

Systems that Model

their Environments

(and Systems that Plan)

Nathaniel Virgo

Earth-Life Science Institute (ELSI),

Tokyo, Japan

funded by the John Templeton Foundation

Introduction

Philosophically motivated:

I research origins of life, and wanted
to know: what is an agent?

i.e. what makes a system try to do something?

As part of that, I'm interested in what it means
for one system to model another, which will
be most of the talk.

In particular, where does the model live in relation
to the system?

Outline

- Background: category-theoretic approach to conjugate priors
- Extension to Bayesian filtering
- How does this arise more abstractly?
 - strongly representable Markov categories
- very brief sketches of approaches to "agents that plan"

Conjugate Priors

reference:

Bart Jacobs 2020

A Channel-Based Perspective
on Conjugate Priors

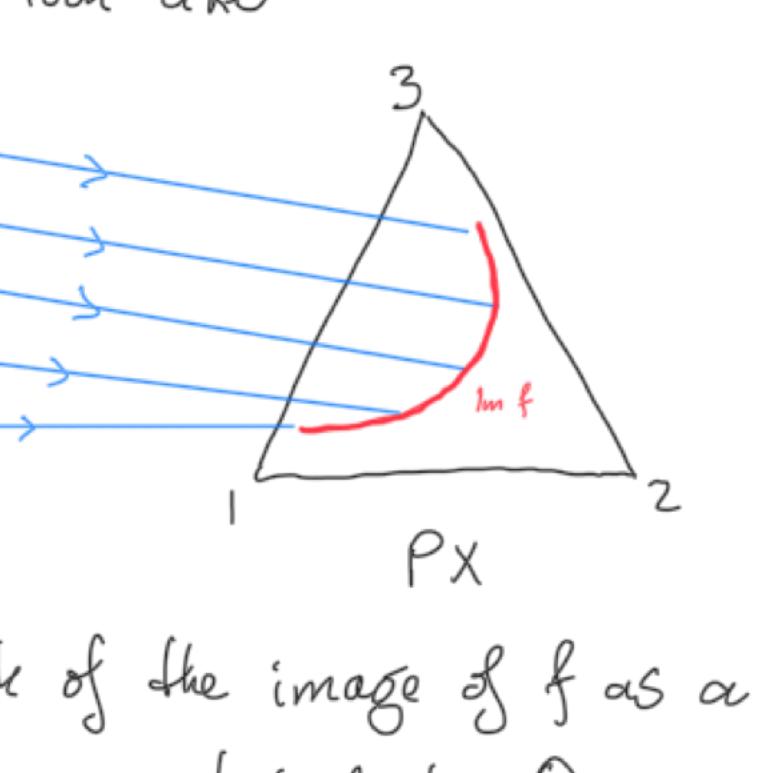
Mathematical Structures in Computer Science 30, p. 44-67

"The main result of this paper... is mathematically trivial." this applies to most of the talk.
But it is not entirely trivial to see that this result is trivial."

A useful perspective came from Jacobs' paper on conjugate priors.

What is a conjugate prior?

Start with a statistical model, which is just a (measurable) function



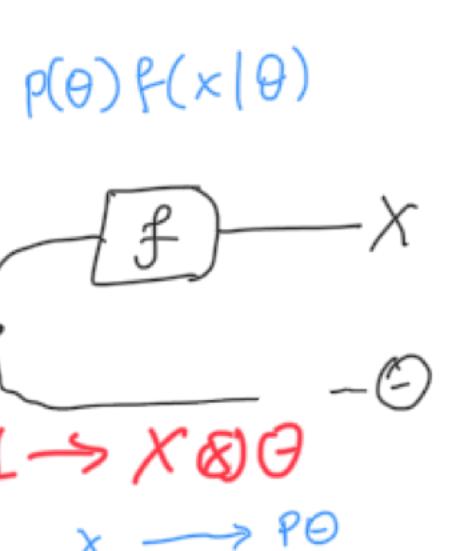
If we work in the Kleisli category of a distribution monad (or any other Markov category) then this is just a morphism

$$\Theta \xrightarrow{f} X$$

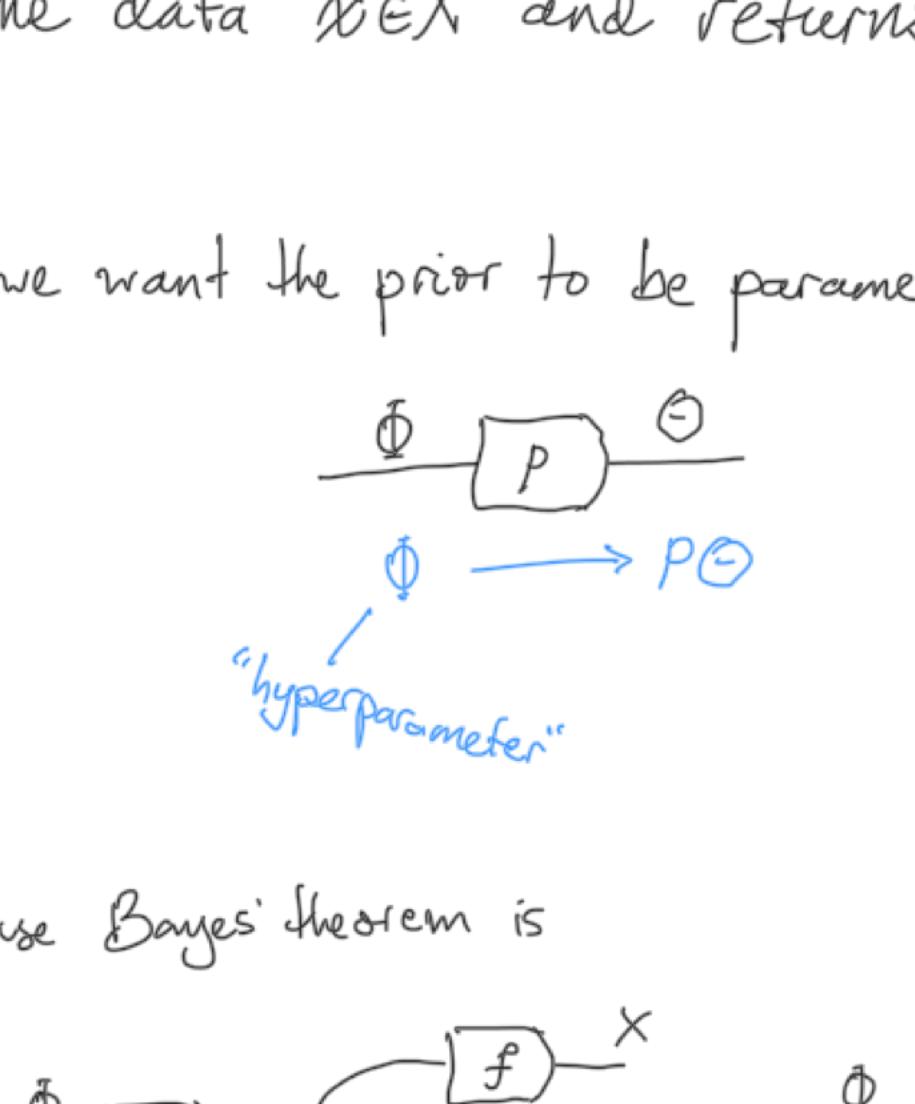
think of this as a family of distributions

e.g. if $X = \{1, 2, 3\}$ and $\Theta = [0, 1]$

then



and f might look like



so we can think of the image of f as a subset of PX , parametrised by Θ .

Typical inference problem:

we have a statistical model $\Theta \xrightarrow{f} X$

we have some prior $p \in P\Theta$ (i.e. the "true" $\theta \in \Theta$ is unknown)

and some data $x \in X$, which is a sample

from $f(\theta)$ the unknown "true" θ

and we want to update to a posterior $\in P\Theta$.

In string diagrams we write Bayes' theorem as

$$P(\theta) p(x|\theta) = P(\theta) \left(\sum_{\theta} p(\theta) f(x|\theta) \right) p(\theta|x)$$

note on how Markov kernels compose; added during the talk:
 $X \xrightarrow{f} Y \xrightarrow{g} Z$
 $\sum_g f(g|x) g(z|y)$

Here, $\xrightarrow{\theta} f_p^\dagger \Theta$ is the "Bayesian inverse" — it takes in the data $x \in X$ and returns the posterior $\in P\Theta$.

Sometimes we want the prior to be parametrised as well.

$$\begin{array}{c} \Phi \\ \square \end{array} \xrightarrow{P} \begin{array}{c} \Theta \\ \square \end{array} \quad \begin{array}{c} \Phi \\ \square \end{array} \xrightarrow{P\Theta} \begin{array}{c} X \\ \square \end{array}$$

$\frac{df}{d\theta}$
 $\xrightarrow{\theta} f_p^\dagger \Theta$

In this case Bayes' theorem is

$$\begin{array}{c} \Phi \\ \square \end{array} \xrightarrow{P} \begin{array}{c} \Theta \\ \square \end{array} = \begin{array}{c} \Phi \\ \square \end{array} \xrightarrow{P} \begin{array}{c} \Theta \\ \square \end{array} \xrightarrow{f} \begin{array}{c} X \\ \square \end{array} \xrightarrow{f_p^\dagger} \begin{array}{c} \Theta \\ \square \end{array}$$

This time

$$\begin{array}{c} X \\ \square \end{array} \xrightarrow{\frac{d}{d\theta}} \begin{array}{c} \Theta \\ \square \end{array}$$

$$X \times \Phi \longrightarrow P\Theta$$

takes the data $x \in X$ and $\phi \in \Phi$ parametrising the prior and returns the posterior in $P\Theta$.

But if we are lucky then it might be that the posterior always lies in the family of distributions defined by P . i.e. there exists $\phi' \in \Phi$ such that the posterior is $p(\phi')$.

This means

$$\begin{array}{c} \Theta \\ \square \end{array} \xrightarrow{P} \begin{array}{c} X \\ \square \end{array}$$

a Bayesian inference interpretation is

a model $\Theta \xrightarrow{f} X$ and an

"interpretation map" $\frac{d}{d\theta} \xrightarrow{\psi} \Theta$

satisfying the same equation as above,

$$\begin{array}{c} \Phi \\ \square \end{array} \xrightarrow{P} \begin{array}{c} \Theta \\ \square \end{array} = \begin{array}{c} \Phi \\ \square \end{array} \xrightarrow{P} \begin{array}{c} \Theta \\ \square \end{array} \xrightarrow{f} \begin{array}{c} X \\ \square \end{array} \xrightarrow{f_p^\dagger} \begin{array}{c} \Theta \\ \square \end{array} \xrightarrow{\frac{d}{d\theta}} \begin{array}{c} \Theta \\ \square \end{array}$$

Given a model f , a conjugate prior is (u, p) satisfying this equation.

(this string diagram equation is one of the main points of Jacobs' paper.)

The point: if you want to implement Bayes on a computer, and if you can find a conjugate prior, then the only thing you need to implement is the function u .

In practice this is often a very simple function.

So we have a situation where the actual physical system is just receiving inputs and updating its state

But we interpret it as updating a Bayesian belief via the map

$$\Theta \xrightarrow{P} X$$

$\Phi \longrightarrow PX$

physical state of machine (exists)

Bayesian beliefs of machine ("doesn't really exist")

Similar ideas exist in neuroscience.

(the "Bayesian brain" hypothesis)

So we switch to a different perspective:

Given a "machine" $\frac{x}{\Phi} \xrightarrow{u} \Phi$,

a Bayesian inference interpretation is

a model $\Theta \xrightarrow{f} X$ and an

"interpretation map" $\frac{d}{d\theta} \xrightarrow{\psi} \Theta$

satisfying the same equation as above,

$$\begin{array}{c} \Phi \\ \square \end{array} \xrightarrow{P} \begin{array}{c} \Theta \\ \square \end{array} = \begin{array}{c} \Phi \\ \square \end{array} \xrightarrow{P} \begin{array}{c} \Theta \\ \square \end{array} \xrightarrow{f} \begin{array}{c} X \\ \square \end{array} \xrightarrow{u} \begin{array}{c} \Phi \\ \square \end{array} \xrightarrow{\frac{d}{d\theta}} \begin{array}{c} \Theta \\ \square \end{array}$$

This says that the expected posterior beliefs have to match the prior beliefs, where the expectation is over the prior beliefs. (cf martingales)

(side point: it matters that the machine is deterministic — this equation isn't strong enough otherwise.)

Note from after the talk: I was a bit confused here — if we take a marginal of this equation, by post-composing both sides with $\frac{d}{d\theta}$, then we get

$$\frac{d}{d\theta} \circ = \frac{d}{d\theta} \frac{d}{d\theta} \circ + \frac{d}{d\theta} \frac{d}{d\theta} \circ$$

This says the expectation of my posterior should equal my prior. (when the expectation is over my beliefs)

The full, un-marginalised equation is a bit more subtle, and says something more along the lines that

given an observation x , my posterior belief should equal my current belief conditioned on X taking that value. (up to an "almost surely" since conditions are only almost surely defined in the first place.)

In general, a machine may have many interpretations.

In particular, if $\frac{\Theta}{\Phi} \xrightarrow{u} \Phi$ is a valid interpretation map

then so is

$$\frac{\Phi}{\Phi} \xrightarrow{u} \frac{\Theta}{\Phi} \xrightarrow{\frac{d}{d\theta}} \frac{\Theta}{\Phi}, \text{ for any } x \in X$$

(assuming $\frac{\Theta}{\Phi} \xrightarrow{\frac{d}{d\theta}} \frac{\Theta}{\Phi}$ has full support).

With the right definitions of morphisms, interpretations form a functor $Machines^{op} \rightarrow \text{Cat}$.

Taking the Grothendieck construction gives us a category of "Bayesian reasoners" = machines equipped with interpretations.

Before returning to this, I want to talk about the extension to Bayesian filtering.

Bayesian Filtering Interpretations

work with Martin Biehl, Simon McGregor

End goal: design (or reason about) agents that model their environments in order to achieve a goal.

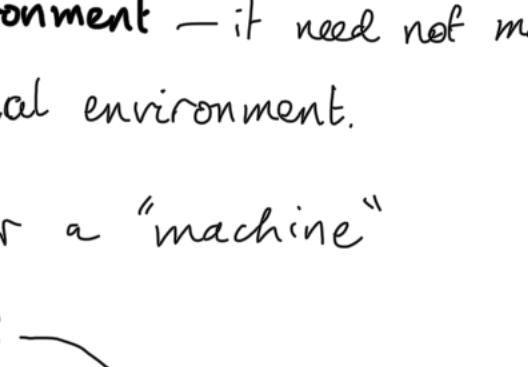
→ keep track of environment state
given data

→ reason about how environment will respond to actions.

(N.B. We're assuming we already have a model
—this is POMDPs rather than reinforcement learning.)

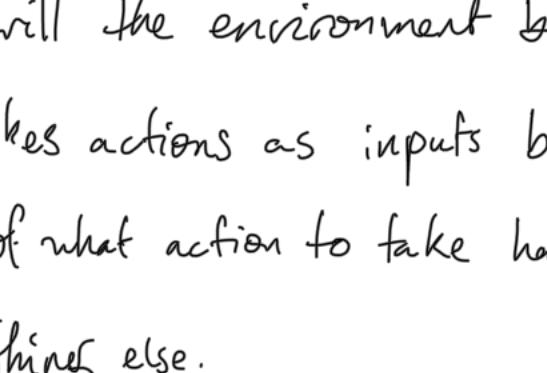
The first point can be achieved by "Bayesian filtering" but I will present a generalised version of filtering that also includes actions.

Suppose we have an environment model that looks like this:



this is meant to be an agent's model of its environment — it need not match the actual environment.

Then consider a "machine"



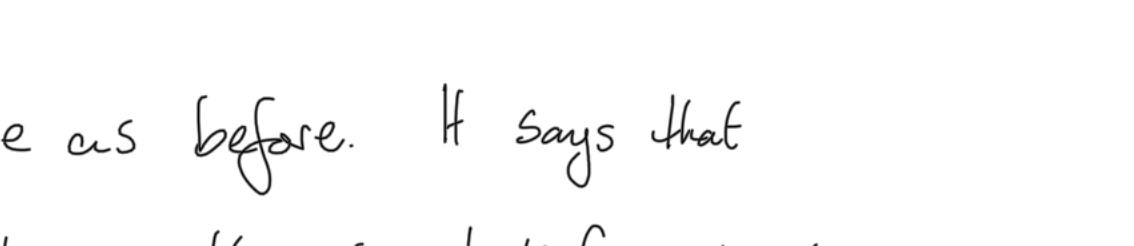
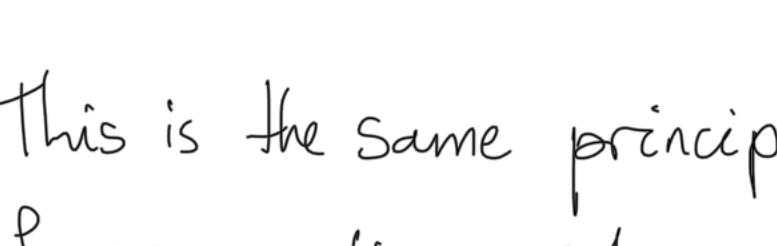
the machine's job is to answer "if I take action $a \in A$ and receive sensor value $s \in S$, then what state will the environment be in?"

H takes actions as inputs because the choice of what action to take has to be taken by something else.

Finally consider an interpretation map



that satisfies



We call (Ψ, Π) a Bayesian filtering interpretation of u .

This is the same principle as before. It says that for every action $a \in A$, the agent's prior beliefs about the next hidden state have to match its posterior beliefs about what will then be the current hidden state, if the data $s \in S$ is sampled from a distribution that matches its current beliefs.

Returning to the more abstract picture

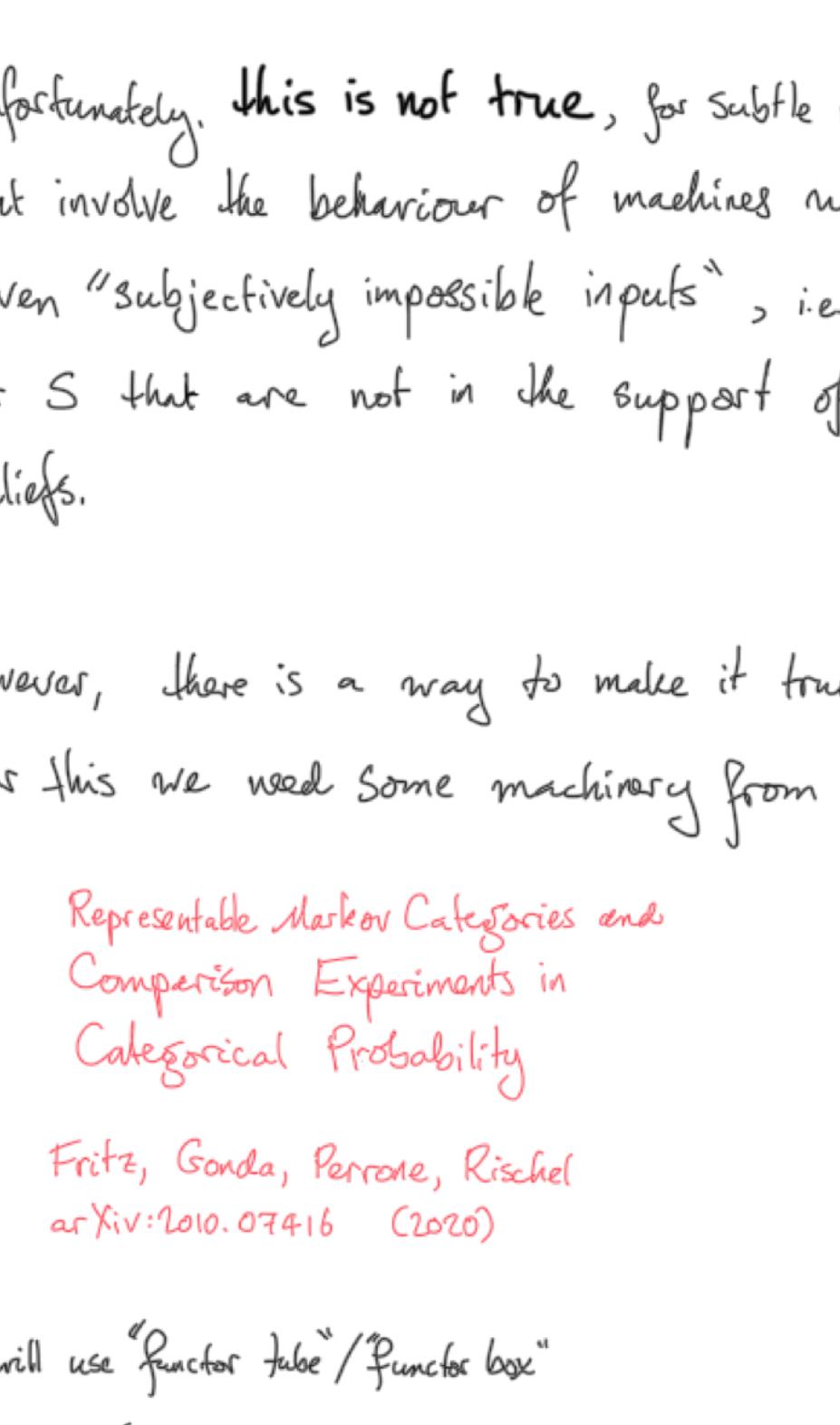
(stuff I figured out last weekend.)

Note: this section is fairly involved, but in the end the only punchline is that Bayesian filtering "pops out" of the machinery of "strongly representable Markov categories". ("...it's not entirely trivial to see that this result is trivial.")

As with conjugate priors, Bayesian filtering interpretations can be made into a functor

$$\text{Machines}^{\text{op}} \rightarrow \text{Cat}$$

that sends each machine to its category of interpretations, and acts on morphisms by "pulling back" the interpretation



We can also fix it when doing this, so that each machine is interpreted as reasoning about the same environment. (If such an interpretation exists.) In that case there is no useful notion of morphism of interpretations, and we just get a functor

$$\text{Machines}^{\text{op}} \rightarrow \text{Set}.$$

Taking the category of elements yields a category of "filtering reasoners" = machines equipped with filtering interpretations, with model given by $\mathbb{P}H$.

What we would like is for this category to have a terminal object. The idea is that the state space of the terminal object would be $\mathbb{P}H$, with the trivial interpretation

$$\mathbb{P}H \xrightarrow{\text{id}_{\mathbb{P}H}} \mathbb{P}H.$$

Then for any other object there would be a unique map into the terminal object given by its interpretation map, seen as a map $\psi: M \rightarrow \mathbb{P}H$.

We interpret it as taking a distribution as input and returning a sample from it, stochastically.

We need the definition of a "strongly representable Markov category" — to define it we need a couple of other definitions.

Definition (Fritz, Gonda, Perrone, Rischel) A synthetic approach to Markov kernels [–] Tobias Fritz, Advances in Mathematics 390 (2020)

A markov category has conditionals if

for every morphism $A \xrightarrow{f} X$,

there exists a morphism $X \xrightarrow{c} Y$

such that

$$A \xrightarrow{f} X = A \xrightarrow{f} Y \xrightarrow{c} X.$$

Conditionals are generally not unique.

This says that if we delete Y then we can recover the joint distribution by multiplying by the conditional. In the discrete case,

$$f(x, y | a) = f(x | a) c(y | x, a)$$

a generalisation of $p(x, y) = p(x) p(y | x)$.

A strongly representable Markov category is one where this establishes a bijection

$$\mathcal{C}_{\text{det}}(A, X \otimes PY) \cong \mathcal{C}(A, X \otimes Y).$$

$$\begin{array}{ccc} A & \xrightarrow{F} & X \otimes PY \\ \downarrow & \lrcorner & \downarrow \\ A & \xrightarrow{F} & X \end{array}$$

Known example: Borel Stoch is strongly representable (Fritz et al.)

Back to machines and interpretations.

Consider an environment model with trivial input, (I haven't worked out the general case yet, but it's probably straightforward.)

$$H \xrightarrow{S} H.$$

$$H \xrightarrow{f} X \xrightarrow{S} H.$$

is deterministic given X if for any choice of conditional,

$$A \xrightarrow{f} X = A \xrightarrow{f} Y \xrightarrow{c} X.$$

such that

$$A \xrightarrow{f} X = A \xrightarrow{f} Y \xrightarrow{c} X.$$

Conditionals are generally not unique.

This says that if we delete Y then we can recover the joint distribution by multiplying

by the conditional. In the discrete case,

$$f(x, y | a) = f(x | a) c(y | x, a)$$

a generalisation of $p(x, y) = p(x) p(y | x)$.

Now pre-compose this with the sampling map

$$H \xrightarrow{S} H.$$

$$A \xrightarrow{f} X \xrightarrow{S} H.$$

Then we can compose with the sampling map

to get

$$A \xrightarrow{f} X \xrightarrow{S} H.$$

of type $A \xrightarrow{f} X \otimes PY$.

A strongly representable Markov category is one where this establishes a bijection

$$\mathcal{C}_{\text{det}}(A, X \otimes PY) \cong \mathcal{C}(A, X \otimes Y).$$

$$\begin{array}{ccc} A & \xrightarrow{F} & X \otimes PY \\ \downarrow & \lrcorner & \downarrow \\ A & \xrightarrow{F} & X \end{array}$$

Known example: Borel Stoch is strongly representable (Fritz et al.)

Back to machines and interpretations.

Consider an environment model with trivial input, (I haven't worked out the general case yet, but it's probably straightforward.)

$$H \xrightarrow{S} H.$$

$$H \xrightarrow{f} X \xrightarrow{S} H.$$

is deterministic given X if for any choice of conditional,

$$A \xrightarrow{f} X = A \xrightarrow{f} Y \xrightarrow{c} X.$$

such that

$$A \xrightarrow{f} X = A \xrightarrow{f} Y \xrightarrow{c} X.$$

Conditionals are generally not unique.

This says that if we delete Y then we can recover the joint distribution by multiplying

by the conditional. In the discrete case,

$$f(x, y | a) = f(x | a) c(y | x, a)$$

a generalisation of $p(x, y) = p(x) p(y | x)$.

Now pre-compose this with the sampling map

$$H \xrightarrow{S} H.$$

$$A \xrightarrow{f} X \xrightarrow{S} H.$$

Then we can use the bijection above to get

the

where c is both almost-surely deterministic in

the sense above, and only defined up to almost

sure equality.

c here is not arbitrary — it is performing

Bayesian updates. Indeed we have

the

which looks a lot like the Bayesian filtering equation above,

the

(This is for the special case

of trivial actions)

This suggests a different category of

machines-with-interpretations, where

- machines' behaviour is only defined up to almost-sure equivalence

- machines are only required to be almost-surely deterministic

In this category the object

is indeed the terminal object, i.e. it is in a sense

the universal model of H .

In the end, the main point here is that Bayesian

filtering "pops out" of strongly representable

Markov categories in a simple way, and the same

machinery might be useful for other similar

questions about Bayesian modelling.

Systems that Plan

I only talked about systems that model their environment, but an obvious next step is systems that use their model to plan actions. Here I will only sketch approaches that we are taking to that.

As with inference, we're interested not in "what does it mean for a system to actually have a plan/goal", but "what does it mean to consistently interpret a system as planning/reaching for a goal."

1. Interpreting systems as solutions to POMDPs

work with Martin Biehl

Idea: we can interpret a system as trying to solve a POMDP task if it is an optimal solution to that task (details worked out in a similar framework to the one I presented)

2. Coalgebraic approach

ongoing work with Simon McGregor, Timo

uses a technically different framework from today's one
main result: if a system is an optimal solution to a certain kind of planning task, then it can be interpreted as doing Bayesian filtering.

(in a sense that looks like what I defined if you squint a bit.)

3 Systems theory

ongoing work with Matteo Cappuccini

→ forget about internal models entirely and figure out what it means for a system to optimally regulate its actual environment, in a very general sense using categorical systems theory.

Conclusions & thanks