

Probabilistic Languages and Semantics

Prakash Panangaden¹

¹School of Computer Science
McGill University

TIFR Mumbai 13th January 2015

Outline

- 1 Introduction
- 2 Conditional probability
- 3 Measures and measurable functions
- 4 Probabilistic relations
- 5 Probabilistic transition systems and probabilistic bisimulation
- 6 Semantics of a language with while loops

What am I trying to do?

- 1 Probability as logic: the central role of conditional probability.
- 2 Describe the key mathematical concepts behind modern probability: measure and integration.
- 3 Probabilistic systems and bisimulation (briefly)
- 4 Semantics of programming languages: part II.

What I am not trying to do

- Drown you in category theory.
- Discuss applications to *e.g.* Bayes nets.
- Discuss metrics or approximation theory.
- Deal with continuous time.
- Prove everything in detail (or anything at all!).

A puzzle

- Imagine a town where every birth is equally likely to give a boy or a girl. $\Pr(\text{boy}) = \Pr(\text{girl}) = \frac{1}{2}$.
- Each birth is an *independent* random event.
- There is a family with two children.
- One of them is a boy (not specified which one), what is the probability that the other one is a boy?
- Since the births are independent, the probability that the other child is a boy should be $\frac{1}{2}$. Right?
- Wrong! Before you are given the additional information that one child is a boy, there are 4 *equally likely* situations: bb, bg, gb, gg.
- The possibility gg is ruled out. So of the three equally likely scenarios: bb, bg, gb, only one has the other child being a boy. The correct answer is $\frac{1}{3}$.
- If I had said, “The *elder* child is a boy”, then the probability that the other child is a boy is indeed $\frac{1}{2}$.

The point of the puzzle

- Conditional probability is tricky!
- Conditional probability/expectation is *the* heart of probabilistic reasoning.
- Conditioning = revising probability (expectation) values in the presence of new information.
- Analogous to *inference* in ordinary logic.

- Sample space: set of possible outcomes; X .
- Event: subset of the sample space; $A, B \subset X$.
- Probability: $\Pr : X \rightarrow [0, 1]$, $\sum_{x \in X} \Pr(x) = 1$.
- Probability of an event A : $\Pr(A) = \sum_{x \in A} \Pr(x)$.
- A, B are independent: $\Pr(A \cap B) = \Pr(A) \cdot \Pr(B)$.
- Subprobability: $\sum_{x \in X} \Pr(x) \leq 1$.

Definition

If A and B are events, the *conditional probability of A given B* , written $\Pr(A \mid B)$, is defined by:

$$\Pr(A \mid B) = \Pr(A \cap B) / \Pr(B).$$

What happens if $\Pr(B) = 0$?

Bayes' Rule

$$\Pr(A | B) = \frac{\Pr(B | A) \cdot \Pr(A)}{\Pr(B)}.$$

- Trivial proof: calculate from the definition.
- Example: Two coins, one fake (two heads) one OK. One coin chosen with equal probability and then tossed to yield a H. What is the probability the coin was fake?
- Answer: $\frac{2}{3}$.
- Bayes' rule shows how to update the *prior* probability of A with the new information that the outcome was B : this gives the *posterior* probability of A given B .

Expectation values

- A *random variable* r is a real-valued function on X .
- The *expectation value* of r is

$$\mathbb{E}[r] = \sum_{x \in X} \Pr(x) r(x).$$

- The *conditional expectation value* of r given A is:

$$\mathbb{E}[r \mid A] = \sum_{x \in X} r(x) \Pr(\{x\} \mid A).$$

- Conditional probability is a special case of conditional expectation.

Kozen's correspondence

Classical logic	Generalization
Truth values $\{0, 1\}$	Probabilities $[0, 1]$
Predicate	Random variable
State	Distribution
The satisfaction relation \models	Integration \int

Model and reason about systems with *continuous* state spaces.

- Hybrid control systems; e.g. flight management systems.
- Telecommunication systems with spatial variation; e.g. mobile (cell) phones.
- Performance modelling.
- Continuous time systems.
- Probabilistic programming languages with recursion.

The Need for Measure Theory

- Basic fact: There are subsets of \mathbf{R} for which no sensible notion of size can be defined.
- More precisely, there is no translation-invariant measure defined on all the subsets of the reals.

- Countability is the key: basic analysis works well with countable summations.
- A σ -algebra Ω on a set X is a family of subsets with the following conditions:
 - 1 $\emptyset, X \in \Omega$
 - 2 $A \in \Omega \Rightarrow A^c \in \Omega$
 - 3 $\{A_i \in \Omega\}_{i \in \mathbb{N}} \Rightarrow \bigcup_i A_i \in \Omega$
- Closure under countable intersections is automatic.
- $A \in \Omega$ and $A \subset B$ or $B \subset A$ does **not** imply $B \in \Omega$.
- A set with a σ -algebra (X, Ω) is called a *measurable space*.

Properties of σ -algebras

- The collection of all subsets of X is always a σ -algebra.
- The intersection of *any* collection of σ -algebras is a σ -algebra.
- Thus, given *any* family \mathcal{F} of subsets of X there is a *least* σ -algebra containing them: $\sigma(\mathcal{F})$; the σ -algebra *generated* by \mathcal{F} .
- For most σ -algebras of interest a “generic” member is hard to describe. We try to work with simpler generating families.
- Because measurable sets are closed under complementation, the character of the subject is very different from topology; *e.g.* closure under limits.

Two Examples

- **R**: the real line. The open intervals do not form a σ -algebra. However, they generate one: the Borel algebra.
- Let \mathcal{A} be an “alphabet” of symbols (say finite) and consider \mathcal{A}^* : words over \mathcal{A} . Let \mathcal{A}^ω be finite and infinite words.
- Let $u \in \mathcal{A}^*$ and let $u \uparrow \stackrel{\text{def}}{=} \{v \in \mathcal{A}^\omega \mid u \leq v\}$.
- A “natural” σ -algebra on \mathcal{A}^ω is the σ -algebra generated by $\{u \uparrow \mid u \in \mathcal{A}^*\}$.

Measurable functions

- $f : (X, \Sigma) \rightarrow (Y, \Omega)$ is *measurable* if for every $B \in \Omega$, $f^{-1}(B) \in \Sigma$.
- Just like the definition of continuous in topology.
- Why is this the definition? Why backwards?
- $x \in f^{-1}(B)$ if and only if $f(x) \in B$.
- No such statement for the forward image.
- Exactly the same reason why we give the Hoare triple for the assignment statement in terms of preconditions.
- Older books (Halmos) give a more general definition that is not compositional.

- If $A \subset X$ is a measurable set, $\mathbf{1}_A(x) = 1$ if $x \in A$ and 0 otherwise is called the *indicator* or *characteristic* function of A and is measurable.
- The sum and product of real-valued measurable functions is measurable.
- If we take *finite* linear combinations of indicators we get *simple* functions: measurable functions with finite range.

Convergence properties

- If $\{f_i : \mathbf{R} \rightarrow \mathbf{R}\}_{i \in \mathbf{N}}$ converges pointwise to f and all the f_i are measurable then so is f .
- Stark difference with continuity.
- If $f : (X, \Sigma) \rightarrow (\mathbb{R}, \mathcal{B})$ is non-negative and measurable then there is a sequence of non-negative *simple* functions s_i such that $s_i \leq s_{i+1} \leq f$ and the s_i converge pointwise to f .
- The secret of integration.

- Want to define a “size” for measurable sets.
- A **measure** on (X, Σ) is a function $\mu : \Sigma \rightarrow [0, \infty]$ or $\mu : \Sigma \rightarrow [0, 1]$ (probability) such that
 - 1 $\mu(\emptyset) = 0$
 - 2 $A \cap B = \emptyset$ implies $\mu(A \cup B) = \mu(A) + \mu(B)$.
 - 3 $A \subset B$ implies $\mu(A) \leq \mu(B)$, follows.
 - 4 $\{A_i\}_{i \in \mathbb{N}} \subset \Sigma$ pairwise disjoint implies $\mu(\bigcup_i A_i) = \sum_i \mu(A_i)$;
subsumes (2).
 - 5 μ is *continuous* with respect to upward and downward chains of sets; follows from (4).
 - 6 Actually, (4) is the only axiom needed.

Examples of measures

- X countable, σ -algebra all subsets of X ; $c(A) =$ number of elements in A . Counting measure; not very useful.
- X any set, σ -algebra $\mathcal{P}(X)$, fix $x_0 \in X$ $\delta_{x_0}(A) = 1$ if $x_0 \in A$, 0 otherwise. Dirac delta “function.”
- $X = \mathbf{R}$, σ -algebra generated by the open (or closed) intervals, the Borel sets \mathcal{B} . $\lambda : \mathcal{B} \rightarrow \mathbf{R}^{\geq 0}$ defined as *the* measure which assigns to intervals their lengths.
- How do we know that such a measure is defined or that it is unique?
- Similarly, we can define measures on \mathbf{R}^n .

- We look for simple “well-structured” families of sets, *e.g.* intervals in \mathbf{R} and define “suitable” functions on them.
- Then we rely on extension theorems to obtain a unique measure on the generated σ -algebra.
- I will skip the “well-structured” conditions on the family of sets and the definition of “suitable” functions.
- A π -system is a family of sets closed under finite intersection.
- If two measures agree on a π -system then they agree on the generated σ -algebra.
- Fantastically useful!!

The Lebesgue integral

- Want to define $\int f d\mu$, where f is measurable and μ is a measure.
- Assume that f is everywhere non-negative and bounded and μ is a probability measure.
- If f is $\mathbf{1}_A$ then we *define* $\int \mathbf{1}_A d\mu = \mu(A)$.
- If f is $r \cdot \mathbf{1}_A$ then we *define* $\int f d\mu = r \cdot \mu(A)$.
- If $f = \sum_{i=1}^k r_i \mathbf{1}_{A_i}$ (simple function) then we define

$$\int f d\mu = \sum_{i=1}^k r_i \cdot \mu(A_i).$$

- Need to check that it does not matter how we write such an f as a simple function.
- There are some subtleties if sets can have infinite measure but these do not arise if we are dealing with probability measures and bounded measurable functions.

The Lebesgue integral

If f is non-negative and measurable and μ a probability measure we define

$$\int f d\mu = \sup \int s d\mu$$

where the *sup* is over all *simple* non-negative functions below f .

- One can define integrals of general functions by splitting them into positive and negative pieces.
- One can prove that the integral is linear and monotone.

The monotone convergence theorem

Let $\{f_n\}$ be a sequence of measurable functions on X such that (1) $\forall x \in X, 0 \leq f_1(x) \leq f_2(x) \leq \dots \leq f_n(x) \leq \dots \leq f(x)$ and (2) $\forall x \in X, \sup_n f_n(x) = f(x)$ then

$$\sup_n \int f_n d\mu = \int f d\mu.$$

- Should remind you of things in domain theory.
- The integral is continuous in an order-theoretic sense.

The monotone convergence mantra

- Want to prove $\int \mathcal{E}(f) d\mu = \int \mathcal{E}'(f) d\nu$.
- Prove it for the special case $f = \mathbf{1}_A$, usually easy.
- Then automatic for simple functions by linearity.
- Then automatic for non-negative bounded measurable functions by the monotone convergence theorem.
- Then clear for general bounded measurable functions.

Ordinary binary relations

- $R : A \rightarrow B$ is just $R \subseteq A \times B$
- Natural converse relation $R^\circ : B \rightarrow A$.
- Composition: $R_1 : A \rightarrow B, R_2 : B \rightarrow C$ then $R_1 \circ R_2 = \{(x, z) \mid \exists y \in B, xR_1y \text{ and } yR_2z\}$.
- Close relation with the powerset construction:
- $\hat{R} : A \rightarrow \mathcal{P}(B)$ is an equivalent description of R .

- A *Markov kernel* on a measurable space (S, Σ) is a function $h : S \times \Sigma \rightarrow [0, 1]$ with (a) $h(s, \cdot) : \Sigma \rightarrow [0, 1]$ a (sub)probability measure and (b) $h(\cdot, A) : X \rightarrow [0, 1]$ a measurable function.
- Though apparently asymmetric, these are the probabilistic analogues of binary relations
- and the uncountable generalization of a matrix.
- They describe transition probabilities in situations where a “point-to-point” approach does not make sense.
- Composition: k “after” h , $(k \circ h)(x, A) = \int k(x', A) dh(x, \cdot)$, where we are integrating the variable x' using the measure $h(x, \cdot)$.
- We construct these things using a major theorem (the Radon-Nikodym theorem).

- Want to define $R : (X, \Sigma) \rightarrow (Y, \Omega)$.
- Define a probabilistic relation R from X to Y to be a Markov kernel of type $R : X \times \Omega \rightarrow [0, 1]$ with the same measurability conditions.
- Given relations $R_1 : (X, \Sigma) \rightarrow (Y, \Omega)$ and $R_2 : (Y, \Omega) \rightarrow (Z, \Lambda)$ we define $R_2 \circ R_1$ ($R_1; R_2$) as
- $(R_2 \circ R_1)(x, C \in \Lambda) = \int R_2(y, C)R_1(x, \cdot)d.$
- Just like the formula for composing ordinary relations with integration for \exists .
- Converse is tricky and requires more machinery and more structure.

The category **SRel**

- Objects: measurable spaces (X, Σ_X)
- Morphisms: $h : (X, \Sigma_X) \rightarrow (Y, \Sigma_Y)$ are Markov kernels $h : X \times \Sigma_Y \rightarrow [0, 1]$.
- Composition: $h : X \rightarrow Y, k : Y \rightarrow Z$ then $\forall x \in X, C \in \Sigma_Z$,
 $(k \circ h)(x, C) = \int_Y k(y, C)h(x, dy)$.
- The identity morphisms: $id : X \rightarrow X$ is $\delta(x, A)$.
- Prove associativity of composition by using the monotone convergence mantra.
- It has countable coproducts; very useful for semantics.
- Unlike **Rel** this category is not self dual.

The Gíry Monad

- Define $\Pi : \mathbf{Mes} \rightarrow \mathbf{Mes}$ by $\Pi((X, \Sigma_X)) = \{\nu \mid \nu : \Sigma_X \rightarrow [0, 1]\}$ where ν is a *subprobability* measure on X .
- Actually, Gíry used probability measures; I made the small change to subprobability measures in order to adapt it to programming language semantics.
- But $\Pi(X)$ has to be a measurable space not just a set.
- For every $A \in \Sigma_X$ we define $\text{ev}_A : \Pi(X) \rightarrow [0, 1]$ by $\text{ev}_A(\nu) = \nu(A)$.
- We define the σ -algebra on $\Pi(X)$ to be the *least* σ -algebra making all the ev_A measurable.
- Given $f : X \rightarrow Y$ define $(\Pi(f)(\nu))(B \in \Sigma_Y) = \nu(f^{-1}(B))$.
- Need natural transformations: $\eta : I \rightarrow \Pi$ and $\mu : \Pi^2 \rightarrow \Pi$.
- $\eta_X(x) = \delta(x, \cdot)$
- $\mu_X(\Omega \in \Pi^2(X)) = \lambda B \in \Sigma_X. \int \text{ev}_B d\Omega_{\Pi(X)}$.

The Kleisli category of Π

- If $T : \mathcal{C} \rightarrow \mathcal{C}$ is a monad, then \mathcal{C}_T has the same objects as \mathcal{C} and the morphisms in \mathcal{C}_T from X to Y are morphisms in \mathcal{C} from X to TY .
- For the powerset monad we get morphisms $X \rightarrow \mathcal{P}(Y)$ which we recognize as just binary relations.
- Here we get $h : X \rightarrow \Pi(Y)$ or $h : X \rightarrow (\Sigma_Y \rightarrow [0, 1])$ or $h : X \times \Sigma_Y \rightarrow [0, 1]$.
- These are exactly the Markov kernels.

Labelled Markov processes

- Labelled Markov processes are probabilistic versions of labelled transition systems. Labelled transition systems where the final state is governed by a probability distribution - no other indeterminacy.
- All probabilistic data is *internal* - no probabilities associated with environment behaviour.
- We observe the interactions - not the internal states.
- **In general, the state space of a labelled Markov process may be a *continuum*.**

Formal Definition of LMPs

- An LMP is a tuple $(S, \Sigma, L, \forall \alpha \in L. \tau_\alpha)$ where $\tau_\alpha : S \times \Sigma \rightarrow [0, 1]$ is a *transition probability* function such that
- $\forall s : S. \lambda A : \Sigma. \tau_\alpha(s, A)$ is a subprobability measure and
- $\forall A : \Sigma. \lambda s : S. \tau_\alpha(s, A)$ is a measurable function.

- Let $\mathcal{S} = (S, \Sigma, \tau)$ be a labelled Markov process. An equivalence relation R on S is a **bisimulation** if whenever sRs' , with $s, s' \in S$, we have that for all $a \in \mathcal{A}$ and every R -closed measurable set $A \in \Sigma$, $\tau_a(s, A) = \tau_a(s', A)$.
- Two states are bisimilar if they are related by a bisimulation relation.
- Can be extended to bisimulation between two different LMPs.



$$\mathcal{L} ::= \top \mid \phi_1 \wedge \phi_2 \mid \langle a \rangle_q \phi$$

- We say $s \models \langle a \rangle_q \phi$ iff

$$\exists A \in \Sigma. (\forall s' \in A. s' \models \phi) \wedge (\tau_a(s, A) > q).$$

- Two systems are bisimilar iff they obey the same formulas of \mathcal{L} .
[DEP 1998 LICS, I and C 2002]

Kozen's Language

$$S ::= x_i := f(\vec{x}) \mid S_1; S_2 \mid \text{if } \mathbf{B} \text{ then } S_1 \text{ else } S_2 \mid \text{while } \mathbf{B} \text{ do } S.$$

- There are a fixed set of variables \vec{x} taking values in a measurable space (X, Σ_X) .
- f is a measurable function.
- B is a measurable subset.

- State transformer semantics: distribution (measure) transformer semantics.
- Meaning of statements: Markov kernels *i.e.* **SRel** morphisms.
- The only subtle part: how to give fixed-point semantics to the while loop?

Partially additive structure

- Back to **SRel** structure.
- Can we “add” **SRel** morphisms?
- Not always, the sum may exceed 1, but we can define *summable families* which may even be countably infinite.
- The homsets of **SRel** form *partially additive monoids*.
- The sums can be rearranged at will (partition-associativity).
- Limit property: If F is a countable family in which every *finite* subfamily is summable then F is summable.
- In the category **SRel**, the sums interact properly with composition.
- If $\{f_i \mid i \in \mathbb{N}\}$ is a countable set of morphisms from X to Y and there is a morphism $f : X \rightarrow (Y + Y + \dots)$ such that when projected onto the X 's we get the f_i , then the family is summable.

Arbib and Manes

Given a partially additive category \mathcal{C} and $f : X \rightarrow X + Y$ we can find a unique pair $f_1 : X \rightarrow X$ and $f_2 : X \rightarrow Y$ such that $f = \iota_1 \circ f_1 + \iota_2 \circ f_2$. Furthermore, there is a morphism $f^* : X \rightarrow Y$ given by

$$f^* = \sum_{n=0}^{\infty} f_2 \circ f_1^n.$$

The theorem says that the family $f_2 \circ f_1^n$ is summable. It is the *iterate* of f .

Semantics of Kozen's Language I

- Statements are **SRel** morphisms of type $(X^n, \Sigma^n) \rightarrow (X^n, \Sigma^n)$.
- **Assignment:** $x := f(\vec{x})$

$$\llbracket x_i := f(\vec{x}) \rrbracket(\vec{x}, \vec{A}) = \delta(x_1, A_1) \dots \delta(x_{i-1}, A_{i-1}) \delta(f(\vec{x}), A_i) \delta(x_{i+1}, A_{i+1}) \dots$$

- **Sequential Composition:** $S_1; S_2$

$$\llbracket S_1; S_2 \rrbracket = \llbracket S_2 \rrbracket \circ \llbracket S_1 \rrbracket$$

where the composition on the right hand side is the composition in **SRel**.

- **Conditionals:** *if* **B** *then* S_1 *else* S_2

$$\llbracket \text{if } \mathbf{B} \text{ then } S_1 \text{ else } S_2 \rrbracket(\vec{x}, \vec{A}) = \delta(\vec{x}, \mathbf{B}) \llbracket S_1 \rrbracket(\vec{x}, \vec{A}) + \delta(\vec{x}, \mathbf{B}^c) \llbracket S_2 \rrbracket(\vec{x}, \vec{A})$$

While Loops: *while* **B** *do* *S*

$$\llbracket \textit{while } \mathbf{B} \textit{ do } S \rrbracket = h^*$$

where we are using the $*$ in **SRel** and the morphism

$$h : (X^n, \Sigma^n) \rightarrow (X^n, \Sigma^n) + (X^n, \Sigma^n)$$

is given by

$$h(\vec{x}, \vec{A}_1 \uplus \vec{A}_2) = \delta(\vec{x}, \mathbf{B}) \llbracket S \rrbracket(\vec{x}, \vec{A}_1) + \delta(\vec{x}, \mathbf{B}^c) \delta(\vec{x}, \vec{A}_2).$$

Weakest precondition semantics

- We can construct a category of probabilistic predicate transformers: **SPT**.
- Objects are measurable spaces.
- Given (X, Σ_X) we can construct the (Banach) space of bounded measurable functions on X (the “predicates”) $\mathcal{F}(X)$.
- A morphism $X \rightarrow Y$ in **SPT** is a bounded (continuous) linear map from $\mathcal{F}(X)$ to $\mathcal{F}(Y)$.



$$\mathbf{SPT} \simeq \mathbf{SRel}^{op}.$$

- This gives us the structure needed for a **wp** semantics.