

# Approximate Inference via Fibrations of Statistical Games

Toby St Clare Smithe

VERSES Research  
&  
Topos Institute



**VERSES**

2023 August 4

## Natural adaptive systems: the story so far

I'm interested in natural adaptive systems: organisms that perceive and act.

We can think of perception as a process of Bayesian inference:

- an organism has sensors, detecting data in  $Y$ ;
- it maintains a belief  $\beta : I \leftrightarrow X$  about the world  $X$ ,
- and about how the world causes its sense data  $c : X \leftrightarrow Y$ .
- Given these, it updates its beliefs via  $c_{\beta}^{\#} : Y \leftrightarrow X$ .

Together,  $c : X \leftrightarrow Y$  and  $c^{\#} : \mathcal{C}(I, X) \rightarrow \mathcal{C}(Y, X)$  are called a **Bayesian lens**  $(c, c^{\#}) : X \leftrightarrow Y$ . (It's a *bidirectional* process!)

- $\mathcal{C}$  is our ambient (concrete) category of 'channels'  $X \leftrightarrow Y$ .
- Arrows  $A \rightarrow B$  are functions.
- Lens composition:  $(d, d^{\#}) \circ (c, c^{\#}) = (d \bullet c, c^{\#} \circ d_c^{\#})$ .

## Exact vs approximate Bayesian lenses

Bayesian lenses: attach ‘state-dependent’ inversions  $c^\sharp$  to channels  $c$ .

This gives a fibration  $\pi_{\text{Lens}} : \mathbf{BayesLens} \rightarrow \mathcal{C}$  by  $(c, c^\sharp) \mapsto c$ .

Bayes’ rule (almost surely [1]) gives a section,  $c \mapsto (c, c^\dagger)$ .

(BUCO: “Bayesian updates compose optically” [2]).

These lenses ‘exactly’ satisfy Bayes’ law.

But there are lots of other Bayesian lenses!

We can think of non-exact lenses as *approximate* inversions:

useful, because exact inference is computationally hard.

This entails a need to measure performance. Enter statistical games!

# Overview of this talk

- 1 Introduction
- 2 Statistical games: the basics
- 3 Relative entropy and copy-composition
- 4 Loss models
- 5 Monoidal structures
- 6 Concluding remarks

## Another motivation: *local* inference

Biological systems are composed of cells.

Moving information around is expensive!

(This is where most energy is spent in machine learning.)

So biology wants systems that can be optimized 'locally'.

Hence: compositional inference, plus compositional 'measurement'.

These measurements will be formalized by **loss functions**.

## Statistical games: attach losses to lenses

A **statistical game**  $X \rightarrow Y$  is a triple  $(c, c^\sharp, L^c)$ ,  
of which  $(c, c^\sharp)$  is a Bayesian lens  $X \leftrightarrow Y$ ,  
and  $L^c$  is an accordingly-typed loss function.

(‘Game’ due to structural similarity to compositional game theory [3, 4].)

This gives a 2<sup>nd</sup> fibration  $\pi_{\text{Loss}} : \mathbf{StatGame} \rightarrow \mathbf{BayesLens}$ .

(We’ll see how loss functions compose later.)

Sections of this fibration capture some important quantities in statistics:

- relative entropy (a.k.a. Kullback-Leibler divergence);
- (log) likelihood;
- free energy (a.k.a. ‘ELBO’).

Let’s see how this works ...

# The relative entropy and its chain rule

Relative entropy measures ‘divergence’ between distributions:

$$D_Y : \mathcal{C}(I, Y) \times \mathcal{C}(I, Y) \rightarrow \mathbb{R}_+, \text{ with } D_Y(\alpha, \alpha') = 0 \iff \alpha = \alpha'.$$

It can be indexed by channels: given parallel  $c, c' : X \rightarrow Y$ , define

$$D_{c,c'} : X \rightarrow \mathbb{R}_+ \text{ by } D_{c,c'}(x) := D_Y(c(x), c'(x)).$$

It satisfies a “chain rule” that we can *almost* write as

$$\begin{aligned} D_{d \bullet c, d' \bullet c'}(x) &= D_{c,c'}(x) + (c^* D_{d,d'})(x) \\ &= D_{c,c'}(x) + \mathbb{E}_{y \sim c(x)} [D_{d,d'}(y)] \end{aligned}$$

... but not quite.

## Problem: a failure of compositionality

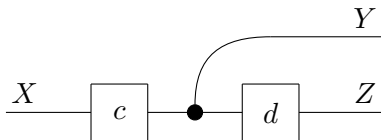
Relative entropy is the expected difference in log probability:

$$D_Y(\alpha, \alpha') = \mathbb{E}_{y \sim \alpha} [\log p_\alpha(y) - \log p_{\alpha'}(y)]$$

But  $p_{d \bullet c}(z|x) = \mathbb{E}_{y \sim c(x)} [p_d(z|y) p_c(y|x)]$ .

And log doesn't commute with expectations!

On the other hand, if we copy the intermediate variable ...





## Solution: copy-composition

... we get a 'copy-composite' channel  $d \bullet^2 c : X \rightarrow Y \otimes Z$   
with density  $p_{d \bullet^2 c}(y, z|x) = p_d(z|y) p_c(y|x)$ .

The expectation has gone away, and so our chain rule is satisfied:

$$D_{d \bullet^2 c, d' \bullet^2 c'}(x) = D_{c, c'}(x) + (c^* D_{d, d'})(x)$$

Challenge: adjust  $\mathcal{C}$  so that composition means *copy-composition*.  
Composites need to carry their intermediaries.

## A bicategory of copy-composite processes

The trouble is that copy-composition is not strictly unital:

$\text{id}_Y \bullet^2 c$  has type  $X \rightarrow Y \otimes Y$ , not  $X \rightarrow Y$ !

But we can think of  $d \bullet^2 c$  as a  $Y$ -coparameterized morphism  $X \xrightarrow[Y]{\bullet} Z$ , and adjust the **Copara** construction [5] accordingly.

This yields a bicategory, **Copara**<sub>2</sub>( $\mathcal{C}$ ).

Horizontal composition is copy-composition.

2-cells are “changes of coparameter”.

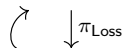
The structure morphisms deal with the book-keeping, introducing and deleting copies where necessary.

\* Still, I'm not sure this is the best way; I've since discovered others! (See bonus slide...)

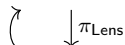
# Fibrations of copy-composite lenses and games

The picture starts to look like this:

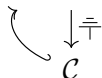
**StatGame**



**BayesLens**



**Copara<sub>2</sub>(C)**



- $\tilde{\tau} : \mathbf{Copara}_2(\mathcal{C}) \rightarrow \mathcal{C}$  discards coparameters.
- There is a canonical (lax) inclusion  $\mathcal{C} \hookrightarrow \mathbf{Copara}_2(\mathcal{C})$ , yielding a section of  $\tilde{\tau}$ .
- Sections of  $\pi_{\text{Lens}}$  are **inference systems**, e.g.  $\dagger$  (since we have a “coparameterized BUCO” result.)
- Sections of  $\pi_{\text{Loss}}$  are **loss models**, e.g. relative entropy.

So what *are* loss functions, in general? And how do they compose?

# Loss functions, externally and internally

'Externally': a loss function is a map  $L^c : \mathcal{C}(I, X) \times Y \rightarrow \mathbb{R}_+$ , associated to a lens  $(c, c')$ .

If  $\mathcal{C}$  has *bilinear effects*<sup>\*</sup>, we can internalize this to a state-dependent effect  $L^c : Y \overset{X}{\dashrightarrow} I$  — i.e., a function  $\mathcal{C}(I, X) \rightarrow \mathcal{C}(Y, I)$ .

Example categories where this works:

- the Kleisli category of a free module monad;
- **sfKrn**, the category of “s-finite kernels”.

**N.B.** Each  $\mathcal{C}(X, I)$  is then a monoidal category: morphisms of effects are ‘differences’; monoidal product is  $+$ . This extends to loss functions (via “state-dependent differences”).

<sup>\*</sup>  $\mathcal{C}(-, I)$  has a commutative monoid structure compatible with copy-composition.

# Examples

**Relative entropy:**  $\text{KL}(c, c^\#)_\beta(y) := D_{c^\#, c^\dagger_\beta}(y) = D_X(c^\#_\beta(y), c^\dagger_\beta(y)).$

**Log likelihood:**  $\text{MLE}(c, c^\#)_\beta(y) := -\log p_{c^\# \bullet \beta}(y).$

**Free energy:**  $\text{FE}(c, c^\#)_\beta(y) := \text{KL}(c, c^\#)_\beta(y) + \text{MLE}(c, c^\#)_\beta(y).$

Alternatively,  $\text{FE} = \text{KL} + \text{MLE}.$

... But how do they compose? (Where does the chain rule enter?)

# Loss models: compositional loss functions

Loss functions compose much like inversions.

Given  $(c, c^\sharp, L^c) : X \rightarrow Y$  and  $(d, d^\sharp, L^d) : Y \rightarrow Z$ ,  $L^d \circ L^c = L_c^d + d_c^{\sharp*} L^c$ .

More explicitly,  $(L^d \circ L^c)_\beta(z) = L_{c^\sharp \bullet \beta}^d(z) + \mathbb{E}_{y \sim d_{c^\sharp \bullet \beta}^\sharp(z)} [L^c(y)]$ .

A loss model is an assignment of loss functions to lenses that is compatible with this chain rule.

Example:  $\text{KL}((d, d^\sharp) \circ (c, c^\sharp)) = \text{KL}(d, d^\sharp)_c + d_c^{\sharp*} \text{KL}(c, c^\sharp)$ .

(This follows from our earlier chain rule.)

However ...

# Laxness

Unlike KL, not all loss models are *strict* sections of  $\pi_{\text{Loss}}$ !

With morphisms between loss functions, we can consider *lax* loss models: the laxators measure the difference from  $L^d \circ L^c$  to  $L^{d \circ c}$ .

Hence, for a given  $\mathcal{C}$ , we have a category of (lax) loss models  $\text{Loss}(\mathcal{C})$ ; morphisms of loss models are 'icons'.

Examples: MLE and thus FE are lax loss models.

But this is not the only kind of laxness around! ...

# Monoidal structures

There are a few monoidal structures here.

- 1  $\mathcal{C}$ ,  $\mathbf{Copara}_2(\mathcal{C})$ , and  $\mathbf{BayesLens}$  are all standardly monoidal; and hence so is  $\mathbf{StatGame}$ .
- 2 We also have the sum of loss functions  $L^c + L^{c'} : Y \xrightarrow{X} I$ .
- 3 And this induces a monoidal structure on  $\mathbf{Loss}(\mathcal{C})$ .  
(We saw this when defining FE as KL + MLE.)
- 4 Typically there are also monoidal structures on loss models themselves!
  - Hence we have  $\mathbf{MonLoss}(\mathcal{C}) \leftrightarrow \mathbf{Loss}(\mathcal{C})$ .
  - These structures are generally lax (due to correlations in the priors).
  - Thus each of KL, MLE, and FE are lax monoidal.
  - (But KL on parallel channels is strict monoidal; see [6, Remark 5.4.11] ...)



## Review of the talk

Statistical games formalize 'local' approximate inference, yielding a fibration over Bayesian lenses.

Using copy-composition, the relative entropy chain rule is witnessed by a (strict) section of this fibration.

Other important statistical quantities give lax sections.

And all of these admit various monoidal structures.

(However, the bicategory  $\mathbf{Copara}_2(\mathcal{C})$  seems inessential ...)

## Bonus: another route to copy-composition

I can see at least two alternatives to  $\mathbf{Copara}_2(\mathcal{C})$ . Here's one.

Instead of channels  $c$  as horizontal 1-cells, use their densities  $p_c$ .  
Construct a decorated cospan double category accordingly:

$$\text{e.g. } A \xrightarrow{p_c} C \xrightarrow{p_d} D \quad := \quad \begin{array}{c} \triangleleft p_c \qquad \triangleleft p_d \\ | \qquad | \qquad \cup \qquad | \qquad | \\ A \quad B \quad C \quad D \end{array}$$

0-cells are sets of objects in  $\mathcal{C}$ .

Horizontal 1-cells are cospans decorated by densities.

Vertical 1-cells are 'reindexings'; 2-cells are compatible channels.

Note that this has copy-composition built in, without hacks.

It's also more faithful to "factor graphs".

# References

1. Braithwaite, D., Hedges, J., and St Clere Smithe, T.: The Compositional Structure of Bayesian Inference. (2023). arXiv: 2305.06112 [cs, math]. <http://arxiv.org/abs/2305.06112> (visited on 05/11/2023). preprint
2. St Clere Smithe, T.: Bayesian Updates Compose Optically. (2020). arXiv: 2006.01631 [math.CT]. preprint
3. Ghani, N., Hedges, J., Winschel, V., and Zahn, P.: Compositional Game Theory. Proceedings of Logic in Computer Science (LiCS) 2018 (2016). arXiv: 1603.04641 [cs.GT]
4. Capucci, M.: Diegetic Representation of Feedback in Open Games. (2022). arXiv: 2206.12338 [cs.GT]. preprint
5. Capucci, M., Gavranović, B., Hedges, J., and Rischel, E.F.: Towards Foundations of Categorical Cybernetics. (2021). arXiv: 2105.06332 [math.CT]. preprint
6. St Clere Smithe, T.: Mathematical Foundations for a Compositional Account of the Bayesian Brain. (2022). arXiv: 2212.12538 [cs, math, q-bio, stat]. <http://arxiv.org/abs/2212.12538> (visited on 01/09/2023). preprint