# August 22: Exploring causal assumptions with string diagrams

**Anonymous Author(s)**
Affiliation
Address
email

## 1 The story at a high level

**Optmizibility:** I make the claim (unproven) that it is possible to find a "universally optimal" decision function if the following identity holds for all decision functions $J : E \rightarrow \Delta(\mathcal{D})$:



$$(1)$$

If the forward direction holds, the reverse direction does not hold - we can take a problem that respects 1 and introduce additional dominated decisions that break 1 without breaking the "universal optimizability" (i.e. decisions we know to be very bad, but exactly how bad depends on the state in a difficult-to-identify manner). It is an open question whether the reverse direction might hold if we exclude such decisions.

**Sufficient conditions for optimizibility:** It is easy to show that 1 holds if there exists some kernel $^*\mu$ such that the following two identities hold:



$$(2)$$



$$(3)$$

The first condition says that $\kappa$ is fixed on the support of $\mu^*\mu$.

The second is that, given some unknown state $\theta$, we can fully characterise the resulting distribution by looking at the output of $\mu$. The diagram says that if we "guess" the underlying state via $\mu^*\mu$ then "simulate" $\mu$ with the resulting guess, the result is the same as if we guessed the underlying state via $\mu^*\mu$ and then simulated $\mu$ using a separate copy of the underlying state.

## 2 Recoverability

A natural assumption suggested by the notion of a CSDP is that of *recoverability* - that a causal theory $\mathcal{T} : E \times D \rightarrow E$ permits some decision function that reproduces the distribution of the observed data. That is, we assume that for every $(\kappa_\theta, \mu_\theta) := \theta \in \mathcal{T}$ there exists $\gamma_\theta \in \Delta(\mathcal{D})$ such that

$$\gamma_\theta \kappa_\theta = \mu_\theta \qquad (4)$$

Suppose also that we have some $\kappa^*$ that, for all $\theta \in \mathcal{T}$, is a Bayesian inversion of $\gamma_\theta$ and $\kappa_\theta$; that is:

$$\tag{5}$$

A sufficient condition for the existence of such a $\kappa^*$ is the assumption that decisions correspond to *variable setting* - that is, there is some variable $\mathsf{X} : E \to X$ such that for all $a \in D$, $\theta \in \mathcal{T}$ we have $\delta_a \kappa_\theta F_\mathsf{X} = \delta_a$ (such an assumption arises in graphical models as hard interventions, and in potential outcomes as "potential-outcome identifiers"). Indeed $F_\mathsf{X}$ is in this case a candidate for $\kappa^*$. It is not necessary that $\kappa^*$ be deterministic, however - suppose every $\kappa$ ignores $D$. Then choose $\gamma_\theta = \gamma$ for arbitrary $\gamma \in \Delta(\mathcal{D})$ and it can be verified that $\kappa^* : b \mapsto \gamma$ satisfies 5.

I believe a weaker sufficient condition for the existence of a universal $\kappa^*$ is that every $\kappa_\theta$ factorises as $\kappa_\theta = h \curlyvee (\mathrm{Id}_F \otimes j_\theta)$ for some fixed $h : D \to \Delta(\mathcal{F})$, but I have not yet shown this.

We will proceed somewhat rashly: suppose that by defining $\gamma : \mathcal{T} \to \Delta(\mathcal{D})$, $\mu : \mathcal{T} \to \Delta(\mathcal{E})$ and $\kappa : \mathcal{T} \times D \to \Delta(\mathcal{E}$ by $\gamma : \theta \to \gamma_\theta$, $\mu : \theta \to \mu_\theta$ and $\kappa : (\theta, d) \to \kappa_\theta(d; \cdot)$ that all resulting objects are Markov kernels, and that $\mathcal{T}$ is a standard measurable space.

By previous assumptions, we have the following properties:

$$\tag{6}$$

$$\tag{7}$$

$$\tag{8}$$

From 7 we also have

$$\tag{9}$$

$$\tag{10}$$

Where 10 follows from 4.

The following assumption is a formalisation of the notion that "we can determine $\mu$ precisely from observation" (alternatively, that we can find an optimal decision for a classical statistical decision problem). Suppose that $\mu$ is characterised by some kernel ${}^*\mu$. That is,

$$\tag{11}$$

An equivalent condition to 11 is that for all $\theta, \theta' \in \mathcal{T}$, $A \in \mathcal{E}$, we have $\mu(\theta; A) = \mu(\theta'; A)$, $\mu^*\mu(\theta; \cdot)$- almost surely. More informally, the support of $\mu^*\mu$ for each input $\theta$ divides $\mathcal{T}$ into equivalence classes such that for all $\theta$ in a given equivalence class, $\mu$ maps to the same probability measure on $\mathcal{E}$.

Note that as a result of 11 we also have $\mu^*\mu\mu = \mu$. This weaker condition is not sufficient for the following result.

43

44 We then have

$$
\begin{array}{ccc}
\text{[diagram: } \mu \ \kappa^* \ \kappa \text{]} & \overset{7}{=} & \text{[diagram: } \mu \ {}^*\mu\mu \ \kappa^* \text{]} 
\end{array} \tag{12}
$$

$$
= \quad \text{[diagram: } \mu^*\mu \ \mu \ \kappa^* \ \kappa \ \kappa^* \text{]} \tag{13}
$$

$$
\overset{7}{=} \quad \text{[diagram: } \mu^*\mu \ \mu \ \kappa^* \ \kappa \text{]} \tag{14}
$$

$$
\overset{34}{=} \quad \text{[diagram: } \mu \ {}^*\mu \ \kappa^* \ \kappa \text{]} \tag{15}
$$

$$
\overset{11}{=} \quad \text{[diagram: } \mu \ {}^*\mu \ \mu\kappa^* \ \kappa \text{]} \tag{16}
$$

45 Equation 16 implies that, given any $\xi \in \Delta(\mathcal{T})$, all distributions of the form

$$
\text{[diagram: } \xi \ \mu \ \kappa^* \ \kappa \rightarrow \mathsf{T}, \ \mathsf{E}, \ \mathsf{D} \text{]} \tag{17}
$$

46 admit both $\kappa := \boxed{\kappa}$ and $\kappa_{\mathrm{fac}} := \boxed{\mu^*\mu} \ \boxed{\kappa}$ as disintegrations from $(\mathsf{D}, \mathsf{T}) \dashrightarrow \mathsf{E}$. Therefore
47 these $\kappa$ and $\kappa_{\mathrm{fac}}$ agree almost surely with respect to the distribution 17 for any prior $\xi$.

48 However, also by assumption 11, we have that for $\theta, \theta' \in \mathcal{T}$ either $\mu(\theta; A) = \mu(\theta'; A)$ for all $A \in \mathcal{E}$,
49 or for any $A \in \mathcal{E} \ \mu(\theta; A) = 0$ or $\mu(\theta'; A) = 0$. That is, any two states either have the same probability
50 measure or probability measures with disjoint support. This is problematic, as the distribution 17
51 then has no support over much of the space $D \times E \times \mathcal{T}$. If $\mu$ were deterministic, for example, and
52 hence associated with some function $f$, while 11 would be guaranteed via a left inverse, 17 would be
53 supported on a subset of $D \times \{(\theta, f(\theta)) | \theta \in \mathcal{T}\}$. In particular, we have no guarantee that the desired
54 equality of $\kappa$ and $\kappa_{\mathrm{fac}}$ holds if we take any decision that doesn't reproduce the observed distribution.
55 This isn't totally trivial: we may live in a world where most actions make things worse, in which case
56 knowing how to keep things the same is valuable.

57 A stronger result can be found if we assume we have an infinite sequence of RVs $\mathsf{X}_i : E \to W$ and
58 $\mathsf{D}_i : D \to V$ such that

59 • $W^{\mathbb{N}} = E$, $V^{\mathbb{N}} = D$ (i.e. the sequence of all $\mathsf{X}_i$'s is identified with $E$ and the sequence of all
60 $\mathsf{D}_i$'s is identified with $D$)

3

- $\mu = \curlyvee \otimes_{i \in \mathbb{N}} \mu F_{\mathsf{X}_i}$ (the $\mathsf{X}_i$'s are "IID conditional on $\theta$")
- There exists $\kappa_0$ such that $\kappa = \curlyvee \otimes_{i \in \mathbb{N}} (F_{\mathsf{D}_i} \otimes \mathrm{Id}_{\mathfrak{T}}) \kappa_0 F_{\mathsf{X}_i}$ ($\kappa$ is "IID conditional on D, $\theta$")

Here we define the "infinite copy map" $\curlyvee \otimes_{i \in \mathbb{N}} \mu F_{\mathsf{X}_i}$ to denote the kernel $\theta \mapsto \nu_\theta$ where $\nu_\theta$ the unique distribution such that for all finite $A \subset \mathbb{N}$ and projections $\pi_A : E \to \Delta(W^{|A|})$, $\nu_\theta \pi_A = \otimes_{i \in A} \mu_\theta F_{\mathsf{X}_i}$. This distribution is unique via the Kolmogorov extension theorem (the symmetry of the copy map guarantees the required consistency conditions) [Tao, 2011].

I assume, for now, that measurability can be worked out in some cases; in particular, that there is a $\sigma$-algebra on infinite sequences that renders the above kernel measurable in the appropriate way.

**Lemma 2.1** ("IID" kernels agree on truncations)**.** *For finite $A \subset \mathbb{N}$, $y, y' \in D$, if $\otimes_{i \in A} \mathsf{X}_i(y) = \otimes_{i \in A} \mathsf{X}_i(y')$ and $\kappa : \mathfrak{T} \times D \to \Delta(\mathcal{E})$ is "IID" in the sense above then for all $\theta \in \mathfrak{T}$, $B \in \mathcal{W}^{|A|}$, $\kappa(\theta, y; B)\pi_A = \kappa(\theta, y'; B)\pi_A$.*

*Proof.* By definition, we have

$$\kappa \pi_A(\theta, y; B) = \otimes_{i \in A} \kappa F_{\mathsf{X}_i}(\theta, \mathsf{D}_i(y); B) \tag{18}$$

$$= \otimes_{i \in A} \kappa F_{\mathsf{X}_i}(\theta, \mathsf{D}_i(y'); B) \tag{19}$$

$$= \kappa \pi_A(\theta, y'; B) \tag{20}$$

$\square$

Suppose both $\mathsf{X}_i$ and $\mathsf{D}_i$ are binary, and that for each $\theta \in \mathfrak{T}$ we have recoverability (Eq. 4) with $\mu_\theta = \gamma_\theta$ (we will conclude that X is "directly controlled" by D, but we will not assume this at the outset). $\kappa^*$ is therefore trivial. For each $\theta$, $\mathsf{X}_i$ are IID Bernoulli variables and so each $\mu_\theta$ is characterised by a single parameter $p$; let $p_\theta$ be the value of this parameter for some given $\theta$. Define $\overline{\mathsf{X}} := \lim_{n \to \infty} \frac{1}{m} \sum_{i \in [n]} \mathsf{X}_i$ and $^*\mu$ to be any kernel $E \to \Delta(\mathfrak{T})$ such that the support of $^*\mu(x; \cdot)$ is a subset of $\{\theta | p_\theta = \overline{\mathsf{X}}(x)\}$. Note that for any $\theta, \theta' \in \mathfrak{T}$ we have either $p_\theta = p_{\theta'}$ and so $\mu(\theta; A) = \mu(\theta'; A)$ for all $A$ or $\theta'$ is not in the support of $\mu^*\mu(\theta; \cdot)$. Thus we have 11, and hence "almost sure" equality of $\kappa$ and $\kappa_{\mathrm{fac}}$.

However with the exception of states where $p_\theta = 0$ or 1, almost sure equality is enough for $\kappa_{\mathrm{fac}} \pi_A(\theta, y; B) = \kappa \pi_A(\theta, y; B)$ for all $y \in D$, finite $A \subset \mathbb{N}$ and $B \in \mathcal{W}^{|A|}$. Then by the Kolmogorov extension theorem, we also have $\kappa_{\mathrm{fac}}(\theta, y; B) = \kappa(\theta, y; B)$ for all $y \in D$ and "almost all" $\theta \in \mathfrak{T}$.

This appears to have similarities to the general case where we are trying to identify a particular function from some set of possible functions and we know the output of that function for a subset of inputs. It still comes down to a question of whether or not the set of functions in question is small enough to be fully characterised by the set of inputs we're allowed to see.

# 3 Notes on category theoretic probability and string diagrams

Category theoretic treatments of probability theory often start with *probability monads* (for a good overview, see [Jacobs, 2018]). A monad on some category $C$ is a functor $T : C \to C$ along with natural transformations called the unit $\eta : 1_C \to T$ and multiplication $\mu : T^2 \to T$. Roughly, functors are maps between categories that preserve identity and composition structure and natural transformations are "maps" between functors that also preserve composition structure. The monad unit is similar to the identity element of a monoid in that application of the identity followed by multiplication yields the identity transformation. The multiplication transformation is also (roughly speaking) associative.

An example of a probability monad is the discrete probability monad given by the functor $\mathcal{D} : \mathbf{Set} \to \mathbf{Set}$ which maps a countable set $X$ to the set of functions from $X \to [0, 1]$ that are probability measures on $X$, denoted $\mathcal{D}(X)$. $\mathcal{D}$ maps a measurable function $f$ to $\mathcal{D}f : X \to \mathcal{D}(X)$ given by $\mathcal{D}f : x \mapsto \delta_{f(x)}$. The unit of this monad is the map $\eta_X : X \to \mathcal{D}(X)$ given by $\eta_X : x \mapsto \delta_x$ (which is equivalent to $\mathcal{D}1_X$) and multiplication is $\mu_X : \mathcal{D}^2(X) \to \mathcal{D}(X)$ where $\mu_X : \Omega \mapsto \sum_\phi \Omega(\phi)\phi$.

For continuous distributions we have the Giry monad on the category **Meas** of mesurable spaces given by the functor $\mathcal{G}$ which maps a measurable space $X$ to the set of probability measures on $X$, denoted $\mathcal{G}(X)$. Other elements of the monad (unit, multiplication and map between morphisms) are the "continuous" version of the above.

Of particular interest is the Kleisli category of the monads above. The Kleisli $C_T$ category of a monad $T$ on category $C$ is the category with the same objects and the morphisms $X \to Y$ in $C_T$ is the set of morphisms $X \to TY$ in $C$. Thus the morphisms $X \to Y$ in the Kleisli category $\mathbf{Set}_{\mathcal{D}}$ are morphisms $X \to \mathcal{D}(Y)$ in **Set**, i.e. stochastic matrices, and in the Kleisli category $\mathbf{Meas}_{\mathcal{G}}$ we have Markov kernels. Composition of arrows in the Kleisli categories correspond to Matrix products and "kernel products" respectively.

Both $\mathcal{D}$ and $\mathcal{G}$ are known to be *commutative* monads, and the Kleisli category of a commutative monad is a symmetric monoidal category.

Diagrams for symmetric monoidal categories consist of wires with arrows, boxes and a couple of special symbols. The identity object (which we identify with the set $\{*\}$) is drawn as nothing at all $\{*\} := \boxed{\phantom{xxx}}$ and identity maps are drawn as bare wires:

$$\mathrm{Id}_X := \quad \uparrow_X \tag{21}$$

We draw Kleisli arrows from the unit (i.e. probability distributions) $\mu : \{*\} \to X$ as triangles and Kleisli arrows $\kappa : X \to Y$ (i.e. Markov kernels $X \to \Delta(\mathcal{Y})$) as boxes. We draw the Kleisli arrow $\mathbb{1}_X : X \to \{*\}$ (which is unique for each $X$) as below

$$\mu := \quad \triangleleft\!\!\!\mu\ \overset{\uparrow X}{} \qquad\qquad \kappa := \quad \boxed{\kappa}\ \overset{\uparrow Y}{} \tag{22}$$

The product of objects in **Meas** is given by $(X, \mathcal{X}) \cdot (Y, \mathcal{Y}) = (X \times Y, \mathcal{X} \otimes \mathcal{Y})$, which we will often write as just $X \times Y$. Horizontal juxtaposition of wires indicates this product, and horizontal juxtaposition also indicates the tensor product of Kleisli arrows. Let $\kappa_1 : X \to W$ and $\kappa_2 : Y \to Z$:

$$(X \times Y, \mathcal{X} \otimes \mathcal{Y}) := \quad \uparrow_X \uparrow_Y \qquad\qquad \kappa_1 \otimes \kappa_2 := \quad \overset{\uparrow W}{\boxed{\kappa_1}}\overset{\uparrow Z}{\boxed{\kappa_2}} \tag{23}$$

Composition of arrows is achieved by "wiring" boxes together. For $\kappa_1 : X \to Y$ and $\kappa_2 : Y \to Z$ we have

$$\kappa_1 \kappa_2(x; A) = \int_Y \kappa_2(y; A)\kappa_1(x; dy) := \quad \begin{matrix}\uparrow Z \\ \boxed{\kappa_2} \\ \boxed{\kappa_1} \\ \uparrow X \end{matrix} \tag{24}$$

Symmetric monoidal categoris have the following coherence theorem[Selinger, 2010]:

**Theorem 3.1** (Coherence (symmetric monoidal)). *A well-formed equation between morphisms in the language of symmetric monoidal categories follows from the axioms of symmetric monoidal categories ifand only if it holds, up to isomorphism of diagrams, in the graphical language.*

Isomorphism of diagrams for symmetric monoidal categories (somewhat informally) is any planar deformation of a diagram including deformations that cause wires to cross. We consider a diagram for a symmetric monoidal category to be well formed only if all wires point upwards.

In fact the Kleisli categories of the probability monads above have (for each object) unique *copy*: $X \to X \times X$ and *erase*: $X \to \{*\}$ maps that satisfy the *commutative comonoid axioms* that (thanks

to the coherence theorem above) can be stated graphically. These differ from the copy and erase maps of *finite product* or *cartesian* categories in that they do not necessarily respect composition of morphisms.

$$\text{Erase} = \mathbb{1}_X := \quad \text{Copy} = x \mapsto \delta_{x,x} := \qquad\qquad (25)$$

$$\qquad = \qquad := \qquad\qquad (26)$$

$$\qquad = \qquad = \qquad\qquad (27)$$

$$\qquad = \qquad\qquad (28)$$

Finally, $\{*\}$ is a terminal object in the Kleisli categories of either probability monad. This means that the map $X \to \{*\}$ is unique for all objects $X$, and as a consequence for all objects $X, Y$ and all $\kappa : X \to Y$ we have

$$\boxed{\kappa} \quad \Big|_X \ = \ \Big|_X \qquad\qquad (29)$$

This is equivalent to requiring for all $x \in X$ $\int_Y \kappa(x; dy) = 1$. In the case of $\mathbf{Set}_{\mathcal{D}}$, this condition is what differentiates a stochastic matrix from a general positive matrix (which live in a larger category than $\mathbf{Set}_{\mathcal{D}}$).

Thus when manipulating diagrams representing Markov kernels in particular (and, importantly, not more general symmetric monoidal categories) diagram isomorphism also includes applications of 26, 27, 28 and 29.

A particular property of the copy map in $\mathbf{Meas}_{\mathcal{G}}$ (and probably $\mathbf{Set}_{\mathcal{D}}$ as well) is that it commutes with Markov kernels iff the markov kernels are deterministic [Fong, 2013].

## 3.1 Disintegration and Bayesian inversion

*Disintegration* is a key operation on probability distributions (equivalently arrows $\{*\} \to X$) in the categories under discussion. It corresponds to "finding the conditional probability" (though conditional probability is usually formalised in a slightly different way).

Given a distribution $\mu : \{*\} \to X \otimes Y$, a disintegration $c : X \to Y$ is a Markov kernel that satisfies

$$X \qquad Y$$
$$\boxed{c}$$
$$X\,Y$$
$$\triangleleft\mu\, \quad = \quad \triangleleft\mu\, \qquad\qquad (30)$$

Disintegrations always exist in $\mathbf{Set}_{\mathcal{D}}$ but not in $\mathbf{Meas}_{\mathcal{G}}$. The do exist in the latter if we restrict ourselves to standard measurable spaces. If $c_1$ and $c_2$ are disintegrations $X \to Y$ of $\mu$, they are equal
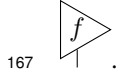
6

156 $\mu$-A.S. In fact, this equality can be strengthened somewhat - they are equal almost surely with respect
157 to any distribution that shares the "$X$-marginal" of $\mu$.

158 Given $\sigma : \{*\} \to X$ and a channel $c : X \to Y$, a Bayesian inversion of $(\sigma, c)$ is a channel $d : Y \to X$
159 such that



$$\tag{31}$$

160 We can obtain disintegrations from Bayesian inversions and vise-versa.

161 Clerc et al. [2017] offer an alternative view of Bayesian inversion which they claim doesn't depend
162 on standard measurability conditions, but there is a step in their proof I didn't follow.

## 3.2 Generalisations

164 Cho and Jacobs [2019] make use of a larger "CD" category by dropping 29. I'm not completely clear
165 whether you end up with arrows being "Markov kernels for general measures" or something else (can
166 we have negative arrows?). This allows for the introduction of "observables" or "effects" of the form

167  .

168 Jacobs et al. [2019] make use of an embedding of $\mathbf{Set}_{\mathcal{D}}$ in $\mathbf{Mat}(\mathbb{R}^+)$ with morphisms all positive
169 matrices (I'm not totally clear on the objects, or how they are self-dual - this doesn't seem to be
170 exactly the same as the category of finite dimensional vector spaces). This latter category is compact
171 closed, which - informally speaking - supports the same diagrams as symmetric monoidal categories
172 with the addition of "upside down" wires.

## 3.3 Key questions for Causal Theories

174 We will first define *labeled diagrams*. Rather than labelling the wires of our diagrams with *spaces* (as is
175 typical [Selinger, 2010]), we assign a unique label to each "wire segment" (with some qualifications).
176 That is, we assign a unique label to each bare wire in the diagram with the following additonal
177 qualifications:

- 178 If we have a box in the diagram representing the identity map, the incoming and outgoing
  179 wires are given the same label
- 180 If we have a wire crossing in the diagram, the diagonally opposite wires are given the same
  181 label
- 182 The input wire and the *two* output wires of the copy map are given the same label

183 Given two diagrams $G_1$ and $G_2$ that are isomorphic under transformations licenced by the axioms of
184 symmetric monoidal categories and commutative comonoid axioms, suppose we have a labelling of
185 $G_1$. We can label $G_2$ using the following translation rule:

- 186 For each box in $G_2$, we can identify a corresponding box in $G_1$ via labels on each box. For
  187 each such pair of boxes, we label the incoming wires of the $G_2$ box with the labels of the
  188 $G_1$ box preserving the left-right order. We do likewise for outgoing wires.

189 These rules will lead to a unique labelling of $G_2$ with all wire segments are labelled. We would like
190 for these rules to yield the following:

- 191 For any sequence of diagram isomorphisms beginning with $G_1$ and ending with $G_2$, we end
  192 up with the same set of labels
- 193 If we label $G_2$ according to the rules above then relabel $G_1$ from $G_2$ according to the same
  194 rules we retrieve the original labels of $G_1$

7

I'm sure one of the papers I read mentioned labeled diagrams, I just couldn't find it when I looked for it

Since writing this, I found Kissinger [2014] as an example of a diagrammatic system with labeled wires, I will follow it up

195 We do not prove these properties here, but motivate them via the following considerations:

- These properties obviously hold for the wire segments into and out of boxes
- The only features a diagram may have apart from boxes and wires are wire crossings, copy maps and erase maps
- The labeling rule for wire crossings respects the symmetry of the swap map
- The labeling rule for copy maps respects the symmetry of the copy map and the property described in Equation 28

202 We will follow the convention whereby "internal" wire labels are omitted from diagrams.

203 Note also that each wire that terminates in a free end can be associated with a random variable.
204 Suppose for $N \in \mathbb{N}$ we have a kernel $\kappa : A \to \Delta(\times_{i \in N} X_i)$. Define by $p_j$ ($j \in [N]$) the projection
205 map $p_j : \times_{i \in N} X_i \to X_j$ defined by $p_j : (x_0, ..., x_N) \mapsto x_j$. $p_j$ is a measurable function, hence
206 a random variable. Define by $\pi_j$ the projection kernel $\mathcal{G}(p_j)$ (that is, $\pi_j : \mathbf{x} \mapsto \delta_{p_j(\mathbf{x})}$). Note that
207 $\kappa(y; p_j^{-1}(A)) = \int_{X_j} \delta_{p_j(\mathbf{x})}(A) \kappa(y; d\mathbf{x}) = \kappa \pi_j$. Diagrammatically, $\pi_j$ is the identity map on the $j$-th
208 wire tensored with the erase map on every other wire. Thus the $j$-th wire carries the distribution
209 associated with the random variable $p_j$. We will therefore consider the labels of the "outgoing" wires
210 of a diagram to denote random varaibles (though there are obviously many random variables not
211 represented by such wires). We will additionally distinguish wire labels from spaces by font - wire
212 labels are sans serif $\mathsf{A}, \mathsf{B}, \mathsf{C}, \mathsf{X}, \mathsf{Y}, \mathsf{Z}$ while spaces are serif $A, B, C, X, Y, Z$.

> Wire labels appear to have a key advantage over random variables: they allow us to "forget" the sample space as the correct typing is handled automatically by composition and erasure of wires

214 **generalised disintegrations**   : Of key importance to our work is generalising the notion of disinte-
215 gration (and possibly Bayesian inversion) to general kernels $X \to Y$ rather than restricting ourselves
216 to probability distributions $\{*\} \to Y$. We will define generalised disintegrations as a straightforward
217 analogy regular disintegrations, but the conditions under which such disintegrations exist are more
218 restrictive than for regular disintegraions.

219 **Definition 3.2** (Label signatures). If a kernel $\kappa : X \to \Delta(Y)$ can be represented by a diagram
220 $G$ with incoming wires $\mathsf{X}_1, ... \mathsf{X}_n$ and outgoing wires $\mathsf{Y}_1, ..., \mathsf{Y}_m$, we can assign the kernel a "label
221 signature" $\kappa : \mathsf{X}_1 \otimes ... \otimes \mathsf{X}_n \dashrightarrow \mathsf{Y}_1 \otimes ... \otimes \mathsf{Y}_m$ or, for short, $\kappa : \mathsf{X}_{[n]} \dashrightarrow \mathsf{Y}_{[m]}$. Note that this
222 signature associates each label with a unique space - the space of $\mathsf{X}_1$ is the space associated with the
223 left-most wire of $G$ and so forth. We will implicitly leverage this correspondence and write with $X_1$
224 the space associated with $\mathsf{X}_1$ and so forth. Note that while $\mathsf{X}_1$ is by construction always different from
225 $\mathsf{X}_2$ (or any other label), the space $X_1$ may coincide with $X_2$ - the fact that labels always maintain
226 distinctions between wires is the fundamental reason for introducing them in the first place.

> There might actually be some sensible way to consider $\kappa$ to be transforming the measurable functions of a type similar to $\otimes_{i \in [n]} X_i$ to functions of a type simlar to $\otimes_{i \in [m]} Y_i$ (or vise versa - perhaps related to Clerc et al. [2017]), but wire labels are all we need at this point

228 **Definition 3.3** (Generalised disintegration). Given a kernel $\kappa : X \to \Delta(Y)$ with label signature
229 $\kappa : \mathsf{X}_{[n]} \dashrightarrow \mathsf{Y}_{[m]}$ and disjoint subsets $S, T \subset [m]$ such that $S \cup T = [m]$, a kernel $c$ is a *g-
230 disintigration from $S$ to $T$* if it's type is compatible with the label signature $c : \mathsf{Y}_\mathsf{S} \dashrightarrow \mathsf{Y}_\mathsf{T}$ and we
231 have the identity (omitting incoming wire labels):



$$\tag{32}$$

8

232

233  In contrast to regular disintegrations, generalised disintegrations "usually" do not exist. Consider
234  $X = \{0, 1\}$, $Y = \{0, 1\}^2$ and $\kappa$ has label signature $\mathsf{X}_1 \dashrightarrow \mathsf{Y}_{\{1,2\}}$ with

$$\kappa : \begin{cases} 1 \mapsto \delta_1 \otimes \delta_1 \\ 0 \mapsto \delta_1 \otimes \delta_0 \end{cases} \tag{33}$$

235  $\kappa$ imposes contradictory requirements for any disintegration $c : \{0, 1\} \to \{0, 1\}$ from $\{1\}$ to $\{2\}$:
236  equality for $\mathsf{X}_1 = 1$ requires $c(1; \cdot) = \delta_1$ while equality for $\mathsf{X}_1 = 0$ requires $c(1; \cdot) = \delta_0$. Subject
237  to some regularity conditions (similar to standard Borel conditions for regular disintegrations),
238  we can define g-disintegrations of a canonically related kernel that do generally exist; intuitively,
239  g-disintegrations exist if they take the "input wires" of $\kappa$ as input wires themselves.

240  **Lemma 3.4.** *Given $\kappa : X \to \Delta(Y)$, a kernel $\kappa^\dagger$ is a right inverse iff we have for all $x \in X$, $A \in \mathcal{X}$,*
241  *$y \in Y$ $\kappa^\dagger(y; A) = \delta_x(A)$, $\kappa(x; \cdot)$-almost surely.*

242  *Proof.* Suppose $\kappa^\dagger$ satisfies the almost sure equality for all $x \in X$. Then for all $x \in X$, $A \in \mathcal{X}$ we
243  have $\kappa\kappa^\dagger(x; A) = \int_Y \kappa^\dagger(y; A)\kappa(x; dy) = \int_Y \delta_x(A)\kappa(x; dy) = \delta_x(A)$; that is, $\kappa\kappa^\dagger = \mathrm{Id}_X$, so $\kappa^\dagger$ is
244  a right inverse of $\kappa$.

245  Suppose we have a right inverse $\kappa^\dagger$. By definition, for all $x \in X$ and $A \in \mathcal{X}$ we have
246  $\int_Y \kappa^\dagger(y; A)\kappa(x; dy) = \delta_x(A)$.

247  Suppose $x \notin A$ and let $B_\epsilon = \kappa_A^{\dagger-1}((\epsilon, 1])$ for some $\epsilon > 0$. We have $\int_Y \kappa^\dagger(y; A)\kappa(x; dy) = 0 \geq$
248  $\epsilon\kappa(x; B_\epsilon)$. Thus for any $\epsilon > 0$ we have $\kappa(x; B_\epsilon) = 0$. Consider the set $B_0 = \kappa_A^{\dagger-1}((0, 1])$. For
249  some sequence $\{\epsilon_i\}_{i \in \mathbb{N}}$ such that $\lim_{i \to \infty} \epsilon_i = 0$ we have $B_0 = \cup_{i \in \mathbb{N}} B_{\epsilon_i}$. By countable additivity,
250  $\kappa(x; B_0) = 0$.

251  Suppose $x \in A$ and let $B^{1-\epsilon} = \kappa_A^{\dagger-1}([0, 1 - \epsilon))$. We have $\int_Y \kappa^\dagger(y; A)\kappa(x; dy) = 1 \leq (1 -$
252  $\epsilon)\kappa(x; B^{1-\epsilon}) + 1 - \kappa(x; B^{1-\epsilon}) = 1 - \epsilon\kappa(x; B^{1-\epsilon})$. Thus $\kappa(x; B^{1-epsilon}) = 0$ for $\epsilon > 0$. By an
253  argument analogous to the above, we also have $\kappa(x; B^1) = 0$. Thus the $\kappa(x; \cdot)$ measure of the set
254  on which $\kappa^\dagger(y; A)$ disagrees with $\delta_x(A)$ is $\kappa(x; B_0) + \kappa(x; B^1) = 0$ and hence $\kappa^\dagger(y; A) = \delta_x(A)$
255  $\kappa(x; \cdot)$-almost surely. $\qquad\square$

256

257

258  **Lemma 3.5.** *Given $\kappa : X \to \Delta(Y)$ and a right inverse $\kappa^\dagger$, we have*



$$\tag{34}$$

9

259 *Proof.* Let the diagram on the left hand side be $L$ and the diagram on the right hand side be $R$.

$$L(x; A \times B) = \int_Y \int_{Y \times Y} \mathrm{Id}_Y \otimes \kappa_S^\dagger(y, y'; A \times B)\delta_{(z,z)}(dy \times dy')\kappa\pi_S(x; dz) \tag{35}$$

$$= \int \mathrm{Id}_Y \otimes \kappa^\dagger(z, z; A \times B)\kappa\pi_S(x; dz) \tag{36}$$

$$= \int \delta_z(A)\kappa_S^\dagger(z; B)\kappa\pi_S(x; dz) \tag{37}$$

$$= \int_A \kappa_S^\dagger(z; B)\kappa\pi_S(x; dz) \tag{38}$$

$$= \delta_x(B)\kappa\pi_S(x; A) \tag{39}$$

260 Where 39 follows from Lemma 3.4.

$$R(x; A \times B) = \int \delta_{(x,x)}(dy \times dy')\kappa\pi_S \otimes \mathrm{Id}_X(y, y'; A \times B) \tag{40}$$

$$= \kappa\pi_S(x; A)\delta_x(B) \qquad\qquad = L \tag{41}$$

261 $\hfill\square$

262 **Theorem 3.6.** *Given countable $X$ and standard measurable $Y$, $n, m \in \mathbb{N}$, $S, T \subset [m]$, $\kappa$ with label*
263 *signature $\mathsf{X}_{[n]} \dashrightarrow \mathsf{Y}_{[m]}$ a g-disintegration exists from $S$ to $T$ if $\kappa\pi_S$ is right-invertible*

264 *via a Markov kernel*

265 *Proof.* In addition, as $R$ is a composition of Markov kernels, and hence a Markov kernel itself, $L$
266 must also be a Markov kernel even if $\kappa^\dagger$ is not.

267 For all $x \in X$ we have a (regular) disintegration $c_x : Y_S \to \Delta(Y_T)$ of $\kappa(x; \cdot)$ by standard mea-
268 surability of $Y$. Define $c : X \otimes Y_S \to \Delta(Y_T)$ by $c : (x, y_S) \mapsto c_x(y_S)$. Clearly, $c(x, y_S)$ is a
269 probability distribution on $Y_T$ for all $(x, y_S) \in X \otimes Y_S$. It remains to show $c(\cdot)^{-1}(B)$ is measurable
270 for all $B \in \mathcal{B}([0, 1])$. But $c(\cdot)^{-1}(B) = \cap_{x \in X} c_y(\cdot)^{-1}(B)$. The right hand side is measurable by
271 measurability of $c_y(\cdot)^{-1}(B)$ countability of $X$, so $c$ is a Markov kernel.

272 By the definition of $c_x$, we have for all $x \in X$



$$\tag{42}$$



$$\tag{43}$$

273 Which implies

$$
\begin{array}{c}
\mathsf{Y}_S\ \mathsf{Y}_T \\
\kappa \\
\end{array}
\quad = \quad
\begin{array}{c}
\mathsf{Y}_S \qquad \mathsf{Y}_T \\
c \\
\ast \\
\kappa \\
\end{array}
\tag{44}
$$

274 Finally, we have

$$
\begin{array}{c}
\mathsf{Y}_S \qquad \mathsf{Y}_T \\
c \\
\kappa_S^{\dagger} \\
\ast \\
\kappa \\
\end{array}
\quad = \quad
\begin{array}{c}
\mathsf{Y}_S \qquad \mathsf{Y}_T \\
c \\
\kappa_S^{\dagger} \\
\ast \\
\kappa \\
\end{array}
\tag{45}
$$

$$
\quad = \quad
\begin{array}{c}
\mathsf{Y}_S \qquad \mathsf{Y}_T \\
c \\
\ast \\
\kappa \\
\end{array}
\tag{46}
$$

275 Where the first line follows from 27 and the second line from 34. If $\kappa_S^{\dagger}$ is a Markov kernel, then
276 $\curlyvee(\mathrm{Id}_{Y_S} \otimes \kappa_S^{\dagger})c$ is a g-disintegration. $\qquad\square$

277 In the reverse direction, suppose $\kappa$ is such that $\kappa\pi_T = \mathrm{Id}_X$; that is, $\pi_T$ is a right inverse of $\kappa$. If
278 $\kappa\pi_S$ is not right invertible then, by definition, there is no $d$ such that $\kappa\pi_S d\pi_T = \mathrm{Id}_X$. However, if a
279 g-disintegration of $\kappa$ exists then there is a $d$ such that $\kappa\pi_S d = \kappa$, a contradiction. Thus if $\kappa\pi_S$ is not
280 right invertible then there is *in general* no g-disintegration from $S$ to $T$.

## 281 References

282 Kenta Cho and Bart Jacobs. Disintegration and Bayesian inversion via string diagrams.
283 *Mathematical Structures in Computer Science*, 29(7):938–971, August 2019. ISSN
284 0960-1295, 1469-8072. doi: 10.1017/S0960129518000488. URL https://www.
285 cambridge.org/core/journals/mathematical-structures-in-computer-science/
286 article/disintegration-and-bayesian-inversion-via-string-diagrams/
287 0581C747DB5793756FE135C70B3B6D51.

288 Florence Clerc, Fredrik Dahlqvist, Vincent Danos, and Ilias Garnier. Pointless learn-
289 ing. *20th International Conference on Foundations of Software Science and Compu-*
290 *tation Structures (FoSSaCS 2017)*, March 2017. doi: 10.1007/978-3-662-54458-7_
291 21. URL https://www.research.ed.ac.uk/portal/en/publications/
292 pointless-learning(694fb610-69c5-469c-9793-825df4f8ddec).html.

293 Brendan Fong. Causal Theories: A Categorical Perspective on Bayesian Networks. *arXiv:1301.6201*
294 *[math]*, January 2013. URL http://arxiv.org/abs/1301.6201. arXiv: 1301.6201.

295 Bart Jacobs. From probability monads to commutative effectuses. *Journal of Logical and*
296 *Algebraic Methods in Programming*, 94:200–237, January 2018. ISSN 2352-2208. doi:
297 10.1016/j.jlamp.2016.11.006. URL http://www.sciencedirect.com/science/article/
298 pii/S2352220816301122.

Bart Jacobs, Aleks Kissinger, and Fabio Zanasi. Causal Inference by String Diagram Surgery. In Mikołaj Bojańczyk and Alex Simpson, editors, *Foundations of Software Science and Computation Structures*, Lecture Notes in Computer Science, pages 313–329. Springer International Publishing, 2019. ISBN 978-3-030-17127-8.

Aleks Kissinger. Abstract Tensor Systems as Monoidal Categories. In Claudia Casadio, Bob Coecke, Michael Moortgat, and Philip Scott, editors, *Categories and Types in Logic, Language, and Physics: Essays Dedicated to Jim Lambek on the Occasion of His 90th Birthday*, Lecture Notes in Computer Science, pages 235–252. Springer Berlin Heidelberg, Berlin, Heidelberg, 2014. ISBN 978-3-642-54789-8. doi: 10.1007/978-3-642-54789-8_13. URL https://doi.org/10.1007/978-3-642-54789-8_13.

Peter Selinger. A survey of graphical languages for monoidal categories. *arXiv:0908.3347 [math]*, 813:289–355, 2010. doi: 10.1007/978-3-642-12821-9_4. URL http://arxiv.org/abs/0908.3347. arXiv: 0908.3347.

Terence Tao. *An Introduction to Measure Theory*. American Mathematical Soc., September 2011. ISBN 978-0-8218-6919-2. Google-Books-ID: HoGDAwAAQBAJ.