# Probabilistic Programming Languages

## Hugo Paquet

*MPRI 25–26*

# Recap

- **Probabilistic programming**: a language for building probabilistic models, including samples and observations.

- Monte Carlo methods for Bayesian inference: importance sampling, particle filters, etc.

- With Monte Carlo methods, the result is always an **approximation.** This is often necessary because the true posterior distribution is generally uncomputable.

# Today

- Semantics for **discrete** probabilistic programs using monads (including a small amount of category theory).

- We add higher-order function types, e.g. $(\mathrm{Nat} \to \mathrm{Nat}) \to \mathrm{Nat}$.

- When the distributions have finite support we can perform **exact inference**.

- Other efficient representations of discrete distributions: probabilistic graphical models.

# 1. Semantics for discrete probabilistic programs using monads

# The need for semantics

- Semantics provides a formal connection between probabilistic programs and the statistical models they represent.

- Using semantics we can reason about program equivalence, e.g. there could be several implementations of the same probability distribution (which may not be equally efficient for inference).
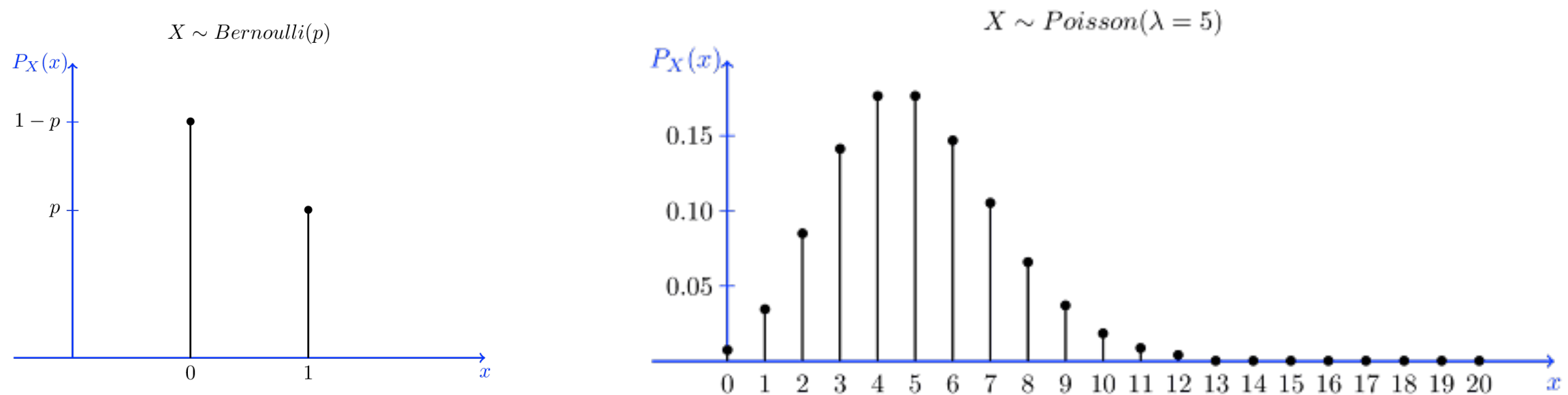
$$
\boxed{\begin{array}{l} \text{let } x = \text{normal 0 2 in} \\ \text{return } x \end{array}} \quad = \quad \boxed{\begin{array}{l} \text{let } y = \text{normal 0 1 in} \\ \text{let } z = \text{normal 0 1 in} \\ \text{return } (y + z) \end{array}}
$$

- **Correctness of inference.** Inference algorithms often approximate. But what are they trying to approximate?

# Discrete probability and measure theory

Discrete probability is based on sets $X$ and distributions functions $X \to [0,1]$.



It is a special case of measure-theoretic probability:

- A measurable space $(X, \Sigma_X)$ is called **discrete** if $\Sigma_X = \mathscr{P}(X)$. Discrete measurable spaces are in correspondence with sets.

- A probability measure $\mu : \Sigma_X \to [0,1]$ on a discrete space is completely determined by its value on the set of **atoms** $\{x \in X \mid \mu\{x\} > 0\}$, and therefore by the a probability mass function $f_\mu : X \to [0,1]$.

- Probability measures have at most countably many atoms. (True even for unnormalized finite measures but not for general measures.)

6

# A probabilistic language

**Types.**

$$A, B ::= \text{Unit} \mid \text{Bool} \mid \text{Nat} \mid \text{Real} \mid A \times B \mid A \to B$$

**Terms.**

*basic operations
+, >, etc.*

$$s, t ::= () \mid b \in \{\text{True}, \text{False}\} \mid r \in \mathbb{R} \mid n \in \mathbb{N} \mid \text{op}$$

$$\langle s, t \rangle \mid \text{fst } s \mid \text{snd } s \mid \text{if } s \text{ then } t \text{ else } t'$$
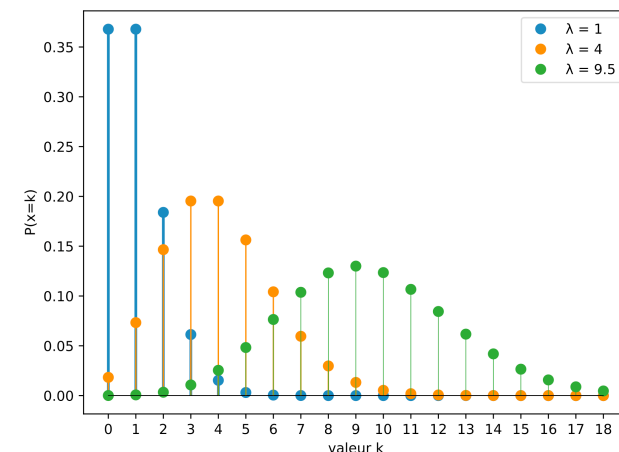
$$\mid x \mid \lambda x.\, s \mid s\ t$$

$$\text{return } s \mid \text{let } x = s \text{ in } t$$

$$\text{bernoulli } s \mid \text{poisson } s$$

poisson($\lambda$) = a distribution on $\mathbb{N}$ with infinite support, defined when $\lambda > 0$.

$$\Pr(X = k) = e^{-\lambda} \frac{\lambda^k}{k!}$$



7

# A higher-order probabilistic language

$$\frac{}{\Gamma, x : A \vdash^D x : A} \qquad \frac{\Gamma, x : A \vdash^P t : B}{\Gamma \vdash^D \lambda x.\, t : A \to B} \qquad \frac{\Gamma \vdash^D t : A \to B \qquad \Gamma \vdash^D s : A}{\Gamma \vdash^P t\, s : B}$$

$$\frac{\Gamma \vdash^D s : A \qquad \Gamma \vdash^D t : B}{\Gamma \vdash^D \langle s, t \rangle : A \times B} \qquad \frac{\Gamma \vdash^D s : A \times B}{\Gamma \vdash^D \mathrm{fst}\ s : A} \qquad \frac{\Gamma \vdash^D s : A \times B}{\Gamma \vdash^D \mathrm{snd}\ s : B}$$

$$\frac{\Gamma \vdash^D s : \mathrm{Bool} \quad \Gamma \vdash^P t : A \quad \Gamma \vdash^P t' : A}{\Gamma \vdash^P \mathrm{if}\ s\ \mathrm{then}\ t\ \mathrm{else}\ t' : A}$$

$$\frac{\Gamma \vdash^D s : A}{\Gamma \vdash^P \mathrm{return}\ s : A} \qquad \frac{\Gamma \vdash^P s : A \qquad \Gamma, x : A \vdash^P t : B}{\Gamma \vdash^P \mathrm{let}\ x = s\ \mathrm{in}\ t : B}$$

$$\frac{\Gamma \vdash^D s : \mathrm{Real}}{\Gamma \vdash^P \mathrm{bernoulli}\ s : \mathrm{Bool}} \qquad \frac{\Gamma \vdash^D s : \mathrm{Real}}{\Gamma \vdash^P \mathrm{poisson}\ s : \mathrm{Nat}}$$

*Ignore out-of-bounds errors for parameters.*

We also assume appropriate typing rules for constants (unit, booleans, integers, reals) and basic operations and functions.

# Expressivity (1)

The type system is very restrictive, for example the program

$\langle$ bernoulli 0.1, bernoulli 0.2 $\rangle$

is not a well-typed program. Instead we must specify the sampling order. The following is a well-typed program of type $\mathrm{Bool} \times \mathrm{Bool}$:

let x = bernoulli 0.1 in

let y = bernoulli 0.2 in

$\langle$ x, y $\rangle$

# Expressivity (2)

The language doesn't have an explicit type constructor $\mathrm{dist}$ for probability distributions.

But we can recover it as:

$$\mathrm{dist}\, A \;:=\; \mathrm{Unit} \to A$$
$$\mathrm{sample}\, d \;:=\; d\,() \qquad \text{(for } d : \mathrm{dist}\, A)$$

In fact, the type $\mathrm{Unit} \to A$ is always a monad in a call-by-value language with function types (prove it.)

# Semantics

We will give formal semantics for this language in terms of sets, functions, and discrete probability distributions.

The semantics of the language consists of semantics for each kind of syntactic element:

| Syntactic elements | Semantics using the probability monad |
|---|---|
| type $A$ | set $[\![A]\!]$ |
| $\Gamma = x_1 : A_1, \ldots, x_n : A_n$ | $[\![\Gamma]\!] = [\![A_1]\!] \times \ldots \times [\![A_n]\!]$ |
| $\Gamma \vdash^{\mathrm{D}} s : A$ | $[\![s]\!] : [\![\Gamma]\!] \to [\![A]\!]$ |
| $\Gamma \vdash^{\mathrm{P}} s : A$ | $[\![s]\!] : [\![\Gamma]\!] \to \mathrm{dist}\,[\![A]\!]$ |

The definition is by induction on types and terms. We will see that it works because $\mathrm{dist}$ is a monad.

# dist : a monad for discrete probability

A **probability distribution** on a set $X$ is a function $d : X \to [0,1]$ such that $\sum_{x \in X} d(x) = 1$.

Let $\mathrm{dist}\, X$ denote the set of all probability distributions on $X$.

The dist operator comes equipped with combinators:

$$\mathrm{return} : X \longrightarrow \mathrm{dist}\, X$$
$$x \longmapsto \delta_x$$

used for the semantics of return

$$>\!\!>\!= \; : \mathrm{dist}\, X \times (X \to \mathrm{dist}\, Y) \longrightarrow \mathrm{dist}\, Y$$
$$(d >\!\!>\!= f) \longmapsto \left( y \mapsto \sum_{x \in X} d(x) f(x)(y) \right)$$

used for the semantics of let … in

They satisfy the three monad equations:

$$d >\!\!>\!= \mathrm{return} \quad = \quad d$$
$$\mathrm{return}(x) >\!\!>\!= f \quad = \quad f(x)$$
$$(d >\!\!>\!= f) >\!\!>\!= g \quad = \quad d >\!\!>\!= \lambda y . (f(y) >\!\!>\!= g)$$

# Semantics of types

Every type defines a **set.**

$$[\![\text{Unit}]\!] = \{ \, * \, \}$$

$$[\![\text{Bool}]\!] = \{\text{True}, \text{False}\}$$

$$[\![\text{Nat}]\!] = \mathbb{N}$$

$$[\![\text{Real}]\!] = \mathbb{R}$$

$$[\![A \times B]\!] = [\![A]\!] \times [\![B]\!]$$

$$[\![A \to B]\!] = [\![A]\!] \to \text{dist} \, [\![B]\!]$$

*$X \to Y :=$ set of all functions from $X$ to $Y$*

- In this language the real numbers are just a set = a discrete measurable space (this is not the usual $\sigma$-algebra on $\mathbb{R}$!).
- But it's ok to have a discrete probability distribution on an uncountable set.
- In fact, even if we remove the type $\text{Real}$, the semantics involves uncountable sets, like $[\![\text{Unit} \to \text{Bool}]\!]$ or $[\![\text{Nat} \to \text{Nat}]\!]$.
- A probability distribution $d$ on an uncountable set $X$ always has countable support, otherwise the sum $\sum_{x \in X} d(x)$ would diverge.

13

# Semantics of terms (1)

$$\frac{}{\Gamma, x : A \vdash^{\mathrm{D}} x : A}$$

$$\Gamma \times A \longrightarrow A$$

is the projection.

---

$$\frac{\Gamma, x : A \vdash^{\mathrm{P}} t : B}{\Gamma \vdash^{\mathrm{D}} \lambda x.\, t : A \to B}$$

Every function
$f : \Gamma \times A \longrightarrow \mathrm{dist}\ B$
determines a function
$$\Gamma \longrightarrow (A \to \mathrm{dist}\ B)$$
$$\gamma \longmapsto f(\gamma, -)$$

(In fact this is a
1-1 correspondence.)

---

$$\frac{\Gamma \vdash^{\mathrm{D}} t : A \to B \qquad \Gamma \vdash^{\mathrm{D}} s : A}{\Gamma \vdash^{\mathrm{P}} t\, s : B}$$

Given two functions
$$g : \Gamma \longrightarrow (A \to \mathrm{dist}\ B)$$
$$h : \Gamma \longrightarrow A$$
we can build the function
$$\Gamma \longrightarrow B$$
$$\gamma \longmapsto g(\gamma)(h(\gamma))$$

# Semantics of terms (2)

$$\frac{\Gamma \vdash^{\mathrm{D}} s : A \qquad \Gamma \vdash^{\mathrm{D}} t : B}{\Gamma \vdash^{\mathrm{D}} \langle s, t \rangle : A \times B}$$

Given two functions
$$g : \Gamma \longrightarrow A$$
$$h : \Gamma \longrightarrow B$$
we can build the function
$$\Gamma \longrightarrow A \times B$$
$$\gamma \longmapsto (g(\gamma), h(\gamma))$$

$$\frac{\Gamma \vdash^{\mathrm{D}} s : A \times B}{\Gamma \vdash^{\mathrm{D}} \mathrm{fst}\ s : A} \qquad \frac{\Gamma \vdash^{\mathrm{D}} s : A \times B}{\Gamma \vdash^{\mathrm{D}} \mathrm{snd}\ s : B}$$

We can post-compose a function
$$g : \Gamma \longrightarrow A \times B$$
with either of the two projection
$$\pi_1 : A \times B \longrightarrow A$$
$$\pi_2 : A \times B \longrightarrow B$$

# Semantics of terms (3)

$$\frac{\Gamma \vdash^D s : \text{Bool} \quad \Gamma \vdash^P t : A \quad \Gamma \vdash^P t' : A}{\Gamma \vdash^P \text{if } s \text{ then } t \text{ else } t' : A}$$

Given functions

$$h : \Gamma \longrightarrow \text{Bool}$$

$$f : \Gamma \longrightarrow \text{dist } A$$

$$g : \Gamma \longrightarrow \text{dist } A$$

we define a function

$$\Gamma \longrightarrow \text{dist } A$$

$$\gamma \longmapsto \begin{cases} f(\gamma) & \text{if } h(\gamma) = \text{True} \\ g(\gamma) & \text{if } h(\gamma) = \text{False} \end{cases}$$

**Remark.** We can easily add a typing rule:

$$\frac{\Gamma \vdash^D s : \text{Bool} \quad \Gamma \vdash^D t : A \quad \Gamma \vdash^D t' : A}{\Gamma \vdash^D \text{if } s \text{ then } t \text{ else } t' : A}$$

# Semantics of terms (4)

$$\frac{\Gamma \vdash^{\mathrm{D}} s : A}{\Gamma \vdash^{\mathrm{P}} \mathrm{return}\ s : A}$$

Given a function
$$g : \Gamma \longrightarrow A$$
we define the function
$$\Gamma \longrightarrow \mathrm{dist}\ A$$
$$\gamma \mapsto \mathrm{return}(g(\gamma))$$

$$\frac{\Gamma \vdash^{\mathrm{P}} s : A \qquad \Gamma, x : A \vdash^{\mathrm{P}} t : B}{\Gamma \vdash^{\mathrm{P}} \mathrm{let}\ x = s\ \mathrm{in}\ t : B}$$

Given functions
$$g : \Gamma \longrightarrow \mathrm{dist}\ A$$
$$h : \Gamma \times A \longrightarrow \mathrm{dist}\ B$$
we define the function
$$\Gamma \longrightarrow \mathrm{dist}\ B$$
$$\gamma \longmapsto g(\gamma) >\!\!>= h(\gamma, -)$$

# Semantics of terms (5)

$$\frac{\Gamma \vdash^D s : \mathrm{Real}}{\Gamma \vdash^P \mathrm{bernoulli}\ s : \mathrm{Bool}} \qquad\qquad \frac{\Gamma \vdash^D s : \mathrm{Real}}{\Gamma \vdash^P \mathrm{poisson}\ s : \mathrm{Nat}}$$

We have functions

$$\mathrm{bernoulli} : [0,1] \longrightarrow \mathrm{dist}\ \mathrm{Bool}$$
$$\mathrm{poisson} : \mathbb{R}_{>0} \longrightarrow \mathrm{dist}\ \mathbb{N}$$

and we need

$$\mathrm{bernoulli?} : \mathbb{R} \longrightarrow \mathrm{dist}\ \mathrm{Bool}$$
$$\mathrm{poisson?} : \mathbb{R} \longrightarrow \mathrm{dist}\ \mathbb{N}$$

The behaviour is undefined but we need to define it.
(Typically, we would use a more precise type system
or add an exception monad.)

# Semantics and enumeration trees

let x = bernoulli 0.2 in
let y = bernoulli 0.3 in
return (x ∧ y)

$$\frac{\dfrac{}{\vdash^{P} \text{bernoulli } 0.2 : \text{Bool}} \qquad \dfrac{x : \text{Bool} \vdash^{P} \text{bernoulli } 0.3 : \text{Bool} \qquad \dfrac{x : \text{Bool}, y : \text{Bool} \vdash^{D} x \wedge y : \text{Bool}}{x : \text{Bool}, y : \text{Bool} \vdash^{P} \text{return } x \wedge y : \text{Bool}}}{x : \text{Bool} \vdash \text{let } y = \text{bernoulli } 0.3 \text{ in return } x \wedge y : \text{Bool}}}{\vdash \text{let } x = \text{bernoulli } 0.2 \text{ in } (\text{let } y = \text{bernoulli } 0.3 \text{ in return } x \wedge y) : \text{Bool}}$$

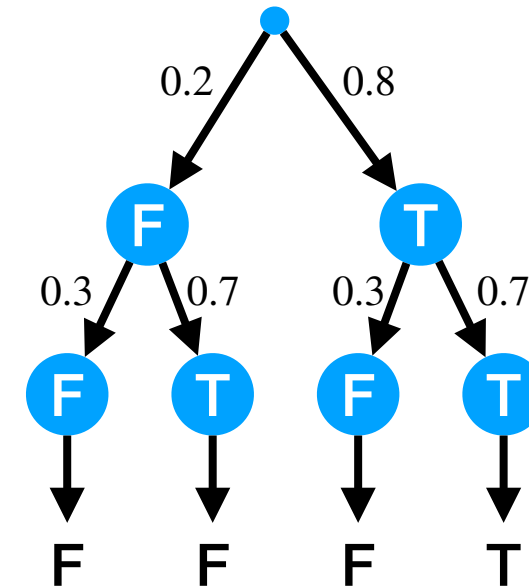We compute the semantics of the program using the basic building blocks:

$$1 \xrightarrow{\text{bernoulli } 0.2} \text{dist Bool} \qquad\qquad 1 \xrightarrow{\text{bernoulli } 0.3} \text{dist Bool} \qquad\qquad \text{Bool} \times \text{Bool} \xrightarrow{\wedge} \text{Bool}$$

# Semantics and enumeration trees

```
let x = bernoulli 0.2 in
let y = bernoulli 0.3 in
return (x ∧ y)
```
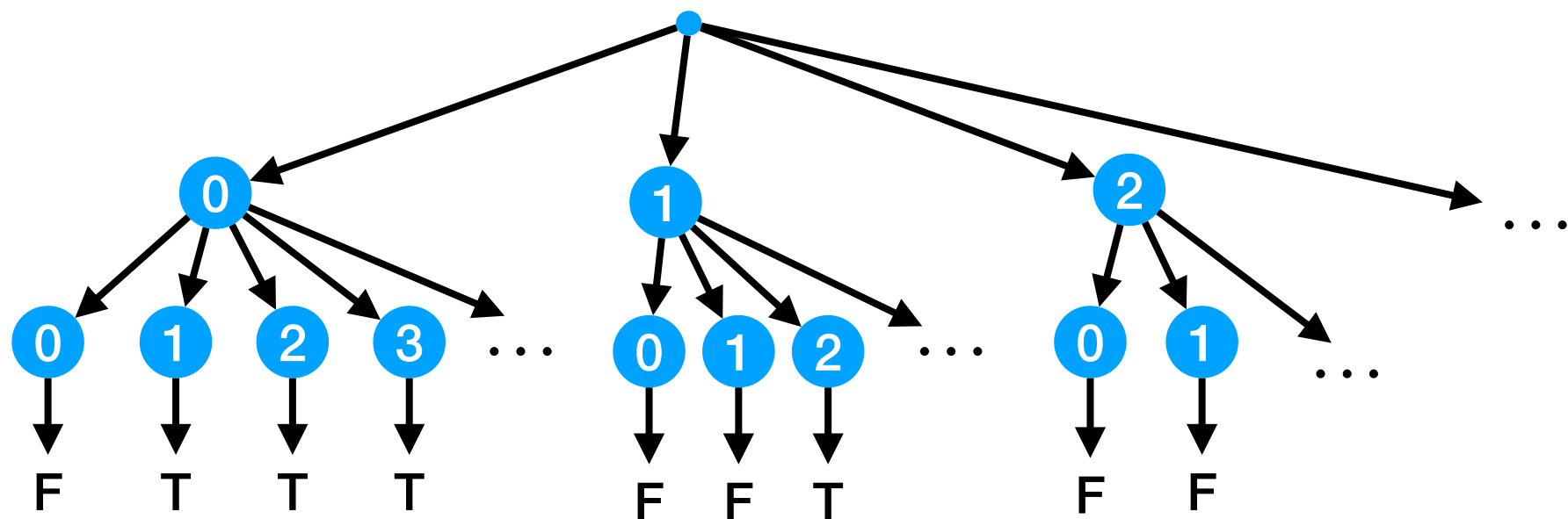
$$\cfrac{\vdash^{P} \text{bernoulli } 0.2 : \text{Bool} \qquad \cfrac{x : \text{Bool} \vdash^{P} \text{bernoulli } 0.3 : \text{Bool} \qquad \cfrac{x : \text{Bool}, y : \text{Bool} \vdash^{D} x \wedge y : \text{Bool}}{x : \text{Bool}, y : \text{Bool} \vdash^{P} \text{return } x \wedge y : \text{Bool}}}{x : \text{Bool} \vdash \text{let } y = \text{bernoulli } 0.3 \text{ in return } x \wedge y : \text{Bool}}}{\vdash \text{let } x = \text{bernoulli } 0.2 \text{ in (let } y = \text{bernoulli } 0.3 \text{ in return } x \wedge y) : \text{Bool}}$$

We compute the semantics of the program using the basic building blocks:

$$1 \xrightarrow{\text{bernoulli } 0.2} \text{dist Bool} \qquad 1 \xrightarrow{\text{bernoulli } 0.3} \text{dist Bool} \qquad \text{Bool} \times \text{Bool} \xrightarrow{\wedge} \text{Bool}$$

let x = poisson 2 in
let y = poisson 3 in
return (x < y)

# An abstract perspective

- A **category** $\mathscr{C}$ has a collection $\operatorname{ob}\mathscr{C}$ of objects and a collection of morphisms $\mathscr{C}(A, B)$ between any two objects $A$ and $B$, with a composition operation and identity morphisms.

- A **functor** $\mathscr{C} \to \mathscr{D}$ consists of a map $F : \operatorname{ob}\mathscr{C} \to \operatorname{ob}\mathscr{D}$ together with a maps $\mathscr{C}(A, B) \longrightarrow \mathscr{D}(FA, FB)$ for every $A, B \in \operatorname{ob}\mathscr{C}$.

> - The category $\mathbf{Set}$ has sets as objects and $\mathbf{Set}(A, B)$ is the set of functions $A \to B$.
>
> - The category $\mathbf{Stoch}$ has sets as objects and $\mathbf{Stoch}(A, B)$ is the set of stochastic matrices, i.e. functions $A \to \operatorname{dist} B$.
>
>   (This is called the *Kleisli category* for the monad dist.)

# An abstract perspective

What structure did we use for the semantics of our language?

- In **Set** we have finite cartesian products $\times$, function spaces $\to$, and *coproducts* $A + A$ (for **Bool** and conditionals).

- In **Stoch** we can also take finite products of objects and morphisms. This is a *monoidal product*, not cartesian because components can be correlated.

- There is an **identity-on-objects** functor $L : \textbf{Set} \to \textbf{Stoch}$ which preserves the products.

- There is a functor $\mathrm{R} : \textbf{Stoch} \to \textbf{Set}$ which sends

$$A \in \textbf{Stoch} \longmapsto \mathrm{dist}\, A \in \textbf{Set} \qquad \text{(action on morphisms?)}$$

- There is a well-behaved (**natural**) correspondence $\textbf{Set}(A, \mathrm{dist}\, B) \cong \textbf{Stoch}(A, B)$ .
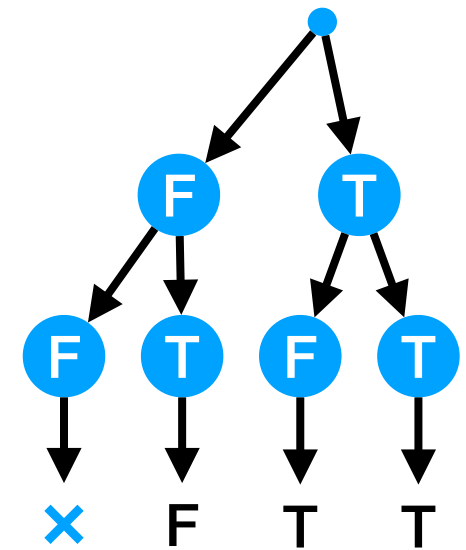
# Extending the language with observations

We add to the language a new syntactic primitive:

$$s, t ::= \dots \mid \mathrm{assume}\ s$$

with typing rule

$$\frac{\Gamma \vdash^{\mathrm{D}} s : \mathrm{Bool}}{\Gamma \vdash^{\mathrm{P}} \mathrm{assume}\ s : \mathrm{Unit}}$$

```
let x = bernoulli 0.5 in
let y = bernoulli 0.5 in
assume (x ‖ y);
return x
```



For the semantics of $\mathrm{assume}$ we need a new monad: $\mathrm{dist}$ is replaced by $\mathrm{dist}(- + 1)$, where $+$ is disjoint union and $1 = \{\ *\ \}$.

**Exercice.** Define the return and bind.

# Extending the language with observations

There is an **isomorphism of monads** $\mathrm{dist}(X + 1) \cong \mathrm{dist}_{\leq}(X)$

where $\mathrm{dist}_{\leq}$ is the monad of sub-probability distributions over $X$. (**Why?**)

The semantics of inference is given by a normalization operator:

$$\mathrm{infer} : \mathrm{dist}_{\leq}(X) \to \mathrm{dist}(X) + 1$$

$$\mu \longmapsto \begin{cases} \mathrm{inl}\dfrac{\mu}{\mu(X)} & \text{if } \mu(X) > 0 \\ \mathrm{inr} * & \text{otherwise} . \end{cases}$$

# Extending the language with observations

**Remark.** If we had a more expressive operation $\mathrm{factor}$ for "soft constraints" then the monad $\mathrm{dist}$ would not be sufficiently expressive for semantics.

Instead we would use either $\mathrm{dist}(X \times \mathbb{R}_+)$ or $\mathrm{meas}$.

# Towards a semantics with random elements

- Suppose that $\Omega$ is an idealised random number generator. Model this as the space $\Omega = [0,1]^{\mathbb{N}}$ with the product of uniform distributions.

- Then we can sample a uniform value ($\Omega \to [0,1]$) and split the infinite supply into two generators ($\Omega \to \Omega \times \Omega$), in a measure-preserving way.

- Then the set of measurable functions $\Omega \to X$ models an implementation of $\mathrm{dist}\, X$. But there are problems in making this precise:

  1. $\Omega \to X$ is not really a measurable space  $\boxed{\text{Fixed in next lecture.}}$

  2. The associativity law of monads is only true up to measure-preserving iso:  $\boxed{\text{Not really a problem.}}$
  $$((d >\!\!>= f) >\!\!>= g) = (d >\!\!>= \lambda y\,.\,(f(y) >\!\!>= g))$$

# Operational semantics

- We can set up a probabilistic version of the usual reduction relation in the $\lambda$-calculus.

- E.g. $\text{bernoulli } 0.2 \xrightarrow{0.2} \text{True}$ and $\text{bernoulli } 0.2 \xrightarrow{0.8} \text{False}$

- Then we can reason about probabilities of terminating or returning a particular value.

- We will not look at operational semantics in this course because normalization in PPLs doesn't have a clear "rewriting" semantics: the operational behaviour of a program depends on the inference algorithm.

*E.g. an operational approach to the Metropolis-Hastings inference algorithm:*

**A lambda-calculus foundation for universal probabilistic programming**

**Authors**: Johannes Borgström, Ugo Dal Lago, Andrew D. Gordon, Marcin Szymczak | Authors Info & Claims

ICFP 2016: Proceedings of the 21st ACM SIGPLAN International Conference on Functional Programming • Pages 33 - 46
https://doi.org/10.1145/2951913.2951942

# 2. TD (discrete semantics and monads)

# 3. Probabilistic graphical models and PPLs

# Key idea

Very simple probabilistic programs can be represented by graphs of random variables. This is connected to "probabilistic graphical models".

```
let b = bernoulli 0.2 in
let c = bernoulli 0.3 in
let n = poisson (if (b ∧ c) then 1 else 2) in
return n
```



- The graphical representation has a formal meaning, and it is often close to the programmer's intuition.

- It allows for efficient inference (we will see this next).

- But the expressivity of graphical models is very limited:

```
let b = bernoulli 0.2 in
if b then
      let n = poisson 1 in
      let y = bernoulli (1/n) in
      return y
else return b
```

# Bayesian networks

The best-known kind of graphical model.

**Informal idea:** a directed graph which encodes conditional independence properties for a joint probability distribution.

An example:



- Sets $X_1, X_2, X_3, X_4$.

- $p \in \mathrm{dist}(X_1 \times X_2 \times X_3 \times X_4)$
  $p(x_1, x_2, x_3, x_4) = p(x_4)\, p(x_1 \mid x_4)\, p(x_3)\, p(x_2 \mid x_1, x_3)$

# Bayesian networks

**Informal idea:** a directed graph which encodes conditional independence properties for a joint probability distribution.

**Definition.** A **Bayesian network** consists of:
- An irreflexive directed acyclic graph $(V, E \subseteq V \times V)$
- A family of sets $\{X_v\}_{v \in V}$
- A distribution

$$P \in \mathrm{dist}\left(\prod_{v \in V} X_v\right)$$

such that for every $\boldsymbol{x} = (x_v)_{v \in V} \in \prod_{v \in V} X_v$,

$$P(\boldsymbol{x}) = \prod_{v \in V} P(x_v \mid \boldsymbol{x}_{\mathrm{pa}(v)}).$$

Equivalently: the distribution of a variable $v$ only depends on the values of its parents $\mathrm{pa}(v)$ in the graph: $P(x_v \mid \boldsymbol{x}_{V \setminus v}) = P(x_v \mid \boldsymbol{x}_{\mathrm{pa}(v)})$

**Remark**. There is no canonical way to represent a given distribution as a graph.

For example, the following graphs express the same independence assumptions (none) and are both compatible with **every** distribution on $X_1 \times X_2$.

Consider the following model with four Boolean variables.

We observe that the grass is wet, and we want to know if it has been raining or if the sprinkler has been turned on.

Our prior distributions depend on whether the sky is cloudy.

$$
C\backslash S \quad \begin{array}{cc} \texttt{t} & \texttt{f} \end{array}
$$
$$
\begin{array}{c} \texttt{t} \\ \texttt{f} \end{array}
\begin{bmatrix} 0.1 & 0.9 \\ 0.5 & 0.5 \end{bmatrix}
$$

$$
S,R\backslash G \quad \begin{array}{cc} \texttt{t} & \texttt{f} \end{array}
$$
$$
\begin{array}{c} \texttt{t,t} \\ \texttt{t,f} \\ \texttt{f,t} \\ \texttt{f,f} \end{array}
\begin{bmatrix} 0.99 & 0.01 \\ 0.9 & 0.1 \\ 0.9 & 0.1 \\ 0 & 1 \end{bmatrix}
$$

$$
\begin{array}{cc} \texttt{t} & \texttt{f} \end{array}
$$
$$
\begin{bmatrix} 3/5 & 2/5 \end{bmatrix}
$$

**sprinkler**

**cloudy**

**grass wet**

**rain**

$$
C\backslash R \quad \begin{array}{cc} \texttt{t} & \texttt{f} \end{array}
$$
$$
\begin{array}{c} \texttt{t} \\ \texttt{f} \end{array}
\begin{bmatrix} 0.8 & 0.2 \\ 0.2 & 0.8 \end{bmatrix}
$$

$$C \backslash S \quad \begin{matrix} t & f \end{matrix}$$
$$\begin{matrix} t \\ f \end{matrix} \begin{bmatrix} 0.1 & 0.9 \\ 0.5 & 0.5 \end{bmatrix}$$

$$S,R \backslash G \quad \begin{matrix} t & f \end{matrix}$$
$$\begin{matrix} t,t \\ t,f \\ f,t \\ f,f \end{matrix} \begin{bmatrix} 0.99 & 0.01 \\ 0.9 & 0.1 \\ 0.9 & 0.1 \\ 0 & 1 \end{bmatrix}$$

$$\begin{matrix} t & f \end{matrix}$$
$$\begin{bmatrix} 3/5 & 2/5 \end{bmatrix}$$

$$C \backslash R \quad \begin{matrix} t & f \end{matrix}$$
$$\begin{matrix} t \\ f \end{matrix} \begin{bmatrix} 0.8 & 0.2 \\ 0.2 & 0.8 \end{bmatrix}$$

**Observation 1.** Efficient representation of the joint distribution.

**Observation 2.** Efficient computation of the marginal distributions, as needed for exact inference. For example:

$$\Pr(G) = \sum_{R,S,C} \Pr(G, R, S, C) \qquad \textit{naive formula}$$

$$= \sum_{R,S,C} \Pr(G \mid S, R) \Pr(R \mid C) \Pr(S \mid C) \Pr(C) \qquad \textit{definition of the joint distribution}$$

$$= \sum_{C} \Pr(C) \sum_{S} \Pr(S \mid C) \left( \sum_{R} \Pr(R \mid C) \Pr(G \mid S, R) \right) \Bigg\} \quad \textit{two possible factorisations}$$

$$= \sum_{R} \sum_{S} \Pr(G \mid S, R) \sum_{C} \Pr(C) \Pr(S \mid C) \Pr(R \mid C)$$

**Algorithm 1** Variable Elimination Algorithm

**input:** $B$ a Baysesian Network, $Q \subset \mathcal{V}$ query variables, $\pi$ an ordering of $\mathcal{V} \setminus Q$

**heuristic**

**output:** $\mathbb{P}(Q)$

$\Phi \leftarrow$ CPTs of network

**for** $X \in \pi$ **do**

    **for** $x \in |X|$ **do**

        $f_x \leftarrow \prod_k f_k$, where $f_k \in S$ and mention variable $X = x$

    **end for**

    $f_X \leftarrow \sum_{x \in |X|} f_x$ and $\Phi \leftarrow f_X \cup \Phi \setminus \{f_k\}$

**end for**

$$\Pr(G) = \sum_{R,S,C} \Pr(G, R, S, C)$$

*naive formula*

$$= \sum_{R,S,C} \Pr(G \mid S, R) \Pr(R \mid C) \Pr(S \mid C) \Pr(C)$$
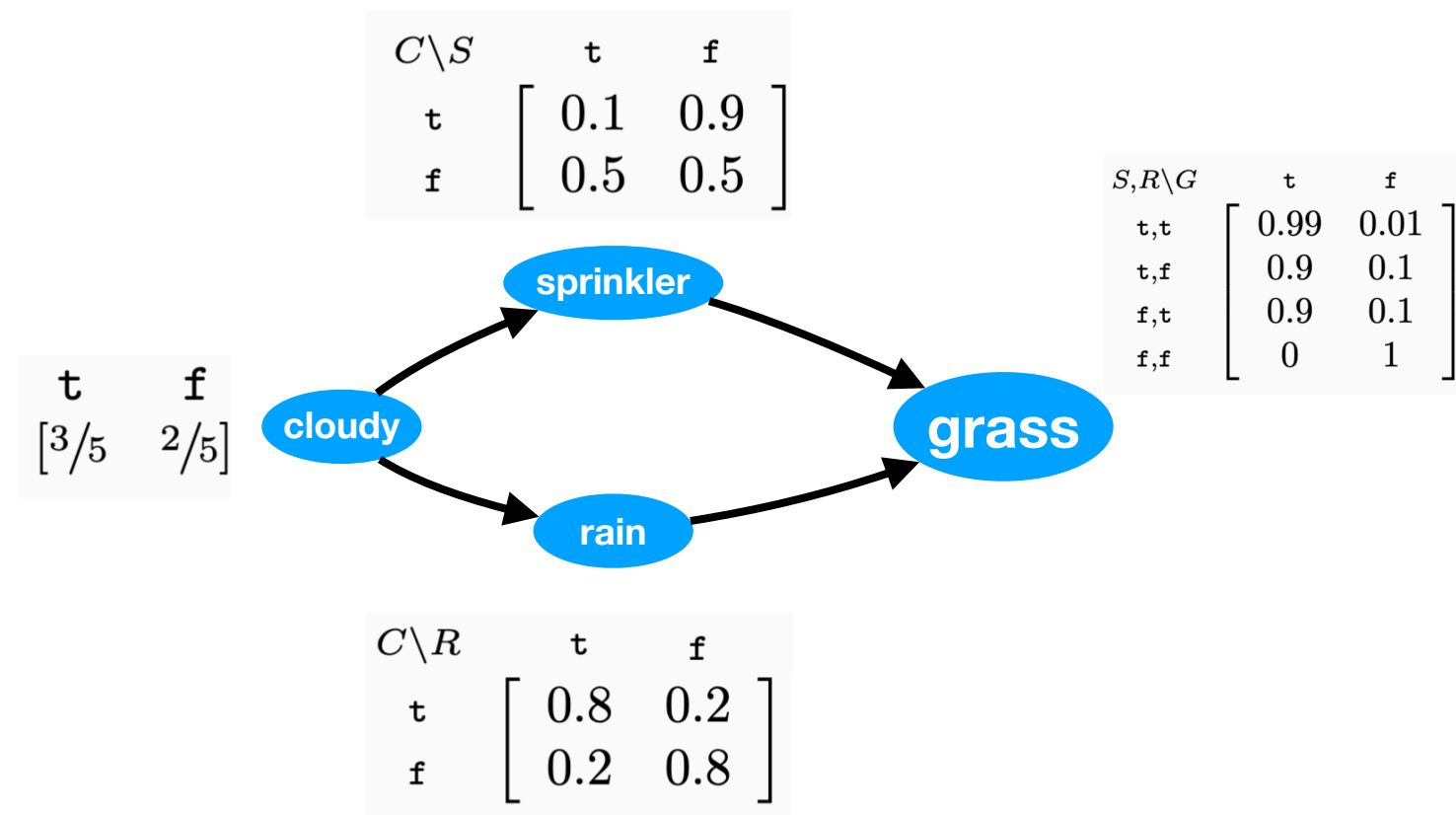
*definition of the joint distribution*

$$= \sum_C \Pr(C) \sum_S \Pr(S \mid C) \left( \sum_R \Pr(R \mid C) \Pr(G \mid S, R) \right)$$

$$= \sum_R \sum_S \Pr(G \mid S, R) \sum_C \Pr(C) \Pr(S \mid C) \Pr(R \mid C)$$

$\Big\}$ *two possible factorisations*

# Bayesian networks and PPLs

$$\begin{array}{c c c}
C\backslash S & \text{t} & \text{f} \\
\text{t} & \begin{bmatrix} 0.1 & 0.9 \\ 0.5 & 0.5 \end{bmatrix} \\
\text{f} &
\end{array}$$

$$\begin{array}{c c c}
S,R\backslash G & \text{t} & \text{f} \\
\text{t,t} & \begin{bmatrix} 0.99 & 0.01 \\ 0.9 & 0.1 \\ 0.9 & 0.1 \\ 0 & 1 \end{bmatrix} \\
\text{t,f} & \\
\text{f,t} & \\
\text{f,f} &
\end{array}$$

$$\begin{array}{c c}
\text{t} & \text{f} \\
[3/5 & 2/5]
\end{array}$$

sprinkler

cloudy

grass

rain

$$\begin{array}{c c c}
C\backslash R & \text{t} & \text{f} \\
\text{t} & \begin{bmatrix} 0.8 & 0.2 \\ 0.2 & 0.8 \end{bmatrix} \\
\text{f} &
\end{array}$$

To implement this in a PPL, we must choose a topological ordering of the DAG.

```
let C = bernoulli (3/5) in
let S = if C then (bernoulli 0.1) else (bernoulli 0.5) in
let R = if C then (bernoulli 0.8) else (bernoulli 0.2) in
let G = if (S ∧ T) then (bernoulli 0.99) else ... in
assume G;
return R
```

**Exercice**: all implementations have the same semantics.

# A categorical perspective

- We can define a Bayesian network as a DAG $G$ together with a structure preserving functor $\mathscr{C}_G \longrightarrow \mathbf{Stoch}$ where $\mathscr{C}_G$ is a category generated from $G$.

- *Informal idea, omitting technical details*: The objects of $\mathscr{C}_G$ are lists of vertices of $G$, and the morphisms are generated from:

  1. "Copy" and "discard" morphisms

  2. Morphisms $v_1, \ldots, v_n \longrightarrow v$ for each vertex $v$, where $\{v_i\}_{i \leq n} = \mathrm{pa}(v)$.

  3. The operations of symmetric monoidal categories.

# Conclusion

- *Sets and probability monads*: semantics for a simple call-by-value PPL. Corresponds to enumeration.

- This method will also extend to continuous probability with measurable spaces. Can also choose an appropriate monad depending on the inference algorithm.

> [Submitted on 9 Nov 2017]
> **Denotational validation of higher-order Bayesian inference**
> Adam Ścibior, Ohad Kammar, Matthijs Vákár, Sam Staton, Hongseok Yang, Yufei Cai, Klaus Ostermann, Sean K. Moss, Chris Heunen, Zoubin Ghahramani

- *Bayesian networks:* an efficient graphical representation of finite multidimensional distributions, with faster inference.

- We looked at the categorical foundations.

> [Submitted on 26 Jan 2013]
> **Causal Theories: A Categorical Perspective on Bayesian Networks**
> Brendan Fong

- Other ongoing research: efficient inference methods for discrete variables in probabilistic programming (cf. Pyro, Stan, Dice, Pluck).