about probability

  - But can't properly express infinities

  - "Spooky interaction at a distance"

- Measure-theoretic probability: global $(\Omega, \Sigma, \mathbb{P})$

  - $\Omega$ : set: all possible "worlds"

  - $X : \Omega \to S$: "getters" or "readers" for worlds

  - $\Sigma$ : set: measurable events (for uncountable $\Omega$)
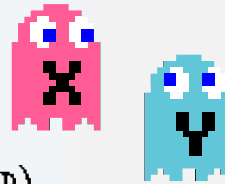
6

# Which Probability Theory?

- Naive/undergraduate/informal probability theory

  ○ How Bayesians tend to think about probability

  ○ But can't properly express infinities

  ○ "Spooky interaction at a distance"

- Measure-theoretic probability: global $(\Omega, \Sigma, \mathbb{P})$

  ○ $\Omega$ : set: all possible "worlds"

  ○ $X : \Omega \to S$: "getters" or "readers" for worlds

  ○ $\Sigma$ : set: measurable events (for uncountable $\Omega$)

  ○ $\mathbb{P} : \Sigma \to [0, 1]$ (or $P : \Omega \to [0, 1]$): probabilities of events
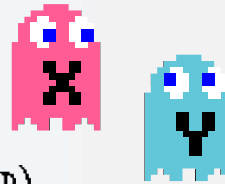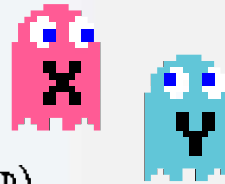
6

# Which Probability Theory?

- Naive/undergraduate/informal probability theory

  - How Bayesians tend to think about probability

  - But can't properly express infinities

  - "Spooky interaction at a distance"

- Measure-theoretic probability: global $(\Omega, \Sigma, \mathbb{P})$

  - $\Omega$ : set: all possible "worlds"

  - $X : \Omega \to S$: "getters" or "readers" for worlds

  - $\Sigma$ : set: measurable events (for uncountable $\Omega$)

  - $\mathbb{P} : \Sigma \to [0, 1]$ (or $P : \Omega \to [0, 1]$): probabilities of events

  - Calculations uncomputable when $\Omega$ infinite

6

# Lambda-ZFC

- Looking for language for compositional measure- theoretic calculations; similarity to Racket a plus

We want to transform notation into measure-theoretic calculations, and eventually approximate the calculations in Racket. For the semantic function's target language, then, we need a call-by-value lambda calculus for expressing uncomputable things that are well-defined in contemporary mathematics.

# Lambda-ZFC

- Looking for language for compositional measure- theoretic calculations; similarity to Racket a plus



$\lambda$ calculus

$$((\lambda x.e, x, e\ e), \alpha, \beta)$$

# Lambda-ZFC

- Looking for language for compositional measure- theoretic calculations; similarity to Racket a plus



$\lambda$ calculus

$((\lambda x.e, x, e\ e), \alpha, \beta)$

Set theory

$(V, \in, =, \{\cdot\}, \bigcup, image, \mathcal{P}, order)$

# Lambda-ZFC

- Looking for language for compositional measure- theoretic calculations; similarity to Racket a plus



$\lambda$ calculus
$((\lambda x.e, x, e\ e), \alpha, \beta)$

$+$

Set theory
$(V, \in, =, \{\cdot\}, \bigcup, image, \mathcal{P}, order)$

$=$

$\lambda_{ZFC}$

7

# Lambda-ZFC

- Looking for language for compositional measure- theoretic calculations; similarity to Racket a plus



$\lambda$ calculus
$((\lambda x.e, x, e\ e), \alpha, \beta)$

Set theory
$(V, \in, =, \{\cdot\}, \bigcup, image, \mathcal{P}, order)$

$\lambda_{ZFC}$

+     =

- "Programming" is doing contemporary math, plus $\lambda x.e$

7

# Lambda-ZFC

- Looking for language for compositional measure- theoretic calculations; similarity to Racket a plus



$\lambda$ calculus

$((\lambda x.e, x, e\ e), \alpha, \beta)$

Set theory

$(V, \in, =, \{\cdot\}, \bigcup, image, \mathcal{P}, order)$

$\lambda_{ZFC}$

- "Programming" is doing contemporary math, plus $\lambda$x.e

- Contains all set-theoretic functions

7

# Lambda-ZFC

- Looking for language for compositional measure- theoretic calculations; similarity to Racket a plus



λ calculus
$((\lambda x.e, x, e\ e), \alpha, \beta)$

Set theory
$(V, \in, =, \{\cdot\}, \bigcup, image, \mathcal{P}, order)$

$\lambda_{ZFC}$

- "Programming" is doing contemporary math, plus λx.e

- Contains all set-theoretic functions

- Can solve any OTM halting problem constructively

To give you an idea of its relative computational power: you can solve any oracle Turing machine halting problem by writing an interpreter in lambda-ZFC. It might seem like too much power, but remember that we want to interpret Bayesian notation exactly. We'll worry about computability when we do the approximations.

# Interpreting Notation

| Syntactic Category | Examples | Computational Structure | Semantic Functions |
|---|---|---|---|

Developing the whole semantics right now would take too much time, so I'm going to give some examples of syntax and talk about the structure of the calculations.

First we have random variable expressions. In the first example, X and Y are random variables, so they're functions of Omega. I've already hinted that you could interpret this by regarding random variables as reader monad computations.

But there's no reason to impose a total order, so we use the corresponding applicative functor, or idiom.

Next, we have statements about random variables. A collection of statements is a probabilistic model. We interpret each statement as transforming the global probability space. The first example, X is distributed Geometric B, extends the probability space. The second example is a `condition,' which asserts that applying the random variable X+Y to any world must yield 4. It *restricts* the global probability space.

A nice way to structure these calculations is with the state monad, with probability-space-valued state.

Last, we have queries. The first example is a `conditional probability query'. It conditions the probability space first, and then asks for the probability that B outputs 1/2. The second example is like the first, but is parameterized on the outputs of B and X+Y. It should return a function, or a distribution, so it's a `distribution query'.

Queries run the probability space monad computation in their own particular way.

8

# Interpreting Notation

| Syntactic Category | Examples | Computational Structure | Semantic Functions |
|---|---|---|---|
| *Expressions* | $X + Y$ <br> $\text{Geometric}(B)$ | | |

Developing the whole semantics right now would take too much time, so I'm going to give some examples of syntax and talk about the structure of the calculations.

First we have random variable expressions. In the first example, X and Y are random variables, so they're functions of Omega. I've already hinted that you could interpret this by regarding random variables as reader monad computations.

But there's no reason to impose a total order, so we use the corresponding applicative functor, or idiom.

Next, we have statements about random variables. A collection of statements is a probabilistic model. We interpret each statement as transforming the global probability space. The first example, X is distributed Geometric B, extends the probability space. The second example is a `condition,' which asserts that applying the random variable X+Y to any world must yield 4. It *restricts* the global probability space.

A nice way to structure these calculations is with the state monad, with probability-space-valued state.

Last, we have queries. The first example is a `conditional probability query'. It conditions the probability space first, and then asks for the probability that B outputs 1/2. The second example is like the first, but is parameterized on the outputs of B and X+Y. It should return a function, or a distribution, so it's a `distribution query'.

Queries run the probability space monad computation in their own particular way.

# Interpreting Notation

| Syntactic Category | Examples | Computational Structure | Semantic Functions |
|---|---|---|---|
| *Expressions* | $X + Y$ <br> $\text{Geometric}(B)$ | Environment idiom: <br> $R\,a = \Omega \to a$ | $\mathcal{R}[\![\cdot]\!]$, $\mathbf{RV}$ |

8

Developing the whole semantics right now would take too much time, so I'm going to give some examples of syntax and talk about the structure of the calculations.

First we have random variable expressions. In the first example, X and Y are random variables, so they're functions of Omega. I've already hinted that you could interpret this by regarding random variables as reader monad computations.

But there's no reason to impose a total order, so we use the corresponding applicative functor, or idiom.

Next, we have statements about random variables. A collection of statements is a probabilistic model. We interpret each statement as transforming the global probability space. The first example, X is distributed Geometric B, extends the probability space. The second example is a `condition,' which asserts that applying the random variable X+Y to any world must yield 4. It *restricts* the global probability space.

A nice way to structure these calculations is with the state monad, with probability-space-valued state.

Last, we have queries. The first example is a `conditional probability query'. It conditions the probability space first, and then asks for the probability that B outputs 1/2. The second example is like the first, but is parameterized on the outputs of B and X+Y. It should return a function, or a distribution, so it's a `distribution query'.

Queries run the probability space monad computation in their own particular way.

# Interpreting Notation

| Syntactic Category | Examples | Computational Structure | Semantic Functions |
|---|---|---|---|
| *Expressions* | $X + Y$<br>$\mathrm{Geometric}(B)$ | Environment idiom:<br>$R\,a = \Omega \to a$ | $\mathcal{R}[\![\cdot]\!]$, $\mathbf{RV}$ |
| *Statements* | $X \sim \mathrm{Geometric}(B)$<br>$X + Y = 4$ | | |

8

# Interpreting Notation

| Syntactic Category | Examples | Computational Structure | Semantic Functions |
|---|---|---|---|
| *Expressions* | $X + Y$ <br> $\mathrm{Geometric}(B)$ | Environment idiom: <br> $R\ a = \Omega \to a$ | $\mathcal{R}[\![\cdot]\!]$, `RV` |
| *Statements* | $X \sim \mathrm{Geometric}(B)$ <br> $X + Y = 4$ | State monad: <br> $M\ b = PS \to (PS, b)$ <br> (usually $b = R\ a$) | $\mathcal{M}[\![\cdot]\!]$, `model` |

8

Developing the whole semantics right now would take too much time, so I'm going to give some examples of syntax and talk about the structure of the calculations.

First we have random variable expressions. In the first example, X and Y are random variables, so they're functions of Omega. I've already hinted that you could interpret this by regarding random variables as reader monad computations.

But there's no reason to impose a total order, so we use the corresponding applicative functor, or idiom.

Next, we have statements about random variables. A collection of statements is a probabilistic model. We interpret each statement as transforming the global probability space. The first example, X is distributed Geometric B, extends the probability space. The second example is a `condition,' which asserts that applying the random variable X+Y to any world must yield 4. It *restricts* the global probability space.

A nice way to structure these calculations is with the state monad, with probability-space-valued state.

Last, we have queries. The first example is a `conditional probability query'. It conditions the probability space first, and then asks for the probability that B outputs 1/2. The second example is like the first, but is parameterized on the outputs of B and X+Y. It should return a function, or a distribution, so it's a `distribution query'.

Queries run the probability space monad computation in their own particular way.

# Interpreting Notation

| Syntactic Category | Examples | Computational Structure | Semantic Functions |
|---|---|---|---|
| *Expressions* | $X + Y$ <br> $\mathrm{Geometric}(B)$ | Environment idiom: <br> $R\ a = \Omega \to a$ | $\mathcal{R}[\![\cdot]\!]$, `RV` |
| *Statements* | $X \sim \mathrm{Geometric}(B)$ <br> $X + Y = 4$ | State monad: <br> $M\ b = PS \to (PS, b)$ <br> (usually $b = R\ a$) | $\mathcal{M}[\![\cdot]\!]$, `model` |
| *Queries* | $\mathbf{P}\big[B = \tfrac{1}{2} \mid X + Y = 4\big]$ <br> $\mathcal{L}[B \mid X + Y]$ | | |

8

Developing the whole semantics right now would take too much time, so I'm going to give some examples of syntax and talk about the structure of the calculations.

First we have random variable expressions. In the first example, X and Y are random variables, so they're functions of Omega. I've already hinted that you could interpret this by regarding random variables as reader monad computations.

But there's no reason to impose a total order, so we use the corresponding applicative functor, or idiom.

Next, we have statements about random variables. A collection of statements is a probabilistic model. We interpret each statement as transforming the global probability space. The first example, X is distributed Geometric B, extends the probability space. The second example is a `condition,' which asserts that applying the random variable X+Y to any world must yield 4. It *restricts* the global probability space.

A nice way to structure these calculations is with the state monad, with probability-space-valued state.

Last, we have queries. The first example is a `conditional probability query'. It conditions the probability space first, and then asks for the probability that B outputs 1/2. The second example is like the first, but is parameterized on the outputs of B and X+Y. It should return a function, or a distribution, so it's a `distribution query'.

Queries run the probability space monad computation in their own particular way.

# Interpreting Notation

| Syntactic Category | Examples | Computational Structure | Semantic Functions |
|---|---|---|---|
| *Expressions* | $X + Y$ <br> $\mathrm{Geometric}(B)$ | Environment idiom: <br> $R\ a = \Omega \to a$ | $\mathcal{R}[\![\cdot]\!]$, RV |
| *Statements* | $X \sim \mathrm{Geometric}(B)$ <br> $X + Y = 4$ | State monad: <br> $M\ b = PS \to (PS, b)$ <br> (usually $b = R\ a$) | $\mathcal{M}[\![\cdot]\!]$, model |
| *Queries* | $\mathbf{P}\left[ B = \tfrac{1}{2} \mid X + Y = 4 \right]$ <br> $\mathcal{L}[B \mid X + Y]$ | State monad run: <br> $b = [0,1]$ or <br> $b = a \to c \to [0,1]$ | $\mathbf{P}[\![\cdot]\!], \mathbf{D}[\![\cdot]\!]$, <br> Pr, Dist |

8

Developing the whole semantics right now would take too much time, so I'm going to give some examples of syntax and talk about the structure of the calculations.

First we have random variable expressions. In the first example, X and Y are random variables, so they're functions of Omega. I've already hinted that you could interpret this by regarding random variables as reader monad computations.

But there's no reason to impose a total order, so we use the corresponding applicative functor, or idiom.

Next, we have statements about random variables. A collection of statements is a probabilistic model. We interpret each statement as transforming the global probability space. The first example, X is distributed Geometric B, extends the probability space. The second example is a `condition,' which asserts that applying the random variable X+Y to any world must yield 4. It *restricts* the global probability space.

A nice way to structure these calculations is with the state monad, with probability-space-valued state.

Last, we have queries. The first example is a `conditional probability query'. It conditions the probability space first, and then asks for the probability that B outputs 1/2. The second example is like the first, but is parameterized on the outputs of B and X+Y. It should return a function, or a distribution, so it's a `distribution query'.

Queries run the probability space monad computation in their own particular way.

# Interpreting Notation

| Syntactic Category | Examples | Computational Structure | Semantic Functions |
|---|---|---|---|
| *Expressions* | $X + Y$ <br> $\mathrm{Geometric}(B)$ | Environment idiom: <br> $R\ a = \Omega \to a$ | $\mathcal{R}[\![\cdot]\!]$, `RV` |
| *Statements* | $X \sim \mathrm{Geometric}(B)$ <br> $X + Y = 4$ | State monad: <br> $M\ b = PS \to (PS, b)$ <br> (usually $b = R\ a$) | $\mathcal{M}[\![\cdot]\!]$, `model` |
| *Queries* | $\mathbf{P}\left[B = \tfrac{1}{2} \mid X + Y = 4\right]$ <br> $\mathcal{L}[B \mid X + Y]$ | State monad run: <br> $b = [0, 1]$ or <br> $b = a \to c \to [0, 1]$ | $\mathbf{P}[\![\cdot]\!]$, $\mathbf{D}[\![\cdot]\!]$, <br> `Pr`, `Dist` |

- Difficult to encode types in most type systems

8

Now, you might think that with all these types and idioms and monads and such, Racket might not be the best implementation language. But Racket's macro system allows us to implement the semantic functions directly. Besides, these types are difficult to encode in most type systems; it seems to require either unityped random variables or dependent types. It's possible in Typed Racket using occurrence typing, but it's a little too much trouble.

# Fun Facts: Semantics

- $\mathcal{R}[\![\cdot]\!] : \lambda_{\mathrm{ZFC}} \to \lambda_{\mathrm{ZFC}}$ interprets anything constructive
  - ○ Uncountable $\Omega$: need to prove measurability conditions

The random variable expression semantic function can turn any lambda-ZFC expression into random variable. When Omega is uncountable, we're going to have to start worrying about something called measurability, but we'll be able to worry about it compositionally.

# Fun Facts: Semantics

- $\mathcal{R}[\![\cdot]\!] : \lambda_{\mathrm{ZFC}} \to \lambda_{\mathrm{ZFC}}$ interprets anything constructive

  - Uncountable $\Omega$: need to prove measurability conditions

- $\mathcal{R}[\![\mathrm{Geometric}(B)]\!] : \Omega \to \mathbb{N} \to [0,1]$ is a discrete *transition kernel*

  - Uncountable $\Omega$ already works: $\mathcal{R}[\![\mathrm{Normal}(M,S)]\!] : \Omega \to \mathcal{B}(\mathbb{R}) \to [0,1]$

Next is a fortunate accident: the random variable semantic function turns notation that denotes conditional distributions into `transition kernels,' which measure-theoretic probability uses to build probability spaces. So the semantics turns Bayesian notation into exactly what measure-theoretic probability requires, and that fact doesn't change when Omega is uncountable. It also allows Bayesians more freedom: any expression with the right type can be a conditional distribution. I'll show an example later.

# Fun Facts: Semantics

- $\mathcal{R}[\![\cdot]\!] : \lambda_{\mathrm{ZFC}} \to \lambda_{\mathrm{ZFC}}$ interprets anything constructive

  - ○ Uncountable $\Omega$: need to prove measurability conditions

- $\mathcal{R}[\![\mathrm{Geometric}(B)]\!] : \Omega \to \mathbb{N} \to [0,1]$ is a discrete *transition kernel*

  - ○ Uncountable $\Omega$ already works: $\mathcal{R}[\![\mathrm{Normal}(M,S)]\!] : \Omega \to \mathcal{B}(\mathbb{R}) \to [0,1]$

- Queries approximate with $\mathbf{finitize}\ (\Omega, P)\ k = (\Omega_k, P_k)$

  - ○ Uncountable $\Omega$: $\mathbf{finitize}\ (\Omega, \Sigma, \mathbb{P})\ k = (\Omega_k, P_k)$ with $\Omega_k$ finite, stochastic

There's a single point of approximation in the approximating semantics: right before a query, `finitize' restricts Omega to a finite subset of size k. Then, as k approaches infinity, the answer to any query approaches the correct value. Finitize also renormalizes P so that it sums to 1.

Approximations for uncountable Omega are going to be tricky, but there are a lot of available approximations. The most efficient ones are randomized algorithms.

# Fun Facts: Implementation

- Almost a transliteration of approximating semantics, except

10

# Fun Facts: Implementation

- Almost a transliteration of approximating semantics, except

  - Lazy lists represent recursively enumerable sets

  - Floats and exact rationals represent probabilities

10

# Fun Facts: Implementation

- Almost a transliteration of approximating semantics, except

  ○ Lazy lists represent recursively enumerable sets

  ○ Floats and exact rationals represent probabilities

- `RV : kstx -> kstx` interprets any Racket expression

10

# Fun Facts: Implementation

- Almost a transliteration of approximating semantics, except

  - Lazy lists represent recursively enumerable sets

  - Floats and exact rationals represent probabilities

- `RV : kstx -> kstx` interprets any Racket expression

- `(define-model name [X ~ ...] ...)` is hygienically referred to by `(with-model name (Pr ... X ...))`

# Duelling Idiots (Paul Nahin)

Let's see how the implementation does on a good, countably infinite probability problem. (By `good,' by the way, I mean that the problem includes gambling and death.) It comes from Paul Nahin's book of puzzlers. Two idiots decide to duel, but they have only one gun, a six-shooter. So they put a bullet in it and take turns spinning the chamber and firing at each other. What's the probability that the player that shoots first wins?

11

# Duelling Idiots (Paul Nahin)

```
(define-model idiot-duel
  [winning-shot ~ (Geometric 1/6)])
```

The trick to answering the query is to recognize that how many shots it takes before the gun finally goes off has a geometric distribution.

# Duelling Idiots (Paul Nahin)

```
(define-model idiot-duel
  [winning-shot ~ (Geometric 1/6)])

(with-model idiot-duel
  (Pr (odd? winning-shot)))
; --> 2/3 as k --> ∞
```

The probability that player one wins is the probability that the winning shot is odd-numbered, and this approaches 2/3 as k approaches infinity. So player one has a much better chance of winning this duel. But suppose the idiots know this, so they come up with a plan to even the odds. Player one takes one shot, then player two takes two shots, player one takes three shots, and so on. What's the probability that player one wins?

11

# Duelling Idiots (Paul Nahin)

```
(define-model idiot-duel
  [winning-shot ~ (Geometric 1/6)])

(with-model idiot-duel
  (Pr (odd? winning-shot)))
; --> 2/3 as k --> ∞

(with-model idiot-duel
  (Pr (p1-fires? winning-shot)))
```

11

# Duelling Idiots (Paul Nahin)

```
(define-model idiot-duel
  [winning-shot ~ (Geometric 1/6)])

(with-model idiot-duel
  (Pr (odd? winning-shot)))
; --> 2/3 as k --> ∞

(with-model idiot-duel
  (Pr (p1-fires? winning-shot)))

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

Designing p1-fires? was the trickiest part of the solution. Don't stare at it too long, though; the point is that it exists and isn't too hard to write.

# Duelling Idiots (Paul Nahin)

```
(define-model idiot-duel
  [winning-shot ~ (Geometric 1/6)])

(with-model idiot-duel
  (Pr (odd? winning-shot)))
; --> 2/3 as k --> ∞

(with-model idiot-duel
  (Pr (p1-fires? winning-shot)))

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

Nahin (MATLAB):        0.5239191275550995247919843
Us (Racket, k=321):    0.5239191275550995247919**8439**

# Duelling Idiot and Half-Wit

But the probelm isn't Bayesian! So suppose that player one is actually a half-wit, and proposes flipping a coin to see whether they will spin the chamber. If they don't spin it, the gun will go off within six shots, and for four of those shots, it will be in player one's hand. But player two is an idiot and agrees to it.

12

# Duelling Idiot and Half-Wit

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

12

# Duelling Idiot and Half-Wit

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

 (with-model half-wit-duel
   (Pr spin? (not (p1-fires? winning-shot))))
```

12

# Duelling Idiot and Half-Wit

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

(with-model half-wit-duel
  (Pr spin? (not (p1-fires? winning-shot))))
```

Answer: about 0.588 (compare `(Pr spin?)` = 0.5)

12

# Observational Equivalence

- Model equivalence: $m \equiv_{\mathbf{D}} m'$ means no query $q$ can distinguish between $m$ and $m'$

# Observational Equivalence

- Model equivalence: $m \equiv_\mathbf{D} m'$ means no query $q$ can distinguish between $m$ and $m'$

- Justifies measure-theoretic optimizations

  - Variable collapse (constant folding for rvs)

  $$X \sim \mathrm{Normal}(0,1); \ Y \sim \mathrm{Normal}(X,1) \ \longrightarrow \ Y \sim \mathrm{Normal}(0,2)$$

  - Propagating conditions (like constraint propagation)

  $$X \sim P_X; \ \ldots; \ X = 3 \ \longrightarrow \ X \sim P_X; \ X = 3; \ \ldots$$

# Observational Equivalence

- Model equivalence: $m \equiv_{\mathbf{D}} m'$ means no query $q$ can distinguish between $m$ and $m'$

- Justifies measure-theoretic optimizations

  ○ Variable collapse (constant folding for rvs)

$$X \sim \mathrm{Normal}(0,1);\ Y \sim \mathrm{Normal}(X,1) \longrightarrow Y \sim \mathrm{Normal}(0,2)$$

  ○ Propagating conditions (like constraint propagation)

$$X \sim P_X;\ \ldots;\ X = 3 \longrightarrow X \sim P_X;\ X = 3;\ \ldots$$

- Justifiable **only** in the **exact semantics**

13

# Observational Equivalence

- Model equivalence: $m \equiv_{\mathbf{D}} m'$ means no query $q$ can distinguish between $m$ and $m'$

- Justifies measure-theoretic optimizations

  ○ Variable collapse (constant folding for rvs)

  $$X \sim \mathrm{Normal}(0,1); \; Y \sim \mathrm{Normal}(X,1) \; \longrightarrow \; Y \sim \mathrm{Normal}(0,2)$$

  ○ Propagating conditions (like constraint propagation)

  $$X \sim P_X; \; \ldots; \; X = 3 \; \longrightarrow \; X \sim P_X; \; X = 3; \; \ldots$$

- Justifiable **only** in the **exact semantics**

  ○ Suppose for $k = 29$, $q \; m = 0.7$ but $q \; m' = 0.2$

# Observational Equivalence

- Model equivalence: $m \equiv_{\mathbf{D}} m'$ means no query $q$ can distinguish between $m$ and $m'$

- Justifies measure-theoretic optimizations

  ○ Variable collapse (constant folding for rvs)

$$X \sim \mathrm{Normal}(0,1); \; Y \sim \mathrm{Normal}(X,1) \; \longrightarrow \; Y \sim \mathrm{Normal}(0,2)$$

  ○ Propagating conditions (like constraint propagation)

$$X \sim P_X; \; \ldots; \; X = 3 \; \longrightarrow \; X \sim P_X; \; X = 3; \; \ldots$$

- Justifiable **only** in the **exact semantics**

  ○ Suppose for $k = 29$, $q\ m = 0.7$ but $q\ m' = 0.2$

  ○ But what if, for $k = 400$, $q\ m = 0.19$?

13

# Observational Equivalence

- Model equivalence: $m \equiv_{\mathbf{D}} m'$ means no query $q$ can distinguish between $m$ and $m'$

- Justifies measure-theoretic optimizations

  - Variable collapse (constant folding for rvs)

$$X \sim \mathrm{Normal}(0,1); \; Y \sim \mathrm{Normal}(X,1) \; \longrightarrow \; Y \sim \mathrm{Normal}(0,2)$$

  - Propagating conditions (like constraint propagation)

$$X \sim P_X; \; \ldots; \; X = 3 \; \longrightarrow \; X \sim P_X; \; X = 3; \; \ldots$$

- Justifiable **only** in the **exact semantics**

  - Suppose for $k = 29$, $q \; m = 0.7$ but $q \; m' = 0.2$

  - But what if, for $k = 400$, $q \; m = 0.19$?