

# Disintegration and Bayesian Inversion via String Diagrams<sup>†</sup>

KENTA CHO and BART JACOBS

*Institute for Computing and Information Sciences, Radboud University*

*P.O.Box 9010, 6500 GL Nijmegen, the Netherlands*

*Email: {K.Cho,bart}@cs.ru.nl*

*Received 14 June 2018*

The notions of disintegration and Bayesian inversion are fundamental in conditional probability theory. They produce channels, as conditional probabilities, from a joint state, or from an already given channel (in opposite direction). These notions exist in the literature, in concrete situations, but are presented here in abstract graphical formulations. The resulting abstract descriptions are used for proving basic results in conditional probability theory. The existence of disintegration and Bayesian inversion is discussed for discrete probability, and also for measure-theoretic probability — via standard Borel spaces and via likelihoods. Finally, the usefulness of disintegration and Bayesian inversion is illustrated in several examples.

## 1. Introduction

The essence of conditional probability can be summarised informally in the following equation about probability distributions:

$$\textit{joint} = \textit{conditional} \cdot \textit{marginal}.$$

A bit more precisely, when we have joint probabilities  $P(x, y)$  for elements  $x, y$  ranging over two sample spaces, the above equation splits into two equations,

$$P(y | x) \cdot P(x) = P(x, y) = P(x | y) \cdot P(y), \quad (1)$$

where  $P(x)$  and  $P(y)$  describe the marginals, which are obtained by discarding variables. We see that conditional probabilities  $P(y | x)$  and  $P(x | y)$  can be constructed in two directions, namely  $y$  given  $x$ , and  $x$  given  $y$ . We also see that we need to copy variables:  $x$  on the left-hand-side of the equations (1), and  $y$  on the right-hand-side.

Conditional probabilities play a crucial role in Bayesian probability theory. They form

<sup>†</sup> The research leading to these results has received funding from the European Research Council under the European Union's Seventh Framework Programme (FP7/2007-2013) / ERC grant agreement nr. 320571.

the nodes of Bayesian networks (Pearl, 1988; Bernardo and Smith, 2000; Barber, 2012), which reflect the conditional independencies of the underlying joint distribution via their graph structure. As part of our approach, we shall capture conditional independence in an abstract manner.

The main notion of this paper is *disintegration*. It is the process of extracting a conditional probability from a joint probability. Disintegration, as we shall formalise it here, gives a structural description of the above equation (1) in terms of *states* and *channels*. In general terms, a *state* is a probability distribution of some sort (discrete, measure-theoretic, or even quantum) and a *channel* is a map or morphism in a probabilistic setting, like  $P(y | x)$  and  $P(x | y)$  as used above. It can take the form of a stochastic matrix, probabilistic transition system, Markov kernel, conditional probability table (in a Bayesian network), or morphism in a Kleisli category of a ‘probability monad’ (Jacobs, 2017). A state is a special kind of channel, with trivial domain. Thus we can work in a monoidal category of channels, where we need discarding and copying — more formally, a comonoid structure on each object — in order to express the above conditional probability equations (1).

In this article we abstract away from interpretation details and will describe disintegration pictorially, in the language of string diagrams. This language can be seen as the internal language of symmetric monoidal categories (Selinger, 2010) — with comonoids in our case. The essence of disintegration becomes: extracting a conditional probability channel from a joint state.

Categorical approaches to Bayesian conditioning have appeared for instance in (Culbertson and Sturtz, 2014; Staton et al., 2016; Clerc et al., 2017) and in (Jacobs et al., 2015; Jacobs and Zanasi, 2016; Jacobs, 2017). The latter references use effectus theory (Jacobs, 2015; Cho et al., 2015), a new comprehensive approach aimed at covering the logic of both quantum theory and probability theory, supported by a Python-based tool *EfProb*, for ‘effectus probability’. This tool is used for the (computationally extensive) examples in this paper.

Disintegration, also known as regular conditional probability, is a notoriously difficult operation in measure-theoretic probability, see *e.g.* (Pollard, 2002; Panangaden, 2009; Chang and Pollard, 1997): it may not exist (Stoyanov, 2014); even if it exists it may be determined only up to negligible sets; and it may not be continuous or computable (Ackerman et al., 2011). Disintegration has been studied using categorical language in (Culbertson and Sturtz, 2014), which focuses on a specific category of probabilistic mappings. Our approach here is more axiomatic.

We thus describe disintegration as going from a joint state to a channel. A closely related concept is *Bayesian inversion*: it turns a channel (with a state) into a channel in opposite direction. We show how Bayesian inversion can be understood and expressed easily in terms of disintegration — and also how, in the other direction, disintegration can be obtained from Bayesian inversion. Bayesian inversion is taken as primitive notion in (Clerc et al., 2017). Here we start from disintegration. The difference is a matter of choice.

Bayesian inversion is crucial for backward inference. We explain it informally: let  $\sigma$  be a state of a domain/type  $X$ , and  $c: X \rightarrow Y$  be a channel; Bayesian inversion yields a

channel  $d: Y \rightarrow X$ . Informally, it produces for an element  $y \in Y$ , seen as singleton/point predicate  $\{y\}$ , the conditioning of the state  $\sigma$  with the pulled back evidence  $c^{-1}(\{y\})$ . A concrete example involving such ‘point observations’ will be described at the end of Section 8. More generally, disintegration and Bayesian inversion are used to structurally organise state updates in the presence of new evidence in probabilistic programming, see *e.g.* (Gordon et al., 2014; Borgström et al., 2013; Staton et al., 2016; Katoen et al., 2015). See also (Shan and Ramsey, 2017), where disintegration is handled via symbolic manipulation.

Disintegration and Bayesian inversion are relatively easy to define in discrete probability theory. The situation is much more difficult in measure-theoretic probability theory, first of all because point predicates  $\{y\}$  do not make much sense there, see also (Chang and Pollard, 1997). A common solution to the problem of the existence of disintegration / Bayesian inversion is to restrict ourselves to standard Borel spaces, as in (Clerc et al., 2017). We take this approach too. There is still an issue that disintegration is determined only up to negligible sets. We address this by defining ‘almost equality’ in our abstract pictorial formulation. This allows us to present a fundamental result from (Clerc et al., 2017) abstractly in our setting, see Section 5.

Another common, more concrete solution is to assume a *likelihood*, that is, a probabilistic relation  $X \times Y \rightarrow \mathbb{R}_{\geq 0}$ . Such a likelihood gives rise to probability density function (pdf), providing a good handle on the situation, see (Pawitan, 2001). The technical core of Section 8 is a generalisation of this likelihood-based approach.

The paper is organised as follows. It starts with a brief introduction to the graphical language that we shall be using, and to the underlying monoidal categories with discarding and copying. Then, Section 3 introduces both disintegration and Bayesian inversion in this graphical language, and relates the two notions. Subsequently, Section 4 contains an elaborated example, namely of naive Bayesian classification. A standard example from the literature (Witten et al., 2011) is redescribed in the current setting: first, channels are extracted via disintegration from a table with given data; next, Bayesian inversion is applied to the combined extracted channels, giving the required classification. This is illustrated in both the discrete and the continuous version of the example.

Next, Section 5 is more technical and elaborates the standard equality notion of ‘equal almost everywhere’ in the current setting. This is used for describing Bayesian inversion in a more formal way, following (Clerc et al., 2017). Section 6 uses our graphical approach to review conditional independence and to prove at an abstract level several known results, namely the equivalence of various formulations of conditional independence, and the ‘graphoid’ axioms from (Verma and Pearl, 1988; Geiger et al., 1990). Section 7 relaxes the requirement that maps are causal, so that ‘effects’ can be used as the duals of states for validity and conditioning. The main result relates conditioning of joint states to forward and backward inference via the extracted channels, in the style of (Jacobs and Zanasi, 2016; Jacobs and Zanasi, 2018); it is illustrated in a concrete example, where a Bayesian network is seen as a graph in a Kleisli category — following (Fong, 2012). Finally, Section 8 gives the likelihood formulation of disintegration and inversion, as briefly described above.

## 2. Graphical language

The basic idea underlying this paper is to describe probability theory in terms of *channels*. A channel  $f: X \rightarrow Y$  is a (stochastic) process from a system of type  $X$  into that of  $Y$ . Concretely, it may be a probability matrix or kernel. Our standing assumption is that types (as objects) and channels (as arrows) form a *symmetric monoidal category*. For the formal definition we refer to (Mac Lane, 1998). We informally summarise that we have the following constructions.

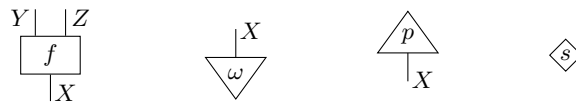
- 1 Sequential composition  $g \circ f: X \rightarrow Z$  for appropriately typed channels  $f: X \rightarrow Y$  and  $g: Y \rightarrow Z$ .
- 2 Parallel composition  $f \otimes g: X \otimes Z \rightarrow Y \otimes W$  for  $f: X \rightarrow Y$  and  $g: Z \rightarrow W$ . This involves composition of types  $X \otimes Z$ .
- 3 Identity channels  $\text{id}_X: X \rightarrow X$ , which ‘do nothing’. Thus  $\text{id} \circ f = f = \text{id} \circ f$ .
- 4 A unit type  $I$ , which represents ‘no system’. Thus  $I \otimes X \cong X \cong X \otimes I$ .
- 5 Swap isomorphisms  $X \otimes Y \cong Y \otimes X$  and associativity isomorphisms  $(X \otimes Y) \otimes Z \cong X \otimes (Y \otimes Z)$ , so the ordering in composed types does not matter.

The representation of such channels in the ordinary ‘formula’ notation easily becomes complex and thus reasoning becomes hard to follow. A graphical language known as *string diagrams* offers a more convenient and intuitive way of reasoning in a symmetric monoidal category.

In string diagrams, types/objects are represented as wires ‘|’, with information flowing bottom to top. The composition of types is depicted by juxtaposition of wires, and the unit type is ‘no diagram’ as below.

$$X \otimes Y = \left| \begin{array}{c} X \\ Y \end{array} \right| \qquad I = \boxed{\phantom{X}}$$

Channels/arrows are represented by boxes with an input wire(s) and an output wire(s), in upward direction. When a box does not have input or output, we write it as a triangle or diamond. For example,  $f: X \rightarrow Y \otimes Z$ ,  $\omega: I \rightarrow X$ ,  $p: X \rightarrow I$ , and  $s: I \rightarrow I$  are respectively depicted as:



The identity channels are represented by ‘no box’, *i.e.* just wires, and the swap isomorphisms are represented by crossing of wires:

$$\boxed{\text{id}} = \left| \begin{array}{c} X \\ X \end{array} \right| \qquad \begin{array}{c} Y \quad X \\ \diagdown \quad \diagup \\ X \quad Y \end{array}$$

Finally, the sequential composition of channels is depicted by connecting the input and

output wires, and the parallel composition is given by juxtaposition, respectively as below:

$$\begin{array}{c} | \\ \boxed{g \circ f} \\ | \end{array} = \begin{array}{c} | \\ \boxed{g} \\ | \\ \boxed{f} \\ | \end{array} \qquad \begin{array}{c} | \\ \boxed{h \otimes k} \\ | \end{array} = \begin{array}{c} | \\ \boxed{h} \\ | \end{array} \begin{array}{c} | \\ \boxed{k} \\ | \end{array}$$

The use of string diagrams is justified by the following ‘coherence’ theorem; see (Joyal and Street, 1991; Selinger, 2010) for details.

**Theorem 2.1.** A well-formed equation between composites of arrows in a symmetric monoidal category follows from the axioms of symmetric monoidal categories if and only if the string diagrams of both sides are equal up to isomorphism of diagrams.  $\square$

We further assume the following structure in our category. For each type  $X$  there are a discarder  $\bar{\dagger}_X: X \rightarrow I$  and a copier  $\Upsilon_X: X \rightarrow X \otimes X$ . They are required to satisfy the following equations:

$$\begin{array}{c} \bar{\dagger} \\ \cup \\ \bullet \\ | \end{array} = | = \begin{array}{c} | \\ \bullet \\ \cup \\ \bar{\dagger} \end{array} \qquad \begin{array}{c} \cup \\ \bullet \\ | \end{array} = \begin{array}{c} \cup \\ \bullet \\ | \end{array} \qquad \begin{array}{c} \cup \\ \cup \\ \bullet \\ | \end{array} = \begin{array}{c} \cup \\ \cup \\ \bullet \\ | \end{array} \qquad \begin{array}{c} \cup \\ \cup \\ \bullet \\ | \end{array} = \begin{array}{c} \cup \\ \bullet \\ | \end{array}$$

This says that  $(\Upsilon_X, \bar{\dagger}_X)$  forms a commutative comonoid on  $X$ . By the associativity we may write:

$$\begin{array}{c} \cup \\ \cup \\ \cup \\ \dots \\ \bullet \\ | \end{array} := \begin{array}{c} \cup \\ \cup \\ \cup \\ \dots \\ \bullet \\ | \end{array}$$

Moreover we assume that the comonoid structures  $(\Upsilon_X, \bar{\dagger}_X)$  are compatible with the monoidal structure  $(\otimes, I)$ , in the sense that the following equations hold.

$$\begin{array}{c} \bar{\dagger} \\ | \\ X \otimes Y \end{array} = \begin{array}{c} \bar{\dagger} \\ | \\ X \end{array} \begin{array}{c} \bar{\dagger} \\ | \\ Y \end{array} \qquad \begin{array}{c} \bar{\dagger} \\ | \\ I \end{array} = \boxed{\phantom{I}} \qquad \begin{array}{c} \cup \\ \cup \\ \bullet \\ | \\ X \otimes Y \end{array} = \begin{array}{c} \cup \\ \cup \\ \bullet \\ | \\ X \end{array} \begin{array}{c} \cup \\ \cup \\ \bullet \\ | \\ Y \end{array} \qquad \begin{array}{c} \cup \\ \bullet \\ | \\ I \end{array} = \boxed{\phantom{I}}$$

Note that we do *not* assume that these maps are natural. Explicitly, we do not necessarily have  $\Upsilon \circ f = (f \otimes f) \circ \Upsilon$  or  $\bar{\dagger} \circ f = \bar{\dagger}$ .

We will use these symmetric monoidal categories throughout in the paper. For convenience, we introduce a term for them.

**Definition 2.2.** A *CD-category* is a symmetric monoidal category  $(\mathbf{C}, \otimes, I)$  with a commutative comonoid  $(\Upsilon_X, \bar{\dagger}_X)$  for each  $X \in \mathbf{C}$ , suitably compatible as described above.

Here ‘CD’ stands for Copy/Discard.

**Definition 2.3.** An arrow  $f: X \rightarrow Y$  in a CD-category is said to be *causal* if

$$\begin{array}{c} \bar{\dagger} \\ | \\ \boxed{f} \\ | \end{array} = \begin{array}{c} \bar{\dagger} \\ | \end{array} .$$

A CD-category is *affine* if all the arrows are causal, or equivalently, the tensor unit  $I$  is a final object.

The term ‘causal’ comes from (categorical) quantum foundation (Coecke and Kissinger, 2017; D’Ariano et al., 2017), and is related to relativistic causality, see *e.g.* (Coecke, 2016).

We reserve the term ‘channel’ for causal arrows. Explicitly, causal arrows  $c: X \rightarrow Y$  in a CD-category are called *channels*. A channel  $\omega: I \rightarrow X$  with input type  $I$  is called a *state* (on  $X$ ). For the time being we will only use channels (and states), and thus only consider affine CD-categories. Non-causal arrows will not appear until Section 7.

Our main examples of affine CD-categories are two Kleisli categories  $\mathcal{Kl}(\mathcal{D})$  and  $\mathcal{Kl}(\mathcal{G})$ , respectively, for discrete probability, and more general, measure-theoretic probability. We explain them in order below. There are more examples, like the Kleisli category of the non-empty powerset monad, or the (opposite of) the category of commutative  $C^*$ -algebras (with positive unital maps), but they are out of scope here.

**Example 2.4.** What we call a *distribution* or a *state* over a set  $X$  is a finite subset  $\{x_1, x_2, \dots, x_n\} \subseteq X$ , called the support where each element  $x_i$  occurs with a multiplicity  $r_i \in [0, 1]$ , such that  $\sum_i r_i = 1$ . Such a convex combination is often written as  $r_1|x_1\rangle + \dots + r_n|x_n\rangle$  with  $r_i \in [0, 1]$ . The ket notation  $|-\rangle$  is meaningless syntactic sugar that is used to distinguish elements  $x \in X$  from occurrences in such formal sums. Notice that a distribution can also be written as a function  $\omega: X \rightarrow [0, 1]$  with finite support  $\text{supp}(\omega) = \{x \in X \mid \omega(x) \neq 0\}$ . We shall write  $\mathcal{D}(X)$  for the set of distributions over  $X$ . This  $\mathcal{D}$  is a monad on the category **Set** of sets and functions.

A function  $f: X \rightarrow \mathcal{D}(Y)$  is called a *Kleisli* map; it forms a channel  $X \rightarrow Y$ . Such maps can be composed as matrices, for which we use special notation  $\circ$ .

$$(g \circ f)(x)(z) = \sum_{y \in Y} f(x)(y) \cdot g(y)(z) \quad \text{for } g: Y \rightarrow \mathcal{D}(Z) \text{ with } x \in X, z \in Z.$$

We write  $1 = \{*\}$  for a singleton set, and  $2 = 1 + 1 = \{0, 1\}$ . Notice that  $\mathcal{D}(1) \cong 1$  and  $\mathcal{D}(2) \cong [0, 1]$ . We can identify a state on  $X$  with a channel  $1 \rightarrow X$ .

The monad  $\mathcal{D}$  is known to be *commutative*. This implies that finite products of sets  $X \times Y$  give rise to a symmetric monoidal structure on the Kleisli category  $\mathcal{Kl}(\mathcal{D})$ . Specifically, for two maps  $f: X \rightarrow \mathcal{D}(Y)$  and  $g: Z \rightarrow \mathcal{D}(W)$ , the tensor product / parallel composition  $f \otimes g: X \times Z \rightarrow \mathcal{D}(Y \times W)$  is given by:

$$(f \otimes g)(x, z)(y, w) = f(x, y) \cdot g(z, w).$$

For each set  $X$  there are a copier  $\Upsilon_X: X \rightarrow \mathcal{D}(X \times X)$  and a discarder  $\bar{\Upsilon}_X: X \rightarrow \mathcal{D}(1)$  given by  $\Upsilon_X(x) = 1|x, x\rangle$  and  $\bar{\Upsilon}_X(x) = 1|*\rangle$ , respectively. They come from the cartesian (finite product) structure of the base category **Set**, through the obvious functor **Set**  $\rightarrow$   $\mathcal{Kl}(\mathcal{D})$ . Therefore  $\mathcal{Kl}(\mathcal{D})$  is a CD-category. It is moreover affine, since the monad is affine in the sense that  $\mathcal{D}(1) \cong 1$ .

**Example 2.5.** Let  $X = (X, \Sigma_X)$  be a measurable space, where  $\Sigma_X$  is a  $\sigma$ -algebra on  $X$ . A *probability measure*, also called a *state*, on  $X$  is a function  $\omega: \Sigma_X \rightarrow [0, 1]$  which is countably additive and satisfies  $\omega(X) = 1$ . We write  $\mathcal{G}(X)$  for the collection of all such probability measures on  $X$ . This set  $\mathcal{G}(X)$  is itself a measurable space with the smallest  $\sigma$ -algebra such that for each  $A \in \Sigma_X$  the ‘evaluation’ map  $\text{ev}_A: \mathcal{G}(X) \rightarrow [0, 1]$ ,  $\text{ev}_A(\omega) = \omega(A)$ , is measurable. Notice that  $\mathcal{G}(X) \cong \mathcal{D}(X)$  when  $X$  is a finite set (as

discrete space). In particular,  $\mathcal{G}(2) \cong \mathcal{D}(2) \cong [0, 1]$ . This  $\mathcal{G}$  is a monad on the category **Meas** of measurable spaces, with measurable functions between them; it is called the Giry monad, after (Giry, 1982).

A Kleisli map, that is, a measurable function  $f: X \rightarrow \mathcal{G}(Y)$  is a channel (or a *probability kernel*, see Example 7.2). These channels can be composed, via Kleisli composition  $\circ$ , using integration<sup>†</sup>:

$$(g \circ f)(x)(C) = \int_Y g(y)(C) f(x)(dy) \quad \text{where } g: Y \rightarrow \mathcal{G}(Z) \text{ and } x \in X, C \in \Sigma_Z.$$

It is well-known that the monad  $\mathcal{G}$  is commutative and affine, see also (Jacobs, 2017). Thus, in a similar manner to the previous example, the Kleisli category  $\mathcal{Kl}(\mathcal{G})$  is an affine CD-category. The parallel composition  $f \otimes g: X \times Z \rightarrow \mathcal{G}(Y \times W)$  for  $f: X \rightarrow \mathcal{G}(Y)$  and  $g: Z \rightarrow \mathcal{G}(W)$  is given as:

$$(f \otimes g)(x, z)(B \times D) = f(x)(B) \cdot g(z)(D),$$

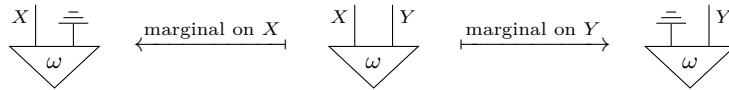
for  $x \in X, z \in Z, B \in \Sigma_Y$ , and  $D \in \Sigma_W$ . Since  $f(x)$  and  $g(z)$  are ( $\sigma$ -)finite measures, this indeed determines a unique measure  $(f \otimes g)(x, z) \in \mathcal{G}(Y \times W)$ , namely the unique product measure of  $f(x)$  and  $g(z)$ . Explicitly, for  $E \in \Sigma_{Y \times W}$ ,

$$(f \otimes g)(x, z)(E) = \int_W \left( \int_Y \mathbf{1}_E(y, w) f(x)(dy) \right) g(z)(dw)$$

where  $\mathbf{1}_E$  is the indicator function.

### 3. Marginalisation, integration and disintegration

Let  $\mathbf{C}$  be an affine CD-category. We think of states  $\omega: I \rightarrow X$  in  $\mathbf{C}$  as abstract (probability) distributions on type  $X$ . States of the form  $\omega: I \rightarrow X \otimes Y$ , often called (bipartite) joint states, are seen as joint distributions on  $X$  and  $Y$ . Later on we shall also consider  $n$ -partite joint states, but for the time being we restrict ourselves to bipartite ones. For a joint distribution  $P(x, y)$  in discrete probability, we can calculate the marginal distribution on  $X$  by summing (or marginalising)  $Y$  out, as  $P(x) = \sum_y P(x, y)$ . The marginal distribution on  $Y$  is also calculated by  $P(y) = \sum_x P(x, y)$ . In our abstract setting, given a joint state  $\omega: I \rightarrow X \otimes Y$ , we can obtain marginal states simply by discarding wires, as in:



In other words, the marginal states are the state  $\omega$  composed with the projection maps  $\pi_1: X \otimes Y \rightarrow X$  and  $\pi_2: X \otimes Y \rightarrow Y$ , as below.

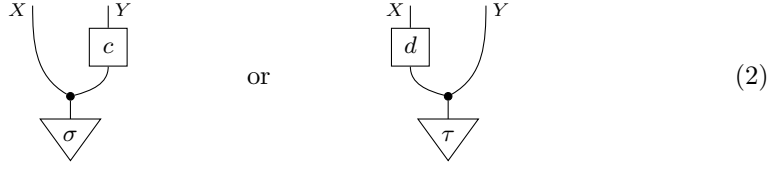
$$\pi_1 := \begin{array}{c} | \\ X \quad | \\ \overline{\overline{\quad}} \end{array} \quad \pi_2 := \begin{array}{c} \overline{\overline{\quad}} \\ | \\ X \quad | \\ Y \end{array}$$

<sup>†</sup> We denote the integral of a function  $f$  with respect to a measure  $\mu$  by  $\int_X f(x) \mu(dx)$ .

**Example 3.1.** For a joint state  $\omega \in \mathcal{D}(X \times Y)$  in  $\mathcal{Kl}(\mathcal{D})$ , the first marginal  $\omega_1 = \pi_1 \circ \omega$  is given by  $\omega_1(x) = \sum_{y \in Y} \omega(x, y)$ , as expected. Similarly the second marginal is given by  $\omega_2(y) = \sum_{x \in X} \omega(x, y)$ .

**Example 3.2.** For a joint state  $\omega \in \mathcal{G}(X \times Y)$  in  $\mathcal{Kl}(\mathcal{G})$ , the first marginal is given by  $\omega_1(A) = \omega(A \times Y)$  for  $A \in \Sigma_X$ , and the second marginal by  $\omega_2(B) = \omega(X \times B)$  for  $B \in \Sigma_Y$ .

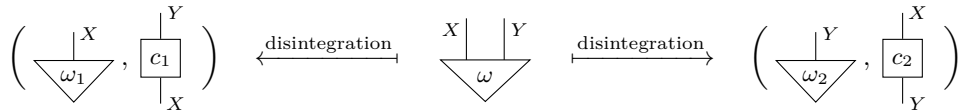
A channel  $c: X \rightarrow Y$  is seen as an abstract *conditional* distribution  $P(y|x)$ . In (discrete) probability theory, we can calculate a joint distribution  $P(x, y)$  from a distribution  $P(x)$  and a conditional distribution  $P(y|x)$  by the formula  $P(x, y) = P(y|x) \cdot P(x)$ , which is often called the *product rule*. Similarly we have  $P(x, y) = P(x|y) \cdot P(y)$ . In our setting, starting from a state  $\sigma: I \rightarrow X$  and a channel  $c: X \rightarrow Y$ , or a state  $\tau: I \rightarrow Y$  and a channel  $d: Y \rightarrow X$ , we can ‘integrate’ them into a joint state on  $X \otimes Y$  as follows, respectively:



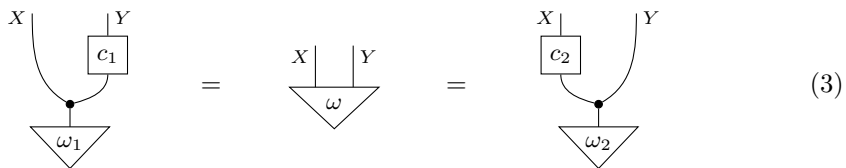
**Example 3.3.** Let  $\sigma \in \mathcal{D}(X)$  and  $c: X \rightarrow \mathcal{D}(Y)$  be a state and a channel in  $\mathcal{Kl}(\mathcal{D})$ . An easy calculation verifies that  $\omega = (\text{id} \otimes c) \circ \Upsilon \circ \sigma$ , the joint state on  $X \times Y$  defined as in (2), satisfies  $\omega(x, y) = c(x)(y) \cdot \sigma(x)$ , as we expect from the product rule.

**Example 3.4.** For a state  $\sigma \in \mathcal{G}(X)$  and a channel  $c: X \rightarrow \mathcal{G}(Y)$  in  $\mathcal{Kl}(\mathcal{G})$ , the joint state  $\omega = (\text{id} \otimes c) \circ \Upsilon \circ \sigma$  is given by  $\omega(A \times B) = \int_A c(x)(B) \sigma(dx)$  for  $A \in \Sigma_X$  and  $B \in \Sigma_Y$ . This ‘integration’ construction of a joint probability measure is standard, see e.g. (Pollard, 2002; Panangaden, 2009).

*Disintegration* is an inverse operation of the ‘integration’ of a state and a channel into a joint state, as in (2). More specifically, it starts from a joint state  $\omega: I \rightarrow X \otimes Y$  and extracts either a state  $\omega_1: I \rightarrow X$  and a channel  $c_1: X \rightarrow Y$ , or a state  $\omega_2: I \rightarrow Y$  and a channel  $c_2: Y \rightarrow X$  as below,



such that the equation on the left or right below holds, respectively.





We immediately see from the equation that  $\omega_1$  and  $\omega_2$  must be marginals of  $\omega$ :

and similarly

Therefore a disintegration of  $\omega$  may be referred to by a channel  $c_i$  only, rather than a pair  $(\omega_i, c_i)$ . This leads to the following definition:

**Definition 3.5.** Let  $\omega: I \rightarrow X \otimes Y$  be a joint state. A channel  $c_1: X \rightarrow Y$  (or  $c_2: Y \rightarrow X$ ) is called a *disintegration* of  $\omega$  if it satisfies the equation (3) with  $\omega_i$  the marginals of  $\omega$ .

Let us look at concrete instances of disintegrations, in the Kleisli categories  $\mathcal{Kl}(\mathcal{D})$  and  $\mathcal{Kl}(\mathcal{G})$  from Examples 2.4 and 2.5.

**Example 3.6.** Let  $\omega \in \mathcal{D}(X \times Y)$  be a joint state in  $\mathcal{Kl}(\mathcal{D})$ . We write  $\omega_1 \in \mathcal{D}(X)$  for the first marginal, given by  $\omega_1(x) = \sum_y \omega(x, y)$ . Then a channel  $c: X \rightarrow \mathcal{D}(Y)$  is a disintegration of  $\omega$  if and only if  $\omega(x, y) = c(x)(y) \cdot \omega_1(x)$  for all  $x \in X$  and  $y \in Y$ . It turns out that there is always such a channel  $c$ . We define a channel  $c$  by:

$$c(x)(y) := \frac{\omega(x, y)}{\omega_1(x)} \quad \text{if } \omega_1(x) \neq 0, \quad (4)$$

and  $c(x) := \tau$  if  $\omega_1(x) = 0$ , for an arbitrary state  $\tau \in \mathcal{D}(Y)$ . (Here  $\mathcal{D}(Y)$  is nonempty, since so is  $\mathcal{D}(X \times Y)$ .) This indeed defines a channel  $c$  satisfying the required equation. Roughly speaking, disintegration in discrete probability is nothing but the ‘definition’ of conditional probability:  $P(y|x) = P(x, y)/P(x)$ . There is still some subtlety — disintegrations need not be unique, when there are  $x \in X$  with  $\omega_1(x) = 0$ . They are, nevertheless, ‘almost surely’ unique; see Section 5.

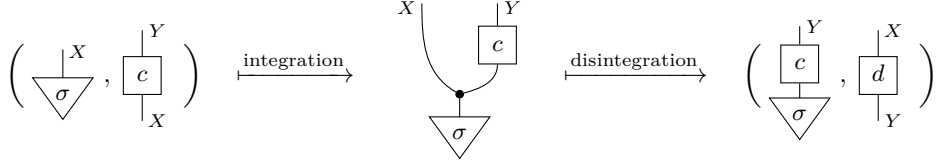
**Example 3.7.** Disintegrations in measure-theoretic probability, in  $\mathcal{Kl}(\mathcal{G})$ , are far more difficult. Let  $\omega \in \mathcal{G}(X \times Y)$  be a joint state, with  $\omega_1 \in \mathcal{G}(X)$  the first marginal. A channel  $c: X \rightarrow \mathcal{G}(Y)$  is a disintegration of  $\omega$  if and only if

$$\omega(A \times B) = \int_A c(x)(B) \omega_1(dx)$$

for all  $A \in \Sigma_X$  and  $B \in \Sigma_Y$ . This is the ordinary notion of *disintegration* (of probability measures), also known as *regular conditional probability*; see *e.g.* (Faden, 1985; Pollard, 2002; Panangaden, 2009). We see that there is no obvious way to obtain a channel  $c$  here, unlike the discrete case. In fact, a disintegration may not exist (Stoyanov, 2014). There are, however, a number of results that guarantee the existence of a disintegration in certain situations. We will come back to this issue later in the section.

*Bayesian inversion* is a special form of disintegration, occurring frequently. We start

from a state  $\sigma: I \rightarrow X$  and a channel  $c: X \rightarrow Y$ . We then integrate them into a joint state on  $X$  and  $Y$ , and disintegrate it in the other direction, as below.



We call the disintegration  $d: Y \rightarrow X$  a *Bayesian inversion* for  $\sigma: I \rightarrow X$  along  $c: X \rightarrow Y$ . By unfolding the definitions, a channel  $d: Y \rightarrow X$  is a Bayesian inversion if and only if

(5)

The composite  $c \circ \sigma$  is also written as  $c_*(\sigma)$ . The operation  $\sigma \mapsto c_*(\sigma)$ , called *state transformation*, is used to explain forward inference in (Jacobs and Zanasi, 2016).

**Example 3.8.** Let  $\sigma \in \mathcal{D}(X)$  and  $c: X \rightarrow \mathcal{D}(Y)$  be a state and a channel in  $\mathcal{Kl}(\mathcal{D})$ . Then a channel  $d: Y \rightarrow \mathcal{D}(X)$  is a Bayesian inversion for  $\sigma$  along  $c$  if and only if  $c(x)(y) \cdot \sigma(x) = d(y)(x) \cdot c_*(\sigma)(y)$ , where  $c_*(\sigma)(y) = \sum_{x'} c(x')(y) \cdot \sigma(x')$ . In a similar manner to Example 3.6, we can obtain such a  $d$  by:

$$d(y)(x) := \frac{c(x)(y) \cdot \sigma(x)}{c_*(\sigma)(y)} = \frac{c(x)(y) \cdot \sigma(x)}{\sum_{x'} c(x')(y) \cdot \sigma(x')}$$

for  $y \in Y$  with  $c_*(\sigma)(y) \neq 0$ . For  $y \in Y$  with  $c_*(\sigma)(y) = 0$ , we may define  $d(y)$  to be an arbitrary state in  $\mathcal{D}(X)$ . We can recognise the above formula as the Bayes formula:

$$P(x|y) = \frac{P(y|x) \cdot P(x)}{P(y)} = \frac{P(y|x) \cdot P(x)}{\sum_{x'} P(y|x') \cdot P(x')} .$$

**Example 3.9.** Let  $\sigma \in \mathcal{G}(X)$  and  $c: X \rightarrow \mathcal{G}(Y)$  be a state and a channel in  $\mathcal{Kl}(\mathcal{G})$ . A channel  $d: Y \rightarrow \mathcal{G}(X)$  is a Bayesian inversion if and only if

$$\int_A c(x)(B) \sigma(dx) = \int_B d(y)(A) c_*(\sigma)(dy)$$

for all  $A \in \Sigma_X$  and  $B \in \Sigma_Y$ . Here  $c_*(\sigma) \in \mathcal{G}(Y)$  is the measure given by  $c_*(\sigma)(B) = \int_X c(x)(B) \sigma(dx)$ . As we see below, Bayesian inversions are in some sense equivalent to disintegrations, and thus, they are as difficult as disintegrations. In particular, a Bayesian inversion need not exist.

In practice, however, the state  $\sigma$  and channel  $c$  are often given via density functions. This setting, so-called (absolutely) continuous probability, makes it easy to compute a Bayesian inversion. Suppose that  $X$  and  $Y$  are subspaces of  $\mathbb{R}$ , and that  $\sigma$  and  $c$  admit

density functions as

$$\sigma(A) = \int_A f(x) dx \qquad c(x)(B) = \int_B \ell(x, y) dy$$

for measurable functions  $f: X \rightarrow \mathbb{R}_{\geq 0}$  and  $\ell: X \times Y \rightarrow \mathbb{R}_{\geq 0}$ . The conditional probability density  $\ell(x, y)$  of  $y$  given  $x$  is often called the *likelihood* of  $x$  given  $y$ . By the familiar Bayes formula for densities — see *e.g.* (Bernardo and Smith, 2000) — the conditional density of  $x$  given  $y$  is:

$$k(y, x) := \frac{\ell(x, y) \cdot f(x)}{\int_X \ell(x', y) \cdot f(x') dx'}$$

This  $k$  then gives a channel  $d: Y \rightarrow \mathcal{G}(X)$  by

$$d(y)(A) = \int_A k(y, x) dx$$

for each  $y \in Y$  such that  $\int_X \ell(x', y) \cdot f(x') dx' \neq 0$ . For the other  $y$ 's we define  $d(y)$  to be some fixed state in  $\mathcal{G}(X)$ . An elementary calculation verifies that  $d$  is indeed a Bayesian inversion for  $\sigma$  along  $c$ . Later, in Section 8, we generalise this calculation into our abstract setting.

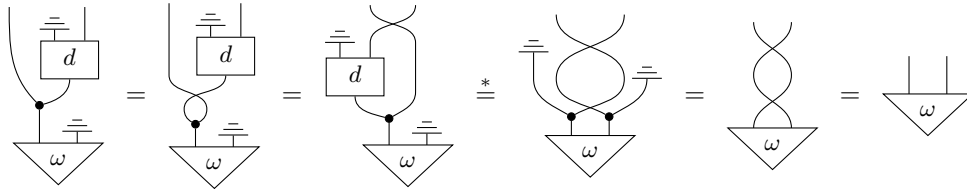
Although Bayesian inversions are a special case of disintegrations, we can conversely obtain disintegrations from Bayesian inversions, as in the proposition below. Therefore, in some sense the two notions are equivalent.

**Proposition 3.10.** Let  $\omega$  be a state on  $X \otimes Y$ . Let  $d: X \rightarrow X \otimes Y$  be a Bayesian inversion for  $\omega$  along the first projection  $\pi_1: X \otimes Y \rightarrow X$  on the left below.

$$\pi_1 = \begin{array}{c} \text{---} \\ | \\ X \text{---} \end{array} \Bigg| \begin{array}{c} \text{---} \\ | \\ Y \text{---} \end{array} \qquad \pi_2 \circ d = \begin{array}{c} \text{---} \\ | \\ \boxed{d} \\ | \\ X \text{---} \end{array} \Bigg| \begin{array}{c} \text{---} \\ | \\ Y \text{---} \end{array}$$

Then the composite  $\pi_2 \circ d: X \rightarrow Y$  shown on the right above is a disintegration of  $\omega$ .

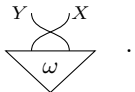
*Proof.* We prove that the first equation in (3) holds for  $c_1 = \pi_2 \circ d$ , as follows.



For the marked equality  $\stackrel{*}{=}$  we used the equation (5) for the Bayesian inversion  $d$ .  $\square$

We say that an affine CD-category  $\mathbf{C}$  admits *disintegration* if for every bipartite state  $\omega: I \rightarrow X \otimes Y$  there exist a disintegration  $c_1: X \rightarrow Y$  of  $\omega$ . Note that in such categories there also exists a disintegration  $c_2: Y \rightarrow X$  of  $\omega$  in the other direction, since it can be

obtained as a disintegration of the following state:



By Proposition 3.10, admitting disintegration is equivalent to admitting Bayesian inversion.

In Example 3.6, we have seen that  $\mathcal{Kl}(\mathcal{D})$  admits disintegration, but that in measure-theoretic probability, in  $\mathcal{Kl}(\mathcal{G})$ , disintegrations may not exist. There are however a number of results that guarantee the existence of disintegrations in specific situations, see *e.g.* (Pachl, 1978; Faden, 1985). We here invoke one of these results and show that there is a subcategory of  $\mathcal{Kl}(\mathcal{G})$  that admits disintegration. A measurable space is called a *standard Borel space* if it is measurably isomorphic to a Polish space with its Borel  $\sigma$ -algebra, or equivalently, if it is measurably isomorphic to a Borel subspace of  $\mathbb{R}$ . Then the following theorem is standard, see *e.g.* (Pollard, 2002, §5.2) or (Faden, 1985, §5).

**Theorem 3.11.** Let  $X$  be any measurable space and  $Y$  be a standard Borel space. Then for any state (*i.e.* a probability measure)  $\omega \in \mathcal{G}(X \times Y)$  in  $\mathcal{Kl}(\mathcal{G})$ , there exists a disintegration  $c_1: X \rightarrow \mathcal{G}(Y)$  of  $\omega$ . □

Let  $\mathbf{pKrn}_{\text{sb}}$  be the full subcategory of  $\mathcal{Kl}(\mathcal{G})$  consisting of standard Borel spaces as objects. It is not hard to see that the product of two standard Borel spaces is a standard Borel space. It follows that  $\mathbf{pKrn}_{\text{sb}}$  is an affine CD-category. Then the previous theorem immediately shows:

**Corollary 3.12.** The category  $\mathbf{pKrn}_{\text{sb}}$  admits disintegration. □

We note that  $\mathbf{pKrn}_{\text{sb}}$  can also be seen as the Kleisli category of the Giry monad restricted on the category of standard Borel spaces, see *e.g.* (Giry, 1982; Doberkat, 2009).

Since there are various ‘existence’ theorems like Theorem 3.11, there may be other subcategories of  $\mathcal{Kl}(\mathcal{G})$  that admit disintegration. A likely candidate is the category of perfect probabilistic mappings in (Culbertson and Sturtz, 2014). We do not go into this question here, since  $\mathbf{pKrn}_{\text{sb}}$  suffices for the present paper.

#### 4. Example: naive Bayesian classifiers via inversion

Bayesian classification is a well-known technique in machine learning that produces a distribution over data classifications, given certain sample data. The distribution describes the probability, for each data (classification) category, that the sample data is in that category. Here we consider an example of ‘naive’ Bayesian classification, where the features are assumed to be independent. We consider a standard classification example from the literature which forms an ideal setting to illustrate the use of both disintegration and Bayesian inversion. Disintegration is used to extract channels from a given table, and subsequently Bayesian inversion is applied to (the tuple of) these channels to obtain the actual classification. The use of channels and disintegration/inversion in this classification setting is new, as far as we know.

For the description of the relevant operations in this example we use notation for

marginalisation and disintegration that we borrowed from the *EfProb* library (Cho and Jacobs, 2017). There are many ways to marginalise an  $n$ -partite state, namely one for each subset of the wires  $\{1, 2, \dots, n\}$ . Such a subset can be described as a *mask*, consisting of a list of  $n$  zero's or one's, where a zero at position  $i$  means that the  $i$ -th wire/component is marginalised out, and a one at position  $i$  means that it remains. Such a mask  $M = [b_1, \dots, b_n]$  with  $b_i \in \{0, 1\}$  is used as a post-fix selection operation in  $\omega.M$  on an  $n$ -partite state  $\omega$ . An example explains it all:

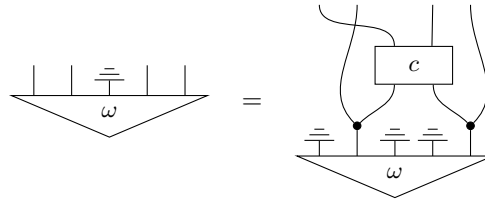


In a similar way one can disintegrate an  $n$ -partite state in  $2^n$  many ways, where a mask of length  $n$  is now used to describe which wires are used as input to the extracted channel and which ones as output. We write  $\omega // M$  for such a disintegration, where  $M$  is a mask, as above. A systematic description will be given in Section 6 below.

In practice it is often useful to be able to marginalise first, and disintegrate next. The general description in  $n$ -ary form is a bit complicated, so we use an example for  $n = 5$ . We shall label the wires with  $x_i$ , as on the left below. We seek the conditional probability written conventionally as  $c = \omega[x_1, x_4 \mid x_2, x_5]$  on the right below.



This channel  $c$  must satisfy:



This picture shows how to obtain the channel  $c$  from  $\omega$ : we first marginalise to restrict to the relevant wires  $x_1, x_2, x_4, x_5$ . This is written as  $\omega.[1, 1, 0, 1, 1]$ . Subsequently we disintegrate with  $x_1, x_4$  as output and  $x_2, x_5$  as input. Hence:

$$c := \omega.[1, 1, 0, 1, 1] // [0, 1, 0, 1]$$

$$=: \omega.[[1, 0, 0, 1, 0] \mid [0, 1, 0, 0, 1]] \quad \text{as we shall write in the sequel.}$$

We see that the latter post-fix  $[[1, 0, 0, 1, 0] \mid [0, 1, 0, 0, 1]]$  is a ‘variable free’ version of the traditional notation  $[x_1, x_4 \mid x_2, x_5]$ , selecting the relevant positions.

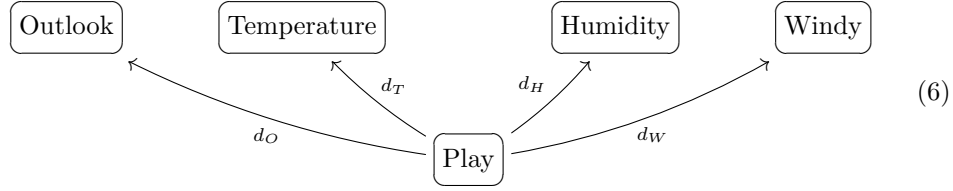
We have now prepared the ground and can turn to the classification example that we announced. It involves the classification of ‘playing’ (yes or no) for certain weather data, used in (Witten et al., 2011). We shall first go through the discrete example in some detail. The relevant data are in the table in Figure 1. The question is: given this table, what can

Outlook	Temperature	Humidity	Windy	Play
Sunny	hot	high	false	no
Sunny	hot	high	true	no
Overcast	hot	high	false	yes
Rainy	mild	high	false	yes
Rainy	cool	normal	false	yes
Rainy	cool	normal	true	no
Overcast	cool	normal	true	yes
Sunny	mild	high	false	no
Sunny	cool	normal	false	yes
Rainy	mild	normal	false	yes
Sunny	mild	normal	true	yes
Overcast	mild	high	true	yes
Overcast	hot	normal	false	yes
Rainy	mild	high	true	no

Fig. 1. Weather and play data, copied from (Witten et al., 2011).

be said about the probability of playing if the outlook is *Sunny*, the temperature is *Cold*, the humidity is *High* and it is Windy?

Our plan is to first organise these table data into four channels  $d_O, d_T, d_H, d_W$  in a network of the form:



We start by extracting the underlying sets for for the categories in the table in Figure 1. We choose abbreviations for the entries in each of the categories.

$$O = \{s, o, r\} \quad W = \{t, f\} \quad T = \{h, m, c\} \quad P = \{y, n\} \quad H = \{h, n\}.$$

These sets are combined into a single product domain:

$$D = O \times T \times H \times W \times P.$$

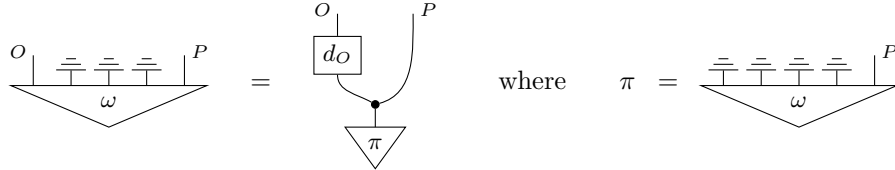
It combines the five columns in Figure 1. The table itself in the figure is represented as a uniform distribution  $\tau \in \mathcal{D}(D)$ . This distribution has 14 entries — like in the table — and looks as follows.

$$\tau = \frac{1}{14}|s, h, h, f, n\rangle + \frac{1}{14}|s, h, h, t, n\rangle + \cdots + \frac{1}{14}|r, m, h, t, n\rangle.$$

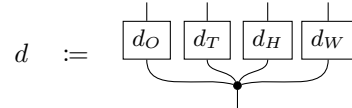
We extract the four channels in Diagram (6) via appropriate disintegrations, from the Play column to the Outlook / Temperature / Humidity / Windy columns.

$$\begin{aligned} d_O &= \tau. [ [1, 0, 0, 0, 0] \mid [0, 0, 0, 0, 1] ] & d_T &= \tau. [ [0, 1, 0, 0, 0] \mid [0, 0, 0, 0, 1] ] \\ d_H &= \tau. [ [0, 0, 1, 0, 0] \mid [0, 0, 0, 0, 1] ] & d_W &= \tau. [ [0, 0, 0, 1, 0] \mid [0, 0, 0, 0, 1] ]. \end{aligned}$$

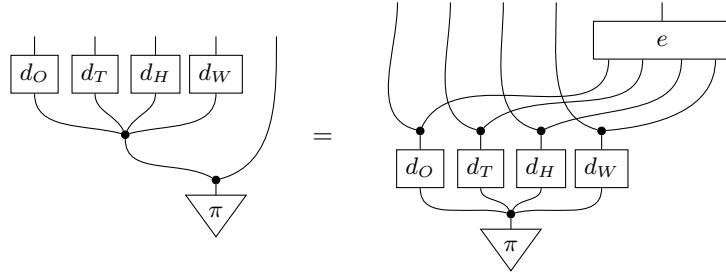
Thus, as described in the beginning of this section, the ‘outlook’ channel  $d_O: P \rightarrow \mathcal{D}(O)$  is extracted by first marginalising the table  $\tau$  to the relevant (first and last) wires, and then disintegrating. Explicitly,  $d_O$  is  $\tau.[1, 0, 0, 0, 1] \parallel [0, 1]$  and satisfies:



In a next step we combine these four channels into a single channel  $d: P \rightarrow O \times T \times H \times W$  via tupling:



The answer that we are looking for will be obtained by Bayesian inversion of this channel  $d$  wrt. the above fifth marginal Play state  $\pi = \tau.[0, 0, 0, 0, 1] = \frac{9}{14}|y\rangle + \frac{5}{14}|n\rangle \in \mathcal{D}(P)$ . We write this inversion as a channel  $e: O \times T \times H \times W \rightarrow P$ . It satisfies, by construction, according to the pattern in (5):



We can finally answer the original classification question. The assumptions — *Sunny* outlook, *Cold* temperature, *High* humidity, *true* windiness — are used as input to the inversion channel  $e$ . This yields the probability of play that we are looking for:

$$e(s, c, H, t) = 0.205|y\rangle + 0.795|n\rangle$$

The resulting classification probability<sup>‡</sup> of 0.205 coincides with the probability of 20.5% that is computed in (Witten et al., 2011) — without the above systematic approach.

One could complain that our approach is ‘too’ abstract, since it remains implicit what these extracted channels do. We elaborate the outlook channel  $d_O: P \rightarrow O$ . For the two elements in  $P = \{y, n\}$  we have:

$$d_O(y) = \frac{2}{9}|s\rangle + \frac{4}{9}|o\rangle + \frac{3}{9}|r\rangle \quad d_O(n) = \frac{3}{5}|s\rangle + 0|o\rangle + \frac{2}{5}|r\rangle.$$

These outcomes arise from the general formula (4). But we can also understand them at a

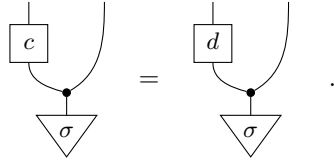
<sup>‡</sup> This outcome has been calculated with the tool *EfProb* (Cho and Jacobs, 2017) that can do both disintegration and inversion.

more concrete level: for the first distribution  $d_O(y)$  we need to concentrate on the 9 lines in Figure 1 for which Play is *yes*; in these lines, in the first Outlook column, 2 out of 9 entries are *Sunny*, 4 out of 9 are *Overcast*, and 3 out of 9 are *Rainy*. This corresponds to the above first distribution  $d_O(y)$ . Similarly, the second distribution  $d_O(n)$  captures the Outlook for the 5 lines where Play is *no*: 3 out of 5 are *Sunny* and 2 out of 5 are *Rainy*.

### 5. Almost equality of channels

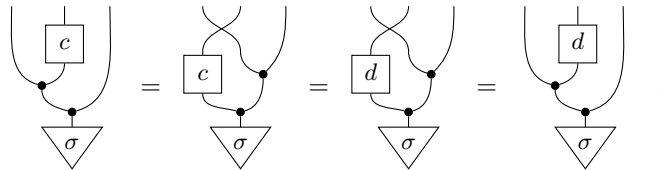
This section explains how the standard notion of ‘equal up to negligible sets’ or ‘equal almost everywhere’ (with respect to a measure) can be expressed abstractly using string diagrams. Via this equality relation Bayesian inversion can be characterised very neatly, following (Clerc et al., 2017). We consider an affine CD-category, continuing in the setting of Section 3.

**Definition 5.1.** Let  $c, d: X \rightarrow Y$  be two parallel channels, and  $\sigma: I \rightarrow X$  be a state on their domain. We say that  $c$  is  $\sigma$ -almost equal to  $d$ , written as  $c \stackrel{\sigma}{\sim} d$  if

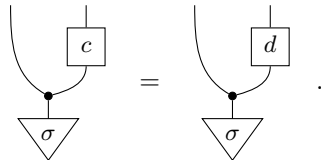


It is obvious that  $\stackrel{\sigma}{\sim}$  is an equivalence relation on channels of type  $X \rightarrow Y$ . When  $S$  is a set of arrows of type  $X \rightarrow Y$ , we write  $S/\sigma$  for the quotient  $S/\stackrel{\sigma}{\sim}$ .

To put it more intuitively, we have  $c \stackrel{\sigma}{\sim} d$  iff  $c$  and  $d$  can be identified whenever the input wires are connected to  $\sigma$ , possibly through copiers. For instance, using the associativity and commutativity of copiers, by  $c \stackrel{\sigma}{\sim} d$  we may reason as:



In particular,  $c \stackrel{\sigma}{\sim} d$  if and only if



Now the following is an obvious consequence from the definition.

**Proposition 5.2.** If both  $c, d: X \rightarrow Y$  are disintegrations of a joint state  $\omega: I \rightarrow X \otimes Y$ , then  $c \stackrel{\omega_1}{\sim} d$ , where  $\omega_1: I \rightarrow X$  is the first marginal of  $\omega$ .  $\square$

For channels  $f, g: X \rightarrow \mathcal{D}(Y)$  and a state  $\sigma \in \mathcal{D}(X)$  in  $\mathcal{Kl}(\mathcal{D})$ , it is easy to see that



$f \stackrel{\sim}{\sim} g$  if and only if  $f(x)(y) \cdot \sigma(x) = g(x)(y) \cdot \sigma(x)$  for all  $x \in X$  and  $y \in Y$  if and only if  $f(x) = g(x)$  for any  $x \in X$  with  $\sigma(x) \neq 0$ . Almost equality in  $\mathcal{Kl}(\mathcal{G})$  is less trivial but characterised in an expected way.

**Proposition 5.3.** Let  $f, g: X \rightarrow \mathcal{G}(Y)$  be channels and  $\mu \in \mathcal{G}(X)$  a state in  $\mathcal{Kl}(\mathcal{G})$ . Then  $f \stackrel{\mu}{\sim} g$  if and only if for any  $B \in \Sigma_Y$ ,  $f(-)(B) = g(-)(B)$   $\mu$ -almost everywhere.

*Proof.* By expanding the definition,  $f \stackrel{\mu}{\sim} g$  if and only if

$$\int_A f(x)(B) \mu(dx) = \int_A g(x)(B) \mu(dx)$$

for all  $A \in \Sigma_X$  and  $B \in \Sigma_Y$ . This is equivalent to  $f(-)(B) = g(-)(B)$   $\mu$ -almost everywhere for all  $B \in \Sigma_Y$ , see (Fremlin, 2000, 131H).  $\square$

Almost-everywhere equality of probability kernels  $f, g: X \rightarrow \mathcal{G}(Y)$  is often formulated by the stronger condition that  $f = g$   $\mu$ -almost everywhere. The next proposition shows that the stronger variant is equivalent under a reasonable assumption (any standard Borel space is countably generated, for example).

**Proposition 5.4.** In the setting of the previous proposition, additionally assume that the measurable space  $Y$  is countably generated. Then  $f \stackrel{\mu}{\sim} g$  if and only if  $f = g$   $\mu$ -almost everywhere.

*Proof.* Let the  $\sigma$ -algebra  $\Sigma_Y$  on  $Y$  be generated by a countable family  $(B_n)_n$ . We may assume that  $(B_n)_n$  is a  $\pi$ -system, *i.e.* a family closed under binary intersections. Let  $A_n = \{x \in X \mid f(x)(B_n) = g(x)(B_n)\}$ , and  $A = \bigcap_n A_n$ . Each  $A_n$  is  $\mu$ -conegligible, and thus  $A$  is  $\mu$ -conegligible. For each  $x \in A$ , we have  $f(x)(B_n) = g(x)(B_n)$  for all  $n$ . By application of the Dynkin  $\pi$ - $\lambda$  theorem, it follows that  $f(x) = g(x)$ . Therefore  $f = g$   $\mu$ -almost everywhere.  $\square$

We can now present a fundamental result from (Clerc et al., 2017, §3.3) in our abstract setting. Let  $\mathbf{C}$  be an affine CD-category and  $(I \downarrow \mathbf{C})$  be the comma (coslice) category. Objects in  $(I \downarrow \mathbf{C})$  are states in  $\mathbf{C}$ , formally pairs  $(X, \sigma)$  of objects  $X \in \mathbf{C}$  and states  $\sigma: I \rightarrow X$ . Arrows from  $(X, \sigma)$  to  $(Y, \tau)$  are *state-preserving* channels, namely  $c: X \rightarrow Y$  in  $\mathbf{C}$  satisfying  $c \circ \sigma = \tau$ . A joint state  $(X \otimes Y, \omega) \in (I \downarrow \mathbf{C})$  is called a *coupling* of two states  $(X, \sigma), (Y, \tau) \in (I \downarrow \mathbf{C})$  if

$$\begin{array}{c} X \\ \text{---} \\ \omega \\ \text{---} \\ \sigma \end{array} = \begin{array}{c} X \\ \text{---} \\ \sigma \end{array} \quad \text{and} \quad \begin{array}{c} \text{---} \\ \omega \\ \text{---} \\ \tau \end{array} = \begin{array}{c} \text{---} \\ \tau \end{array} .$$

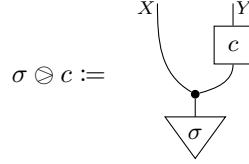
We write  $\text{Coupl}((X, \sigma), (Y, \tau))$  for the set of couplings of  $(X, \sigma)$  and  $(Y, \tau)$ .

**Theorem 5.5.** Let  $\mathbf{C}$  be an affine CD-category that admits disintegration. For each pair of states  $(X, \sigma), (Y, \tau) \in (I \downarrow \mathbf{C})$ , there is the following bijection:

$$(I \downarrow \mathbf{C})((X, \sigma), (Y, \tau)) / \sigma \cong \text{Coupl}((X, \sigma), (Y, \tau))$$

*Proof.* For each  $c \in (I \downarrow \mathbf{C})((X, \sigma), (Y, \tau))$ , we define a joint state  $I \rightarrow X \otimes Y$  to be

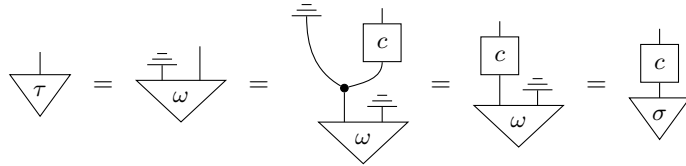
the ‘integration’ of  $\sigma$  and  $c$  as below, for which we use the following *ad hoc* notation:



It is easy to check that  $\sigma \circledast c$  is a coupling of  $\sigma$  and  $\tau$ . For two channels  $c, d: X \rightarrow Y$ , we have  $\sigma \circledast c = \sigma \circledast d$  if and only if  $c \overset{\sigma}{\sim} d$ , by the definition of  $\overset{\sigma}{\sim}$ . This means the mapping

$$c \longmapsto \sigma \circledast c, \quad (I \downarrow \mathbf{C})((X, \sigma), (Y, \tau)) / \sigma \longrightarrow \text{Coupl}((X, \sigma), (Y, \tau))$$

is well-defined and injective. To prove the surjectivity let  $(X \otimes Y, \omega) \in \text{Coupl}((X, \sigma), (Y, \tau))$ . Let  $c: X \rightarrow Y$  be a disintegration of  $\omega$ . Then  $c$  is state-preserving since



Moreover we have  $\sigma \circledast c = \omega$ , as desired. □

Via the symmetry  $X \otimes Y \xrightarrow{\cong} Y \otimes X$  we have the obvious bijection  $\text{Coupl}((X, \sigma), (Y, \tau)) \cong \text{Coupl}((Y, \tau), (X, \sigma))$ . This immediately gives the following corollary.

**Corollary 5.6.** Let  $\mathbf{C}$  be an affine CD-category that admits disintegration. For any states  $(X, \sigma), (Y, \tau) \in (I \downarrow \mathbf{C})$  we have

$$(I \downarrow \mathbf{C})((X, \sigma), (Y, \tau)) / \sigma \cong (I \downarrow \mathbf{C})((Y, \tau), (X, \sigma)) / \tau$$

The bijection sends a channel  $c: X \rightarrow Y$  to a Bayesian inversion  $d: Y \rightarrow X$  for  $\sigma$  along  $c$ . □

Theorem 2 of (Clerc et al., 2017) is obtained as an instance, for the category  $\mathbf{pKrn}_{\text{sb}}$ . This bijective correspondence yields a ‘dagger’  $(-)^{\dagger}$  functor on (a suitable quotient of) the comma category  $(I \downarrow \mathbf{C})$  — as noted by the authors of (Clerc et al., 2017).

### Strong almost equality

Unfortunately, the almost equality defined above is not the most useful notion for equational reasoning between string diagrams. Indeed, later in the proofs of Proposition 6.4 and Theorem 8.3 we encounter situations where we need a stronger notion of almost equality. We define the stronger one as follows.

**Definition 5.7.** Let  $c, d: X \rightarrow Y$  be channels and  $\sigma: I \rightarrow X$  be a state. We say that  $c$  is *strongly  $\sigma$ -almost equal to  $d$*  if

$$\begin{array}{c} \boxed{c} \\ | \\ \triangle \omega \end{array} = \begin{array}{c} \boxed{d} \\ | \\ \triangle \omega \end{array}$$

holds for any  $\omega: I \rightarrow X \otimes Z$  such that  $\sigma$  is the first marginal of  $\omega$ , that is:

$$\begin{array}{c} | \\ \triangle \sigma \end{array} = \begin{array}{c} | \\ \triangle \omega \\ \hline \sigma \end{array}.$$

Notice that strong almost quality implies almost quality, via the following  $\omega$ :

$$\begin{array}{c} | \\ \triangle \omega \end{array} := \begin{array}{c} \cup \\ | \\ \triangle \sigma \end{array}.$$

The good news is that the converse often holds too. We say that an affine CD-category admits *equality strengthening* if  $c$  is strongly  $\sigma$ -almost equal to  $d$  whenever  $c \overset{\sigma}{\sim} d$ , i.e.  $c$  is  $\sigma$ -almost equal to  $d$ . We present two propositions that guarantee this property.

**Proposition 5.8.** If an affine CD-category admits disintegration, then it admits equality strengthening.

*Proof.* Suppose that  $c \overset{\sigma}{\sim} d$  holds for a state  $\sigma: I \rightarrow X$  and for channels  $c, d: X \rightarrow Y$ . Let  $\omega: I \rightarrow X \otimes Z$  be a state whose first marginal is  $\sigma$ . We can disintegrate  $\omega$  as in:

$$\begin{array}{c} | \\ \triangle \omega \end{array} = \begin{array}{c} \cup \\ \boxed{e} \\ | \\ \triangle \sigma \end{array}.$$

Then we have

$$\begin{array}{c} \boxed{c} \\ | \\ \triangle \omega \end{array} = \begin{array}{c} \boxed{c} \quad \boxed{e} \\ | \quad | \\ \cup \\ | \\ \triangle \sigma \end{array} = \begin{array}{c} \boxed{d} \quad \boxed{e} \\ | \quad | \\ \cup \\ | \\ \triangle \sigma \end{array} = \begin{array}{c} \boxed{d} \\ | \\ \triangle \omega \end{array}.$$

Therefore  $c$  is strongly  $\sigma$ -almost equal to  $d$ . □

Recall that  $\mathcal{Kl}(\mathcal{G})$  does not admit disintegration. Nevertheless, it admits equality strengthening.

**Proposition 5.9.** The category  $\mathcal{Kl}(\mathcal{G})$  admits equality strengthening.

*Proof.* Assume  $c \overset{\sigma}{\sim} d$  for  $\sigma \in \mathcal{G}(X)$  and  $c, d: X \rightarrow \mathcal{G}(Y)$ . Let  $\omega \in \mathcal{G}(X \otimes Z)$  be a probability measure whose first marginal is  $\sigma$ , i.e.  $(\pi_1)_*(\omega) = \sigma$ . We need to prove

$(c \otimes \eta_Z) \circ \omega = (d \otimes \eta_Z) \circ \omega$ , which is equivalent to:

$$\int_{X \times Z} c(x)(B) \mathbf{1}_C(z) \omega(d(x, z)) = \int_{X \times Z} d(x)(B) \mathbf{1}_C(z) \omega(d(x, z)) \quad (7)$$

for all  $B \in \Sigma_Y$  and  $C \in \Sigma_Z$ . Using

$$|c(x)(B) \mathbf{1}_C(z) - d(x)(B) \mathbf{1}_C(z)| = |c(x)(B) - d(x)(B)| \mathbf{1}_C(z) \leq |c(x)(B) - d(x)(B)| ,$$

we have

$$\begin{aligned} & \left| \int_{X \times Z} (c(x)(B) \mathbf{1}_C(z) - d(x)(B) \mathbf{1}_C(z)) \omega(d(x, z)) \right| \\ & \leq \int_{X \times Z} |c(x)(B) \mathbf{1}_C(z) - d(x)(B) \mathbf{1}_C(z)| \omega(d(x, z)) \\ & \leq \int_{X \times Z} |c(x)(B) - d(x)(B)| \omega(d(x, z)) \\ & = \int_X |c(x)(B) - d(x)(B)| (\pi_1)_*(\omega)(dx) \\ & = \int_X |c(x)(B) - d(x)(B)| \sigma(dx) \\ & = 0 , \end{aligned}$$

where the last equality holds since  $c(-)(B) = d(-)(B)$   $\sigma$ -almost everywhere by Proposition 5.3. Therefore the desired equality (7) holds.  $\square$

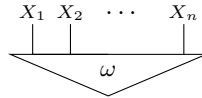
## 6. Conditional independence

Throughout this section, we consider an affine CD-category that admits disintegration.

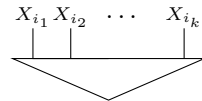
### 6.1. Disintegration of multipartite states

So far we have concentrated on bipartite states — except in the classification example in Section 4. In order to deal with a general  $n$ -partite state  $\omega: I \rightarrow X_1 \otimes \cdots \otimes X_n$ , we will introduce several notations and conventions in Definitions 6.1, 6.2 and 6.3 below; they are in line with standard practice in probability theory.

In the conventions, an  $n$ -partite state, as below, is fixed, and used implicitly.

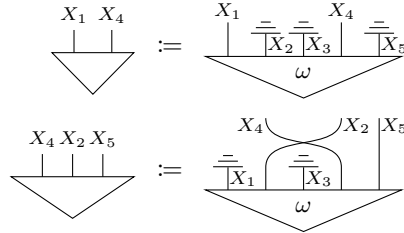


**Definition 6.1.** When we write



where  $i_1, \dots, i_k$  are distinct, it denotes the state  $I \rightarrow X_{i_1} \otimes \cdots \otimes X_{i_k}$  obtained from  $\omega$  by

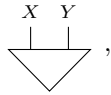
marginalisation and permutation of wires (if necessary). Let us give a couple of examples, for  $n = 5$ .



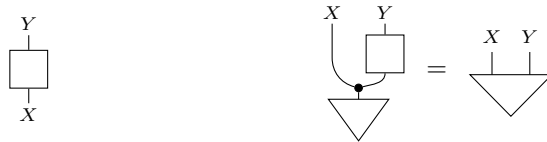
We permute wires via a combination of crossing. This is unambiguous by the coherence theorem.

Below we will use symbols  $X, Y, Z, W, \dots$  to denote not only a single wire  $X_i$  but also multiple wires  $X_i \otimes X_j \otimes \dots$ . Disintegrations more general than in the bipartite case are now introduced as follows.

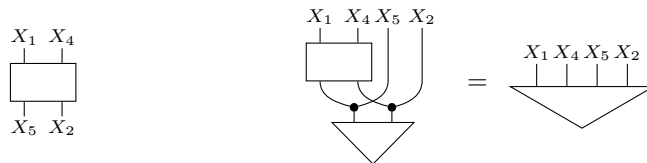
**Definition 6.2.** For  $X = X_{i_1} \otimes \dots \otimes X_{i_k}$  and  $Y = X_{j_1} \otimes \dots \otimes X_{j_l}$ , with all  $i_1, \dots, i_k, j_1, \dots, j_l$  distinct, a disintegration  $X \rightarrow Y$  is defined to be a disintegration of



the marginal state given by the previous convention. We denote the disintegration simply as on the left below,

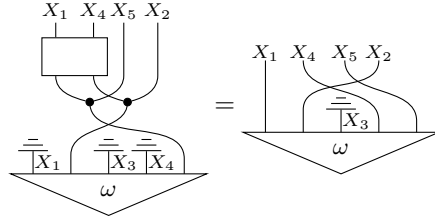


By definition, it must satisfy the equation on the right above. Let us give an example. The disintegration  $X_5 \otimes X_2 \rightarrow X_1 \otimes X_4$  on the left below is defined by the equation on the right.



More specifically, assuming  $n = 5$  and expanding the notation for marginals, the equation

is:



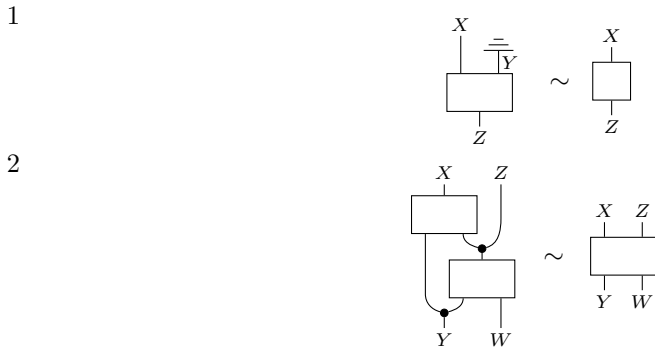
Note that disintegrations need not be unique. Thus when we write  $\begin{array}{c} \square \\ \downarrow \\ X \end{array}$ , we in fact *choose* one of them. Nevertheless, such disintegrations are unique up to almost-equality with respect to  $\begin{array}{c} \square \\ \downarrow \\ X \end{array}$ , which is good enough for our purpose.

Finally we make a convention about almost equality (Definition 5.1).

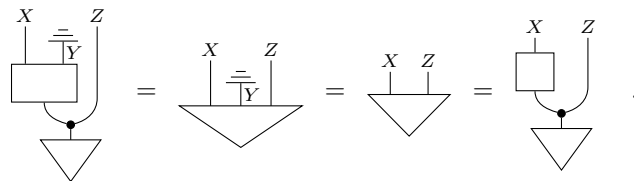
**Definition 6.3.** Let  $S$  and  $T$  be string diagrams of type  $X \rightarrow Y$  that are made from marginals and disintegrations of  $\omega$  as defined in Definitions 6.1 and 6.2. When we say  $S$  is almost equal to  $T$  (or write  $S \sim T$ ) without reference to a state, it means that  $S$  is almost equal to  $T$  with respect to the state  $\begin{array}{c} \square \\ \downarrow \\ X \end{array}$ .

We shall make use of the following auxiliary equations involving discarding and composition of disintegrations.

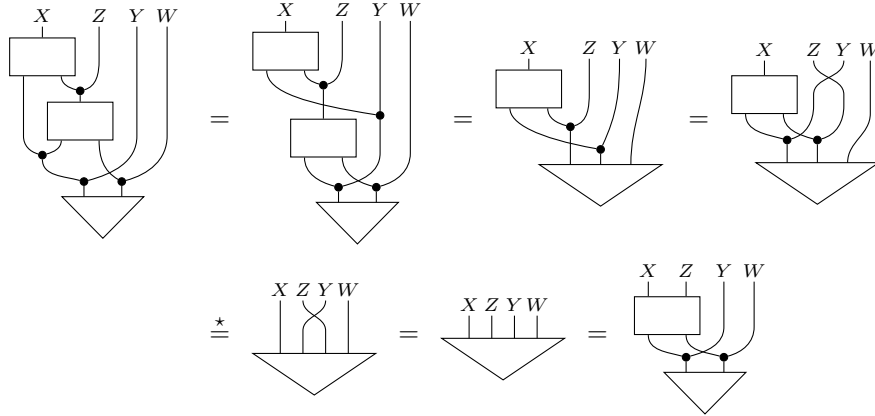
**Proposition 6.4.** In the conventions and notations above, the following hold.



*Proof.* By the definition of almost equality, 1 is proved by:



Similarly, we prove 2 as follows.



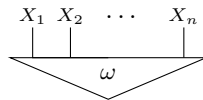
The marked equality  $\stackrel{*}{=}$  holds by strong almost equality, see Proposition 5.8. □

The equations correspond respectively to  $\sum_y P(x, y|z) = P(x|z)$  and  $P(x|y, z) \cdot P(z|y, w) = P(x, z|y, w)$  in discrete probability.

**Remark 6.5.** We here keep our notation somewhat informal, e.g. using symbols  $X, Y, Z, \dots$  as a sort of meta-variables. We refer to (Joyal and Street, 1991; Selinger, 2010) for more formal aspects of string diagrams.

### 6.2. Conditional independence

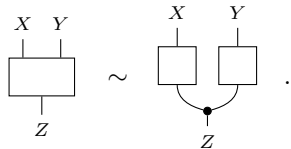
We continue using the notations in the previous subsection. Recall that we fix an  $n$ -partite state



and use symbols  $X, Y, Z, W, \dots$  to denote a wire  $X_i$  or multiple wires  $X_i \otimes X_j \otimes \dots$ .

We now introduce the notion of conditional independence. Although it is defined with respect to the underlying state  $\omega$ , we leave the state  $\omega$  implicit, like an underlying probability space  $\Omega$  in conventional probability theory.

**Definition 6.6.** Let  $X, Y, Z$  denote distinct wires. Then we say  $X$  and  $Y$  are conditionally independent given  $Z$ , written as  $X \perp\!\!\!\perp Y \mid Z$ , if



The definition is analogous to the condition  $P(x, y|z) = P(x|z)P(y|z)$  in discrete

probability. Indeed our definition coincides with the usual conditional independence, as explained below.

**Example 6.7.** In  $\mathcal{Kl}(\mathcal{D})$ , let  $c_{X|Z}: Z \rightarrow \mathcal{D}(X)$ ,  $c_{Y|Z}: Z \rightarrow \mathcal{D}(Y)$ ,  $c_{XY|Z}: Z \rightarrow \mathcal{D}(X \times Y)$  be disintegrations of some joint state, say  $\omega \in \mathcal{D}(X \times Y \times Z)$ . Let  $\omega_Z \in \mathcal{D}(Z)$  be the marginal on  $Z$ . Then  $X \perp\!\!\!\perp Y \mid Z$  if and only if

$$c_{XY|Z}(z)(x, y) = c_{X|Z}(z)(x) \cdot c_{Y|Z}(z)(y) \quad \text{whenever } \omega_Z(z) \neq 0$$

for all  $x \in X$ ,  $y \in Y$  and  $z \in Z$ . If we write  $P(x, y|z) = c_{XY|Z}(z)(x, y)$ ,  $P(x|z) = c_{X|Z}(z)(x)$ ,  $P(y|z) = c_{Y|Z}(z)(y)$ , and  $P(z) = \omega_Z(z)$ , then the condition will look more familiar:

$$P(x, y|z) = P(x|z) \cdot P(y|z) \quad \text{whenever } P(z) \neq 0 .$$

**Example 6.8.** Similarly, in  $\mathcal{Kl}(\mathcal{G})$ , let  $c_{X|Z}: Z \rightarrow \mathcal{G}(X)$ ,  $c_{Y|Z}: Z \rightarrow \mathcal{G}(Y)$ ,  $c_{XY|Z}: Z \rightarrow \mathcal{G}(X \times Y)$ , and  $\omega_Z \in \mathcal{G}(Z)$  be appropriate disintegrations and a marginal of some joint probability measure  $\omega$ . Then  $X \perp\!\!\!\perp Y \mid Z$  if and only if

$$c_{XY|Z}(z)(A \times B) = c_{X|Z}(z)(A) \cdot c_{Y|Z}(z)(B) \quad \text{for } \omega_Z\text{-almost all } z \in Z$$

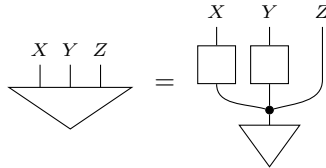
for all  $A \in \Sigma_X$  and  $B \in \Sigma_Y$ .

The equivalences in the next result are well-known in conditional probability. Our contribution is that we formulate and prove them at an abstract, graphical level.

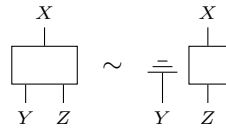
**Proposition 6.9.** The following are equivalent.

1  $X \perp\!\!\!\perp Y \mid Z$

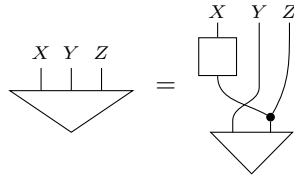
2



3



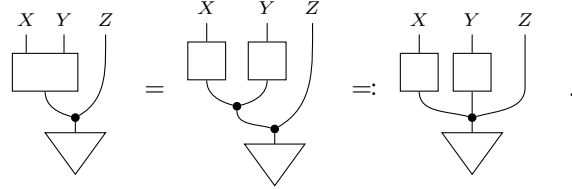
4



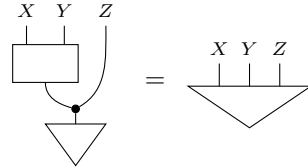
As we will see below, conditional independence  $X \perp\!\!\!\perp Y \mid Z$  is symmetric in  $X$  and  $Y$ . Therefore the obvious symmetric counterparts of 3 and 4 are also equivalent to them.



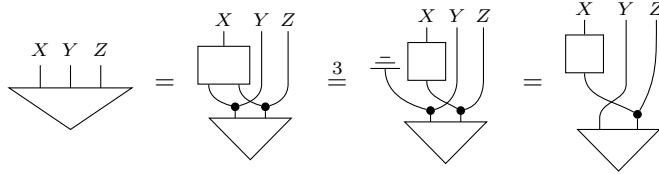
*Proof.* By definition of almost equality, 1 is equivalent to



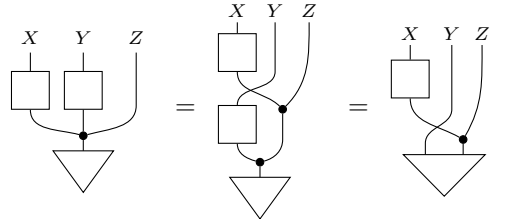
We then have  $1 \Leftrightarrow 2$ , since the identity below holds by the definition of disintegration.



Next, assuming 3, we obtain 4 as follows.



We prove  $4 \Rightarrow 3$  similarly. Finally note that the following equation holds.



From this  $2 \Leftrightarrow 4$  is immediate. □

Note that the condition 3 of the proposition is an analogue of  $P(x|y, z) = P(x|z)$ . The other conditions 2 and 4 say that the joint state can be factorised in certain ways, corresponding to the following equations:

$$P(x, y, z) = P(x|z) P(y|z) P(z) = P(x|z) P(y, z).$$

The proposition below shows that our abstract formulation of conditional independence does satisfy the basic ‘rules’ of conditional independence, which are known as *(semi-)graphoids axioms* (Verma and Pearl, 1988; Geiger et al., 1990).

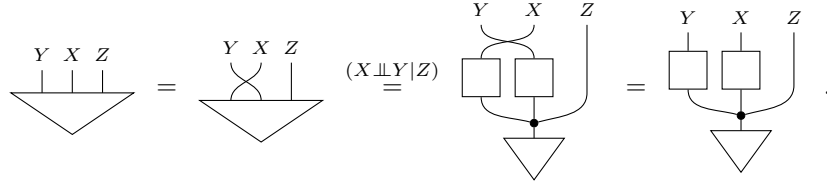
**Proposition 6.10.** Conditional independence  $(-) \perp\!\!\!\perp (-) \mid (-)$  satisfies:

- 1 (Symmetry)  $X \perp\!\!\!\perp Y \mid Z$  if and only if  $Y \perp\!\!\!\perp X \mid Z$ .
- 2 (Decomposition)  $X \perp\!\!\!\perp Y \otimes Z \mid W$  implies  $X \perp\!\!\!\perp Y \mid W$  and  $X \perp\!\!\!\perp Z \mid W$ .
- 3 (Weak union)  $X \perp\!\!\!\perp Y \otimes Z \mid W$  implies  $X \perp\!\!\!\perp Y \mid Z \otimes W$ .

4 (Contraction)  $X \perp\!\!\!\perp Z \mid W$  and  $X \perp\!\!\!\perp Y \mid Z \otimes W$  imply  $X \perp\!\!\!\perp Y \otimes Z \mid W$ .

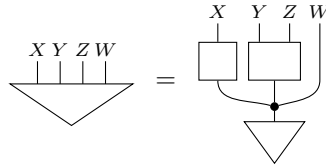
*Proof.* We will freely use Proposition 6.9.

(1) Suppose  $X \perp\!\!\!\perp Y \mid Z$ . Then

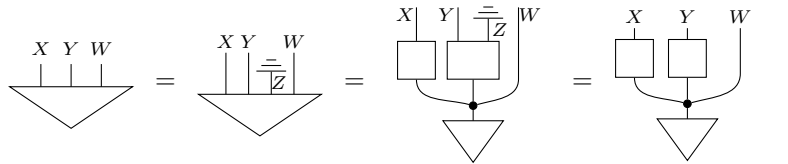


This means  $Y \perp\!\!\!\perp X \mid Z$ .

(2) Suppose  $X \perp\!\!\!\perp Y \otimes Z \mid W$ , namely:

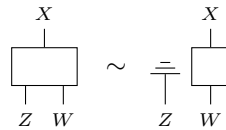


Marginalising  $Z$ , we obtain

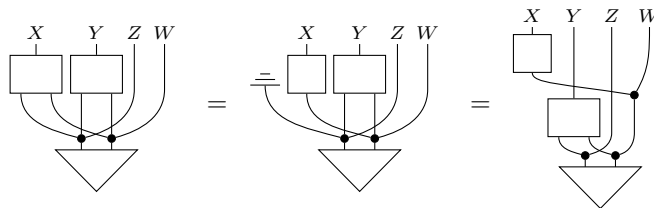


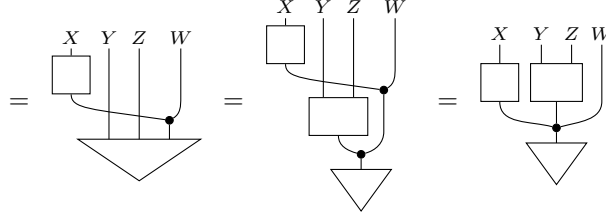
by Proposition 6.4.1. Thus  $X \perp\!\!\!\perp Y \mid W$ . Similarly we prove  $X \perp\!\!\!\perp Z \mid W$ .

Finally, we prove 3 and 4 at the same time. Note that  $X \perp\!\!\!\perp Y \otimes Z \mid W$  implies  $X \perp\!\!\!\perp Z \mid W$ , as shown above. Therefore what we need to prove is that  $X \perp\!\!\!\perp Y \otimes Z \mid W$  if and only if  $X \perp\!\!\!\perp Y \mid Z \otimes W$ , under  $X \perp\!\!\!\perp Z \mid W$ . Assume  $X \perp\!\!\!\perp Z \mid W$ , so we have



Then





This proves  $X \perp\!\!\!\perp Y \otimes Z \mid W$  if and only if  $X \perp\!\!\!\perp Y \mid Z \otimes W$ . □

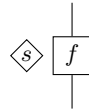
The four properties from the graphoid axioms are essential in reasoning of conditional independence with DAGs or Bayesian networks (Verma and Pearl, 1988; Geiger et al., 1990).

Conditional independence has been studied categorically in (Simpson, 2017). There, a categorical notion of (*conditional*) *independence structure* is introduced, generalising algebraic axiomatisations of conditional independence, such as with graphoids and separoids (Dawid, 2001). We leave it to future work to precisely relate our approach to conditional probability to Simpson’s categorical framework.

### 7. Beyond causal channels

All CD-categories  $\mathbf{C}$  that we have considered so far are affine in the sense that all arrows  $f: X \rightarrow Y$  are causal:  $\bar{\tau} \circ f = \bar{\tau}$ . We now drop the affineness, in order to enlarge our category to include ‘non-causal’ arrows, which enables us to have new notions such as scalars and effects. Essentially, we lose nothing by this change: all the arguments so far can still be applied to the subcategory  $\text{Caus}(\mathbf{C}) \subseteq \mathbf{C}$  containing all the objects and causal arrows. The category  $\text{Caus}(\mathbf{C})$  is an affine CD-category, inheriting the monoidal structure  $(\otimes, I)$  and the comonoid structures  $(\Upsilon, \bar{\tau})$  from  $\mathbf{C}$ .

Recall that *channels* in  $\mathbf{C}$  are causal arrows, *i.e.* arrows in  $\text{Caus}(\mathbf{C})$ . *States* are channels of the form  $\sigma: I \rightarrow X$ . We call endomaps  $I \rightarrow I$  on the tensor unit *scalars*. The set  $\mathbf{C}(I, I)$  of scalars forms a monoid via the composition  $s \cdot t = s \circ t$  and  $1 = \text{id}_I$ . The monoid of scalars is always commutative — in fact, this is the case for any monoidal category, see *e.g.* (Abramsky and Coecke, 2009, §3.2). In string diagram scalars are written as  $\diamond s$  or simply as  $s$ . We can multiply scalars  $s$  to any arrows  $f: X \rightarrow Y$  by the parallel composition, or diagrammatically by juxtaposition:



We call an arrow  $\sigma: I \rightarrow X$  *normalisable* if the scalar  $\bar{\tau} \circ \sigma: I \rightarrow I$  is (multiplicatively) invertible. In that case we can normalise  $\sigma$  into a proper state as follows.

$$\text{norm}(\sigma) := \left( \bar{\tau} \circ \sigma \right)^{-1} \circ \sigma$$

*Effects* in  $\mathbf{C}$  are arrows of the form  $p: X \rightarrow I$ ; they correspond to observables, with

predicates as special case. Diagrammatically they are written as on the left below.

$$\begin{array}{c} \triangle \\ \uparrow \\ p \end{array} \qquad \sigma \models p := \begin{array}{c} \triangle \\ \hline p \\ \hline \sigma \\ \triangle \end{array}$$

On the right the *validity*  $\sigma \models p$  of a state  $\sigma: I \rightarrow X$  and a effect  $p: X \rightarrow I$  is defined. It is the scalar given by composition. Note that effects are not causal in general; by definition, only discards  $\bar{\tau}$  are causal ones. States  $\sigma: I \rightarrow X$  can be *conditioned* by effects  $p: X \rightarrow I$  via normalisation, as follows.

$$\sigma|_p := \text{norm} \left( \begin{array}{c} \triangle \\ \uparrow \\ p \\ \bullet \\ \downarrow \\ \sigma \\ \triangle \end{array} \right) = \left( \begin{array}{c} \triangle \\ \hline p \\ \hline \sigma \\ \triangle \end{array} \right)^{-1} \begin{array}{c} \triangle \\ \uparrow \\ p \\ \bullet \\ \downarrow \\ \sigma \\ \triangle \end{array}$$

The conditional state  $\sigma|_p$  is defined if the validity  $\sigma \models p$  is invertible. Conditioning  $\sigma|_p$  generalises conditional probability  $P(A|B)$  given *event*  $B$ , see Example 7.2 below. At the end of the section we will explain how conditioning and disintegration are related.

Recall, from Examples 2.4 and 2.5, that our previous examples  $\mathcal{Kl}(\mathcal{D})$  and  $\mathcal{Kl}(\mathcal{G})$  are both affine. We give two non-affine CD-categories that have  $\mathcal{Kl}(\mathcal{D})$  and  $\mathcal{Kl}(\mathcal{G})$  as subcategories, respectively.

**Example 7.1.** For discrete probability, we use *multisets* (or *unnormalised distributions*) over nonnegative real numbers  $\mathbb{R}_{\geq 0} = [0, \infty)$ , such as

$$1|x\rangle + 0.5|y\rangle + 3|z\rangle \quad \text{on a set } X = \{x, y, z, \dots\}$$

We denote by  $\mathcal{M}(X)$  the set of multisets over  $\mathbb{R}_{\geq 0}$  on  $X$ . More formally:

$$\mathcal{M}(X) = \{\phi: X \rightarrow \mathbb{R}_{\geq 0} \mid \phi \text{ has finite support}\} .$$

It extends to a commutative monad  $\mathcal{M}: \mathbf{Set} \rightarrow \mathbf{Set}$ , see (Coumans and Jacobs, 2013). In a similar way to the distribution monad  $\mathcal{D}$ , we can check that the Kleisli category  $\mathcal{Kl}(\mathcal{M})$  is a CD-category. For a Kleisli map  $f: X \rightarrow \mathcal{M}(Y)$ , causality  $\bar{\tau} \circ f = \bar{\tau}$  amounts to the condition  $\sum_y f(x)(y) = 1$  for all  $x \in X$ . It is thus easy to see that  $\text{Caus}(\mathcal{Kl}(\mathcal{M})) \cong \mathcal{Kl}(\mathcal{D})$ . In fact, the distribution monad  $\mathcal{D}$  can be obtained from  $\mathcal{M}$  as its *affine submonad*, see (Jacobs, 2017).

An effect  $p: X \rightarrow 1$  in  $\mathcal{Kl}(\mathcal{M})$  is a function  $p: X \rightarrow \mathbb{R}_{\geq 0}$ . Its validity  $\sigma \models p$  in a state  $\omega$  is given by the expected value  $\sum_x \sigma(x) \cdot p(x)$ . The state  $\sigma|_p$  updated with ‘evidence’  $p$  is defined as  $\sigma|_p(x) = \frac{\sigma(x) \cdot p(x)}{\sigma \models p}$ .

**Example 7.2.** For general, measure-theoretic probability, we use *s-finite kernels* between measurable spaces (Kallenberg, 2017; Staton, 2017). Let  $X$  and  $Y$  be measurable spaces. A function  $f: X \times \Sigma_Y \rightarrow [0, \infty]$  is called a *kernel* from  $X$  to  $Y$  if

- $f(x, -): \Sigma_Y \rightarrow [0, \infty]$  is a measure for each  $X$ ; and
- $f(-, B): X \rightarrow [0, \infty]$  is measurable<sup>§</sup> for each  $B \in \Sigma_Y$ .

<sup>§</sup> The  $\sigma$ -algebra on  $[0, \infty]$  is the standard one generated by  $\{\infty\}$  and measurable subsets of  $\mathbb{R}_{\geq 0}$ .

We write  $f: X \rightsquigarrow Y$  when  $f$  is a kernel from  $X$  to  $Y$ . A *probability kernel* is a kernel  $f: X \rightsquigarrow Y$  with  $f(x, Y) = 1$  for all  $x \in X$ . A kernel  $f: X \rightsquigarrow Y$  is *finite* if there exists  $r \in [0, \infty)$  such that for all  $x \in X$ ,  $f(x, Y) \leq r$ . (Note that it must be ‘uniformly’ finite.) A kernel  $f: X \rightsquigarrow Y$  is *s-finite* if  $f = \sum_n f_n$  for some countable family  $(f_n: X \rightarrow Y)_{n \in \mathbb{N}}$  of finite kernels.

For two s-finite kernels  $f: X \rightsquigarrow Y$  and  $g: Y \rightsquigarrow Z$ , we define the (sequential) composite  $g \circ f: X \rightsquigarrow Z$  by

$$(g \circ f)(x, C) = \int_Y g(y, C) f(x, dy)$$

for  $x \in X$  and  $C \in \Sigma_Z$ . There are identity kernels  $\eta_X: X \rightsquigarrow X$  given by  $\eta_X(x, A) = \mathbf{1}_A(x)$ . With these data, measurable spaces and s-finite kernels form a category, which we denote by **sfKrn**. There is a monoidal structure on **sfKrn**. For measurable spaces  $X, Y$  we define the tensor product  $X \otimes Y = X \times Y$  to be the cartesian product of measurable spaces. The tensor unit  $I = 1$  is the singleton space. For s-finite kernels  $f: X \rightsquigarrow Y$  and  $g: Z \rightsquigarrow W$ , we define  $f \otimes g: X \times Z \rightsquigarrow Y \times W$  by

$$\begin{aligned} (f \otimes g)((x, z), E) &= \int_Y \left( \int_W \mathbf{1}_E(y, w) g(z, dw) \right) f(x, dy) \\ &= \int_W \left( \int_Y \mathbf{1}_E(y, w) f(x, dy) \right) g(z, dw) \end{aligned}$$

for  $x \in X, z \in Z, E \in \Sigma_{Y \times W}$ . The latter equality holds by the Fubini-Tonelli theorem for s-finite measures. These make the category **sfKrn** symmetric monoidal. Finally, for each measurable space  $X$  there is a ‘copier’  $\Upsilon: X \rightsquigarrow X \times X$  and a ‘discarder’  $\ddagger: X \rightsquigarrow 1$ , given by  $\Upsilon(x, E) = \mathbf{1}_E(x, x)$  and  $\ddagger(x, 1) = 1$ , so that **sfKrn** is a CD-category. For more technical details we refer to (Kallenberg, 2017; Staton, 2017).

Note that an s-finite kernel  $f: X \rightsquigarrow Y$  is causal if and only if it is a probability kernel, which is nothing but a Kleisli map  $X \rightarrow \mathcal{G}(Y)$  for the Giry monad. Therefore the causal subcategory of **sfKrn** is the Kleisli category of the Giry monad:  $\text{Caus}(\mathbf{sfKrn}) \cong \mathcal{Kl}(\mathcal{G})$ . In particular, states in **sfKrn** are probability measures  $\sigma \in \mathcal{G}(X)$ .

An effect  $p: X \rightsquigarrow 1$  in **sfKrn**, *i.e.* an s-finite kernel  $p: X \times \Sigma_1 \rightarrow [0, \infty]$ , can be identified with a measurable function  $p: X \rightarrow [0, \infty]$ . The validity  $\sigma \models p$  is then the integral  $\int_X p(x) \sigma(dx)$ , defined in  $[0, \infty]$ . The conditional state  $\sigma|_p \in \mathcal{G}(X)$  is defined by:

$$\sigma|_p(A) = \frac{\int_A p(x) \sigma(dx)}{\sigma \models p}$$

for  $A \in \Sigma_X$ , when the validity  $\sigma \models p$  is neither 0 nor  $\infty$ . In particular, for any ‘event’  $B \in \Sigma_X$ , the obvious indicator function  $\mathbf{1}_B: X \rightarrow [0, \infty]$  is an effect. Then conditioning yields a new state on  $X$ :

$$\sigma|_{\mathbf{1}_B}(A) = \frac{\int_A \mathbf{1}_B(x) \sigma(dx)}{\sigma \models \mathbf{1}_B} = \frac{\sigma(A \cap B)}{\sigma(B)} \quad \text{for } A \in \Sigma_X.$$

This amounts to conditional probability  $P(A|B) = P(A, B)/P(B)$  given event  $B$ .

We need to generalise definitions from the previous sections in the non-affine/causal setting. Here we make only a minimal generalisation that is required in the next section.

For example, we still restrict ourselves to disintegration of (causal) states, although in the literature, the notion of disintegration exists also for non-probability (i.e. non-causal) measures and even for non-finite measures, see e.g. (Chang and Pollard, 1997, Definition 1). We leave such generalisations to future work.

Let  $\sigma: I \rightarrow X$  be a state. We define  $\sigma$ -almost equality  $f \overset{\sigma}{\sim} g$  between arbitrary arrows  $f, g: X \rightarrow Y$  in the same way as Definition 5.1. Similarly, strong  $\sigma$ -almost equality between arbitrary arrows is defined as in Definition 5.7 (here  $\omega$  still ranges over states).

Disintegrations of a joint state  $\omega: I \rightarrow X \otimes Y$  are defined as in Definition 3.5, except that an arrow  $c_1: X \rightarrow Y$  (or  $c_2: Y \rightarrow X$ ) is not necessarily causal. Nevertheless, disintegrations of a state are *almost causal* in the following sense.

**Definition 7.3.** Let  $\sigma: I \rightarrow X$  be a state. We say that an arrow  $c: X \rightarrow Y$  is  $\sigma$ -almost causal if  $\overline{\overline{\sigma}} \circ c \overset{\sigma}{\sim} \overline{\overline{\sigma}}$ .

Then the following is immediate from the definition.

**Proposition 7.4.** Let  $c_1: X \rightarrow Y$  be a disintegration of a state  $\omega: I \rightarrow X \otimes Y$ . Then  $c_1$  is  $\omega_1$ -almost causal, where  $\omega_1: I \rightarrow X$  is the first marginal of  $\omega$ .  $\square$

**Example 7.5.** In **sfKrn**, a s-finite kernel  $f: X \rightsquigarrow Y$  is  $\sigma$ -almost causal if and only if  $f(x, Y) = 1$  for  $\sigma$ -almost all  $x \in X$ . In that case, we can find a probability kernel  $f': X \rightsquigarrow Y$  such that  $f \overset{\sigma}{\sim} f'$ , by tweaking  $f$  in the obvious way. From this it follows that a disintegration of a joint probability measure  $\omega \in \mathcal{G}(X \times Y)$  exists in **sfKrn** if and only if it exists in  $\mathcal{Kl}(\mathcal{G})$ .

By Proposition 5.9, we can strengthen almost equality between channels in **sfKrn**. It is easy to see that equality strengthening is valid also for almost causal maps: if  $f, g: X \rightsquigarrow Y$  are  $\sigma$ -almost causal maps with  $f \overset{\sigma}{\sim} g$ , then  $f, g$  are strongly  $\sigma$ -almost equal. (It is not clear whether equality strengthening is valid for arbitrary maps in **sfKrn**, but we will not need such a general result in this paper.)

In the remainder of the section, we present a basic relationship between disintegration and conditioning  $\sigma|_p$ , introduced above. We assume a joint state  $\omega$  with its two disintegrations  $c_1$  and  $c_2$  in:

$$\begin{array}{c} X \\ \curvearrowright \\ \bullet \\ \overline{\overline{\sigma}} \\ \square c_1 \\ \overline{\overline{\sigma}} \\ \triangleleft \omega \end{array} = \begin{array}{c} X \quad Y \\ | \quad | \\ \triangleleft \omega \end{array} = \begin{array}{c} X_1 \quad Y \\ \square c_2 \\ \bullet \\ \overline{\overline{\sigma}} \\ \triangleleft \omega \end{array} \tag{8}$$

We write  $\omega_1$  and  $\omega_2$  for the first and second marginals of  $\omega$ . (Thus the equations (8) above are the same as (3).)

Let  $q$  be an effect on  $Y$ . It can be extended to an effect  $\mathbf{1} \otimes q$  on  $X \otimes Y$ , where:

$$\mathbf{1} \otimes q \quad := \quad \overline{\overline{\sigma}} \quad \triangleleft \begin{array}{c} q \\ \triangleleft \\ Y \end{array}$$

Then we can form the conditioned state  $\omega|_{\mathbf{1} \otimes q}$ . In a next step we take its first marginal,

written as  $(\omega|_{\mathbf{1} \otimes q})_1$ . It turns out that, in general, this first marginal is different from the original first marginal  $\omega_1$ , even though the effect  $q$  only applies to the second coordinate. This is called ‘crossover influence’ in (Jacobs and Zanasi, 2017). It happens when the state  $\omega$  is ‘entwined’, that is, when its two coordinates are correlated.

A fundamental result in this context is that this crossover influence can also be captured via the channels  $c_1, c_2$  that are extracted from  $\omega$  via disintegrations. This works via effect transformation  $c^*(p) := p \circ c$  and state transformation  $c_*(\sigma) := c \circ \sigma$  along a channel.

**Theorem 7.6.** In the above setting, assuming that the relevant conditioned states exist, there are equalities of states:

$$\omega_1|_{c_1^*(q)} = (\omega|_{\mathbf{1} \otimes q})_1 = (c_2)_*(\omega_2|_q). \quad (9)$$

Following (Jacobs and Zanasi, 2016) we can say that the expression on the left in (9) uses *backward* inference, and the one on the right uses *forward* inference.

*Proof.* We first note that the state in the middle of (9) is the first marginal of:

$$\left( \begin{array}{c} \triangle q \\ \text{---} \\ \omega \end{array} \right)^{-1} \begin{array}{c} \text{---} \\ \triangle q \\ \omega \end{array}$$

Hence:

$$(\omega|_{\mathbf{1} \otimes q})_1 = \left( \begin{array}{c} \triangle q \\ \text{---} \\ \omega \end{array} \right)^{-1} \begin{array}{c} \triangle q \\ \omega \end{array} \quad (10)$$

We note that the above scalar (that is inverted) can also be obtained as:

$$\begin{array}{c} \triangle q \\ \square c_1 \\ \omega \end{array} = \begin{array}{c} \text{---} \\ \triangle q \\ \square c_1 \\ \omega \end{array} \stackrel{(8)}{=} \begin{array}{c} \triangle q \\ \text{---} \\ \omega \end{array} \quad (11)$$

Hence we can prove the equation on the left in (9):

$$\omega_1|_{c_1^*(q)} = \left( \begin{array}{c} \triangle q \\ \square c_1 \\ \omega \end{array} \right)^{-1} \begin{array}{c} \triangle q \\ \omega \end{array} \stackrel{(8),(11)}{=} \left( \begin{array}{c} \triangle q \\ \text{---} \\ \omega \end{array} \right)^{-1} \begin{array}{c} \triangle q \\ \omega \end{array}$$

In a similar way we prove the equation on the right in (9), since  $(c_2)_*(\omega_2|_q)$  equals:

$$\left( \begin{array}{c} \triangle q \\ \omega \end{array} \right)^{-1} \begin{array}{c} \triangle q \\ \omega \end{array} = \left( \begin{array}{c} \triangle q \\ \omega \end{array} \right)^{-1} \begin{array}{c} \triangle q \\ \omega \end{array} \stackrel{(8)}{=} \left( \begin{array}{c} \triangle q \\ \omega \end{array} \right)^{-1} \begin{array}{c} \triangle q \\ \omega \end{array}$$

□

The two equations in Theorem 9 will be illustrated in the ‘disease and mood’ example below, where a particular state (probability) will be calculated in three different ways. For a gentle introduction to Bayesian inference with channels, and many more examples, we refer to (Jacobs and Zanasi, 2018).

**Remark 7.7.** Another way of relating conditioning to disintegration / Bayesian inversion was pointed out by a reviewer of the paper: for a state  $\sigma: I \rightarrow X$  and an effect  $p: X \rightarrow I$ , a conditional state  $\sigma|_p: I \rightarrow X$  is a Bayesian inversion for  $\sigma$  along  $p$ . This, however, involves disintegration of non-causal maps, which we leave to future work.

*Disease and mood example*

We describe an example of probabilistic (Bayesian) reasoning. The setting is the following. We consider a joint state about the occurrence and non-occurrence of a disease, written as a two-element set  $D = \{d, d^\perp\}$ , jointly with the occurrence and non-occurrence of a good mood, written as  $M = \{m, m^\perp\}$ , where  $m^\perp$  stands for a bad (not good) mood. The joint distribution, called  $\omega \in \mathcal{D}(M \times D)$ , that we start from is of the form:

$$\omega = 0.05|m, d\rangle + 0.4|m, d^\perp\rangle + 0.5|m^\perp, d\rangle + 0.05|m^\perp, d^\perp\rangle. \tag{12}$$

Suppose there is a test for the disease, which is positive in 90% of all cases of people having the disease, and still 5% positive for people without the disease. Suppose the disease comes out positive. What is then the mood? It is expected that the mood will deteriorate, since the disease and the mood are ‘entwined’ (correlated) in the above joint state (12): a high likelihood of disease corresponds to a low mood.

The sensitivity is expressed as a channel  $s: D \rightarrow 2$ , where  $2 = \{t, f\}$ , with:

$$c(d) = \frac{9}{10}|t\rangle + \frac{1}{10}|f\rangle \quad c(d^\perp) = \frac{1}{20}|t\rangle + \frac{19}{20}|f\rangle.$$

We calculate the first and second marginal of  $\omega$ :

$$\omega_1 := \omega.[1, 0] = 0.45|m\rangle + 0.55|m^\perp\rangle \quad \omega_2 := \omega.[0, 1] = 0.55|d\rangle + 0.45|d^\perp\rangle.$$

The a priori probability of a positive test is obtained by state transformation, applied to the second marginal:

$$s_*(\omega_2) = 0.518|t\rangle + 0.482|f\rangle.$$

As explained above, we are interested in the mood after a positive test. We write  $tt$



for the predicate on  $2 = \{t, f\}$  given by  $tt(t) = 1$  and  $tt(f) = 0$ . We can transform it to a predicate  $s^*(tt)$  on  $D = \{d, d^\perp\}$ , namely  $s^*(tt)(d) = \frac{9}{10}$  and  $s^*(tt)(d^\perp) = \frac{1}{20}$ . We now describe three equivalent ways to calculate the posterior mood, as in Theorem 7.6, with predicate  $q = s^*(tt)$ .

- 1 First we use the truth predicate  $\mathbf{1}$  on  $M = \{m, m^\perp\}$ , which is always 1, to extend (weaken) the predicate  $s^*(tt)$  on  $D$  to  $\mathbf{1} \otimes s^*(tt)$  on  $M \times D$ . The latter predicate can be used to update the joint state  $\omega \in \mathcal{D}(M \times D)$ . If we then take the first marginal we obtain the updated mood probability:

$$(\omega|_{\mathbf{1} \otimes s^*(tt)}) \cdot [1, 0] = 0.126|m\rangle + 0.874|m^\perp\rangle.$$

Clearly, a positive test leads to a lower mood: a reduction from 0.45 to 0.126.

- 2 Next we extract channels  $c_1: M \rightarrow D$  and  $c_2: D \rightarrow M$  from  $\omega$  via disintegration as in (8). For instance,  $c_1$  is defined as:

$$\begin{aligned} c_1(m) &= \frac{\omega(m,d)}{\omega_1(m)}|d\rangle + \frac{\omega(m,d^\perp)}{\omega_1(m)}|d^\perp\rangle = \frac{1}{9}|d\rangle + \frac{8}{9}|d^\perp\rangle \\ c_1(m^\perp) &= \frac{\omega(m^\perp,d)}{\omega_1(m^\perp)}|d\rangle + \frac{\omega(m^\perp,d^\perp)}{\omega_1(m^\perp)}|d^\perp\rangle = \frac{10}{11}|d\rangle + \frac{1}{11}|d^\perp\rangle. \end{aligned}$$

We can now transform the predicate  $s^*(tt)$  on  $D$  along  $c_1$  to get the predicate  $c_1^*(s^*(tt)) = (s_1 \circ c_1)^*(tt)$  on  $M$  given by:  $d \mapsto 13/90$  and  $d^\perp \mapsto 181/220$ . We can now use this predicate to update the a priori ‘mood’ marginal  $\omega_1 \in \mathcal{D}(M)$ . This gives the same a posteriori outcome as before:

$$\omega_1|_{c_1^*(s^*(tt))} = 0.126|m\rangle + 0.874|m^\perp\rangle.$$

- 3 We can also update the second ‘disease’ marginal  $\omega_2 \in \mathcal{D}(D)$  directly with the predicate  $s^*(tt)$  and then do state transformation along the channel  $c_2: D \rightarrow M$ . This gives the same updated mood:

$$(c_2)_*(\omega_2|_{s^*(tt)}) = 0.126|m\rangle + 0.874|m^\perp\rangle.$$

## 8. Disintegration via likelihoods

We continue in the setting of Section 7 in a CD-category that is not necessarily affine. The goal of this section is to present Theorem 8.3, which generalises a construction of Bayesian inversions using densities/likelihoods shown in Example 3.9.

We first introduce ‘likelihoods’ in our setting.

**Definition 8.1.** We say a channel  $c: X \rightarrow Y$  is represented by a effect  $\ell$  on  $X \otimes Y$  with respect to an arrow  $\nu: I \rightarrow Y$  if

$$\begin{array}{c} \boxed{c} \\ \text{---} \end{array} = \begin{array}{c} \triangle \text{---} \ell \\ \text{---} \text{---} \nu \\ \text{---} \end{array} \quad (13)$$

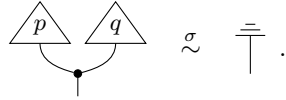
We call  $\ell$  a *likelihood relation* for the channel  $c$  with respect to  $\nu$ .

Interpreted in the category **sfKrn**, the definition says: a kernel  $c: X \rightsquigarrow Y$  satisfies

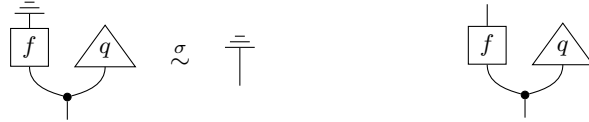
$$c(x, B) = \int_B \ell(x, y) \nu(dy)$$

for a kernel  $\ell: X \times Y \rightsquigarrow 1$  (identified with a measurable function  $\ell: X \times Y \rightarrow [0, \infty]$ ) and a measure  $\nu: 1 \rightsquigarrow Y$ . This is basically the same as what we have in Example 3.9, but here  $\nu$  is not necessarily the Lebesgue measure.

**Definition 8.2.** Let  $\sigma: I \rightarrow X$  be a state. An effect  $p: X \rightarrow I$  is  $\sigma$ -almost invertible if there is an effect  $q: X \rightarrow I$  such that:



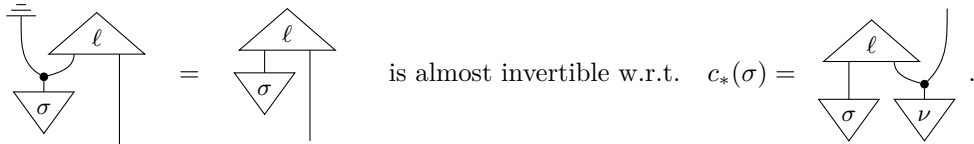
The definition allows us to normalise an arrow  $f: X \rightarrow Y$  into an almost causal one, as follows. If an effect  $\bar{\tau} \circ f$  is  $\sigma$ -almost inverted by  $q$ , as on the left below,



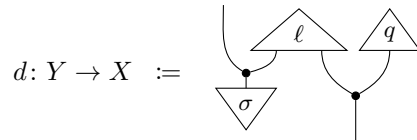
then clearly the arrow  $X \rightarrow Y$  on the right is  $\sigma$ -almost causal.

We can now formulate and prove our main technical result.

**Theorem 8.3.** Let  $\sigma$  be a state on  $X$ , and  $c: X \rightarrow Y$  be a channel represented by a likelihood relation  $\ell$  with respect to  $\nu$  as in (13) above. Assume that the category admits equality strengthening for almost causal maps, and that the effect

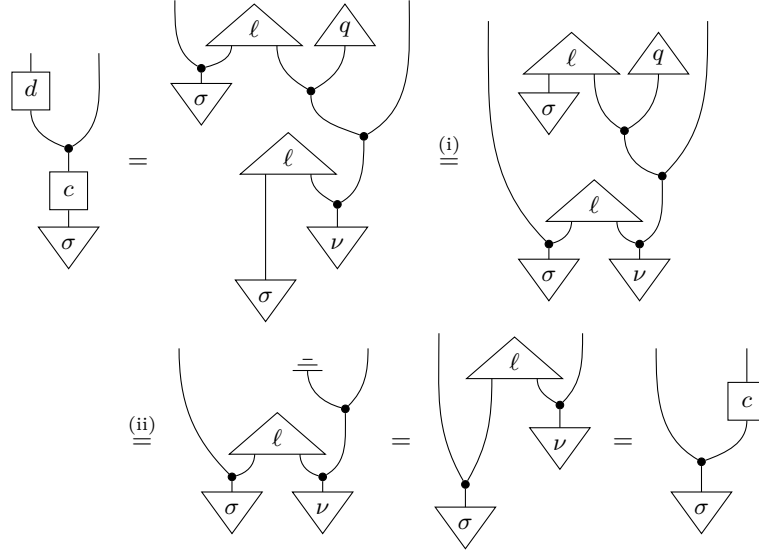


Then, writing  $q: Y \rightarrow I$  for an almost inverse to the effect, the channel

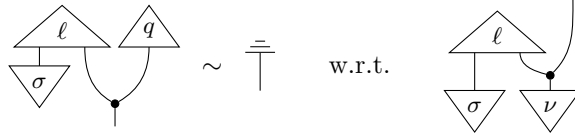


is a Bayesian inversion for  $\sigma$  along  $c: X \rightarrow Y$ . Namely, together they satisfy the equation (5).

*Proof.* We reason as follows.



For the equality  $\stackrel{(i)}{=}$  we use associativity and commutativity of copiers  $\Upsilon$ . The equality  $\stackrel{(ii)}{=}$  follows by:



via equality strengthening. □

**Example 8.4.** We instantiate the Theorem 8.3 in **sfKrn**. Let  $c: X \rightsquigarrow Y$  be a probability kernel represented by a likelihood relation  $\ell: X \times Y \rightsquigarrow 1$  with respect to  $\nu: 1 \rightsquigarrow Y$ . The relation  $\ell$  is identified with a measurable function  $\ell: X \times Y \rightarrow [0, \infty]$  and  $\nu$  with a measure  $\nu: \Sigma_Y \rightarrow [0, \infty]$ . The equation (13) amounts to

$$c(x, B) = \int_B \ell(x, y) \nu(dy) .$$

In particular, each  $\ell(x, -)$  is a probability density function, satisfying  $\int_Y \ell(x, y) \nu(dy) = 1$ . Typically, we use the Lebesgue measure as  $\nu$ , with  $Y$  a subspace of  $\mathbb{R}$ . Let  $\sigma: 1 \rightsquigarrow X$  be a probability measure. Then  $c_*(\sigma): 1 \rightsquigarrow Y$  is given as:

$$c_*(\sigma)(B) = \int_X c(x, B) \sigma(dx) = \int_X \int_B \ell(x, y) \nu(dy) \sigma(dx)$$

The effect

$$p: Y \rightsquigarrow 1 = \begin{array}{c} \triangle \ell \\ | \\ \triangle \sigma \end{array} \quad \text{is given as:} \quad p(y) = \int_X \ell(x, y) \sigma(dx) .$$

To define an inverse of  $p$ , we claim that  $0 < p < \infty$ ,  $c_*(\sigma)$ -almost everywhere. We prove

that  $p^{-1}(\{0, \infty\}) = p^{-1}(0) \cup p^{-1}(\infty)$  is  $c_*(\sigma)$ -negligible, as:

$$\begin{aligned} c_*(\sigma)(p^{-1}(0)) &= \int_X \int_{p^{-1}(0)} \ell(x, y) \nu(dy) \sigma(dx) \\ &= \int_{p^{-1}(0)} p(x) \nu(dy) \\ &= \int_{p^{-1}(0)} 0 \nu(dy) = 0 \end{aligned}$$

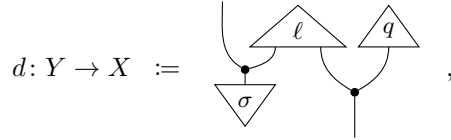
and, similarly we have

$$\int_{p^{-1}(\infty)} \infty \nu(dy) = \int_{p^{-1}(\infty)} p(x) \nu(dy) = c_*(\sigma)(p^{-1}(\infty)) \leq 1$$

but this is possible only when  $\nu(p^{-1}(\infty)) = 0$ , hence  $c_*(\sigma)(p^{-1}(\infty)) = \int_{p^{-1}(\infty)} p(x) \nu(dy) = 0$ . Now define an effect  $q: Y \rightsquigarrow 1$  by

$$q(y) = \begin{cases} p(y)^{-1} & \text{if } 0 < p(y) < \infty \\ 0 & \text{otherwise.} \end{cases}$$

Then  $p$  is  $c_*(\sigma)$ -almost inverted by  $q$ . By Theorem 8.3, the Bayesian inversion for  $\sigma$  along  $c$  is given by



namely,

$$\begin{aligned} d(y, A) &= q(y) \int_A \ell(x, y) \sigma(dx) \\ &= \frac{\int_A \ell(x, y) \sigma(dx)}{\int_X \ell(x, y) \sigma(dx)} \quad \text{whenever } 0 < \int_X \ell(x, y) \sigma(dx) < \infty \end{aligned} \tag{14}$$

This may be seen as a variant of the Bayes formula. The calculation in Example 3.9 is reproduced when  $\sigma$  is also given via a density function.

In the end we reconsider the naive Bayesian classification example from Section 4. There we only considered the discrete version. The original source (Witten et al., 2011) also contains a ‘hybrid’ version, combining discrete and continuous probability.

In that hybrid form the Temperature and Humidity columns in Figure 1 are different, and are given by numerical values. What these values are does not matter too much here, since they are only used to calculate *mean* ( $\mu$ ) and *standard deviation* ( $\sigma$ ) values. This is done separately for the cases where Play is *yes* or *no*. It results in the following table.

	Temperature	Humidity
Play = yes	$\mu = 73, \sigma = 6.2$	$\mu = 79.1, \sigma = 10.2$
Play = no	$\mu = 74.6, \sigma = 7.9$	$\mu = 86.2, \sigma = 9.7$

These  $\mu, \sigma$  values are used to define two functions  $c_T: P \rightarrow \mathcal{G}(\mathbb{R})$  and  $c_H: P \rightarrow \mathcal{G}(\mathbb{R})$ ,

with normal distributions  $\mathcal{N}$  as pdf. To be precise, these channels are defined as follows, where  $A \subseteq \mathbb{R}$  is a measurable subset.

$$\begin{aligned} c_T(y)(A) &= \int_A \mathcal{N}(73, 6.2)(x) \, dx & c_H(y)(A) &= \int_A \mathcal{N}(79.1, 10.2)(x) \, dx \\ c_T(n)(A) &= \int_A \mathcal{N}(74.6, 7.9)(x) \, dx & c_H(y)(A) &= \int_A \mathcal{N}(86.2, 9.7)(x) \, dx. \end{aligned}$$

These functions  $c_T$  and  $c_H$  form channels for the Giry monad  $\mathcal{G}$ .

The two columns Outlook and Windy in Figure 1 remain the same. The corresponding maps  $d_O: P \rightarrow \mathcal{D}(O)$  and  $d_W: P \rightarrow \mathcal{D}(W)$  for the discrete probability monad  $\mathcal{D}$  — in Diagram (6) — are now written as maps  $c_O: P \rightarrow \mathcal{G}(O)$  and  $c_W: P \rightarrow \mathcal{G}(W)$  for the monad  $\mathcal{G}$ , using the obvious inclusion  $\mathcal{D} \hookrightarrow \mathcal{G}$ .

We now have four channels  $c_O: P \rightarrow O$ ,  $c_T: P \rightarrow \mathbb{R}$ ,  $c_H: P \rightarrow \mathbb{R}$  and  $c_W: P \rightarrow W$ . We combine them, as before, into a single channel  $c: P \rightarrow T \times \mathbb{R} \times \mathbb{R} \times W$ . In combination with the Play marginal distribution  $\pi$  from Section 4, we can compute  $c$ 's Bayesian inversion  $f: T \times \mathbb{R} \times \mathbb{R} \times W \rightarrow \mathcal{G}(P)$ . We apply it to the input data used in (Witten et al., 2011), and get the following Play distribution.

$$f(s, 66, 90, t) = 0.207|y\rangle + 0.793|n\rangle.$$

The latter inversion computation produces the probability of 0.207 for playing when the outlook is *Sunny*, the temperature is *66* (Fahrenheit), the humidity is *90%* and the windiness is *true*. The value computed in (Witten et al., 2011) is 20.8%. The minor difference of 0.001 with our outcome can be attributed to (intermediate) rounding errors.

Our computation is done via *EfProb*, using the formula (14).

## References

- Abramsky, S. and Coecke, B. (2009). Categorical quantum mechanics. In *Handbook of Quantum Logic and Quantum Structures: Quantum Logic*, pages 261–323. Elsevier.
- Ackerman, N. L., Freer, C. E., and Roy, D. M. (2011). Noncomputable conditional distributions. In *Logic in Computer Science*, pages 107–116. IEEE.
- Barber, D. (2012). *Bayesian Reasoning and Machine Learning*. Cambridge Univ. Press.
- Bernardo, J. and Smith, A. (2000). *Bayesian Theory*. John Wiley & Sons.
- Borgström, J., Gordon, A., Greenberg, M., Margetson, J., and Gael, J. V. (2013). Measure transformer semantics for Bayesian machine learning. *Logical Methods in Comp. Sci.*, 9(3):1–39.
- Chang, J. T. and Pollard, D. (1997). Conditioning as disintegration. *Statistica Neerlandica*, 51(3):287–317.
- Cho, K. and Jacobs, B. (2017). The EfProb library for probabilistic calculations. In *Conference on Algebra and Coalgebra in Computer Science*, volume 72 of *LIPICs*. Schloss Dagstuhl.
- Cho, K., Jacobs, B., Westerbaan, A., and Westerbaan, B. (2015). An introduction to effectus theory. Preprint. arXiv:1512.05813 [cs.LO].
- Clerc, F., Danos, V., Dahlqvist, F., and Garnier, I. (2017). Pointless learning. In *Foundations of Software Science and Computation Structures*, volume 10203 of *Lect. Notes Comp. Sci.*, pages 355–369. Springer.
- Coecke, B. (2016). Terminality implies no-signalling ...and much more than that. *New Generation Computing*, 34(1):69–85.

- Coecke, B. and Kissinger, A. (2017). *Picturing Quantum Processes: A First Course in Quantum Theory and Diagrammatic Reasoning*. Cambridge University Press.
- Coumans, D. and Jacobs, B. (2013). Scalars, monads and categories. In Heunen, C., Sadrzadeh, M., and Grefenstette, E., editors, *Quantum Physics and Linguistics. A Compositional, Diagrammatic Discourse*, pages 184–216. Oxford Univ. Press.
- Culbertson, J. and Sturtz, K. (2014). A categorical foundation for Bayesian probability. *Applied Categorical Structures*, 22(4):647–662.
- D’Ariano, G. M., Chiribella, G., and Perinotti, P. (2017). *Quantum Theory from First Principles: An Informational Approach*. Cambridge University Press.
- Dawid, A. (2001). Separoids: A mathematical framework for conditional independence and irrelevance. *Annals of Mathematics and Artificial Intelligence*, 32(1):335–372.
- Doberkat, E.-E. (2009). *Stochastic Coalgebraic Logic*. Monographs in Theoretical Computer Science. Springer.
- Faden, A. M. (1985). The existence of regular conditional probabilities: Necessary and sufficient conditions. *The Annals of Probability*, 13(1):288–298.
- Fong, B. (2012). Causal theories: A categorical perspective on Bayesian networks. Master’s thesis, Univ. of Oxford. arXiv:1301.6201 [math.PR].
- Fremlin, D. H. (2000). *Measure Theory (5 volumes)*. Torres Fremlin.
- Geiger, D., Verma, T., and Pearl, J. (1990). Identifying independence in Bayesian networks. *Networks*, 20:507–534.
- Giry, M. (1982). A categorical approach to probability theory. In *Categorical Aspects of Topology and Analysis*, volume 915 of *Lecture Notes in Mathematics*, pages 68–85. Springer.
- Gordon, A., Henzinger, T., Nori, A., and Rajamani, S. (2014). Probabilistic programming. In *Future of Software Engineering*, pages 167–181. ACM.
- Jacobs, B. (2015). New directions in categorical logic, for classical, probabilistic and quantum logic. *Logical Methods in Comp. Sci.*, 11(3):1–76.
- Jacobs, B. (2017). From probability monads to commutative effectuses. *Journ. of Logical and Algebraic Methods in Programming*, 156.
- Jacobs, B., Westerbaan, B., and Westerbaan, A. (2015). States of convex sets. In *Foundations of Software Science and Computation Structures*, volume 9034 of *Lect. Notes Comp. Sci.*, pages 87–101. Springer.
- Jacobs, B. and Zanasi, F. (2016). A predicate/state transformer semantics for Bayesian learning. In *Math. Found. of Programming Semantics*, volume 325 of *Elect. Notes in Theor. Comp. Sci.*, pages 185–200. Elsevier.
- Jacobs, B. and Zanasi, F. (2017). A formal semantics of influence in Bayesian reasoning. In *Math. Found. of Computer Science*, volume 83 of *LIPICs*. Schloss Dagstuhl.
- Jacobs, B. and Zanasi, F. (2018). The logical essentials of Bayesian reasoning. See [arxiv.org/abs/1804.01193](https://arxiv.org/abs/1804.01193).
- Joyal, A. and Street, R. (1991). The geometry of tensor calculus, I. *Advances in Mathematics*, 88(1):55–112.
- Kallenberg, O. (2017). *Random Measures, Theory and Applications*. Springer.
- Katoen, J.-P., Gretz, F., Jansen, N., Lucien Kaminski, B., and Olmedo, F. (2015). Understanding probabilistic programs. In *Correct System Design*, volume 9360 of *Lect. Notes Comp. Sci.*, pages 15–32. Springer.
- Mac Lane, S. (1998). *Categories for the Working Mathematician*. Springer, second edition.
- Pachl, J. K. (1978). Disintegration and compact measures. *Mathematica Scandinavica*, 43:157–168.
- Panangaden, P. (2009). *Labelled Markov Processes*. Imperial College Press.

- Pawitan, Y. (2001). *In All Likelihood. Statistical Modelling and Inference Using Likelihood*. Clarendon Press.
- Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems*. Morgan Kaufmann.
- Pollard, D. (2002). *A User's Guide to Measure Theoretic Probability*. Cambridge University Press.
- Selinger, P. (2010). A survey of graphical languages for monoidal categories. In *New Structures for Physics*, volume 813 of *Lecture Notes in Physics*, pages 289–355. Springer.
- Shan, C.-c. and Ramsey, N. (2017). Exact Bayesian inference by symbolic disintegration. In *Princ. of Programming Languages*, pages 130–144. ACM.
- Simpson, A. (2017). Category-theoretic structure for independence and conditional independence. In *Math. Found. of Programming Semantics*. To appear.
- Staton, S. (2017). Commutative semantics for probabilistic programming. In *European Symp. on Programming*, volume 10201 of *Lect. Notes Comp. Sci.*, pages 855–879. Springer.
- Staton, S., Yang, H., Heunen, C., Kammar, O., and Wood, F. (2016). Semantics for probabilistic programming: higher-order functions, continuous distributions, and soft constraints. In *Logic in Computer Science*, pages 525–534. ACM.
- Stoyanov, J. M. (2014). *Counterexamples in Probability*. Dover, third edition.
- Verma, T. and Pearl, J. (1988). Causal networks: Semantics and expressiveness. In *Uncertainty in Artificial Intelligence*.
- Witten, I., Frank, E., and Hall, M. (2011). *Data Mining – Practical Machine Learning Tools and Techniques*. Elsevier, Amsterdam.