

Semantics of probabilistic programming

Fredrik Dahlqvist

Work in progress with Dexter Kozen and help from Vincent Danos, Ilias Garnier and Alexandra Silva

London, 15 November 2018

Warning



What is probabilistic programming?

```
(defquery example
  (let [x (sample (normal 0 1))]
    (observe (normal x 1) 0.5)
    (> x 1)))
```

Denotational and Operational semantics

Operational semantics

Denotational semantics

Denotational and Operational semantics

Operational semantics

- Step-by-step execution of the program

Denotational semantics

Denotational and Operational semantics

Operational semantics

- Step-by-step execution of the program
- Sampling *actually* occurs

Denotational semantics

Denotational and Operational semantics

Operational semantics

- Step-by-step execution of the program
- Sampling *actually* occurs
- Statistical properties emerge after many 'runs'

Denotational semantics

Denotational and Operational semantics

Operational semantics

- Step-by-step execution of the program
- Sampling *actually* occurs
- Statistical properties emerge after many 'runs'

Denotational semantics

- Mathematical meaning of the program

Denotational and Operational semantics

Operational semantics

- Step-by-step execution of the program
- Sampling *actually* occurs
- Statistical properties emerge after many 'runs'

Denotational semantics

- Mathematical meaning of the program
- Sampling = distribution

Denotational and Operational semantics

Operational semantics

- Step-by-step execution of the program
- Sampling *actually* occurs
- Statistical properties emerge after many 'runs'

Denotational semantics

- Mathematical meaning of the program
- Sampling = distribution
- Statistical properties immediately available

Denotational and Operational semantics

Operational semantics

- Step-by-step execution of the program
- Sampling *actually* occurs
- Statistical properties emerge after many 'runs'

Denotational semantics

- Mathematical meaning of the program
- Sampling = distribution
- Statistical properties immediately available



Probabilistic Adequacy

Denotational Semantics

Denotational semantics [Kozen '81]: the main ideas

- Probabilistic programs transform probabilities.

Denotational semantics [Kozen '81]: the main ideas

- Probabilistic programs transform probabilities.
- Probabilities are measures. Measures form vector spaces

$$\mathcal{M}(X) = \text{all 'finite' measures}$$

Denotational semantics [Kozen '81]: the main ideas

- Probabilistic programs transform probabilities.
- Probabilities are measures. Measures form vector spaces

$$\mathcal{M}(X) = \text{all 'finite' measures}$$

- Probabilistic programs transform measures *linearly*

$$\mathcal{M}(X) \xrightarrow{P} \mathcal{M}(X)$$

Denotational semantics [Kozen '81]: the main ideas

- Probabilistic programs transform probabilities.
- Probabilities are measures. Measures form vector spaces

$$\mathcal{M}(X) = \text{all 'finite' measures}$$

- Probabilistic programs transform measures *linearly*

$$\mathcal{M}(X) \xrightarrow{P} \mathcal{M}(X)$$

- Measures are normed (Banach space) \Rightarrow can study convergence.

Denotational semantics [Kozen '81]: the main ideas

- Probabilistic programs transform probabilities.
- Probabilities are measures. Measures form vector spaces

$$\mathcal{M}(X) = \text{all 'finite' measures}$$

- Probabilistic programs transform measures *linearly*

$$\mathcal{M}(X) \xrightarrow{P} \mathcal{M}(X)$$

- Measures are normed (Banach space) \Rightarrow can study convergence.
- Measures are (partially) ordered \Rightarrow can study fixpoints (while loops)

Denotational semantics [Kozen '81]: the main ideas

- Probabilistic programs transform probabilities.
- Probabilities are measures. Measures form vector spaces

$$\mathcal{M}(X) = \text{all 'finite' measures}$$

- Probabilistic programs transform measures *linearly*

$$\mathcal{M}(X) \xrightarrow{P} \mathcal{M}(X)$$

- Measures are normed (Banach space) \Rightarrow can study convergence.
- Measures are (partially) ordered \Rightarrow can study fixpoints (while loops)
- Measures belong to a *monoidal closed category* \Rightarrow higher-order.

Assignments

```
{  
  x := 0.5  
}
```

Assignments

```
{  
  x := 0.5  
}
```

$$\frac{\vdash 0.5 : \text{real}}{x : \text{real} \vdash x := 0.5 : \text{real}}$$

Assignments

```
{
  x := 0.5
}
```

$$\frac{\vdash 0.5 : \text{real}}{x : \text{real} \vdash x := 0.5 : \text{real}}$$

$$\llbracket \text{ } \rrbracket = \mathbb{R} \xrightarrow{\llbracket 0.5 \rrbracket = 1 \mapsto \delta_{0.5}} \llbracket \text{real} \rrbracket = \mathcal{MR}$$

$$\llbracket \text{real} \rrbracket = \mathcal{MR} \xrightarrow{\mu \mapsto \mu(\mathbb{R}) \delta_{0.5}} \llbracket \text{real} \rrbracket = \mathcal{MR}$$

Some initial observations

Some initial observations

- Ground types T (e.g. `real`) are interpreted by $\mathcal{M}[[T]]$

Some initial observations

- Ground types T (e.g. `real`) are interpreted by $\mathcal{M}[[T]]$
- Very similar to 'flat domains' in Scott's denotational semantics

Some initial observations

- Ground types T (e.g. `real`) are interpreted by $\mathcal{M}[[T]]$
- Very similar to 'flat domains' in Scott's denotational semantics
- \mathcal{M} turns any partial map into a total linear operator

Some initial observations

- Ground types T (e.g. `real`) are interpreted by $\mathcal{M}[[T]]$
- Very similar to ‘flat domains’ in Scott’s denotational semantics
- \mathcal{M} turns any partial map into a total linear operator
- Again, similar to monotone maps between ‘flat domains’ and

Some initial observations

- Ground types T (e.g. `real`) are interpreted by $\mathcal{M}[[T]]$
- Very similar to ‘flat domains’ in Scott’s denotational semantics
- \mathcal{M} turns any partial map into a total linear operator
- Again, similar to monotone maps between ‘flat domains’ and
- Rule of thumb: denotational semantics will be one \mathcal{M} -level up

Sampling

```
{  
  x := sample(normal(0,1))  
}
```

Sampling

```
{  
  x := sample(normal(0,1))  
}
```

$$\frac{\frac{\vdash \text{normal}(0,1) : \text{M real}}{\vdash \text{sample}(\text{normal}(0,1)) : \text{real}}}{x : \text{real} \vdash x := \text{sample}(\text{normal}(0,1)) : \text{real}}$$

Sampling

```
{
  x := sample(normal(0, 1))
}
```

$$\frac{\frac{\vdash \text{normal}(0, 1) : \text{M real}}{\vdash \text{sample}(\text{normal}(0, 1)) : \text{real}}}{x : \text{real} \vdash x := \text{sample}(\text{normal}(0, 1)) : \text{real}}$$

$$\begin{array}{ccc}
 \llbracket \] = \mathbb{R} & \xrightarrow{1 \mapsto \delta_{\mathcal{N}(0,1)}} & \llbracket \text{M real} \rrbracket = \mathcal{M}^2\mathbb{R} \\
 & \searrow & \downarrow \mu \mapsto \int_{B^+(\mathcal{M}\mathbb{R})} x \, d\mu(x) \\
 \llbracket \text{sample}(\text{normal}(0,1)) \rrbracket & & \llbracket \text{real} \rrbracket = \mathcal{M}\mathbb{R}
 \end{array}$$

Some observation about sampling

Some observation about sampling

- The transformation $\mathcal{M}V \rightarrow V$ is completely generic: it is given by the *Bochner integral* $\mu \mapsto \int_{B^+(V)} x \, d\mu(x)$

Some observation about sampling

- The transformation $\mathcal{M}V \rightarrow V$ is completely generic: it is given by the *Bochner integral* $\mu \mapsto \int_{B^+(V)} x \, d\mu(x)$
- Denotationally $\llbracket \text{sample}(\text{normal}(0, 1)) \rrbracket$ is proportional to $\mathcal{N}(0, 1)$ as expected.

Some observation about sampling

- The transformation $\mathcal{M}V \rightarrow V$ is completely generic: it is given by the *Bochner integral* $\mu \mapsto \int_{B^+(V)} x \, d\mu(x)$
- Denotationally $\llbracket \text{sample}(\text{normal}(0, 1)) \rrbracket$ is proportional to $\mathcal{N}(0, 1)$ as expected.
- Bochner integrals are an essential part of the mathematical universe allowing higher-order functions.

Higher-order functions

```
{  
  fn x. normal(x,y)  
}
```

Higher-order functions

```
{  
  fn x. normal(x,y)  
}
```

$$\frac{x : \text{real}, y : \text{real} \vdash \text{normal}(x, y) : \text{M real}}{y : \text{real} \vdash \text{fn } x. \text{normal}(x, y) : \text{real} \rightarrow \text{M real}}$$

Higher-order functions

```
{
  fn x. normal(x,y)
}
```

$$\frac{x : \text{real}, y : \text{real} \vdash \text{normal}(x, y) : \text{M real}}{y : \text{real} \vdash \text{fn } x. \text{normal}(x, y) : \text{real} \rightarrow \text{M real}}$$

$$\llbracket \text{real} \rrbracket \otimes \llbracket \text{real} \rrbracket = \mathcal{M}\mathbb{R} \otimes \mathcal{M}\mathbb{R} \longrightarrow \llbracket \text{M real} \rrbracket = \mathcal{M}^2\mathbb{R}$$

$$\mathcal{M}\mathbb{R} \xrightarrow{\llbracket \text{fn } x. \text{normal}(x, y) \rrbracket} \mathcal{L}_r(\mathcal{M}\mathbb{R}, \mathcal{M}^2\mathbb{R})$$

A computer science perspective on tensor products

A computer science perspective on tensor products

- Given a map f in two arguments U, V into W , we want to *curry*

$$U \rightarrow [V, W] \quad V \rightarrow [U, W]$$

A computer science perspective on tensor products

- Given a map f in two arguments U, V into W , we want to *curry*

$$U \rightarrow [V, W] \quad V \rightarrow [U, W]$$

- These morphisms must be linear transformations

A computer science perspective on tensor products

- Given a map f in two arguments U, V into W , we want to *curry*

$$U \rightarrow [V, W] \quad V \rightarrow [U, W]$$

- These morphisms must be linear transformations
- This means our map f is linear in U and V *separately* (bilinear)

$$f(\lambda u + \lambda' u', v) = \lambda f(u, v) + \lambda' f(u', v)$$

A computer science perspective on tensor products

- Given a map f in two arguments U, V into W , we want to *curry*

$$U \rightarrow [V, W] \quad V \rightarrow [U, W]$$

- These morphisms must be linear transformations
- This means our map f is linear in U and V *separately* (bilinear)

$$f(\lambda u + \lambda' u', v) = \lambda f(u, v) + \lambda' f(u', v)$$

- This does not mean that it is *jointly* linear:

$$f(\lambda(u, v)) = f(\lambda u, \lambda v) = \lambda f(u, \lambda v) = \lambda^2 f(u, v)$$

A computer science perspective on tensor products

- Given a map f in two arguments U, V into W , we want to *curry*

$$U \rightarrow [V, W] \quad V \rightarrow [U, W]$$

- These morphisms must be linear transformations
- This means our map f is linear in U and V *separately* (bilinear)

$$f(\lambda u + \lambda' u', v) = \lambda f(u, v) + \lambda' f(u', v)$$

- This does not mean that it is *jointly* linear:

$$f(\lambda(u, v)) = f(\lambda u, \lambda v) = \lambda f(u, \lambda v) = \lambda^2 f(u, v)$$

- So $f : U \times V \rightarrow W$ is *not* linear!

A computer science perspective on tensor products

- Given a map f in two arguments U, V into W , we want to *curry*

$$U \rightarrow [V, W] \quad V \rightarrow [U, W]$$

- These morphisms must be linear transformations
- This means our map f is linear in U and V *separately* (bilinear)

$$f(\lambda u + \lambda' u', v) = \lambda f(u, v) + \lambda' f(u', v)$$

- This does not mean that it is *jointly* linear:

$$f(\lambda(u, v)) = f(\lambda u, \lambda v) = \lambda f(u, \lambda v) = \lambda^2 f(u, v)$$

- So $f : U \times V \rightarrow W$ is *not* linear!
- BUT: $\hat{f} : U \otimes V \rightarrow W$ is.

Denotational semantics: conclusion

- Typed language accommodating many important classical and probabilistic constructs

Denotational semantics: conclusion

- Typed language accommodating many important classical and probabilistic constructs
- Very powerful semantics in terms of ordered Banach space

Denotational semantics: conclusion

- Typed language accommodating many important classical and probabilistic constructs
- Very powerful semantics in terms of ordered Banach space
- Advanced but completely mainstream mathematics

Denotational semantics: conclusion

- Typed language accommodating many important classical and probabilistic constructs
- Very powerful semantics in terms of ordered Banach space
- Advanced but completely mainstream mathematics
- Many 'moral' similarities with Scott's semantics

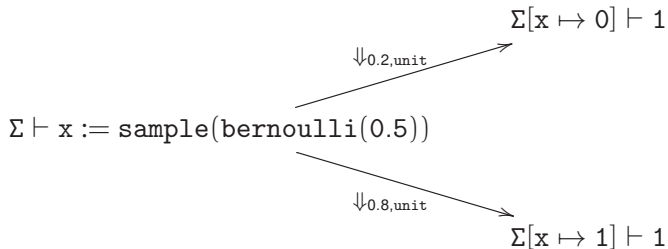
Operational Semantics

Operational semantics: discrete case

```
{  
  x=sample(bernoulli(0.2))  
}
```

Operational semantics: discrete case

```
{  
  x=sample(bernoulli(0.2))  
}
```

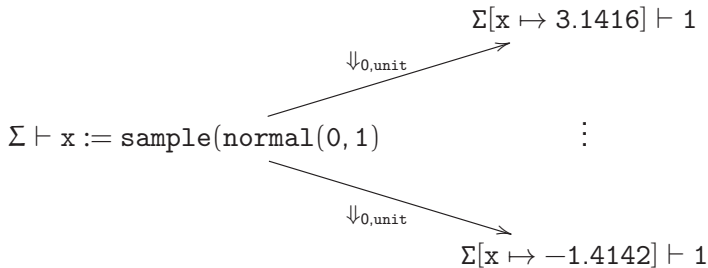


Operational semantics: continuous case

```
{  
  x=sample(normal(0,1))  
}
```

Operational semantics: continuous case

```
{
  x = sample(normal(0, 1))
}
```



Operational semantics: continuous case

```
{  
  x=sample(normal(0,1))  
}
```

Operational semantics: continuous case

```
{  
  x = sample(normal(0, 1))  
}
```

$$\begin{aligned}(\Sigma, \text{seed}) \vdash x = \text{sample}(\text{normal}(0, 1)) \Downarrow_{\text{unit}} \\ (\Sigma[x \mapsto \llbracket \text{normal}(0, 1) \rrbracket(\text{seed})], \text{seed} + 1) \vdash 1\end{aligned}$$

where

$$\llbracket \text{normal}(0, 1) \rrbracket : \mathbb{N} \rightarrow \mathbb{R}$$

with certain properties

Probabilistic Adequacy

Probabilistic adequacy: general aim

Probabilistic adequacy: general aim

- Ideally, the step-by-step execution of the model 'agrees' with its intended mathematical meaning

Probabilistic adequacy: general aim

- Ideally, the step-by-step execution of the model 'agrees' with its intended mathematical meaning
- But what does it mean for a probabilistic program?

Probabilistic adequacy: general aim

- Ideally, the step-by-step execution of the model ‘agrees’ with its intended mathematical meaning
- But what does it mean for a probabilistic program?
- Operationally: program = empirical process.
- Empirical distribution for $A \subseteq \llbracket \Sigma \rrbracket$

$$\mathbb{P}_n(A) = \frac{1}{n} \sum_{i=1}^n I_A(X_i)$$

Probabilistic adequacy: general aim

- Ideally, the step-by-step execution of the model ‘agrees’ with its intended mathematical meaning
- But what does it mean for a probabilistic program?
- Operationally: program = empirical process.
- Empirical distribution for $A \subseteq \llbracket \Sigma \rrbracket$

$$\mathbb{P}_n(A) = \frac{1}{n} \sum_{i=1}^n I_A(X_i)$$

- Probabilistic adequacy:

Does the empirical distribution converge to the denotational semantics? If yes, how fast?

Concentration of measure – an example

- Consider the hamming cube $\{0, 1\}^n$

Concentration of measure – an example

- Consider the hamming cube $\{0, 1\}^n$
- Metric space with $d((x_1, \dots, x_n), (y_1, \dots, y_n)) = \frac{1}{n} \sum_i (x_i + y_i \bmod 2)$

Concentration of measure – an example

- Consider the hamming cube $\{0, 1\}^n$
- Metric space with $d((x_1, \dots, x_n), (y_1, \dots, y_n)) = \frac{1}{n} \sum_i (x_i + y_i \bmod 2)$
- Measured space with counting measure $\mu(A) = \frac{\#A}{2^n}$

Concentration of measure – an example

- Consider the hamming cube $\{0, 1\}^n$
- Metric space with $d((x_1, \dots, x_n), (y_1, \dots, y_n)) = \frac{1}{n} \sum_i (x_i + y_i \bmod 2)$
- Measured space with counting measure $\mu(A) = \frac{\#A}{2^n}$
- Consider the function $f : \{0, 1\}^n \rightarrow \mathbb{R}, (x_1, \dots, x_n) \mapsto \frac{1}{n} \sum_i x_i$

Concentration of measure – an example

- Consider the hamming cube $\{0, 1\}^n$
- Metric space with $d((x_1, \dots, x_n), (y_1, \dots, y_n)) = \frac{1}{n} \sum_i (x_i + y_i \bmod 2)$
- Measured space with counting measure $\mu(A) = \frac{\#A}{2^n}$
- Consider the function $f : \{0, 1\}^n \rightarrow \mathbb{R}, (x_1, \dots, x_n) \mapsto \frac{1}{n} \sum_i x_i$
- The median of f is $\frac{1}{2}$: $\mu\{x \mid f(x) \leq \frac{1}{2}\} = \mu\{x \mid f(x) \geq \frac{1}{2}\}$

Concentration of measure – an example

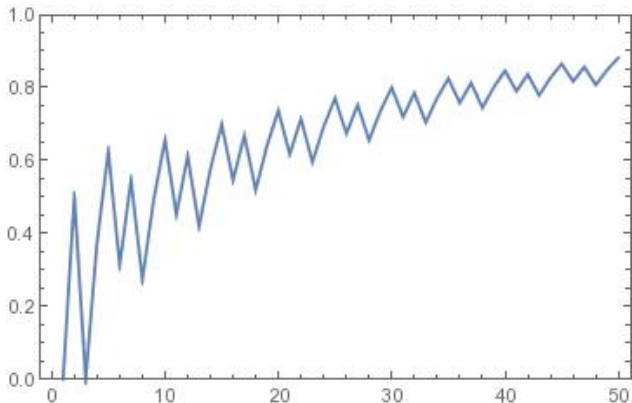
- Consider the hamming cube $\{0, 1\}^n$
- Metric space with $d((x_1, \dots, x_n), (y_1, \dots, y_n)) = \frac{1}{n} \sum_i (x_i + y_i \bmod 2)$
- Measured space with counting measure $\mu(A) = \frac{\#A}{2^n}$
- Consider the function $f : \{0, 1\}^n \rightarrow \mathbb{R}$, $(x_1, \dots, x_n) \mapsto \frac{1}{n} \sum_i x_i$
- The median of f is $\frac{1}{2}$: $\mu\{x \mid f(x) \leq \frac{1}{2}\} = \mu\{x \mid f(x) \geq \frac{1}{2}\}$
- How far are we away from the median on average?

$$A_f(\varepsilon, n) := \left\{ (x_1, \dots, x_n) \mid \left| f(x_1, \dots, x_n) - \frac{1}{2} \right| < \varepsilon \right\}$$

$$\mu(A_f(\varepsilon, n)) = \frac{1}{2^n} \sum_{k=\lceil n\varepsilon \rceil}^{\lfloor n\varepsilon \rfloor} \binom{n}{k}$$

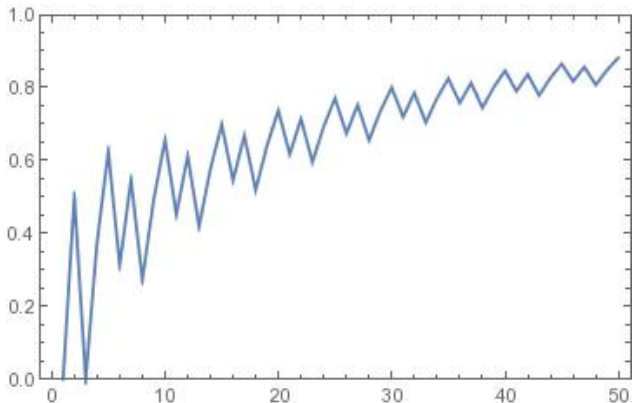
Concentration of measure – an example

- For $\varepsilon = \frac{1}{10}$ this is what $\mu(A_f(\varepsilon), n)$ varies as n increases:



Concentration of measure – an example

- For $\varepsilon = \frac{1}{10}$ this is what $\mu(A_f(\varepsilon), n)$ varies as n increases:



- This is a general pattern in *metric measured spaces*

Probabilistic adequacy and concentration of measure

```
{  
  x=sample(bernoulli(0.5))  
}
```

Probabilistic adequacy and concentration of measure

```
{  
  x=sample(bernoulli(0.5))  
}
```

- Denotationally $\llbracket x = \text{sample}(\text{bernoulli}(0.5)) \rrbracket (\mu) \propto \text{Bern}(0.5)$

Probabilistic adequacy and concentration of measure

```
{  
  x=sample(bernoulli(0.5))  
}
```

- Denotationally $\llbracket x = \text{sample}(\text{bernoulli}(0.5)) \rrbracket (\mu) \propto \text{Bern}(0.5)$
- The function f computes the empirical probability of $\Sigma = [x \mapsto 1]$

Probabilistic adequacy and concentration of measure

```
{  
  x=sample(bernoulli(0.5))  
}
```

- Denotationally $\llbracket x = \text{sample}(\text{bernoulli}(0.5)) \rrbracket (\mu) \propto \text{Bern}(0.5)$
- The function f computes the empirical probability of $\Sigma = [x \mapsto 1]$
- The convergence of $\mu(A_f)(\varepsilon, n)$ given above shows that the empirical probability (multiple runs of the program) converges with the denotational semantics.

Probabilistic adequacy and concentration of measure

```
{  
  x=sample(bernoulli(0.5))  
}
```

- Denotationally $\llbracket x = \text{sample}(\text{bernoulli}(0.5)) \rrbracket (\mu) \propto \text{Bern}(0.5)$
- The function f computes the empirical probability of $\Sigma = [x \mapsto 1]$
- The convergence of $\mu(A_f)(\varepsilon, n)$ given above shows that the empirical probability (multiple runs of the program) converges with the denotational semantics.
- Moreover, the rate of convergence can also be bounded (\sqrt{n})

Thank you.