

A ProbProg Language Taxonomy

Ohad Kammar
36th International Conference on
Mathematical Foundations
of
Programming Semantics
(MFPS'20)



THE UNIVERSITY of EDINBURGH

informatics ifcs

Laboratory for Foundations
of Computer Science



supported by:



THE ROYAL
SOCIETY

The
Alan Turing
Institute

Facebook Research

Probabilistic programming languages everywhere

Languages: Anglican, BLOGS, BUGS, Church
Edward, Gen, HackPPL, Hecara,
Monad-Bytes, Pro Stan, Turing,
Venture, WebPPL, ...

Semantics: Boolean-valued models, Boolean spaces
probabilistic coherence spaces,
event structures, Σ -finite kernels ...

Communities: MFPS, LICS, POPL, NeurIPS, ICFP, PLDI, ...

How to organise ProbProg Languages?

Proposed Taxonomy:

Conditioning Sampling	density	distribution
graded	Stan, Pyro	Haskell Tree-Typed Gen
non-graded	Anglican Monad-Byes	Talk structure ?

ProbProg Basics

Dataset: (x, y)

Model M (Bayesian Linear Regression)

$$a \sim \mathcal{N}(0, 2)$$

$$1.1 \sim \mathcal{N}(a \times 1, \frac{1}{4})$$

$$1.9 \sim \mathcal{N}(a \times 2, \frac{1}{4})$$

$$2.7 \sim \mathcal{N}(a \times 3, \frac{1}{4})$$

$$(1, 1.1)$$

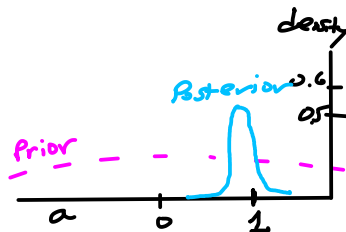
$$(2, 1.9)$$

$$(3, 2.7)$$

$$y = a \cdot x$$

$$P_{a \sim M} [l \leq a \leq u] \propto$$

$$\int_l^u dx e^{-\frac{x^2}{2}} \cdot e^{-\frac{4(1.1-x)^2}{2}} \cdot e^{-\frac{4(1.9-x \cdot 2)^2}{2}} \cdot e^{-\frac{4(2.7-x \cdot 3)^2}{2}}$$



Ingredients of a ProbProg Language

([Razey-Shar'17])
notation)

Sampling

$x \sim \mu$ in
 N

Conditioning

$M \rightsquigarrow \mu$

Sequencing

let $x = M$ in N

+ other PL features (H-O Functions, algebraic
data types, state, ...
Not in this talk)

$A, B ::=$

$\{l_1, \dots, l_n\}$

$| \mathbb{N}$

$| [a, b) \mid [a, b] \mid \dots$

$| A \times B$

Syntax for spaces

finite discrete

countable discrete

continuous
($a, b \in [-\infty, \infty]$)

products

$\Gamma ::= x_1 : A_1, \dots, x_n : A_n$

contexts


 $\mathcal{P} ::=$
 $\text{cat } \{l_1, l_2, \dots, l_n\}$
 $| \text{Counting}$
 $| \text{Lebesgue } [a, b]$
 $| \text{Lebesgue}$
 $\Omega ::= \mathcal{P}_1, \dots, \mathcal{P}_n$
 $\mu, \nu ::=$
 $\{l_1: \mu_1, \dots, l_n: \mu_n\}$
 $| \text{Geometric } M$
 $| \text{Uniform}$
 $| \mathcal{N}(M_{\text{mean}}, M_{\text{stdv}})$

Stock measures

categorical/discrete
countable counting

bounded Lebesgue
($a, b \in (-\infty, \infty)$)

Unbounded Lebesgue

Sample spaces

built-in
probability
distributions

Type system

Term judgements

$$\Gamma \mid \Omega \vdash M : A$$

→ graded type
system

mostly standard, e.g.:

$$\Gamma \mid \Omega_1 \vdash M : A \quad \Gamma, x:A \mid \Omega_2 \vdash N : B$$

$$\Gamma \mid \Omega_1, \Omega_2 \vdash \text{let } x = M \text{ in } N : B$$

$$\Gamma \mid \Omega_1 \vdash \text{mean} : \mathbb{R} \quad \Gamma \mid \Omega_2 \vdash \text{stdv} : (0, \infty)$$

$$\Gamma \mid \Omega_1, \Omega_2 \vdash \mathcal{N}(\text{mean}, \text{stdv}) \ll \text{Lebesgue}$$

Distribution judgements

$$\Gamma \mid \Omega \vdash \mu \ll \rho$$

Typing judgements

Sampling

e.g. Lebesgue := \mathbb{R}
etc.

$$\frac{\Gamma \mid \Omega_1 \vdash \mu \ll \rho \quad \Gamma, x:\underline{P} \mid \Omega_2 \vdash N:A}{\Gamma \mid \Omega_1, \rho, \Omega_2 \vdash \gg \mu \text{ in } N : A}$$

Conditioning

$$\frac{\Gamma \mid \Omega_1 \vdash M:\underline{P} \quad \Gamma \mid \Omega_2 \vdash \mu \ll \rho}{\Gamma \mid \Omega_1, \Omega_2 \vdash M \rightsquigarrow \mu : \underline{1}}$$

$\underline{c} := \{x\}$

Syntax

- Space \mathcal{T} or A
- Stochastic measure space \mathcal{P}

denotes

Standard Borel space (SBS)

SBS + σ -finite distribution
(counting or Lebesgue)

e.g.
$$\int \mathbb{I} \text{cat}(l, 1 \dots 1 l_n) \mathbb{I}_{\text{dist}}(dl) f(l) = \sum_{i=1}^n f(l_i)$$

$$\int \mathbb{I} \text{Lebesgue} \mathbb{I}_{\text{dist}} dx f(x) = \int dx f(x)$$

etc.

$$\mathbb{I}\Omega\mathbb{I} = \left(\prod_{\mathcal{P} \in \Omega} \mathbb{I}\mathcal{P}\mathbb{I}, \bigotimes_{\mathcal{P} \in \Omega} \mathbb{I}\mathcal{P}\mathbb{I}_{\text{dist}} \right)$$

product measure

Semantics (ctd)

Syntax

- terms

$\llbracket \Gamma \mid \Omega \vdash M : A \rrbracket : \llbracket \Gamma \rrbracket$

- distributions

$\llbracket \Gamma \mid \Omega \vdash \mu \ll \rho \rrbracket : \llbracket \Gamma \rrbracket$

denotes

density + random variable

$(\llbracket M \rrbracket_{\text{den}}, \llbracket M \rrbracket_{\text{val}}) \xrightarrow{\llbracket \Omega \rrbracket, \llbracket \Omega \rrbracket} \mathcal{W} \times \llbracket A \rrbracket$

(parameterised) density functions
 $\mathcal{W}^{\underline{\Omega} \times \underline{\rho}}$

$\longrightarrow \mathcal{W}^{\underline{\Omega} \times \underline{\rho}}$

E.g.:

"Bayes' Theorem"

$$\llbracket a \leftarrow \mu \text{ in } N \rrbracket_{\text{den}}^{\mathcal{M}} := \lambda(\omega_1, a, \omega_2).$$

$$\llbracket \mu \rrbracket_{\text{den}}(\mathcal{M}; \omega_1) \times$$

$$\llbracket N \rrbracket_{\text{den}}(\mathcal{M}[a \mapsto a]; \omega_2)$$

$$\llbracket M \rightsquigarrow \mu \rrbracket_{\text{den}}^{\mathcal{M}} := \lambda(\omega_1, \omega_2).$$

$$\llbracket M \rrbracket_{\text{den}}(\mathcal{M}; \omega_1) \times \llbracket \mu \rrbracket_{\text{den}}(\mathcal{M}; \omega_2, \llbracket M \rrbracket_{\text{val}}(\mathcal{M}; \omega_1))$$

etc.

A graded monad

Rest is standard [Katsumata'14]:

$$T_{\Omega} S := W^{\Omega} \times S^{\Omega} \cong (W \times S)^W$$

$(W, ; 1)$ monoid \uparrow

Each model $\Gamma / \Omega \vdash M : A$ gives:

- A **kernel**

$$\llbracket \Gamma \rrbracket \xrightarrow{\llbracket M \rrbracket_{\text{dist}}} \llbracket A \rrbracket$$

$$\llbracket M \rrbracket_{\text{dist}}(\Gamma, U) := \int \llbracket \Omega \rrbracket d\omega$$

$$\llbracket M \rrbracket_{\text{den}}(\Gamma; \omega) \times \left[\llbracket M \rrbracket_{\text{val}}(\Gamma; \omega) \in U \right]$$

- model evidence:**

used to evaluate/
diagnose
sit

$$\llbracket M \rrbracket_{\text{ev}}^{\Gamma} := \llbracket M \rrbracket_{\text{dist}}(\Gamma, \llbracket A \rrbracket)$$

Foundations

$\mathbb{W} \stackrel{?}{=} \text{not a measurable space}$
eg. $[0, \infty]^{\mathbb{R}}$

(Aumann's Thm)

I used quasi-Borel spaces: [Hansen, Staton, Leinster, May 16, 17]
 $(\underline{X}, \mathcal{M}_X)$ $\xrightarrow{\text{subset}} \mathcal{M}_X \subseteq \mathcal{P}(X^{\mathbb{R}})$

+ closure axioms.

E.g. S measurable space $\Rightarrow \mathcal{M}_S := \text{Meas}(\mathbb{R}, S)$

Foundations (ctd)

Qbs as a category:

$f: X \rightarrow Y$ is

Function $\underline{X} \xrightarrow{f} \underline{Y}$
s.t. $\forall \alpha \in M_X$.

$$\mathbb{R} \xrightarrow{\alpha} \underline{X} \xrightarrow{f} \underline{Y} \\ \in M_Y$$

Thm [HSK17]:

For SBS S, T :

$$\text{Qbs}(S, T) = \text{Meas}(S, T)$$

So Qbs is a conservative extension of SBS.

Distributions

For distributions, a noned $\text{Dist: obs} \rightarrow \text{obs}$ [Scibior et al. 17]

$S \subseteq \text{Obs}$, $\text{Dist } S = S\text{-finite measures on } S$

$\mathcal{M}_{\text{Dist } S} = S\text{-finite kernels } \mathcal{R} \mapsto S$



[Staton'17, Kalenbay'17]

$\kappa \text{ } S\text{-finite} = \sigma\text{-affine combination of prob. kernels:}$

$$\kappa = \sum_{n=0}^{\infty} w_n \cdot \kappa_n \quad w_n \in [0, \infty]$$



Programming with Ω explicitly is tedious
Prefer:

and semantically:

$$\Gamma \vdash M : A \qquad \Gamma \vdash \mu \ll \text{Dist } P$$

$$\llbracket \Gamma \rrbracket \xrightarrow{\llbracket M \rrbracket_{\text{dist}}} \text{Dist } \llbracket A \rrbracket \qquad \llbracket \Gamma \rrbracket \xrightarrow{\llbracket \mu \rrbracket_{\text{dist}}} W^{\llbracket P \rrbracket}$$

[Staton'17]

We can now include any primitive probability distributions
for sampling
Conditioning still requires density:

$$\frac{\Gamma \vdash \mu \ll P \quad \Gamma \vdash M : \underline{P}}{\Gamma \vdash M \rightsquigarrow \mu : \underline{1}}$$



- Inferene benefits from graded information
(e.g. Stan's HMC, ADVI)
- Programming is easier without grading

We can^o pay the price (either way)

- use static analysis to compile
non-graded \rightarrow graded

(e.g., SlicStan [Gorinovan, Gordon, Sutton'19]
TreeTypes [Lew et al. '20])



(ongoing!)

$\text{Dist}_{\ll P} =$ distributions with density
w.r.t. $\ll P \gg$

$$\begin{aligned} \text{So: } \ll \Gamma \mid \Omega \vdash M : A \gg : \ll \Gamma \gg &\xrightarrow{(\ll M \gg_{\text{let}}, \ll M \gg_{\text{val}})} \text{Dist}_{\ll \ll \Omega \gg \times \ll A \gg}^{\ll \Omega \gg} \\ \ll \Gamma \mid \Omega \vdash \mu \ll P \gg : \ll \Gamma \gg &\xrightarrow{(\ll \mu \gg_{\text{let}}, \ll \mu \gg_{\text{val}})} \text{Dist}_{\ll \ll \Omega \gg \times \text{Dist}_{\ll P}^{\ll P \gg}}^{\ll \Omega \gg} \end{aligned}$$

$$\ll M \rightsquigarrow \mu \gg_{\text{r}} := \varphi \circ \ll M \gg_{\text{let}}^{\text{r}} \quad \text{where}$$

$$\varphi : \ll \Omega \gg \rightarrow \mathcal{W}$$

$$\varphi := \frac{\Delta \alpha_{\text{r}}^{\text{r}} \ll \mu \gg_{\text{dist}}^{\text{r}}}{\Delta \ll \Omega \gg}$$

and

$$\begin{array}{ccc} & \xrightarrow{\alpha} & \\ \ll \Omega \gg & \text{disintegrate} & \ll P \gg \\ & \text{w.r.t. } \Omega & \\ & \xleftarrow{\alpha^{\text{r}, \Omega}} & \end{array}$$

(Ongoing!)

- build on/relate to
 - Hakem [Ramsy-Sham '17, Narayanan '19]
 - Ong-Mattinson [unpub 13, 16]
 - Bayesian inversion [Dahlquist et al. '16-'18, '20]
- Foundational hazards: disintegrating σ -finite measures
- non-grabed [c.f. Vakar-Ong '17]

Summary

- Proposed Taxonomy:

Conditioning Sampling	density	distribution
graded	Stan, Pyro	Hakan Tree-Typed Gen
non- graded	Anglica Monad-Pyres	?

- quasi-formal species as a convenient notation-theory