

# Semantic Foundations for Probabilistic Programming

Chris Heunen



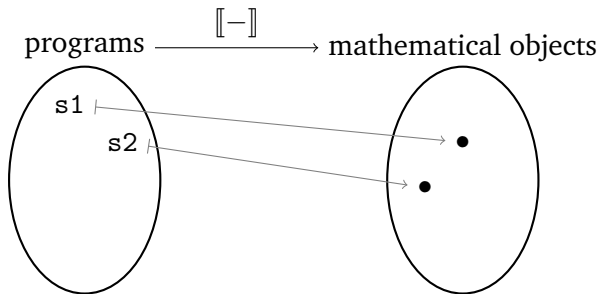
THE UNIVERSITY of EDINBURGH  
**informatics**

Ohad Kammar, Sam Staton,  
Frank Wood, Hongseok Yang

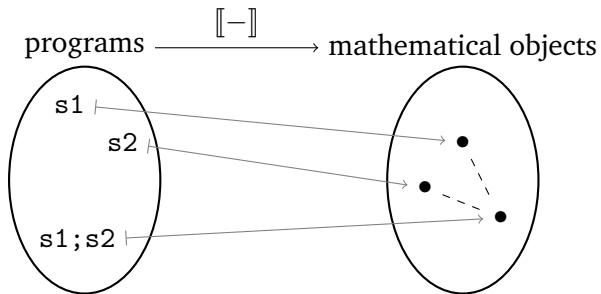


DEPARTMENT OF  
**COMPUTER  
SCIENCE**

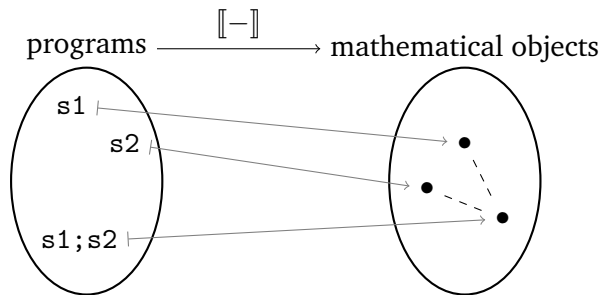
# Semantic foundations



# Semantic foundations

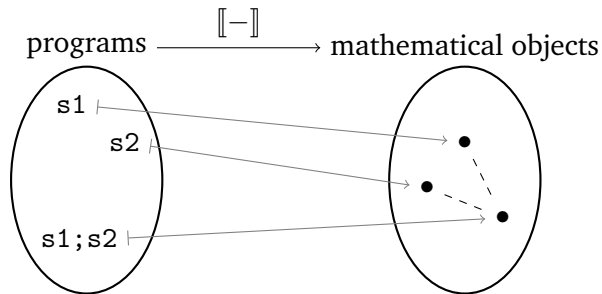


# Semantic foundations



- ▶ *Operational*: remember implementation details (efficiency)
- ▶ *Denotational*: see what program does conceptually (correctness)

# Semantic foundations

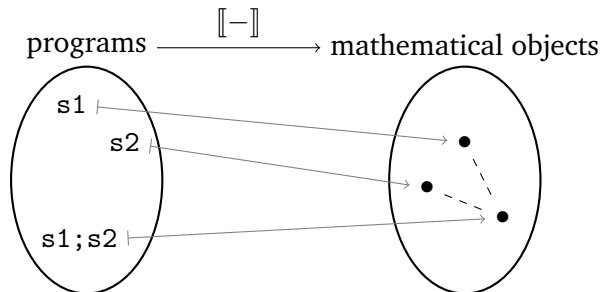


- ▶ *Operational*: remember implementation details (efficiency)
- ▶ *Denotational*: see what program does conceptually (correctness)

## Motivation:

- ▶ Ground programmer's unspoken intuitions
- ▶ Justify/refute/suggest program transformations
- ▶ Understand programming through mathematics

# Semantic foundations



- ▶ *Operational*: remember implementation details (efficiency)
- ▶ *Denotational*: see what program does conceptually (correctness)

## Motivation:

- ▶ Ground programmer's unspoken intuitions
- ▶ Justify/refute/suggest program transformations
- ▶ Understand *probability* through *program equations*

# Probabilistic programming

$$P(A | B) = \frac{P(B | A) \times P(A)}{P(B)}$$

# Probabilistic programming

$$P(A | B) \propto P(B | A) \times P(A)$$



# Probabilistic programming

$$P(A | B) \propto P(B | A) \times P(A)$$

posterior  $\propto$  likelihood  $\times$  prior

# Probabilistic programming

$$P(A | B) \propto P(B | A) \times P(A)$$

posterior  $\propto$  likelihood  $\times$  prior

idealized Anglican = functional programming +  
**normalize**    **observe**    **sample**

# Overview

- ▶ Interpret types as measurable spaces e.g.  $\llbracket \text{real} \rrbracket = \mathbb{R}$
- ▶ Interpret (open) terms as kernels
- ▶ Interpret closed terms as measures
- ▶ Inference **normalizes** measures posterior  $\propto$  **likelihood**  $\times$  **prior**

[Kozen, “Semantics of probabilistic programs”, J Comp Syst Sci, 1981]

# Overview

- ▶ Interpret types as measurable spaces e.g.  $\llbracket \text{real} \rrbracket = \mathbb{R}$
- ▶ Interpret (open) terms as kernels
- ▶ Interpret closed terms as measures
- ▶ Inference **normalizes** measures posterior  $\propto$  **likelihood**  $\times$  **prior**

But:

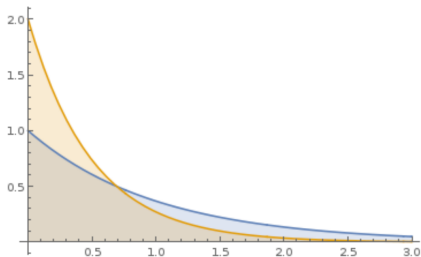
- ▶ Commutativity? Fubini not true for all kernels
- ▶ Higher order functions?  $\mathbb{R} \rightarrow \mathbb{R}$  not a measurable space
- ▶ Extensionality?
- ▶ Recursion?

[Kozen, “Semantics of probabilistic programs”, J Comp Syst Sci, 1981]

[Aumann, “Borel structures for function spaces”, Ill J Math, 1961]

## Example

1. Toss a fair coin to get outcome  $x$
2. Set up exponential decay with rate  $r$  depending on  $x$
3. Observe immediate decay
4. What is the outcome  $x$ ?



## Example

1. Toss a fair coin to get outcome  $x$
2. Set up exponential decay with rate  $r$  depending on  $x$
3. Observe immediate decay
4. What is the outcome  $x$ ?

```
let x = sample(bern(0.5)) in
let r = if x then 2.0 else 1.0
observe(0.0 from exp(r));
return x
```

## Example

1. Toss a fair coin to get outcome  $x$
2. Set up exponential decay with rate  $r$  depending on  $x$
3. Observe immediate decay
4. What is the outcome  $x$ ?

```
let x = sample(bern(0.5)) in
let r = if x then 2.0 else 1.0
observe(0.0 from exp(r));
return x
```

two traces:

	0.5	0.5
x=true		x=false

## Example

1. Toss a fair coin to get outcome  $x$
2. Set up exponential decay with rate  $r$  depending on  $x$
3. Observe immediate decay
4. What is the outcome  $x$ ?

two traces:

```
let x = sample(bern(0.5)) in
let r = if x then 2.0 else 1.0
observe(0.0 from exp(r));
return x
```

```
0.5          0.5
x=true       x=false
r=2.0
score 2
return true
```



## Example

1. Toss a fair coin to get outcome  $x$
2. Set up exponential decay with rate  $r$  depending on  $x$
3. Observe immediate decay
4. What is the outcome  $x$ ?

```
let x = sample(bern(0.5)) in
let r = if x then 2.0 else 1.0
observe(0.0 from exp(r));
return x
```

two traces:

0.5	0.5
x=true	x=false
r=2.0	r=1.0
score 2	score 1
return true	return false

## Example

1. Toss a fair coin to get outcome  $x$
2. Set up exponential decay with rate  $r$  depending on  $x$
3. Observe immediate decay
4. What is the outcome  $x$ ?

```
let x = sample(bern(0.5)) in
let r = if x then 2.0 else 1.0
observe(0.0 from exp(r));
return x
```

two traces:

0.5	0.5
x=true	x=false
r=2.0	r=1.0
score 2	score 1
return true	return false

posterior  $\propto$  likelihood  $\times$  prior

$2 \times 0.5$ : true  
 $1 \times 0.5$ : false

## Example

1. Toss a fair coin to get outcome  $x$
2. Set up exponential decay with rate  $r$  depending on  $x$
3. Observe immediate decay
4. What is the outcome  $x$ ?

$$P(\text{true}) = 1, P(\text{false}) = 0.5$$

two traces:

```
let x = sample(bern(0.5)) in
let r = if x then 2.0 else 1.0
observe(0.0 from exp(r));
return x
```

0.5	0.5
x=true	x=false
r=2.0	r=1.0
score 2	score 1
return true	return false

posterior  $\propto$  likelihood  $\times$  prior

$2 \times 0.5$ : true  
 $1 \times 0.5$ : false

## Example

1. Toss a fair coin to get outcome  $x$
2. Set up exponential decay with rate  $r$  depending on  $x$
3. Observe immediate decay model evidence (score): 1.5
4. What is the outcome  $x$ ?  $P(\text{true}) = 66\%$ ,  $P(\text{false}) = 33\%$

two traces:

	0.5	0.5
let $x = \text{sample}(\text{bern}(0.5))$ in	$x = \text{true}$	$x = \text{false}$
let $r = \text{if } x \text{ then } 2.0 \text{ else } 1.0$	$r = 2.0$	$r = 1.0$
$\text{observe}(0.0 \text{ from } \text{exp}(r));$	score 2	score 1
return $x$	return true	return false

posterior  $\propto$  likelihood  $\times$  prior

$2 \times 0.5$ : true  
 $1 \times 0.5$ : false

## Example

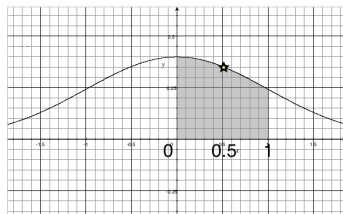
1. **Toss** a fair coin to get outcome  $x$
2. Set up exponential decay with rate  $r$  depending on  $x$
3. **Observe** immediate decay model evidence (score): 1.5
4. **What** is the outcome  $x$ ?  $P(\text{true}) = 66\%$ ,  $P(\text{false}) = 33\%$

Programs may also sample continuous distributions  
so have to deal with uncountable number of traces:

```
let y = sample(gauss(7,2))
```

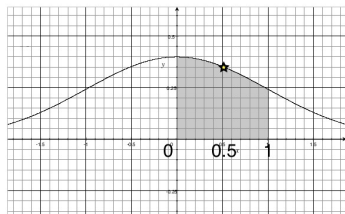
# Measure theory

Impossible to sample 0.5 from standard normal distribution  
But sample in interval  $(0, 1)$  with probability around 0.34



# Measure theory

Impossible to sample 0.5 from standard normal distribution  
But sample in interval  $(0, 1)$  with probability around 0.34



A **measurable space** is a set  $X$  with a family  $\Sigma_X$  of subsets that is closed under countable unions and complements

A **(probability) measure** on  $X$  is a function  $p: \Sigma_X \rightarrow [0, \infty]$  that satisfies  $p(\sum U_n) = \sum p(U_n)$  (and has  $p(X) = 1$ )

# First order language

► Types:  $\mathbb{A}, \mathbb{B} ::= \mathbb{R} \mid \mathbf{P}(\mathbb{A}) \mid \mathbf{1} \mid \mathbb{A} \times \mathbb{B} \mid \sum_{i \in I} \mathbb{A}_i$

real numbers

finite products

distributions over  $\mathbb{A}$

countable sums

`bool := 1 + 1`

`nat :=  $\sum_{i \in \mathbb{N}} 1$`



# First order language

► **Types:**  $\mathbb{A}, \mathbb{B} ::= \mathbb{R} \mid \mathbf{P}(\mathbb{A}) \mid \mathbf{1} \mid \mathbb{A} \times \mathbb{B} \mid \sum_{i \in I} \mathbb{A}_i$

► **Deterministic terms** may not sample:

► variables

$x, y, z$

► constructors for sums and products

case,  $\text{in}_i$ , if, false, true

► measurable functions

bern, exp, gauss, dirac

$\vdash_d 42.0 : \mathbb{R}$

$\vdash_d \text{gauss}(2.0, 7.0) : \mathbf{P}(\mathbb{R})$

$x : \mathbb{R}, y : \mathbb{R} \vdash_d x + y : \mathbb{R}$

$x : \mathbb{R}, y : \mathbb{R} \vdash_d x < y : \text{bool}$

# First order language

► **Types:**  $\mathbb{A}, \mathbb{B} ::= \mathbb{R} \mid \mathbb{P}(\mathbb{A}) \mid \mathbf{1} \mid \mathbb{A} \times \mathbb{B} \mid \sum_{i \in I} \mathbb{A}_i$

► **Deterministic terms** may not sample:

- variables x, y, z
- constructors for sums and products case, in<sub>i</sub>, if, false, true
- measurable functions bern, exp, gauss, dirac

► **Probabilistic terms** may sample:

- sequencing return, let
- constraints score
- priors sample

$$\frac{\Gamma \vdash_d t : \mathbb{A}}{\Gamma \vdash_p \text{return}(t) : \mathbb{A}}$$

$$\frac{\Gamma \vdash_p t : \mathbb{A} \quad x : \mathbb{A} \vdash_p u : \mathbb{B}}{\Gamma \vdash_p \text{let } x = t \text{ in } u : \mathbb{B}}$$

$$\frac{\Gamma \vdash_d t : \mathbb{R}}{\Gamma \vdash_p \text{score}(t) : \mathbf{1}}$$

$$\frac{\Gamma \vdash_d t : \mathbb{P}(\mathbb{A})}{\Gamma \vdash_p \text{sample}(t) : \mathbb{A}}$$

# First order language

► **Types:**  $\mathbb{A}, \mathbb{B} ::= \mathbb{R} \mid \mathbb{P}(\mathbb{A}) \mid \mathbf{1} \mid \mathbb{A} \times \mathbb{B} \mid \sum_{i \in I} \mathbb{A}_i$

► **Deterministic terms** may not sample:

- variables x, y, z
- constructors for sums and products case, in<sub>i</sub>, if, false, true
- measurable functions bern, exp, gauss, dirac
- inference norm

► **Probabilistic terms** may sample:

- sequencing return, let
- constraints score
- priors sample

$$\frac{\Gamma \vdash_d t : \mathbb{A}}{\Gamma \vdash_p \text{return}(t) : \mathbb{A}}$$

$$\frac{\Gamma \vdash_p t : \mathbb{A} \quad x : \mathbb{A} \vdash_p u : \mathbb{B}}{\Gamma \vdash_p \text{let } x = t \text{ in } u : \mathbb{B}}$$

$$\frac{\Gamma \vdash_d t : \mathbb{R}}{\Gamma \vdash_p \text{score}(t) : \mathbf{1}}$$

$$\frac{\Gamma \vdash_d t : \mathbb{P}(\mathbb{A})}{\Gamma \vdash_p \text{sample}(t) : \mathbb{A}}$$

# First order semantics

## Interpret

- ▶ **type**  $\mathbb{A}$  as measurable space  $[[\mathbb{A}]]$
- ▶ **deterministic term**  $\Gamma \vdash_d t : \mathbb{A}$  as measurable function  $[[\Gamma]] \rightarrow [[\mathbb{A}]]$
- ▶ **probabilistic term**  $\Gamma \vdash_p t : \mathbb{A}$  as kernel  $[[t]] : [[\Gamma]] \times \Sigma_{[[\mathbb{A}]]} \rightarrow [0, \infty]$   
fixing first argument: measure, fixing second argument: measurable

# First order semantics

## Interpret

- ▶ **type**  $\mathbb{A}$  as measurable space  $[[\mathbb{A}]]$
- ▶ **deterministic term**  $\Gamma \vdash_d t : \mathbb{A}$  as measurable function  $[[\Gamma]] \rightarrow [[\mathbb{A}]]$
- ▶ **probabilistic term**  $\Gamma \vdash_p t : \mathbb{A}$  as kernel  $[[t]] : [[\Gamma]] \times \Sigma_{[[\mathbb{A}]]} \rightarrow [0, \infty]$   
fixing first argument: measure, fixing second argument: measurable

$$\frac{\Gamma \vdash_d t : \mathbb{R}}{\Gamma \vdash_p \text{score}(t) : 1}$$

$$[[\text{score}(t)]](\gamma, *) = [[t]](\gamma)$$

# First order semantics

## Interpret

- ▶ **type**  $\mathbb{A}$  as measurable space  $\llbracket \mathbb{A} \rrbracket$
- ▶ **deterministic term**  $\Gamma \vdash_d t : \mathbb{A}$  as measurable function  $\llbracket \Gamma \rrbracket \rightarrow \llbracket \mathbb{A} \rrbracket$
- ▶ **probabilistic term**  $\Gamma \vdash_p t : \mathbb{A}$  as kernel  $\llbracket t \rrbracket : \llbracket \Gamma \rrbracket \times \Sigma_{\llbracket \mathbb{A} \rrbracket} \rightarrow [0, \infty]$   
fixing first argument: measure, fixing second argument: measurable

$$\frac{\Gamma \vdash_d t : \mathbb{R}}{\Gamma \vdash_p \text{score}(t) : 1}$$

$$\llbracket \text{score}(t) \rrbracket(\gamma, *) = \llbracket t \rrbracket(\gamma)$$

$$\frac{\Gamma \vdash_d t : \text{P}(\mathbb{A})}{\Gamma \vdash_p \text{sample}(t) : \mathbb{A}}$$

$$\llbracket \text{sample}(t) \rrbracket(\gamma, U) = (\llbracket t \rrbracket(\gamma))(U)$$

# First order semantics

## Interpret

- ▶ **type**  $\mathbb{A}$  as measurable space  $\llbracket \mathbb{A} \rrbracket$
- ▶ **deterministic term**  $\Gamma \vdash_d t : \mathbb{A}$  as measurable function  $\llbracket \Gamma \rrbracket \rightarrow \llbracket \mathbb{A} \rrbracket$
- ▶ **probabilistic term**  $\Gamma \vdash_p t : \mathbb{A}$  as kernel  $\llbracket t \rrbracket : \llbracket \Gamma \rrbracket \times \Sigma_{\llbracket \mathbb{A} \rrbracket} \rightarrow [0, \infty]$   
fixing first argument: measure, fixing second argument: measurable

$$\frac{\Gamma \vdash_d t : \mathbb{R}}{\Gamma \vdash_p \text{score}(t) : 1}$$

$$\llbracket \text{score}(t) \rrbracket(\gamma, *) = \llbracket t \rrbracket(\gamma)$$

$$\frac{\Gamma \vdash_d t : \mathbb{P}(\mathbb{A})}{\Gamma \vdash_p \text{sample}(t) : \mathbb{A}}$$

$$\llbracket \text{sample}(t) \rrbracket(\gamma, U) = (\llbracket t \rrbracket(\gamma))(U)$$

$$\frac{\Gamma \vdash_p t : \mathbb{A} \quad x : \mathbb{A} \vdash_p u : \mathbb{B}}{\Gamma \vdash_p \text{let } x = t \text{ in } u : \mathbb{B}}$$

$$\begin{aligned} & \llbracket \text{let } x = t \text{ in } u \rrbracket(\gamma, U) \\ &= \int_{\llbracket \mathbb{A} \rrbracket} \llbracket u \rrbracket(\gamma, x, U) \llbracket t \rrbracket(\gamma, dx) \end{aligned}$$

# First order semantics

## Interpret

- ▶ **type**  $\mathbb{A}$  as measurable space  $\llbracket \mathbb{A} \rrbracket$
- ▶ **deterministic term**  $\Gamma \vdash_d t : \mathbb{A}$  as measurable function  $\llbracket \Gamma \rrbracket \rightarrow \llbracket \mathbb{A} \rrbracket$
- ▶ **probabilistic term**  $\Gamma \vdash_p t : \mathbb{A}$  as kernel  $\llbracket t \rrbracket : \llbracket \Gamma \rrbracket \times \Sigma_{\llbracket \mathbb{A} \rrbracket} \rightarrow [0, \infty]$   
fixing first argument: measure, fixing second argument: measurable

$$\frac{\Gamma \vdash_d t : \mathbb{R}}{\Gamma \vdash_p \text{score}(t) : 1}$$

$$\llbracket \text{score}(t) \rrbracket(\gamma, *) = \llbracket t \rrbracket(\gamma)$$

$$\frac{\Gamma \vdash_d t : \mathbb{P}(\mathbb{A})}{\Gamma \vdash_p \text{sample}(t) : \mathbb{A}}$$

$$\llbracket \text{sample}(t) \rrbracket(\gamma, U) = (\llbracket t \rrbracket(\gamma))(U)$$

$$\frac{\Gamma \vdash_p t : \mathbb{A} \quad x : \mathbb{A} \vdash_p u : \mathbb{B}}{\Gamma \vdash_p \text{let } x = t \text{ in } u : \mathbb{B}}$$

$$\llbracket \text{let } x = t \text{ in } u \rrbracket = \int_{\llbracket \mathbb{A} \rrbracket} \llbracket u \rrbracket d\llbracket t \rrbracket$$



## Example

```
[[ let x = sample(bern(0.5)) in
   let r = if x then 2.0 else 1.0
   observe(0.0 from exp(r));
   return x ]]
```

The meaning of a program returning values in  $X$  is a measure on  $X$

$\emptyset$	has measure	0.0
{true}	has measure	$1.0 = 0.5 \times 2.0$
{false}	has measure	$0.5 = 0.5 \times 1.0$
{true, false}	has measure	1.5

Normalization: **posterior**  $\propto$  **likelihood**  $\times$  **prior**

$$\frac{\Gamma \vdash_p t: \mathbb{A}}{\Gamma \vdash_d \text{norm}(t): \mathbb{R} \times \mathbb{P}(\mathbb{A}) + \underbrace{1 + 1}_{\text{errors}}}$$

model evidence

normalized posterior

errors

Normalization: **posterior**  $\propto$  **likelihood**  $\times$  **prior**

$$\frac{\Gamma \vdash_p t: \mathbb{A}}{\Gamma \vdash_d \text{norm}(t): \mathbb{R} \times \mathbb{P}(\mathbb{A}) + \underbrace{1 + 1}_{\text{errors}}}$$

model evidence

normalized posterior

errors

Interpretation of probabilistic term is kernel  $\llbracket \Gamma \rrbracket \times \Sigma_{\llbracket \mathbb{A} \rrbracket} \rightarrow [0, \infty]$  so  
fixing first argument gives measure

$$\frac{\llbracket t \rrbracket(\gamma, -)}{\llbracket t \rrbracket(\gamma, \llbracket \mathbb{A} \rrbracket)}$$

is normalized probability measure

Normalization: **posterior**  $\propto$  **likelihood**  $\times$  **prior**

$$\frac{\Gamma \vdash_p t: \mathbb{A}}{\Gamma \vdash_d \text{norm}(t): \mathbb{R} \times \mathbb{P}(\mathbb{A}) + \underbrace{1 + 1}_{\text{errors (constant is 0 or } \infty)}}}$$

model evidence
normalized posterior

Interpretation of probabilistic term is kernel  $\llbracket \Gamma \rrbracket \times \Sigma_{\llbracket \mathbb{A} \rrbracket} \rightarrow [0, \infty]$  so  
fixing first argument gives measure

$$\frac{\llbracket t \rrbracket(\gamma, -)}{\llbracket t \rrbracket(\gamma, \llbracket \mathbb{A} \rrbracket)}$$

is normalized probability measure  
normalizing constant is model evidence

Normalization: **posterior**  $\propto$  **likelihood**  $\times$  **prior**

$$\frac{\Gamma \vdash_p t : \mathbb{A}}{\Gamma \vdash_d \text{norm}(t) : \mathbb{R} \times \mathbb{P}(\mathbb{A}) + \underbrace{1 + 1}_{\text{errors (constant is 0 or } \infty)}}}$$

model evidence

normalized posterior

```
[[ let x = sample(bern(0.5)) in
  let r = if x then 2.0 else 1.0
  observe(0.0 from exp(r));
  return x ]]
```

$$= \begin{pmatrix} \text{true} : & 2.0 \times 0.5 \\ \text{false} : & 1.0 \times 0.5 \end{pmatrix}$$

Normalization: **posterior**  $\propto$  **likelihood**  $\times$  **prior**

$$\frac{\Gamma \vdash_p t: \mathbb{A}}{\Gamma \vdash_d \text{norm}(t): \mathbb{R} \times \mathbb{P}(\mathbb{A}) + \underbrace{1 + 1}_{\text{errors (constant is 0 or } \infty)}}}$$

model evidence

normalized posterior

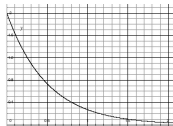
```
norm(  
  let x = sample(bern(0.5)) in  
  let r = if x then 2.0 else 1.0  
  observe(0.0 from exp(r));  
  return x ) = in1(1.5, bern(0.66))
```

## Example: sequential Monte Carlo

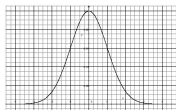
$$\left[ \begin{array}{l} \text{norm}(\text{let } x=t \\ \text{in } u) \end{array} \right] = \left[ \begin{array}{l} \text{norm}(\text{let } (e,d) = \text{norm}(t) \text{ in} \\ \text{score}(e); \text{let } x=\text{sample}(d) \\ \text{in } u) \end{array} \right]$$

## Example: importance sampling

```
[[ sample(exp(2)) ]]
```



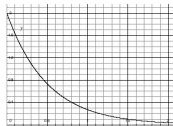
```
= [[ let x = sample(gauss(0,1))  
   score(exp-pdf(2,x) / gauss-pdf(0,1,x));  
   return x ]]
```



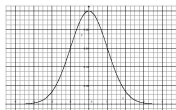


## Example: importance sampling

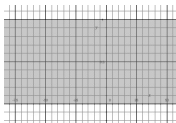
```
[[ sample(exp(2)) ]]
```



```
= [[ let x = sample(gauss(0,1))  
    score(exp-pdf(2,x) / gauss-pdf(0,1,x));  
    return x ]]
```

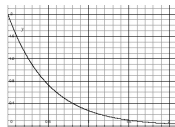


```
= [[ let x = sample(gauss(0,1))  
    score(1 / gauss-pdf(0,1,x));  
    score(exp-pdf(2,x));  
    return x ]]
```

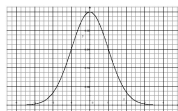


# Example: importance sampling

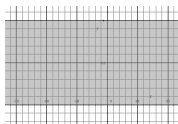
```
norm( sample(exp(2)) )
```



```
= norm(  
  let x = sample(gauss(0,1))  
  score(exp-pdf(2,x) / gauss-pdf(0,1,x));  
  return x )
```



```
≠ norm( norm(  
  let x = sample(gauss(0,1))  
  score(1 / gauss-pdf(0,1,x)); )  
  score(exp-pdf(2,x));  
  return x )
```



Don't **normalize** as you go

## Commutativity

Reordering lines is very useful program transformation

$$\left[ \begin{array}{l} \text{let } x=t \text{ in} \\ \text{let } y=u \text{ in} \\ v \end{array} \right] = \left[ \begin{array}{l} \text{let } y=u \text{ in} \\ \text{let } x=t \text{ in} \\ v \end{array} \right]$$

## Commutativity

Reordering lines is very useful program transformation

$$\left[ \begin{array}{l} \text{let } x=t \text{ in} \\ \text{let } y=u \text{ in} \\ v \end{array} \right] = \left[ \begin{array}{l} \text{let } y=u \text{ in} \\ \text{let } x=t \text{ in} \\ v \end{array} \right]$$

amounts to Fubini's theorem

$$\int_{[A]} \int_{[B]} [v] d[u] d[t] = \int_{[B]} \int_{[A]} [v] d[t] d[u]$$

## Commutativity

Reordering lines is very useful program transformation

$$\left[ \left[ \begin{array}{l} \text{let } x=t \text{ in} \\ \text{let } y=u \text{ in} \end{array} \right] \right]_{\mathbf{v}} = \left[ \left[ \begin{array}{l} \text{let } y=u \text{ in} \\ \text{let } x=t \text{ in} \end{array} \right] \right]_{\mathbf{v}}$$

amounts to Fubini's theorem

$$\int_{\mathbb{A}} \int_{\mathbb{B}} \llbracket v \rrbracket d\llbracket u \rrbracket d\llbracket t \rrbracket = \int_{\mathbb{B}} \int_{\mathbb{A}} \llbracket v \rrbracket d\llbracket t \rrbracket d\llbracket u \rrbracket$$



Not true for arbitrary kernels, only for **s-finite kernels**

kernel is s-finite when countable sum of bounded ones

$k: \llbracket \Gamma \rrbracket \times \Sigma_{\llbracket \mathbb{A} \rrbracket} \rightarrow [0, \infty]$  bounded if  $\exists n \forall \gamma \forall U: k(\gamma, U) < n$

# Commutativity

Reordering lines is very useful program transformation

$$\left[ \left[ \text{let } x=t \text{ in} \right] \left[ \text{let } y=u \text{ in} \right] \right]_{\mathbb{V}} = \left[ \left[ \text{let } y=u \text{ in} \right] \left[ \text{let } x=t \text{ in} \right] \right]_{\mathbb{V}}$$

amounts to Fubini's theorem

$$\int_{\mathbb{A}} \int_{\mathbb{B}} \llbracket v \rrbracket d\llbracket u \rrbracket d\llbracket t \rrbracket = \int_{\mathbb{B}} \int_{\mathbb{A}} \llbracket v \rrbracket d\llbracket t \rrbracket d\llbracket u \rrbracket$$



Not true for arbitrary kernels, only for **s-finite kernels**

kernel is s-finite when countable sum of bounded ones

$k: \llbracket \Gamma \rrbracket \times \Sigma_{\llbracket \mathbb{A} \rrbracket} \rightarrow [0, \infty]$  bounded if  $\exists n \forall \gamma \forall U: k(\gamma, U) < n$

- ▶ kernel  $k$  is s-finite iff it can be built from sub-probability distributions, **score**, and *binding*

$$k \gg= l \text{ is } (\gamma, V) \mapsto \int_{\llbracket \mathbb{A} \rrbracket} l(\gamma, x, V) k(\gamma, dx)$$

- ▶ measurable spaces and s-finite kernels form distributive symmetric monoidal category

# Commutativity

Reordering lines is very useful program transformation

$$\left[ \left[ \begin{array}{l} \text{let } x=t \text{ in} \\ \text{let } y=u \text{ in} \end{array} \right] \right]_{\mathbf{v}} = \left[ \left[ \begin{array}{l} \text{let } y=u \text{ in} \\ \text{let } x=t \text{ in} \end{array} \right] \right]_{\mathbf{v}}$$

amounts to Fubini's theorem

$$\int_{[\mathbf{A}]} \int_{[\mathbf{B}]} \llbracket v \rrbracket d[u] d[t] = \int_{[\mathbf{B}]} \int_{[\mathbf{A}]} \llbracket v \rrbracket d[t] d[u]$$



Not true for arbitrary kernels, only for **s-finite kernels**

kernel is s-finite when countable sum of bounded ones

$k: [\Gamma] \times \Sigma_{[\mathbf{A}]} \rightarrow [0, \infty]$  bounded if  $\exists n \forall \gamma \forall U: k(\gamma, U) < n$

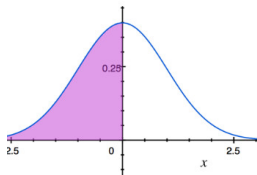
Interpret terms as s-finite kernels

## Example: facts about distributions

```
[[ let x = sample(gauss(0.0,1.0))  
  in return (x<0) ]]
```

 = 

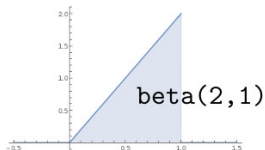
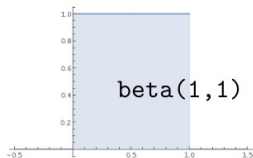
```
[[ sample(bern(0.5)) ]]
```





## Example: conjugate priors

```
[[let x = sample(beta(1,1))  
  in observe(bern(x), true);  
  return x]] = [[observe(bern(0.5), true);  
  let x = sample(beta(2,1))  
  in return x]]
```



# Higher order functions

Allow probabilistic terms as input/output for other terms

```
(define (ibp-stick-breaking-process concentration base-measure)
  (let ((sticks (mem (lambda j (random-beta 1.0 concentration))))
        (atoms (mem (lambda j (base-measure)))))
    (lambda ()
      (let loop ((j 1) (dualstick (sticks 1)))
        (append (if (flip dualstick)                ;; with prob. dualstick
                    (atoms j)                       ;; add feature j
                    '())                             ;; otherwise, next stick
                (loop (+ j 1) (* dualstick (sticks (+ j 1))))))))))
```

[Roy et al, “A stochastic programming perspective on nonparametric Bayes”, ICML 2008]

## Higher order functions

Allow probabilistic terms as input/output for other terms

```
(define (ibp-stick-breaking-process concentration base-measure)
  (let ((sticks (mem (lambda j (random-beta 1.0 concentration))))
        (atoms (mem (lambda j (base-measure))))))
    (lambda ()
      (let loop ((j 1) (dualstick (sticks 1)))
        (append (if (flip dualstick) ;; with prob. dualstick
                    (atoms j) ;; add feature j
                    '()) ;; otherwise, next stick
                (loop (+ j 1) (* dualstick (sticks (+ j 1))))))))))
```

$$\mathbb{A}, \mathbb{B} ::= \mathbb{R} \mid \mathbb{P}(\mathbb{A}) \mid \mathbf{1} \mid \mathbb{A} \times \mathbb{B} \mid \sum_{i \in I} \mathbb{A}_i \mid \mathbb{A} \rightarrow \mathbb{B}$$

[Roy et al, “A stochastic programming perspective on nonparametric Bayes”, ICML 2008]

# Higher order functions

Allow probabilistic terms as input/output for other terms

```
(define (ibp-stick-breaking-process concentration base-measure)
  (let ((sticks (mem (lambda j (random-beta 1.0 concentration))))
        (atoms (mem (lambda j (base-measure))))))
    (lambda ()
      (let loop ((j 1) (dualstick (sticks 1)))
        (append (if (flip dualstick) ;; with prob. dualstick
                    (atoms j) ;; add feature j
                    '()) ;; otherwise, next stick
                (loop (+ j 1) (* dualstick (sticks (+ j 1))))))))))
```

$$\mathbb{A}, \mathbb{B} ::= \mathbb{R} \mid \mathcal{P}(\mathbb{A}) \mid \mathbf{1} \mid \mathbb{A} \times \mathbb{B} \mid \sum_{i \in I} \mathbb{A}_i \mid \mathbb{A} \rightarrow \mathbb{B}$$



$\mathbb{R} \rightarrow \mathbb{R}$  is not a measurable space

[Roy et al, “A stochastic programming perspective on nonparametric Bayes”, ICML 2008]

[Aumann, “Borel structures for function spaces”, Ill J Math, 1961]

# Higher order functions

Allow probabilistic terms as input/output for other terms

```
(define (ibp-stick-breaking-process concentration base-measure)
  (let ((sticks (mem (lambda j (random-beta 1.0 concentration))))
        (atoms (mem (lambda j (base-measure))))))
    (lambda ()
      (let loop ((j 1) (dualstick (sticks 1)))
        (append (if (flip dualstick) ;; with prob. dualstick
                    (atoms j) ;; add feature j
                    '()) ;; otherwise, next stick
                (loop (+ j 1) (* dualstick (sticks (+ j 1))))))))))
```

$$\mathbb{A}, \mathbb{B} ::= \mathbb{R} \mid \mathcal{P}(\mathbb{A}) \mid \mathbf{1} \mid \mathbb{A} \times \mathbb{B} \mid \sum_{i \in I} \mathbb{A}_i \mid \mathbb{A} \rightarrow \mathbb{B}$$



$\mathbb{R} \rightarrow \mathbb{R}$  is not a measurable space

Easy to handle operationally.

What to do denotationally?

[Roy et al, “A stochastic programming perspective on nonparametric Bayes”, ICML 2008]

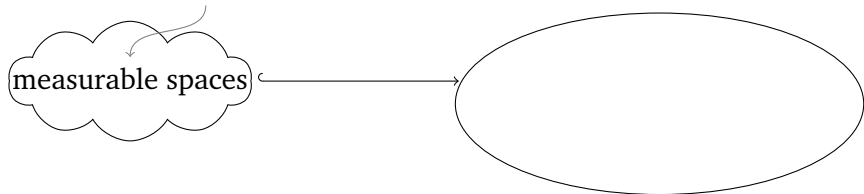
[Aumann, “Borel structures for function spaces”, Ill J Math, 1961]

[Borgström et al, “Measure transformer semantics for Bayesian machine learning”, ESOP2011]

# Higher order semantics

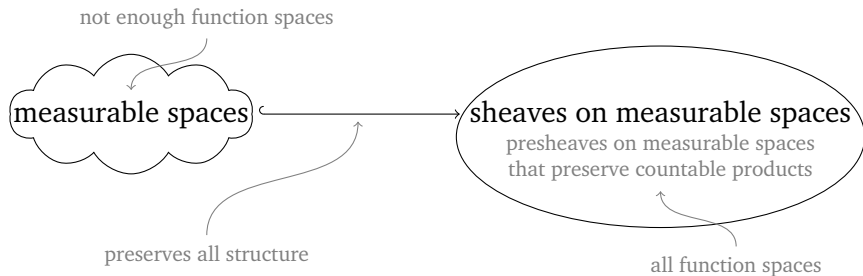
Use category theory to extend measure theory

not enough function spaces



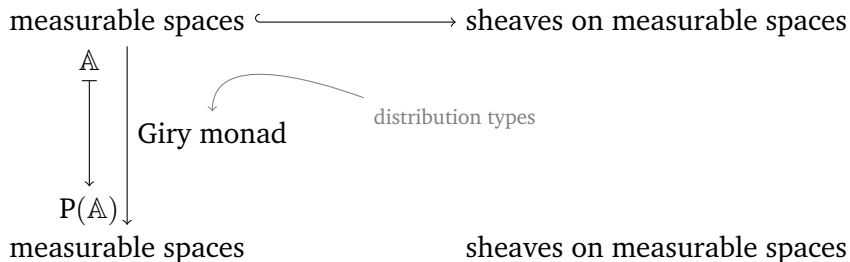
# Higher order semantics

Use category theory to extend measure theory



# Higher order semantics

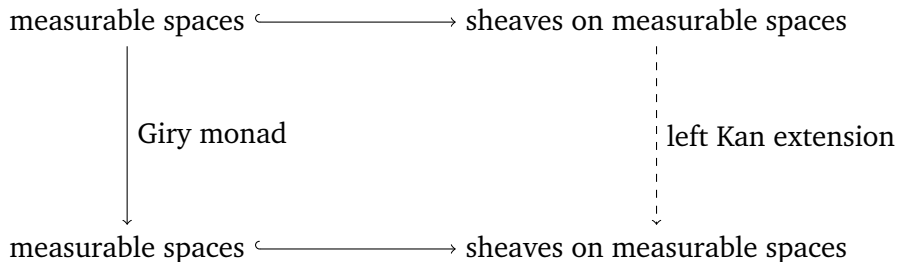
Use category theory to extend measure theory





# Higher order semantics

Use category theory to extend measure theory



# Higher order semantics

Use category theory to extend measure theory

measurable spaces  $\hookrightarrow$  sheaves on measurable spaces

- ▶  $\llbracket 1 \rightarrow (\mathbb{R} \rightarrow \mathbb{R}) \rrbracket$  consists of **random functions** measurable  $\Omega \times \mathbb{R} \rightarrow \mathbb{R}$
- ▶ All definable functions  $\mathbb{R} \rightarrow \mathbb{R}$  are measurable “Church-Turing”
- ▶ Denotational and operational semantics match soundness & adequacy

# Extensionality



Not **extensional**:  $1 \xrightarrow{p} \mathbb{A} \begin{matrix} \xrightarrow{f} \\ \xrightarrow{g} \end{matrix} \mathbb{B}$  for all  $p \not\Rightarrow f = g$

Solution: restrict to subcategory that is extensional

# Extensionality



Not **extensional**:  $1 \xrightarrow{p} \mathbb{A} \begin{matrix} \xrightarrow{f} \\ \xrightarrow{g} \end{matrix} \mathbb{B}$  for all  $p \not\Rightarrow f = g$

Solution: restrict to subcategory that is extensional

A **quasi-measurable space** is a set  $X$  with  $M_X \subseteq [\mathbb{R} \rightarrow X]$  satisfying

- ▶ if  $f: \mathbb{R} \rightarrow \mathbb{R}$  is measurable and  $g \in M$ , then  $gf \in M$
- ▶ if  $f: \mathbb{R} \rightarrow X$  is constant, then  $f \in M$
- ▶ if  $f: \mathbb{R} \rightarrow \mathbb{N}$  is measurable and  $g_n \in M$ , then  $[g_n]f \in M \quad t \mapsto g_{f(t)}(t)$

*morphisms* are functions  $f: X \rightarrow Y$  with  $g \in M_X \Rightarrow fg \in M_Y$

# Extensionality



Not **extensional**:  $1 \xrightarrow{p} \mathbb{A} \xrightleftharpoons[g]{f} \mathbb{B}$  for all  $p \not\Rightarrow f = g$

Solution: restrict to subcategory that is extensional

A **quasi-measurable space** is a set  $X$  with  $M_X \subseteq [\mathbb{R} \rightarrow X]$  satisfying

- ▶ if  $f: \mathbb{R} \rightarrow \mathbb{R}$  is measurable and  $g \in M$ , then  $gf \in M$
- ▶ if  $f: \mathbb{R} \rightarrow X$  is constant, then  $f \in M$
- ▶ if  $f: \mathbb{R} \rightarrow \mathbb{N}$  is measurable and  $g_n \in M$ , then  $[g_n]f \in M$   $t \mapsto g_{f(t)}(t)$

*morphisms* are functions  $f: X \rightarrow Y$  with  $g \in M_X \Rightarrow fg \in M_Y$

*Example*:  $X$  measurable space,  $M_X$  measurable functions  $\mathbb{R} \rightarrow X$   
morphism  $X \rightarrow Y$  is measurable function

# Extensionality



Not **extensional**:  $1 \xrightarrow{p} \mathbb{A} \xrightleftharpoons[g]{f} \mathbb{B}$  for all  $p \not\Rightarrow f = g$

Solution: restrict to subcategory that is extensional

A **quasi-measurable space** is a set  $X$  with  $M_X \subseteq [\mathbb{R} \rightarrow X]$  satisfying

- ▶ if  $f: \mathbb{R} \rightarrow \mathbb{R}$  is measurable and  $g \in M$ , then  $gf \in M$
- ▶ if  $f: \mathbb{R} \rightarrow X$  is constant, then  $f \in M$
- ▶ if  $f: \mathbb{R} \rightarrow \mathbb{N}$  is measurable and  $g_n \in M$ , then  $[g_n]f \in M \quad t \mapsto g_{f(t)}(t)$

*morphisms* are functions  $f: X \rightarrow Y$  with  $g \in M_X \Rightarrow fg \in M_Y$

*Example*:  $X$  measurable space,  $M_X$  measurable functions  $\mathbb{R} \rightarrow X$   
morphism  $X \rightarrow Y$  is measurable function

*Theorem*: this gives cartesian closed category with countable sums

*Corollary*: if term  $t$  has first order type, then  $\llbracket t \rrbracket$  is measurable

even if  $t$  involves higher order functions

# Extensionality



Not **extensional**:  $1 \xrightarrow{p} \mathbb{A} \xrightleftharpoons[g]{f} \mathbb{B}$  for all  $p \not\Rightarrow f = g$

Solution: restrict to subcategory that is extensional

A **quasi-measurable space** is a set  $X$  with  $M_X \subseteq [\mathbb{R} \rightarrow X]$  satisfying

- ▶ if  $f: \mathbb{R} \rightarrow \mathbb{R}$  is measurable and  $g \in M$ , then  $gf \in M$
- ▶ if  $f: \mathbb{R} \rightarrow X$  is constant, then  $f \in M$
- ▶ if  $f: \mathbb{R} \rightarrow \mathbb{N}$  is measurable and  $g_n \in M$ , then  $[g_n]f \in M \quad t \mapsto g_{f(t)}(t)$

A **measure** on  $(X, M_X)$  is a measure  $\mu$  on  $\mathbb{R}$  with a function  $f \in M$

*Proposition:* measures on  $[X \rightarrow Y]$  are random functions

measurable map  $\mathbb{R} \times X \rightarrow Y$  modulo measure on  $\mathbb{R}$

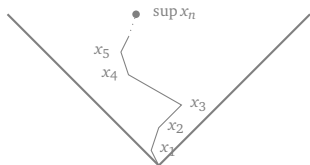
# Recursion



No recursion / least fixed points

Idea: restrict to presheaves over domains

An  $\omega$ -complete partial order has suprema of increasing sequences  
morphisms preserve suprema of increasing sequences and infima



A quasi-measurable space is **ordered** when  $X$  is an  $\omega$ cpo and  $M$  is closed under pointwise increasing suprema

*Example:* Any  $\omega$ cpo, e.g.  $[0,1]$  take  $M$  all measurable functions  $\mathbb{R} \rightarrow X$   
where  $X$  has the Borel  $\sigma$ -algebra on the Lawson topology

*Theorem:* this gives a cartesian closed category with countable sums



## Example: von Neumann's trick

```
[[let g = bern(0.66) in
  letrec f() = (let x = sample(g)
                let y = sample(g)
                if x=y then f()
                else return x)
  in f()]]  $\stackrel{?}{=} \llbracket \text{sample}(\text{bern}(0.5)) \rrbracket$ 
```

# Conclusion

Foundational semantics for probabilistic programming:

- ▶ continuous distributions
- ▶ soft constraints
- ▶ commutativity
- ▶ higher order functions
- ▶ recursion

can verify/suggest program transformations. Approximations?