
August 22: Exploring causal assumptions with string diagrams

Anonymous Author(s)

Affiliation

Address

email

1 The story at a high level

2 Take a causal theory \mathcal{T} where we label each pair $\theta := (\kappa_\theta, \mu_\theta) \in \mathcal{T}$. Define the kernels $\kappa : \mathcal{T} \times D \rightarrow \mathcal{E}$
 3 and $\mu : \mathcal{T} \rightarrow \mathcal{E}$.

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

$$\text{---} \boxed{\mu} \boxed{J} \boxed{\kappa} \text{---} \triangleright = \text{---} \boxed{\mu} \boxed{\beta} \text{---} \boxed{J} \boxed{\kappa} \text{---} \triangleright \quad \text{with a curved arrow from } \beta \text{ to } \kappa \text{ labeled } \alpha \quad (1)$$

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

11 **Sufficient conditions for optimizability:** It is easy to show that ?? holds if there exists some kernel
 12 μ^* such that the following two identities hold:

$$\text{---} \boxed{\mu^* \mu} \boxed{\kappa} \text{---} = \text{---} \boxed{\kappa} \text{---} \quad (2)$$

$$\text{---} \boxed{\mu} \boxed{\mu^*} \text{---} \text{---} \boxed{\mu} \text{---} = \text{---} \boxed{\mu} \boxed{\mu^*} \text{---} \text{---} \boxed{\mu} \text{---} \quad (3)$$

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

14 The second is less obvious. It implies that if we “guess” the underlying state via $\mu^* \mu$ this is as good
 15 as having the actual underlying state for the purposes of determining the output of μ , but it is stronger
 16 than this. In particular, the *joint distribution* between the “guess” and the observations must be the
 17 same whether we use the guess or the true underlying state as input to μ .

18 Two sufficient conditions for ?? to obtain are 1) when μ is deterministic (as μ then has a left inverse)
 19 and 2) if observations are an infinite sequence of binary random variables where each μ_θ corresponds
 20 to a Bernoulli distribution for a particular parameter p_θ (via a μ^* that witnesses the strong law of
 21 large numbers).

22 A more general sufficient graphical condition is available, but it is not presently clear if it is also a
 23 necessary one.

24 These conditions are not necessary for ??; observations may be “too informative”. For example,
 25 if \mathcal{T} contains many different μ_θ but only one κ_θ , then we can always perform ??, while we do not
 26 generally have ??.

27 Below, I document additional assumptions that, along with ?? yield ??.

28 I’m not sure how interesting the assumptions themselves are. One interesting point about the big
 29 picture story is that from one point of view the assumptions boil down to:

- 30 • We can characterise the input-output behaviour of κ for any given state and a small subset
 31 of available decisions
- 32 • κ is sufficiently regular that its behaviour on said subset of decisions characterises its
 33 complete behaviour

34 2 Recoverability

35 A natural assumption suggested by the notion of a CSDP is that of *recoverability* - that a causal theory
 36 $\mathcal{T} : E \times D \rightarrow E$ permits some decision function that reproduces the distribution of the observed data.
 37 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

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

38 “Traditional” causal inference doesn’t have a strict equivalent of this assumption, though it corresponds
 39 roughly to the “easy” cases (for example, it is satisfied by a CBN where there are no backdoor paths
 40 between the “intervened” variable and the “target” variable). One reason I think it’s interesting is
 41 that *randomised recoverability* may be quite a general assumption - that is, there is “in principle” a
 42 stochastic decision that recovers the observed distribution, but we are practically limited to taking
 43 mixed decisions that cannot necessarily accomplish this.

44 Suppose also that we have some κ^* that, for all $\theta \in \mathcal{T}$, is a Bayesian inversion of γ_θ and κ_θ ; that is:

$$\begin{array}{c} \boxed{\kappa_\theta} \text{ E} \\ \swarrow \quad \searrow \\ \triangle \gamma_\theta \\ \swarrow \quad \searrow \\ \text{D} \end{array} = \begin{array}{c} \text{E} \\ \swarrow \quad \searrow \\ \triangle \gamma_\theta \\ \swarrow \quad \searrow \\ \boxed{\kappa^*} \text{ D} \end{array} \quad (5)$$

45 A sufficient condition for the existence of such a κ^* is the assumption that decisions correspond to
 46 *variable setting* - that is, there is some variable $X : E \rightarrow X$ such that for all $a \in D$, $\theta \in \mathcal{T}$ we have
 47 $\delta_a \kappa_\theta F_X = \delta_a$ (such an assumption arises in graphical models as hard interventions, and in potential
 48 outcomes as “potential-outcome identifiers”). Indeed F_X is in this case a candidate for κ^* . It is not
 49 necessary that κ^* be deterministic, however - suppose every κ ignores D . Then choose $\gamma_\theta = \gamma$ for
 50 arbitrary $\gamma \in \Delta(\mathcal{D})$ and it can be verified that $\kappa^* : b \mapsto \gamma$ satisfies ??.

51 I believe a weaker sufficient condition for the existence of a universal κ^* is that every κ_θ factorises as
 52 $\kappa_\theta = h \curlywedge (\text{Id}_F \otimes j_\theta)$ for some fixed $h : D \rightarrow \Delta(\mathcal{F})$, but I have not yet shown this.

53 We will proceed somewhat rashly: suppose that by defining $\gamma : \mathcal{T} \rightarrow \Delta(\mathcal{D})$, $\mu : \mathcal{T} \rightarrow \Delta(\mathcal{E})$ and
 54 $\kappa : \mathcal{T} \times D \rightarrow \Delta(\mathcal{E})$ by $\gamma : \theta \rightarrow \gamma_\theta$, $\mu : \theta \rightarrow \mu_\theta$ and $\kappa : (\theta, d) \rightarrow \kappa_\theta(d; \cdot)$ that all resulting objects are
 55 Markov kernels, and that \mathcal{T} is a standard measurable space.

$$\text{Diagram (12)} \quad ?? \quad (12)$$

$$\text{Diagram (13)} \quad = \quad (13)$$

$$\text{Diagram (14)} \quad ?? \quad (14)$$

$$\text{Diagram (15)} \quad ?? \quad (15)$$

$$\text{Diagram (16)} \quad ?? \quad (16)$$

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

$$\text{Diagram (17)} \quad (17)$$

admit both $\kappa := \text{---} \boxed{\kappa} \text{---}$ and $\kappa_{\text{fac}} := \text{---} \boxed{\mu^* \mu} \boxed{\kappa} \text{---}$ as disintegrations from $(D, T) \dashrightarrow E$. Therefore these κ and κ_{fac} agree almost surely with respect to the distribution ?? for any prior ξ .

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

A stronger result can be found if we assume we have an infinite sequence of RVs $X_i : E \rightarrow W$ and $D_i : D \rightarrow V$ such that

- $W^{\mathbb{N}} = E, V^{\mathbb{N}} = D$ (i.e. the sequence of all X_i 's is identified with E and the sequence of all D_i 's is identified with D)
- $\mu = \bigvee \otimes_{i \in \mathbb{N}} \mu F_{X_i}$ (the X_i 's are "IID conditional on θ ")
- There exists κ_0 such that $\kappa = \bigvee \otimes_{i \in \mathbb{N}} (F_{D_i} \otimes \text{Id}_{\mathcal{T}}) \kappa_0 F_{X_i}$ (κ is "IID conditional on D, θ ")

this might be closely related to exchangeability via de Finetti?

87 Here we define the “infinite copy map” $\vee \otimes_{i \in \mathbb{N}} \mu F_{X_i}$ to denote the kernel $\theta \mapsto \nu_\theta$ where ν_θ the unique
 88 distribution such that for all finite $A \subset \mathbb{N}$ and projections $\pi_A : E \rightarrow \Delta(W^{|A|})$, $\nu_\theta \pi_A = \otimes_{i \in A} \mu_\theta F_{X_i}$.
 89 This distribution is unique via the Kolmogorov extension theorem (the symmetry of the copy map
 90 guarantees the required consistency conditions) [?].

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

91
 92 **Lemma 2.1** (“IID” kernels agree on truncations). *For finite $A \subset \mathbb{N}$, $y, y' \in D$, if $\otimes_{i \in A} X_i(y) =$
 93 $\otimes_{i \in A} X_i(y')$ and $\kappa : \mathcal{T} \times D \rightarrow \Delta(\mathcal{E})$ is “IID” in the sense above then for all $\theta \in \mathcal{T}$, $B \in \mathcal{W}^{|A|}$,*
 94 $\kappa(\theta, y; B) \pi_A = \kappa(\theta, y'; B) \pi_A$.

95 *Proof.* By definition, we have

$$\kappa \pi_A(\theta, y; B) = \otimes_{i \in A} \kappa F_{X_i}(\theta, D_i(y); B) \quad (18)$$

$$= \otimes_{i \in A} \kappa F_{X_i}(\theta, D_i(y'); B) \quad (19)$$

$$= \kappa \pi_A(\theta, y'; B) \quad (20)$$

96

□

97 Suppose both X_i and D_i are binary, and that for each $\theta \in \mathcal{T}$ we have recoverability (Eq. ??) with
 98 $\mu_\theta = \gamma_\theta$ (we will conclude that X is “directly controlled” by D , but we will not assume this at
 99 the outset). κ^* is therefore trivial. For each θ , X_i are IID Bernoulli variables and so each μ_θ
 100 is characterised by a single parameter p ; let p_θ be the value of this parameter for some given θ .
 101 Define $\bar{X} := \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i \in [n]} X_i$ and $^*\mu$ to be any kernel $E \rightarrow \Delta(\mathcal{T})$ such that the support of
 102 $^*\mu(x; \cdot)$ is a subset of $\{\theta | p_\theta = \bar{X}(x)\}$. Note that for any $\theta, \theta' \in \mathcal{T}$ we have either $p_\theta = p_{\theta'}$ and so
 103 $\mu(\theta; A) = \mu(\theta'; A)$ for all A or θ' is not in the support of $^*\mu(\theta; \cdot)$. Thus we have ??, and hence
 104 “almost sure” equality of κ and κ_{fac} .

105 However with the exception of states where $p_\theta = 0$ or 1, almost sure equality is enough for
 106 $\kappa_{\text{fac}} \pi_A(\theta, y; B) = \kappa \pi_A(\theta, y; B)$ for all $y \in D$, finite $A \subset \mathbb{N}$ and $B \in \mathcal{W}^{|A|}$. Then by the
 107 Kolmogorov extension theorem, we also have $\kappa_{\text{fac}}(\theta, y; B) = \kappa(\theta, y; B)$ for all $y \in D$ and “almost
 108 all” $\theta \in \mathcal{T}$.

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

113 3 Notes on category theoretic probability and string diagrams

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

122 An example of a probability monad is the discrete probability monad given by the functor $\mathcal{D} : \mathbf{Set} \rightarrow$
 123 \mathbf{Set} which maps a countable set X to the set of functions from $X \rightarrow [0, 1]$ that are probability
 124 measures on X , denoted $\mathcal{D}(X)$. \mathcal{D} maps a measurable function f to $\mathcal{D}f : X \rightarrow \mathcal{D}(X)$ given by
 125 $\mathcal{D}f : x \mapsto \delta_{f(x)}$. The unit of this monad is the map $\eta_X : X \rightarrow \mathcal{D}(X)$ given by $\eta_X : x \mapsto \delta_x$ (which
 126 is equivalent to $\mathcal{D}1_X$) and multiplication is $\mu_X : \mathcal{D}^2(X) \rightarrow \mathcal{D}(X)$ where $\mu_X : \Omega \mapsto \sum_{\phi} \Omega(\phi) \phi$.

127 For continuous distributions we have the Giry monad on the category **Meas** of measurable spaces
 128 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 $X \rightarrow Y$ in C_T is the set of morphisms $X \rightarrow TY$ in C . Thus the morphisms $X \rightarrow Y$ in the Kleisli category $\mathbf{Set}_{\mathcal{D}}$ are morphisms $X \rightarrow \mathcal{D}(Y)$ in \mathbf{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 $\{*\} := \square$ and identity maps are drawn as bare wires:

$$\text{Id}_X := \uparrow_X \quad (21)$$

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

$$\mu := \triangleup_X \quad \kappa := \boxed{\kappa}_Y \quad (22)$$

The product of objects in \mathbf{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 $\kappa_1 : X \rightarrow W$ and $\kappa_2 : Y \rightarrow Z$:

$$(X \times Y, \mathcal{X} \otimes \mathcal{Y}) := \uparrow_X \uparrow_Y \quad \kappa_1 \otimes \kappa_2 := \begin{array}{c} \uparrow_W \quad \uparrow_Z \\ \boxed{\kappa_1} \quad \boxed{\kappa_2} \\ \downarrow_X \quad \downarrow_Y \end{array} \quad (23)$$

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

$$\kappa_1 \kappa_2(x; A) = \int_Y \kappa_2(y; A) \kappa_1(x; dy) := \begin{array}{c} \uparrow_Z \\ \boxed{\kappa_2} \\ \downarrow_Y \\ \boxed{\kappa_1} \\ \downarrow_X \end{array} \quad (24)$$

Symmetric monoidal categories 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 if and 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*: $X \rightarrow X \times X$ and *erase*: $X \rightarrow \{*\}$ 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

160 maps of *finite product* or *cartesian* categories in that they do not necessarily respect composition of
 161 morphisms.

$$\text{Erase} = \mathbb{1}_X := \begin{array}{c} * \\ | \end{array} \quad \text{Copy} = x \mapsto \delta_{x,x} := \begin{array}{c} \swarrow \quad \searrow \\ | \end{array} \quad (25)$$

$$\begin{array}{c} \swarrow \quad \searrow \\ | \end{array} = \begin{array}{c} \swarrow \quad \searrow \\ | \end{array} := \begin{array}{c} \swarrow \quad \searrow \\ | \end{array} \quad (26)$$

$$\begin{array}{c} * \\ | \end{array} = \begin{array}{c} \swarrow \quad \searrow \\ | \end{array} = \begin{array}{c} \uparrow \end{array} \quad (27)$$

$$\begin{array}{c} \swarrow \quad \searrow \\ \cup \end{array} = \begin{array}{c} \swarrow \quad \searrow \\ | \end{array} \quad (28)$$

162 Finally, $\{*\}$ is a terminal object in the Kleisli categories of either probability monad. This means
 163 that the map $X \rightarrow \{*\}$ is unique for all objects X , and as a consequence for all objects X, Y and all
 164 $\kappa : X \rightarrow Y$ we have

$$\begin{array}{c} * \\ \boxed{\kappa} \\ | \end{array} X = \begin{array}{c} * \\ | \end{array} X \quad (29)$$

165 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
 166 what differentiates a stochastic matrix from a general positive matrix (which live in a larger category
 167 than $\mathbf{Set}_{\mathcal{D}}$).

168 Thus when manipulating diagrams representing Markov kernels in particular (and, importantly, not
 169 more general symmetric monoidal categories) diagram isomorphism also includes applications of 6,
 170 7, 8 and 9.

171 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
 172 Markov kernels iff the markov kernels are deterministic [Fong, 2013].

173 3.1 Disintegration and Bayesian inversion

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

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

$$\begin{array}{c} X \quad Y \\ \swarrow \quad \searrow \\ \mu \end{array} = \begin{array}{c} X \quad Y \\ \swarrow \quad \searrow \\ \mu \end{array} \quad (30)$$

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

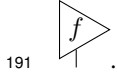
180 μ -A.S. In fact, this equality can be strengthened somewhat - they are equal almost surely with respect
 181 to any distribution that shares the “ X -marginal” of μ .
 182 Given $\sigma : \{*\} \rightarrow X$ and a channel $c : X \rightarrow Y$, a Bayesian inversion of (σ, c) is a channel $d : Y \rightarrow X$
 183 such that

$$(31)$$

184 We can obtain disintegrations from Bayesian inversions and vise-versa.
 185 Clerc et al. [2017] offer an alternative view of Bayesian inversion which they claim doesn’t depend
 186 on standard measurability conditions, but there is a step in their proof I didn’t follow.

187 3.2 Generalisations

188 Cho and Jacobs [2019] make use of a larger “CD” category by dropping 9. I’m not completely clear
 189 whether you end up with arrows being “Markov kernels for general measures” or something else (can
 190 we have negative arrows?). This allows for the introduction of “observables” or “effects” of the form



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

197 3.3 Key questions for Causal Theories

198 We will first define *labeled diagrams*. Rather than labelling the wires of our diagrams with *spaces* (as is
 199 typical [Selinger, 2010]), we assign a unique label to each “wire segment” (with some qualifications).
 200 That is, we assign a unique label to each bare wire in the diagram with the following additional
 201 qualifications:

- 202 • If we have a box in the diagram representing the identity map, the incoming and outgoing
 203 wires are given the same label
- 204 • If we have a wire crossing in the diagram, the diagonally opposite wires are given the same
 205 label
- 206 • The input wire and the *two* output wires of the copy map are given the same label

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

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

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

- 215 • For any sequence of diagram isomorphisms beginning with G_1 and ending with G_2 , we end
 216 up with the same set of labels
- 217 • If we label G_2 according to the rules above then relabel G_1 from G_2 according to the same
 218 rules we retrieve the original labels of G_1

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 ? as an example of a diagrammatic system with labeled wires, I will follow it up

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

- 220 • These properties obviously hold for the wire segments into and out of boxes
- 221 • The only features a diagram may have apart from boxes and wires are wire crossings, copy
- 222 maps and erase maps
- 223 • The labeling rule for wire crossings respects the symmetry of the swap map
- 224 • The labeling rule for copy maps respects the symmetry of the copy map and the property
- 225 described in Equation 8

226 We will follow the convention whereby “internal” wire labels are omitted from diagrams.

227 Note also that each wire that terminates in a free end can be associated with a random variable.
 228 Suppose for $N \in \mathbb{N}$ we have a kernel $\kappa : A \rightarrow \Delta(\times_{i \in N} X_i)$. Define by p_j ($j \in [N]$) the projection
 229 map $p_j : \times_{i \in N} X_i \rightarrow X_j$ defined by $p_j : (x_0, \dots, x_N) \mapsto x_j$. p_j is a measurable function, hence
 230 a random variable. Define by π_j the projection kernel $\mathcal{G}(\pi_j)$ (that is, $\pi_j : \mathbf{x} \mapsto \delta_{p_j(\mathbf{x})}$). Note that
 231 $\kappa(y; p_j^{-1}(A)) = \int_{X_j} \delta_{p_j(\mathbf{x})}(A) \kappa(y; d\mathbf{x}) = \kappa \pi_j$. Diagrammatically, π_j is the identity map on the j -th
 232 wire tensored with the erase map on every other wire. Thus the j -th wire carries the distribution
 233 associated with the random variable p_j . We will therefore consider the labels of the “outgoing” wires
 234 of a diagram to denote random variables (though there are obviously many random variables not
 235 represented by such wires). We will additionally distinguish wire labels from spaces by font - wire
 236 labels are sans serif A, B, C, X, Y, 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

237
 238 **generalised disintegrations** : Of key importance to our work is generalising the notion of disinte-
 239 gration (and possibly Bayesian inversion) to general kernels $X \rightarrow Y$ rather than restricting ourselves
 240 to probability distributions $\{*\} \rightarrow Y$. We will define generalised disintegrations as a straightforward
 241 analogy regular disintegrations, but the conditions under which such disintegrations exist are more
 242 restrictive than for regular disintegrations.

243 **Definition 3.2** (Label signatures). If a kernel $\kappa : X \rightarrow \Delta(Y)$ can be represented by a diagram
 244 G with incoming wires X_1, \dots, X_n and outgoing wires Y_1, \dots, Y_m , we can assign the kernel a “label
 245 signature” $\kappa : X_1 \otimes \dots \otimes X_n \dashrightarrow Y_1 \otimes \dots \otimes Y_m$ or, for short, $\kappa : X_{[n]} \dashrightarrow Y_{[m]}$. Note that this
 246 signature associates each label with a unique space - the space of X_1 is the space associated with the
 247 left-most wire of G and so forth. We will implicitly leverage this correspondence and write with X_1
 248 the space associated with X_1 and so forth. Note that while X_1 is by construction always different from
 249 X_2 (or any other label), the space X_1 may coincide with X_2 - the fact that labels always maintain
 250 distinctions between wires is the fundamental reason for introducing them in the first place.

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

251
 252 **Definition 3.3** (Generalised disintegration). Given a kernel $\kappa : X \rightarrow \Delta(Y)$ with label signature
 253 $\kappa : X_{[n]} \dashrightarrow Y_{[m]}$ and disjoint subsets $S, T \subset [m]$ such that $S \cup T = [m]$, a kernel c is a *g-*
 254 *disintegration from S to T* if it’s type is compatible with the label signature $c : Y_S \dashrightarrow Y_T$ and we
 255 have the identity (omitting incoming wire labels):

$$\begin{array}{c} Y_S \quad Y_T \\ \hline \boxed{\kappa} \end{array} = \begin{array}{c} Y_S \quad Y_T \\ \swarrow \quad \searrow \\ \boxed{C} \\ \swarrow \quad \searrow \\ \boxed{\kappa}^* \end{array} \quad (32)$$

I have introduced without definition additional labeling operations here: first, each label has a particular space associated with it (in order to license the notion of “type compatible with label signature”), and we have supposed labels can be “bundled”.

256

257 In contrast to regular disintegrations, generalised disintegrations “usually” do not exist. Consider
 258 $X = \{0, 1\}$, $Y = \{0, 1\}^2$ and κ has label signature $X_1 \dashrightarrow Y_{\{1,2\}}$ with

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

259 κ imposes contradictory requirements for any disintegration $c : \{0, 1\} \rightarrow \{0, 1\}$ from $\{1\}$ to $\{2\}$:
 260 equality for $X_1 = 1$ requires $c(1; \cdot) = \delta_1$ while equality for $X_1 = 0$ requires $c(1; \cdot) = \delta_0$. Subject
 261 to some regularity conditions (similar to standard Borel conditions for regular disintegrations),
 262 we can define g-disintegrations of a canonically related kernel that do generally exist; intuitively,
 263 g-disintegrations exist if they take the “input wires” of κ as input wires themselves.

264 **Lemma 3.4.** *Given $\kappa : X \rightarrow \Delta(Y)$, a kernel κ^\dagger is a right inverse iff we have for all $x \in X$, $A \in \mathcal{X}$,
 265 $y \in Y$ $\kappa^\dagger(y; A) = \delta_x(A)$, $\kappa(x; \cdot)$ -almost surely.*

266 *Proof.* Suppose κ^\dagger satisfies the almost sure equality for all $x \in X$. Then for all $x \in X$, $A \in \mathcal{X}$ we
 267 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, $\kappa\kappa^\dagger = \text{Id}_X$, so κ^\dagger is
 268 a right inverse of κ .

269 Suppose we have a right inverse κ^\dagger . By definition, for all $x \in X$ and $A \in \mathcal{X}$ we have
 270 $\int_Y \kappa^\dagger(y; A)\kappa(x; dy) = \delta_x(A)$.

271 Suppose $x \notin A$ and let $B_\epsilon = \kappa_A^{\dagger-1}((\epsilon, 1])$ for some $\epsilon > 0$. We have $\int_Y \kappa^\dagger(y; A)\kappa(x; dy) = 0 \geq$
 272 $\epsilon\kappa(x; B_\epsilon)$. Thus for any $\epsilon > 0$ we have $\kappa(x; B_\epsilon) = 0$. Consider the set $B_0 = \kappa_A^{\dagger-1}((0, 1])$. For
 273 some sequence $\{\epsilon_i\}_{i \in \mathbb{N}}$ such that $\lim_{i \rightarrow \infty} \epsilon_i = 0$ we have $B_0 = \cup_{i \in \mathbb{N}} B_{\epsilon_i}$. By countable additivity,
 274 $\kappa(x; B_0) = 0$.

275 Suppose $x \in A$ and let $B^{1-\epsilon} = \kappa_A^{\dagger-1}([0, 1 - \epsilon))$. We have $\int_Y \kappa^\dagger(y; A)\kappa(x; dy) = 1 \leq (1 -$
 276 $\epsilon)\kappa(x; B^{1-\epsilon}) + 1 - \kappa(x; B^{1-\epsilon}) = 1 - \epsilon\kappa(x; B^{1-\epsilon})$. Thus $\kappa(x; B^{1-\epsilon}) = 0$ for $\epsilon > 0$. By an
 277 argument analogous to the above, we also have $\kappa(x; B^1) = 0$. Thus the $\kappa(x; \cdot)$ measure of the set
 278 on which $\kappa^\dagger(y; A)$ disagrees with $\delta_x(A)$ is $\kappa(x; B_0) + \kappa(x; B^1) = 0$ and hence $\kappa^\dagger(y; A) = \delta_x(A)$
 279 $\kappa(x; \cdot)$ -almost surely. \square

I haven't shown that any map inverting κ implies the existence of a Markov kernel that does so

280

I am using countable sets below to get my general argument in order without getting too hung up on measurability; I will try to lift it to standard measurable once it's all there

281

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

$$\begin{array}{c} \text{X} \quad \text{Y} \\ \downarrow \quad \uparrow \\ \boxed{\kappa^\dagger} \\ \downarrow \quad \uparrow \\ \boxed{\kappa} \\ \text{X} \end{array} = \begin{array}{c} \text{X} \quad \text{Y} \\ \downarrow \quad \uparrow \\ \boxed{\kappa} \\ \downarrow \quad \uparrow \\ \boxed{\kappa^\dagger} \\ \text{X} \end{array} \quad (34)$$

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

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

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

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

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

$$= \delta_x(B) \kappa \pi_S(x; A) \quad (39)$$

284 Where ?? follows from Lemma ??.

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

$$= \kappa \pi_S(x; A) \delta_x(B) = L \quad (41)$$

285

□

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

288 via a Markov kernel

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

291 For all $x \in X$ we have a (regular) disintegration $c_x : Y_S \rightarrow \Delta(Y_T)$ of $\kappa(x; \cdot)$ by standard mea-
292 surability of Y . Define $c : X \otimes Y_S \rightarrow \Delta(Y_T)$ by $c : (x, y_S) \mapsto c_x(y_S)$. Clearly, $c(x, y_S)$ is a
293 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
294 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
295 measurability of $c_y(\cdot)^{-1}(B)$ countability of X , so c is a Markov kernel.

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

$$\quad (42)$$

$$\quad (43)$$

297 Which implies

$$\begin{array}{c} Y_S \quad Y_T \\ \downarrow \quad \downarrow \\ \boxed{\kappa} \end{array} = \begin{array}{c} Y_S \quad Y_T \\ \downarrow \quad \downarrow \\ \boxed{C} \\ \downarrow \\ \boxed{\kappa^*} \end{array} \quad (44)$$

298 Finally, we have

$$\begin{array}{c} Y_S \quad Y_T \\ \downarrow \quad \downarrow \\ \boxed{C} \\ \downarrow \\ \boxed{\kappa_S^\dagger} \\ \downarrow \\ \boxed{\kappa^*} \end{array} = \begin{array}{c} Y_S \quad Y_T \\ \downarrow \quad \downarrow \\ \boxed{C} \\ \downarrow \\ \boxed{\kappa_S^\dagger} \\ \downarrow \\ \boxed{\kappa^*} \end{array} \quad (45)$$

$$\begin{array}{c} Y_S \quad Y_T \\ \downarrow \quad \downarrow \\ \boxed{C} \\ \downarrow \\ \boxed{\kappa^*} \end{array} = \begin{array}{c} Y_S \quad Y_T \\ \downarrow \quad \downarrow \\ \boxed{C} \\ \downarrow \\ \boxed{\kappa^*} \end{array} \quad (46)$$

299 Where the first line follows from 7 and the second line from ?? . If κ_S^\dagger is a Markov kernel, then
 300 $\Upsilon(\text{Id}_{Y_S} \otimes \kappa_S^\dagger)c$ is a g-disintegration. \square

301 In the reverse direction, suppose κ is such that $\kappa\pi_T = \text{Id}_X$; that is, π_T is a right inverse of κ . If
 302 $\kappa\pi_S$ is not right invertible then, by definition, there is no d such that $\kappa\pi_S d\pi_T = \text{Id}_X$. However, if a
 303 g-disintegration of κ exists then there is a d such that $\kappa\pi_S d = \kappa$, a contradiction. Thus if $\kappa\pi_S$ is not
 304 right invertible then there is *in general* no g-disintegration from S to T .

305 References

- 306 Kenta Cho and Bart Jacobs. Disintegration and Bayesian inversion via string diagrams.
 307 *Mathematical Structures in Computer Science*, 29(7):938–971, August 2019. ISSN
 308 0960-1295, 1469-8072. doi: 10.1017/S0960129518000488. URL [https://www.](https://www.cambridge.org/core/journals/mathematical-structures-in-computer-science/article/disintegration-and-bayesian-inversion-via-string-diagrams/0581C747DB5793756FE135C70B3B6D51)
 309 [cambridge.org/core/journals/mathematical-structures-in-computer-science/](https://www.cambridge.org/core/journals/mathematical-structures-in-computer-science/article/disintegration-and-bayesian-inversion-via-string-diagrams/0581C747DB5793756FE135C70B3B6D51)
 310 [article/disintegration-and-bayesian-inversion-via-string-diagrams/](https://www.cambridge.org/core/journals/mathematical-structures-in-computer-science/article/disintegration-and-bayesian-inversion-via-string-diagrams/0581C747DB5793756FE135C70B3B6D51)
 311 [0581C747DB5793756FE135C70B3B6D51](https://www.cambridge.org/core/journals/mathematical-structures-in-computer-science/article/disintegration-and-bayesian-inversion-via-string-diagrams/0581C747DB5793756FE135C70B3B6D51).
- 312 Florence Clerc, Fredrik Dahlqvist, Vincent Danos, and Ilias Garnier. Pointless learn-
 313 ing. *20th International Conference on Foundations of Software Science and Compu-*
 314 *tation Structures (FoSSaCS 2017)*, March 2017. doi: 10.1007/978-3-662-54458-7_
 315 21. URL [https://www.research.ed.ac.uk/portal/en/publications/](https://www.research.ed.ac.uk/portal/en/publications/pointless-learning(694fb610-69c5-469c-9793-825df4f8ddec).html)
 316 [pointless-learning\(694fb610-69c5-469c-9793-825df4f8ddec\).html](https://www.research.ed.ac.uk/portal/en/publications/pointless-learning(694fb610-69c5-469c-9793-825df4f8ddec).html).
- 317 Brendan Fong. Causal Theories: A Categorical Perspective on Bayesian Networks. *arXiv:1301.6201*
 318 *[math]*, January 2013. URL <http://arxiv.org/abs/1301.6201>. arXiv: 1301.6201.
- 319 Bart Jacobs. From probability monads to commutative effectuses. *Journal of Logical and*
 320 *Algebraic Methods in Programming*, 94:200–237, January 2018. ISSN 2352-2208. doi:
 321 10.1016/j.jlamp.2016.11.006. URL [http://www.sciencedirect.com/science/article/](http://www.sciencedirect.com/science/article/pii/S2352220816301122)
 322 [pii/S2352220816301122](http://www.sciencedirect.com/science/article/pii/S2352220816301122).

- 323 Bart Jacobs, Aleks Kissinger, and Fabio Zanasi. Causal Inference by String Diagram Surgery. In
324 Mikołaj Bojańczyk and Alex Simpson, editors, *Foundations of Software Science and Computation*
325 *Structures*, Lecture Notes in Computer Science, pages 313–329. Springer International Publishing,
326 2019. ISBN 978-3-030-17127-8.
- 327 Peter Selinger. A survey of graphical languages for monoidal categories. *arXiv:0908.3347 [math]*,
328 813:289–355, 2010. doi: 10.1007/978-3-642-12821-9_4. URL [http://arxiv.org/abs/0908.](http://arxiv.org/abs/0908.3347)
329 3347. arXiv: 0908.3347.