

PDF hosted at the Radboud Repository of the Radboud University Nijmegen

The following full text is a publisher's version.

For additional information about this publication click this link.

<https://hdl.handle.net/2066/216753>

Please be advised that this information was generated on 2020-09-18 and may be subject to change.

PAPER

A channel-based perspective on conjugate priors[†]

B. Jacobs*

Institute for Computing and Information Sciences, Radboud University, P.O. Box 9010, 6500 GL Nijmegen, The Netherlands

*Email: bart@cs.ru.nl

(Received 1 July 2017; revised 5 March 2019; accepted 14 April 2019)

Abstract

A desired closure property in Bayesian probability is that an updated posterior distribution be in the same class of distributions – say Gaussians – as the prior distribution. When the updating takes place via a statistical model, one calls the class of prior distributions the ‘conjugate priors’ of the model. This paper gives (1) an abstract formulation of this notion of conjugate prior, using channels, in a graphical language, (2) a simple abstract proof that such conjugate priors yield Bayesian inversions and (3) an extension to multiple updates. The theory is illustrated with several standard examples.

1. Introduction

The main result of this paper, Theorem 6.3, is mathematically trivial. But it is not entirely trivial to see that this result is trivial. The effort and contribution of this paper lie in setting up a framework – using the abstract language of channels, Kleisli maps and string diagrams for probability theory – to define the notion of conjugate prior in such a way that there is a trivial proof of the main statement, saying that conjugate priors yield Bayesian inversions. This is indeed what conjugate priors are meant to be.

Conjugate priors form a fundamental topic in Bayesian theory. They are commonly described via a closure property of a class of prior distributions, namely as being closed under certain Bayesian updates. Conjugate priors are especially useful because they do not only involve a closure *property*, but also a particular *structure*, namely an explicit function that performs an analytical computation of posterior distributions via updates of the parameters. This greatly simplifies Bayesian analysis. For instance, the Beta distribution is conjugate prior to the Bernoulli (or ‘flip’) distribution, and also to the binomial distribution: updating a $\text{Beta}(\alpha, \beta)$ prior via a Bernoulli/binomial statistical model yields a new $\text{Beta}(\alpha', \beta')$ prior, with adapted parameters α', β' that can be computed explicitly from α, β and the observation at hand. Despite this importance, the descriptions in the literature of what it means to be a conjugate prior are remarkably informal. One does find several lists of classes of distributions, for instance at Wikipedia,¹ together with formulas about how to recompute parameters. The topic has a long and rich history in statistics (see e.g. Bishop (2006)), with much emphasis on exponential families (Diaconis and Ylvisaker 1979), but a precise, general definition is hard to find.

We briefly review some common approaches, without any pretension to be complete: the definition in Alpaydin (2010) is rather short, based on an example, and just says: ‘We see that the

[†]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.

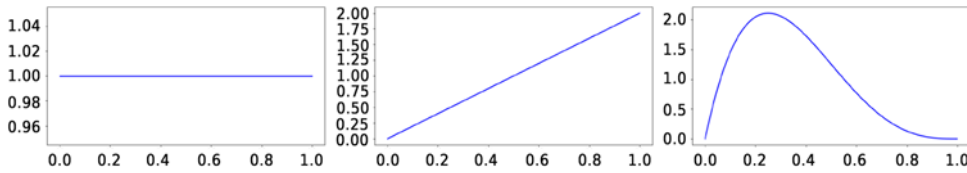


Figure 1. Uniform prior, and two posterior probability density functions on $[0, 1]$, after observing head, and after observing head–tail–tail–tail. These functions correspond respectively to Beta(1, 1), Beta(2, 1), and Beta(2, 4). These three plots can also be obtained via suitable Bayesian updates (inversions).

posterior has the same form as the prior and we call such a prior a *conjugate prior*. Also Russell and Norvig (2003) mention the term ‘conjugate prior’ only in relation to an example. There is a separate section in Bishop (2006) about conjugate priors, but no precise definition. Instead, there is the informal description ‘... the posterior distribution has the same functional form as the prior’. The most precise definition (known to the author) is in Bernardo and Smith (2000, Section 5.2), where the conjugate family with respect to a statistical model, assuming a ‘sufficient statistic’, is described. It comes close to our channel-based description, since it explicitly mentions the conjugate family as a conditional probability distribution with (recomputed) parameters. The approach is rather concrete however, and the high level of mathematical abstraction that we seek here is missing in Bernardo and Smith (2000).

This paper presents a novel systematic perspective for precisely defining what conjugate priorship means. It uses the notion of ‘channel’ as starting point. The basis of this approach lies in category theory, especially effectus theory (Cho et al. 2015; Jacobs 2015). However, we try to make this paper accessible to non-category theorists, by using the term ‘channel’ instead of morphism in a Kleisli category of a suitable monad. Moreover, a graphical language is used for channels, which hopefully makes the approach more intuitive. Thus, the emphasis of the paper is on *what it means* to have conjugate priors. It does not offer new perspectives on how to find/obtain them.

The paper is organised as follows. It starts in Section 2 with a high-level description of the main ideas, without going into technical details. Preparatory definitions are provided in Sections 3 and 4, dealing with channels in probabilistic computation, with a diagrammatic language for channels, and with Bayesian inversion. Then, Section 5 contains the novel channel-based definition of conjugate priorship; it also illustrates how several standard examples fit in this new setting. Section 6 establishes the (expected) close relationship between conjugate priors and Bayesian inversions (Cho and Jacobs, 2019; Clerc et al., 2017). Section 7 illustrates that multiple updates are handled basically in the same way as single updates, and also how the notion of ‘sufficient statistic’ fits in.

2. Main Ideas

This section gives an informal description of the main ideas underlying this paper. It starts with a standard example, and then proceeds with a step-by-step introduction to the essentials of the perspective of this paper.

A well-known example in Bayesian reasoning is inferring the (unknown) bias of a coin from a sequence of consecutive head/tail observations. The bias is a number $r \in [0, 1]$ in the unit interval, giving the ‘Bernoulli’ or ‘flip’ probability r for head, and $1 - r$ for tail. Initially we assume a uniform distribution for r , as described by the constant probability density function (pdf) on the left in Figure 1. After observing one head, this pdf changes to the second picture. After observing head–tail–tail–tail we get the third pdf. These pictures are obtained by Bayesian inversion, see Section 4.

It is a well-known fact that all the resulting distributions are instances of the Beta(α, β) family of distributions, for different parameters α, β . After each observation, one can recompute the entire updated distribution, via Bayesian inversion. But in fact there is a much more efficient way

to obtain the revised distribution, namely by computing the new parameter values: increment α by one, for head, and increment β by one for tail, see Example 5.3 for details. The family of distributions $\text{Beta}(\alpha, \beta)$, indexed by parameters α, β , is thus suitably closed under updates with Bernoulli. It is the essence of the statement that Beta is conjugate prior to Bernoulli. This will be made precise later on.

Let $X = (X, \Sigma)$ be a measurable space, where $\Sigma \subseteq \mathcal{P}(X)$ is a σ -algebra of measurable subsets. We shall write $\mathcal{G}(X)$ for the set of probability distributions on X . Elements $\omega \in \mathcal{G}(X)$ are thus countably additive functions $\omega: \Sigma \rightarrow [0, 1]$ with $\omega(X) = 1$.

Idea 1: A family of distributions on X , indexed by a measurable space P of parameters, is a (measurable) function $P \rightarrow \mathcal{G}(X)$. Categorically, such a function is a Kleisli map for \mathcal{G} , considered as monad on the category of measurable spaces (see Section 3). These Kleisli maps are also called channels and will be written simply as arrows $P \rightarrow X$, or diagrammatically as boxes $\begin{array}{c} \square \\ \downarrow \\ P \end{array}^X$ where we imagine that information is flowing upwards.

The study of families of distributions goes back a long way, e.g. as ‘experiments’ (Blackwell, 1951).

Along these lines we shall describe the family of Beta distributions as a channel with $P = \mathbb{R}_{>0} \times \mathbb{R}_{>0}$ and $X = [0, 1]$, namely as function:

$$\mathbb{R}_{>0} \times \mathbb{R}_{>0} \xrightarrow{\text{Beta}} \mathcal{G}([0, 1]) \tag{1}$$

For $(\alpha, \beta) \in \mathbb{R}_{>0} \times \mathbb{R}_{>0}$ there is the probability distribution $\text{Beta}(\alpha, \beta) \in \mathcal{G}([0, 1])$ determined by its value on a measurable subset $M \subseteq [0, 1]$, which is obtained via integration:

$$\text{Beta}(\alpha, \beta)(M) = \int_M \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)} dx, \tag{2}$$

where $B(\alpha, \beta) = \int_{[0,1]} x^{\alpha-1}(1-x)^{\beta-1} dx$ is a normalisation constant.

A conjugate prior relationship involves a family of distributions $P \rightarrow \mathcal{G}(X)$ which is closed wrt. updates based on observations (or: data) from a separate domain O . Each ‘parameter’ element $x \in X$ gives rise to a separate distribution on O . This is what is usually called a *statistical* or *parametric* model. We shall also describe it as a channel.

Idea 2: The observations for a family $P \rightarrow \mathcal{G}(X)$ arise via another ‘Kleisli’ map $X \rightarrow \mathcal{G}(O)$ representing the statistical model. Conjugate priorship will be defined for two such composable channels $P \rightarrow X \rightarrow O$, where O is the space of observations.

In the above coin example, the space O of observations is the two-element set $2 = \{0, 1\}$ where 0 is for tail and 1 for head. The Bernoulli channel is written as $\text{Flip}: [0, 1] \rightarrow \mathcal{G}(2)$. A probability $r \in [0, 1]$ determines a Bernoulli/flip/coin probability distribution $\text{Flip}(r) \in \mathcal{G}(2)$ on 2, formally sending the subset $\{1\}$ to r and $\{0\}$ to $1 - r$.

Idea 3: A channel $c: P \rightarrow X$ is a conjugate prior to a channel $d: X \rightarrow O$ if there is a *parameter translation function* $h: P \times O \rightarrow P$ satisfying a suitable equation.

The idea is that $c(p)$ is a prior, for $p \in P$, which gets updated via the statistical model (channel) d , in the light of observation $y \in O$. The revised, updated distribution is $c(h(p, y))$. The model d is usually written as a conditional probability $d(y | \theta)$.

In the coin example we have $h: \mathbb{R}_{>0} \times \mathbb{R}_{>0} \times 2 \rightarrow \mathbb{R}_{>0} \times \mathbb{R}_{>0}$ given by $h(\alpha, \beta, 1) = (\alpha + 1, \beta)$ and $h(\alpha, \beta, 0) = (\alpha, \beta + 1)$, see Example 5.3 below for more information.

What has been left unexplained is the ‘suitable’ equation that the parameter translation function $h: P \times O \rightarrow P$ should satisfy. It is not entirely trivial, because it is an equation between channels in what is called the Kleisli category $\mathcal{Kl}(\mathcal{G})$ of the Girly monad \mathcal{G} . At this stage we need to move to a more categorical description. The equation, which will appear in Definition 5.1, bears similarities with the notion of Bayesian inversion, which will be introduced in Section 4.

3. Channels and conditional probabilities

This section will describe conditional probabilities as arrows and will show how to compose them. Thereby we are entering the world of category theory. We aim to suppress the underlying categorical machinery and make this work accessible to readers without such background. For those with categorical background knowledge: we will be working in the Kleisli categories of the distribution monad \mathcal{D} for discrete probability, and of the Girly monad \mathcal{G} for continuous probability, see e.g. Girly (1982), Panangaden (2009) and Jacobs (2018). Discrete distributions may be seen as a special case of continuous distributions, via a suitable inclusion map $\mathcal{D} \rightarrow \mathcal{G}$. Hence one could give one account, using \mathcal{G} only. However, in computer science, unlike for instance in statistics, discrete distributions are so often used that they merit separate treatment.

We thus start with discrete probability. We write a (finite, discrete) distribution on a set X as a formal convex sum $r_1|x_1\rangle + \dots + r_n|x_n\rangle$ of elements $x_i \in X$ and probabilities $r_i \in [0, 1]$ with $\sum_i r_i = 1$. The ‘ket’ notation $|-\rangle$ is syntactic sugar, used to distinguish elements of x from their occurrence $|x\rangle$ in such formal convex sums.² A distribution as above can be identified with a ‘probability mass’ function $\omega: X \rightarrow [0, 1]$ which is r_i on x_i and 0 elsewhere. We often implicitly identify distributions with such functions. We shall write $\mathcal{D}(X)$ for the set of distributions on X .

We shall focus on functions of the form $c: X \rightarrow \mathcal{D}(Y)$. They give for each element $x \in X$ a distribution $c(x)$ on Y . Hence such functions form an X -indexed collection $(c(x))_{x \in X}$ of distributions $c(x)$ on Y . They can be understood as conditional probabilities $P(y|x) = r$, if $c(x)$ is of the form $\dots r|y\rangle \dots$, with weight $r = c(x)(y) \in [0, 1]$ for $y \in Y$. Thus, by construction, $\sum_y P(y|x) = 1$, for each $x \in X$. Moreover, if the sets X and Y are finite, we can describe $c: X \rightarrow \mathcal{D}(Y)$ as a stochastic matrix, with entries $P(y|x)$, adding up to one – per row or column, depending on the chosen orientation of the matrix.

We shall often write functions $X \rightarrow \mathcal{D}(Y)$ simply as arrows $X \rightarrow Y$, call them ‘channels’, and write them as ‘boxes’ in diagrams. This arrow notation is justified, because there is a natural way to compose channels, as we shall see shortly. But first we describe *state transformation*, also called *prediction*. Given a channel $c: X \rightarrow \mathcal{D}(Y)$ and a state $\omega \in \mathcal{D}(X)$, we can form a new state, written as $c \gg \omega$, on Y . It is defined as:

$$c \gg \omega := \sum_y \left(\sum_x \omega(x) \cdot c(x)(y) \right) |y\rangle. \tag{3}$$

The outer sum \sum_y is a formal convex sum, whereas the inner sum \sum_x is an actual sum in the unit interval $[0, 1]$. Using state transformation \gg it is easy to define composition of channels: given functions $c: X \rightarrow \mathcal{D}(Y)$ and $d: Y \rightarrow \mathcal{D}(Z)$, we use the ordinary composition symbol \circ to form a composite channel $d \circ c: X \rightarrow \mathcal{D}(Z)$, where:

$$(d \circ c)(x) := d \gg c(x) = \sum_{z \in Z} \left(\sum_y c(x)(y) \cdot d(y)(z) \right) |z\rangle. \tag{4}$$

Essentially, this is matrix composition for stochastic matrices. Channel composition \circ is associative and also has a neutral element, namely the identity channel $\eta: X \rightarrow X$ given by the ‘Dirac’ function $\eta(x) = 1|x\rangle$. It is not hard to see that $(d \circ c) \gg \omega = d \gg (c \gg \omega)$, see e.g. Jacobs and Zanasi (2019) for more details.

We turn to channels in continuous probability. As already mentioned in Section 2, we write $\mathcal{G}(X)$ for the set of probability distributions $\omega: \Sigma_X \rightarrow [0, 1]$, where $X = (X, \Sigma_X)$ is a measurable space. These probability distributions are (also) called states. The set $\mathcal{G}(X)$ carries a σ -algebra itself, but that does not play an important role here. Each element $x \in X$ yields a probability measure $\eta(x) \in \mathcal{G}(X)$, with $\eta(x)(M) = \mathbf{1}_M(x)$, which is 1 if $x \in M$ and 0 otherwise. This map $\mathbf{1}_M: X \rightarrow [0, 1]$ is called the indicator function for the subset $M \in \Sigma_X$.

For a state/measure $\omega \in \mathcal{G}(X)$ and a measurable function $f: X \rightarrow \mathbb{R}_{\geq 0}$ we write $\int f d\omega$ for the Lebesgue integral, if it exists. We follow the notation of Jacobs (2013) and refer there for details, or alternatively, to Panangaden (2009). We recall that an integral $\int_M f d\omega$ over a measurable subset $M \subseteq X$ of the domain of f is defined as $\int \mathbf{1}_M \cdot f d\omega$, and that $\int \mathbf{1}_M d\omega = \omega(M)$. Moreover, $\int f d\eta(x) = f(x)$.

For a measurable function $g: X \rightarrow Y$ between measurable spaces X, Y there is the ‘push forward’ function $\mathcal{G}(g): \mathcal{G}(X) \rightarrow \mathcal{G}(Y)$, given by $\mathcal{G}(g)(\omega)(N) = \omega(g^{-1}(N))$. It satisfies:

$$\int f d\mathcal{G}(g)(\omega) = \int f \circ g d\omega. \tag{5}$$

Often, the measurable space X is a subset $X \subseteq \mathbb{R}$ of the real numbers and a probability distribution ω on X is given by a probability density function (pdf), that is, by a measurable function $f: X \rightarrow \mathbb{R}_{\geq 0}$ with $\int_X f(x) dx = 1$. Such a pdf f gives rise to a state $\omega \in \mathcal{G}(X)$, namely:

$$\omega(M) = \int_M f(x) dx. \tag{6}$$

We then write $\omega = \int f$. In such cases where the measure is not explicitly mentioned, one should assume that it is the Lebesgue measure on \mathbb{R} .

In this continuous context a channel is a measurable function $c: X \rightarrow \mathcal{G}(Y)$, for measurable spaces X, Y . Like in the discrete case, it gives an X -indexed collection $(c(x))_{x \in X}$ of probability distributions on Y . The channel c can transform a state $\omega \in \mathcal{G}(X)$ on X into a state $c \gg \omega \in \mathcal{G}(Y)$ on Y , given on a measurable subset $N \subseteq Y$ as:

$$(c \gg \omega)(N) = \int c(-)(N) d\omega. \tag{7}$$

For another channel $d: Y \rightarrow \mathcal{G}(Z)$ there is a composite channel $d \circ c: X \rightarrow \mathcal{G}(Z)$, via integration:

$$(d \circ c)(x)(K) := (d \gg c(x))(K) = \int d(-)(K) dc(x) \tag{8}$$

In many situations a channel $c: X \rightarrow \mathcal{G}(Y)$ is given by an indexed probability density function (pdf) $u: X \times Y \rightarrow \mathbb{R}_{\geq 0}$, with $\int u(x, y) dy = 1$ for each $x \in X$. The associated channel c is:

$$c(x)(N) = \int_N u(x, y) dy. \tag{9}$$

In that case we simply write $c = \int u$ and call c a pdf-channel. We have already seen such a description of the Beta distribution as a pdf-channel in (2).

(In these pdf-channels $X \rightarrow Y$, we use a collection of pdf’s $u(x, -)$ which are all dominated by the Lebesgue measure. This domination happens via the relationship \ll of absolute continuity, using the Radon–Nikodym Theorem, see e.g. Panangaden (2009).)

Various additional computation rules for integrals are given in the Appendix.

4. Bayesian inversion in string diagrams

In this paper, we make superficial use of string diagrams to graphically represent sequential and parallel composition of channels, mainly in order to provide an intuitive visual overview. We refer to Selinger (2011) for mathematical details and mention here only the essentials.

A channel $X \rightarrow Y$, for instance of the sort discussed in the previous section, can be written as a box $\begin{array}{c} \square \\ \downarrow \\ X \end{array}$ with information flowing upwards, from the wire labelled with X to the wire labelled with Y . Composition of channels, as in (4) or (8), simply involves connecting wires (of the same type). The identity channel is just a wire. We use a triangle notation ∇^x for a state on X . It is special case of a channel, namely of the form $1 \rightarrow X$ with trivial singleton domain 1.

In the present (probabilistic) setting, we allow copying of wires, written diagrammatically as Υ . We briefly describe such copy channels for discrete and continuous probability:

$$\begin{array}{ccc} X \xrightarrow{\Upsilon} \mathcal{D}(X \times X) & & X \xrightarrow{\Upsilon} \mathcal{G}(X \times X) \\ x \mapsto 1|x, x & & x \mapsto (M \times N \mapsto \mathbf{1}_{M \cap N}(x)) \end{array}$$

After such a copy we can use parallel channels. We briefly describe how this works, first in the discrete case. For channels $c: X \rightarrow \mathcal{D}(Y)$ and $d: A \rightarrow \mathcal{D}(B)$, we have a channel $c \otimes d: X \times A \rightarrow \mathcal{D}(Y \times B)$ given by:

$$(c \otimes d)(x, a) = \sum_{y,b} c(x)(y) \cdot d(a)(b)|y, b\rangle.$$

Similarly, in the continuous case, for channels $c: X \rightarrow \mathcal{G}(Y)$ and $d: A \rightarrow \mathcal{G}(B)$ we get $c \otimes d: X \times A \rightarrow \mathcal{G}(Y \times B)$ given by:

$$(c \otimes d)(x, a)(M \times N) = c(x)(M) \cdot d(a)(N).$$

Recall that the product σ -algebra on $Y \times B$ is generated by measurable rectangles of the form $M \times N$, for $M \in \Sigma_Y$ and $N \in \Sigma_B$. Hence measures are uniquely determined by their actions on such rectangles.

We can now formulate what Bayesian inversion is. The definition is couched in purely diagrammatic language, but is applied only to probabilistic interpretations in this paper.

Definition 4.1. *The Bayesian inversion of a channel $c: X \rightarrow Y$ with respect to a state ω of type X , if it exists, is a channel in the opposite direction, written as $c^\dagger_\omega: Y \rightarrow X$, such that the following equation holds.*

$$\begin{array}{c} \text{Diagram 1: } \begin{array}{c} \square \\ \downarrow \\ \omega \end{array} \quad = \quad \begin{array}{c} \square \\ \downarrow \\ \square \\ \downarrow \\ \omega \end{array} \end{array} \tag{10}$$

The dagger notation c^\dagger_ω is copied from Clerc et al. (2017), see also Cho and Jacobs (2019). There the state ω is left implicit, via a restriction to a certain comma category of kernels. In that setting the operation $(-)^\dagger$ is functorial and forms a dagger category (see e.g. Abramsky and Coecke (2009), Selinger (2007) for definitions). In particular, it preserves composition and identities of channels. Equation (10) can also be written as: $\langle \text{id}, c \rangle \gg \omega = \langle c^\dagger_\omega, \text{id} \rangle \gg (c \gg \omega)$. Alternatively, in the discrete case, with variables explicit, it says: $c(x)(y) \cdot \omega(x) = c^\dagger_\omega(y)(x) \cdot (c \gg \omega)(y)$. This comes close to the ‘adjointness’ formulations that are typical for daggers.

It is relatively easy to define Bayesian inversion in discrete probability theory: for a channel $c: X \rightarrow \mathcal{D}(Y)$ and a state/distribution $\omega \in \mathcal{D}(X)$ one can define a channel $c^\dagger_\omega: Y \rightarrow \mathcal{D}(X)$ as:

$$c^\dagger_\omega(y)(x) = \frac{\omega(x) \cdot c(x)(y)}{(c \gg \omega)(y)} = \frac{\omega(x) \cdot c(x)(y)}{\sum_z \omega(z) \cdot c(z)(y)}, \tag{11}$$

assuming that the denominator is non-zero. This corresponds to the familiar formula $P(B | A) = P(A,B)/P(A)$ for conditional probability. The state $c^\dagger_\omega(y)$ can alternatively be defined via updating the state ω with the point predicate $\{y\}$, transformed via c into a predicate $c \ll \mathbf{1}_{\{y\}}$ on X , see Cho and Jacobs (2019) for details.

The situation is much more difficult in continuous probability theory, since Bayesian inversions may not exist (Ackerman et al. 2011; Stoyanov 2014) or may be determined only up to measure zero. But when restricted to e.g. standard Borel spaces, as in Clerc et al. (2017), existence is ensured, see also Culbertson and Sturtz (2014) and Faden (1985). Another common solution is to assume that we have a pdf-channel: there is a map $u: X \times Y \rightarrow \mathbb{R}_{\geq 0}$ that defines a channel $c: X \rightarrow \mathcal{G}(Y)$, like in (9), as $c(x)(N) = \int_N u(x, y) dy$. Then, for a distribution $\omega \in \mathcal{G}(X)$ we can take as Bayesian inversion:

$$\begin{aligned} c^\dagger_\omega(y)(M) &= \frac{\int_M u(-, y) d\omega}{\int_X u(-, y) d\omega} \\ &= \frac{\int_M f(x) \cdot u(x, y) dx}{\int_X f(x) \cdot u(x, y) dx} \quad \text{when } \omega = \int f(x) dx. \end{aligned} \tag{12}$$

We prove that this definition satisfies the inversion Equation (10), using the calculation rules from the Appendix.

$$\begin{aligned} (\langle c^\dagger_\omega, \text{id} \rangle \gg (c \gg \omega))(M \times N) &\stackrel{(7)}{=} \int \langle c^\dagger_\omega, \text{id} \rangle(-)(M \times N) d(c \gg \omega) \\ &\stackrel{(A.2,A.3)}{=} \int \left(\int f(x) \cdot u(x, y) dx \right) \cdot \langle c^\dagger_\omega, \text{id} \rangle(y)(M \times N) dy \\ &\stackrel{(A.6)}{=} \int \left(\int f(x) \cdot u(x, y) dx \right) \cdot c^\dagger_\omega(y)(M) \cdot \mathbf{1}_N(y) dy \\ &\stackrel{(12)}{=} \int_N \left(\int f(x) \cdot u(x, y) dx \right) \cdot \frac{\int_M f(x) \cdot u(x, y) dx}{\int f(x) \cdot u(x, y) dx} dy \\ &= \int_N \int_M f(x) \cdot u(x, y) dx dy \\ &\stackrel{(A.7)}{=} (\langle \text{id}, c \rangle \gg \omega)(M \times N). \end{aligned}$$

5. Conjugate priors

We now come to the core of this paper. As described in the introduction, the informal definition says that a class of distributions is conjugate prior to a statistical model if the associated posteriors are *in the same class* of distributions. The posteriors can be computed via Bayesian inversion (12) of the statistical model.

This definition of ‘conjugate prior’ is a bit vague, since it loosely talks about ‘classes of distributions’, without further specification. As described in ‘Idea 1’ in Section 2, we interpret ‘class of states on X ’ as a channel $P \rightarrow X$, where P is the type of parameters of the class.

We have already seen this channel-based description for the class Beta distributions, in (1), as channel Beta: $\mathbb{R}_{>0} \times \mathbb{R}_{>0} \rightarrow [0, 1]$. This works more generally, for instance for Gaussian

(normal) distributions $\text{Norm}(\mu, \sigma)$, where μ is the mean parameter and σ is the standard deviation parameter, giving a channel of the form:

$$\mathbb{R} \times \mathbb{R}_{>0} \xrightarrow{\text{Norm}} \mathcal{G}(\mathbb{R}) \tag{13}$$

It is determined by its value on a measurable subset $M \subseteq \mathbb{R}$ as the standard integral:

$$\text{Norm}(\mu, \sigma)(M) = \int_M \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx \tag{14}$$

Given a channel $c: P \rightarrow X$, we shall look at states $c(p)$, for parameters $p \in P$, as priors. The statistical model, for which these $c(p)$'s will be described as conjugate priors, goes from X to some other object O of 'observations'. Thus our starting point is a statistical model consisting of a pair of (composable) channels of the form:

$$P \xrightarrow{c} X \xrightarrow{d} O \quad \text{or, as diagram,} \quad \begin{array}{c} \boxed{d} \\ \boxed{c} \end{array} \tag{15}$$

Such a pair of composable channels may be seen as a two-stage hierarchical Bayesian model, in which we standardly use $O =$ observables, $X =$ parameters of the statistical model, and $P =$ parameters over the parameters, i.e. hyperparameters, see e.g. Bernardo and Smith (2000). There, esp. in Definition 5.6 of conjugate priorship, one can also distinguish two channels, written as $p(\theta | \tau)$ and $p(x | \theta)$, corresponding respectively to our channels c and d . The τ form the hyperparameters.

In this setting we come to our main definition that formulates the notion of conjugate prior in an abstract manner, avoiding classes of distributions. It contains the crucial equation that was missing in the informal description in Section 2.

All our examples of (conjugate prior) channels are maps in the Kleisli category of the Giry monad, but the formulation applies more generally. In fact, abstraction purifies the situation and shows the essentials. The definition below speaks of 'deterministic' channels, between brackets. This part will be explained later on, in the beginning of Section 6. It can be ignored for now.

Definition 5.1. *In the situation (15) we call channel c a conjugate prior to channel d if there is a (deterministic) channel $h: P \times O \rightarrow P$ for which the following equation holds:*

$$\begin{array}{c} \boxed{d} \\ \bullet \\ \boxed{c} \end{array} = \begin{array}{c} \boxed{h} \\ \bullet \\ \boxed{d} \\ \boxed{c} \end{array} \tag{16}$$

Equivalently, in equational form: $\langle \text{id}, d \rangle \circ c = ((c \circ h) \otimes \text{id}) \circ \langle \text{id}, \Upsilon \circ d \circ c \rangle$.

The idea is that the map $h: P \times O \rightarrow P$ translates parameters, with an observation from O as additional argument. Informally, one gets a posterior state $c(h(p, y))$ from the prior state $c(p)$, given the observation $y \in O$. The power of this 'analytic' approach is that it involves simple recomputation of parameters, instead of more complicated updating of entire states. This will be illustrated in several standard examples below.

The above Equation (16) is formulated in an abstract manner – which is its main strength. We will derive an alternative formulation of Equation (16) for pdf-channels. It greatly simplifies the calculations in examples.

Lemma 5.2. Consider composable channels $P \xrightarrow{c} X \xrightarrow{d} O$, as in (15), for the Giry monad \mathcal{G} , where $c: P \rightarrow \mathcal{G}(X)$ and $d: X \rightarrow \mathcal{G}(O)$ are given by pdf's $u: P \times X \rightarrow \mathbb{R}_{\geq 0}$ and $v: X \times O \rightarrow \mathbb{R}_{\geq 0}$, as pdf-channels $c = \int u$ and $d = \int v$. Let c be conjugate prior to d , via a measurable function $h: P \times O \rightarrow P$.

Equation (16) then amounts to, for an element $p \in P$ and for measurable subsets $M \subseteq X$ and $N \subseteq O$,

$$\begin{aligned} & \int_N \int_M u(p, x) \cdot v(x, y) \, dx \, dy \\ &= \int_N \left(\int u(p, x) \cdot v(x, y) \, dx \right) \cdot \left(\int_M u(h(p, y), x) \, dx \right) \, dy. \end{aligned} \tag{17}$$

In order to prove this equation, it suffices to prove that the two functions under the outer integral \int_N are equal, that is, it suffices to prove for each $y \in O$,

$$\int_M u(p, x) \cdot v(x, y) \, dx = \left(\int u(p, x) \cdot v(x, y) \, dx \right) \cdot \left(\int_M u(h(p, y), x) \, dx \right). \tag{18}$$

This formulation will be used in the examples below.

Proof. We extensively use the equations for integration from Section 3 and from the Appendix, in order to prove (17). The left-hand side of Equation (16) gives the left-hand side of (17):

$$((\text{id}, d) \circ c)(p)(M \times N) \stackrel{(8)}{=} ((\text{id}, d) \gg c(p))(M \times N) \stackrel{(A.7)}{=} \int_N \int_M u(p, x) \cdot v(x, y) \, dx \, dy.$$

Unravelling the right-hand side of (16) is a bit more work:

$$\begin{aligned} & ((c \circ h) \otimes \text{id}) \circ (\text{id}, \Upsilon \circ d \circ c)(p)(M \times N) \\ & \stackrel{(8)}{=} \int (c \circ h) \otimes \text{id}(-)(M \times N) \, d(\text{id}, \Upsilon \circ d \circ c)(p) \\ & \stackrel{(A.6)}{=} \int ((c \circ h) \otimes \text{id})(-)(M \times N) \, d(\eta(p) \otimes (\Upsilon \circ d \circ c)(p)) \\ & \stackrel{(A.1)}{=} \int \int ((c \circ h) \otimes \text{id})(-, -)(M \times N) \, d\eta(p) \, d(\Upsilon \circ d \circ c)(p) \\ &= \int ((c \circ h) \otimes \text{id})(p, -)(M \times N) \, d\mathcal{G}(\Upsilon)(d \gg c(p)) \\ & \stackrel{(5)}{=} \int ((c \circ h) \otimes \text{id})(p, \Upsilon(-))(M \times N) \, d(d \gg c(p)) \\ & \stackrel{(A.2,A.4)}{=} \int \left(\int u(p, x) \cdot v(x, y) \, dx \right) \cdot ((c \circ h) \otimes \text{id})(p, y, y)(M \times N) \, dy \\ &= \int \left(\int u(p, x) \cdot v(x, y) \, dx \right) \cdot c(h(p, y))(M) \cdot \mathbf{1}_N(y) \, dy \\ &= \int_N \left(\int u(p, x) \cdot v(x, y) \, dx \right) \cdot \left(\int_M u(h(p, y), x) \, dx \right) \, dy. \end{aligned}$$

By combining this outcome with the earlier one, we get the desired equation (17). □

One can reorganise Equation (18) as a normalisation fraction:

$$\int_M u(h(p, y), x) \, dx = \frac{\int_M u(p, x) \cdot v(x, y) \, dx}{\int u(p, x) \cdot v(x, y) \, dx}. \tag{19}$$

It now strongly resembles Equation (12) for Bayesian inversion. This connection will be established more generally in Theorem 6.3. Essentially, the above normalisation fraction (19) occurs in Bernardo and Smith (2000, Definition 5.6).

We are now ready to review some standard examples. The first one describes the structure underlying the coin example in Section 2.

Example 5.3. It is well known that the beta distributions are conjugate prior to the Bernoulli ‘flip’ likelihood function. We shall reformulate this fact following the pattern of Definition 5.1, with two composable channels, as in (15), namely:

$$\mathbb{N}_{>0} \times \mathbb{N}_{>0} \xrightarrow{\text{Beta}} [0, 1] \xrightarrow{\text{Flip}} 2 \quad \text{where } 2 = \{0, 1\}.$$

The Beta channel is as in (1), but now restricted to the non-negative natural numbers $\mathbb{N}_{>0}$. We recall that the normalisation constant $B(\alpha, \beta)$ is $\int_{[0,1]} x^{\alpha-1}(1-x)^{\beta-1} \, dx$.

The Flip channel sends a probability $r \in [0, 1]$ to the Bernoulli(r) distribution, which can also be written as a discrete distribution $\text{Flip}(r) = r|1\rangle + (1-r)|0\rangle$. More formally, as a Kleisli map $[0, 1] \rightarrow \mathcal{G}(2)$ it is, for a subset $N \subseteq 2$,

$$\text{Flip}(r)(N) = \int_N r^i \cdot (1-r)^{1-i} \, di = \sum_{i \in N} r^i \cdot (1-r)^{1-i} = \begin{cases} 0 & \text{if } N = \emptyset \\ r & \text{if } N = \{1\} \\ 1-r & \text{if } N = \{0\} \\ 1 & \text{if } N = \{0, 1\}. \end{cases}$$

The i in di refers here to the counting measure.

In order to show that Beta is a conjugate prior of Flip, we have to produce a parameter translation function $h: \mathbb{N}_{>0} \times \mathbb{N}_{>0} \times 2 \rightarrow \mathbb{N}_{>0} \times \mathbb{N}_{>0}$. It is defined by distinguishing the elements in $2 = \{0, 1\}$

$$h(\alpha, \beta, 1) = (\alpha + 1, \beta) \quad \text{and} \quad h(\alpha, \beta, 0) = (\alpha, \beta + 1). \tag{20}$$

Thus, in one formula, $h(\alpha, \beta, i) = (\alpha + i, \beta + (1 - i))$.

We prove Equation (18) for $c = \text{Beta} = \int u$ and $d = \text{Flip} = \int v$. We start from its right-hand side, for an arbitrary $i \in 2$,

$$\begin{aligned} & \left(\int u(\alpha, \beta, x) \cdot v(x, i) \, dx \right) \cdot \left(\int_M u(h(\alpha, \beta, i), x) \, dx \right) \\ &= \left(\int \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)} \cdot x^i \cdot (1-x)^{1-i} \, dx \right) \cdot \left(\int_M \frac{x^{\alpha+i-1}(1-x)^{\beta+(1-i)-1}}{B(\alpha+i, \beta+(1-i))} \, dx \right) \\ &= \left(\frac{\int x^{\alpha+i-1}(1-x)^{\beta+(1-i)-1} \, dx}{B(\alpha, \beta)} \right) \cdot \left(\int_M \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha+i, \beta+(1-i))} \cdot x^i \cdot (1-x)^{1-i} \, dx \right) \\ &= \left(\frac{B(\alpha+i, \beta+(1-i))}{B(\alpha, \beta)} \right) \cdot \left(\int_M \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha+i, \beta+(1-i))} \cdot x^i \cdot (1-x)^{1-i} \, dx \right) \\ &= \int_M \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)} \cdot x^i \cdot (1-x)^{1-i} \, dx \\ &= \int_M u(\alpha, \beta, x) \cdot v(x, i) \, dx. \end{aligned}$$

The latter expression is the left-hand side of (18). We see that the essence of the verification of the conjugate prior equation is the shifting of functions and normalisation factors. This is a general pattern.

In the remainder of this section, we review how some of the standard examples fit in the current setting, without each time doing the entire ‘big’ calculation to verify Equation (18).

Example 5.4. In a similar way, one verifies that the Beta channel is a conjugate prior to the binomial channel. For the latter we fix a natural number $n > 0$, and consider the two channels:

$$\mathbb{N}_{>0} \times \mathbb{N}_{>0} \xrightarrow{\text{Beta}} [0, 1] \xrightarrow{\text{Binom}_n} \{0, 1, \dots, n\}$$

The binomial channel Binom_n is defined for $r \in [0, 1]$ and $M \subseteq \{0, 1, \dots, n\}$ as:

$$\text{Binom}_n(r)(M) = \int_M \binom{n}{i} \cdot r^i \cdot (1-r)^{n-i} \, di = \sum_{i \in M} \binom{n}{i} \cdot r^i \cdot (1-r)^{n-i}.$$

The conjugate prior property requires in this situation a parameter translation function $h: \mathbb{N}_{>0} \times \mathbb{N}_{>0} \times \{0, 1, \dots, n\} \rightarrow \mathbb{N}_{>0} \times \mathbb{N}_{>0}$, given by: $h(\alpha, \beta, i) = (\alpha + i, \beta + n - i)$.

Here is another well-known conjugate prior relationship, namely between Dirichlet and ‘multinomial’ distributions. The latter are simply called discrete distributions in the present context.

Example 5.5. Here we shall identify a number $n \in \mathbb{N}$ with the n -element set $\{0, 1, \dots, n - 1\}$. We then write $\mathcal{D}_*(n)$ for the set of n -tuples $(x_0, \dots, x_{n-1}) \in (\mathbb{R}_{>0})^n$ with $\sum_i x_i = 1$.

For a fixed $n > 0$, let $O = \{y_0, \dots, y_{n-1}\}$ be a set of ‘observations’. We consider the following two channels.

$$(\mathbb{N}_{>0})^n \xrightarrow{\text{Dir}_n} \mathcal{D}_*(n) \xrightarrow{\text{Mult}} O$$

The multinomial channel is defined as $\text{Mult}(x_0, \dots, x_{n-1}) = x_0 | y_0 + \dots + x_{n-1} | y_{n-1}$. The Dirichlet channel Dir_n is more complicated: for an n -tuple $\vec{\alpha} = (\alpha_0, \dots, \alpha_{n-1})$ it is given via pdf’s d_n , in:

$$\text{Dir}_n(\vec{\alpha}) = \int d_n(\vec{\alpha}) \quad \text{where} \quad d_n(\vec{\alpha})(x_0, \dots, x_{n-1}) = \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \cdot \prod_i x_i^{\alpha_i - 1},$$

for $(x_0, \dots, x_{n-1}) \in \mathcal{D}_*(n)$. The operation Γ is the ‘Gamma’ function, which is defined on natural numbers $k > 1$ as $\Gamma(k) = (k - 1)!$.

The parameter translation function $h: (\mathbb{N}_{>0})^n \times O \rightarrow (\mathbb{N}_{>0})^n$ is:

$$h(\alpha_0, \dots, \alpha_{n-1}, y) = (\alpha_0, \dots, \alpha_i + 1, \dots, \alpha_{n-1}) \quad \text{if } y = y_i.$$

Example 5.6. The γ distribution is conjugate prior to the Poisson distribution, in:

$$\mathbb{R}_{>0} \times \mathbb{R}_{>0} \xrightarrow{\text{Gam}} \mathbb{R}_{>0} \xrightarrow{\text{Pois}} \mathbb{N}$$

where:

$$\text{Gam}(\alpha, \beta)(M) = \int_M \frac{\beta^\alpha \cdot x^{\alpha-1} \cdot e^{-\beta x}}{\Gamma(\alpha)} \, dx \quad \text{and} \quad \text{Pois}(x) = \sum_{k=0}^{k=\infty} e^{-x} \cdot \frac{x^k}{k!} |k\rangle.$$

The parameter translation function $h: \mathbb{R}_{>0} \times \mathbb{R}_{>0} \times \mathbb{N} \rightarrow \mathbb{R}_{>0} \times \mathbb{R}_{>0}$ is: $h(\alpha, \beta, k) = (\alpha + k, \beta + 1)$.

We include one more example, illustrating that normal channels are conjugate priors to themselves. This fact is also well known.

Example 5.7. Consider the following two normal channels.

$$\mathbb{R} \times \mathbb{R}_{>0} \xrightarrow{\text{Norm}} \mathbb{R} \xrightarrow{\text{Norm}(-, \nu)} \mathbb{R}_{>0}$$

The channel Norm is described in (13); it is used twice here, the second time with a fixed standard deviation ν , for ‘noise’. This second channel is typically used for observation, like in Kalman filtering, for which a fixed noise level can be assumed. In this situation, the first normal channel Norm is a conjugate prior to the second channel Norm($-$, ν), via the parameter translation function $h: \mathbb{R} \times \mathbb{R}_{>0} \times \mathbb{R}_{>0} \rightarrow \mathbb{R} \times \mathbb{R}_{>0}$ given by:

$$h(\mu, \sigma, y) = \left(\frac{\mu \cdot \nu^2 + y \cdot \sigma^2}{\nu^2 + \sigma^2}, \frac{\nu \cdot \sigma}{\sqrt{\nu^2 + \sigma^2}} \right)$$

6. Conjugate priors form bayesian inversions

This section connects the main two notions of this paper, by showing that conjugate priors give rise to Bayesian inversion. The argument is a very simple example of diagrammatic reasoning. Before we come to it, we have to clarify an issue that was left open earlier, regarding ‘deterministic’ channels, see Definition 5.1.

Definition 6.1. A channel c is called deterministic if it commutes with copiers, that is, if it satisfies the equation on the left below.



As a special case, a state ω is called deterministic if it satisfies the equation on the right, above.

The state description is a special case of the channel description since a state on X is a channel $1 \rightarrow X$ and copying on the trivial (final) object 1 does nothing, up to isomorphism.

Few channels (or states) are deterministic. In deterministic and continuous computation, the ordinary functions $f: X \rightarrow Y$ are deterministic, when considered as a channel $\eta \circ f$. We check this explicitly for point states, since this is what we need later on.

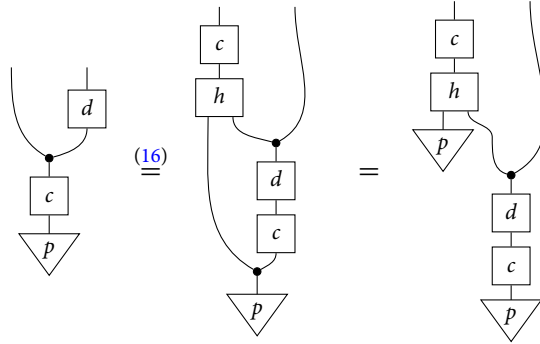
Example 6.2. Let x be an element of a measurable space X . The associated point state $\eta(x) \in \mathcal{G}(X)$ is deterministic, where $\eta(x)(M) = \mathbf{1}_M(x)$. We check the equation on the right in Definition 6.1:

$$\begin{aligned} (\Upsilon \circ \eta(x))(M \times N) &= \eta(x, x)(M \times N) = \mathbf{1}_{M \times N}(x, x) = \mathbf{1}_M(x) \cdot \mathbf{1}_N(x) \\ &= \eta(x)(M) \cdot \eta(x)(N) = (\eta(x) \otimes \eta(x))(M \times N). \end{aligned}$$

We now come to the main result.

Theorem 6.3. Let $P \xrightarrow{c} X \xrightarrow{d} O$ be channels, where c is conjugate prior to d , say via $h: P \times O \rightarrow P$. Then for each deterministic (copyable) state p , the map $c \circ h(p, -): O \rightarrow X$ is a Bayesian inversion of d , wrt. the transformed state $c \gg p$.

Proof. We have to prove Equation (10), for channel d and state $c \gg p$, with the channel $c \circ h(p, -)$ playing the role of Bayesian inversion $d_{c \gg p}^\dagger$. This is easiest to see graphically, using that the state p is deterministic and thus commutes with copiers Υ , see the equation on the right in Definition 6.1.



This is it. □

When we specialise to Giry-channels, we get an ‘if-and-only-if’ statement, since there we can reason elementwise.

Corollary 6.4. *Let $P \xrightarrow{c} X \xrightarrow{d} O$ be two channels in $\mathcal{K}(\mathcal{G})$, and let $h: P \times O \rightarrow P$ be a measurable function. The following two points are equivalent:*

- (i) *c is a conjugate prior to d, via h;*
- (ii) *$c(h(p, -)): O \rightarrow \mathcal{G}(X)$ is a Bayesian inverse for channel d with state $c(p)$, i.e. is $d_{c(p)}^\dagger$, for each parameter $p \in P$.* □

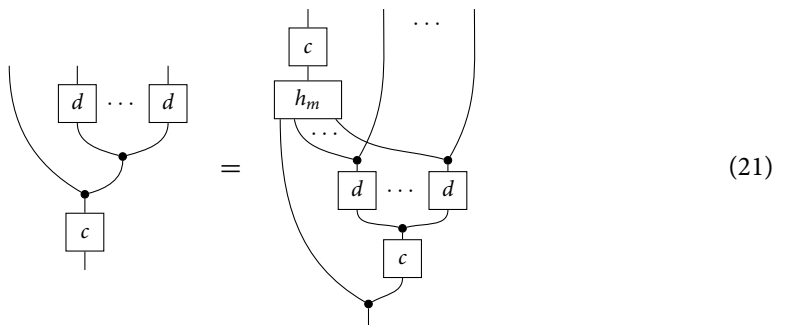
7. Multiple updates

So far we have dealt with the situation where there is a single observation $y \in O$ that leads to an update of a prior distribution. In this final section we briefly look at how to handle multiple observations $y_1, \dots, y_m \in O$. It will lead to the notion of *sufficient statistic*.

Multiple updates fit in our set-up $P \xrightarrow{c} X \xrightarrow{d} O$ by using a tuple-channel as output:

$$P \xrightarrow{c} X \xrightarrow{\langle d, \dots, d \rangle} O^m \quad \text{where} \quad \langle d, \dots, d \rangle = \begin{array}{c} \boxed{d} \quad \dots \quad \boxed{d} \\ \vdots \\ \bullet \end{array}$$

Following Definition 5.1 the channel c is conjugate prior to the m-tuple $\langle d, \dots, d \rangle$ if there is a (deterministic) parameter translation function $h_m: P \times O^m \rightarrow P$ with:



This construction involves in essence a Bayesian inversion, as in Section 6.

We shall illustrate this for the Beta–Flip relationship from Example 5.3. Suppose we have multiple head/tail observations $y_1, \dots, y_m \in 2 = \{0, 1\}$ which we wish to incorporate into a prior

distribution $\text{Beta}(\alpha, \beta)$. We use $\sum_i y_i$ as the number of 1's among the observations y_i and $\sum_i (1 - y_i)$ as the number of 0's, and then define as (multiple) parameter update function:

$$\mathbb{N}_{>0} \times \mathbb{N}_{>0} \times 2^m \xrightarrow{h_m} \mathbb{N}_{>0} \times \mathbb{N}_{>0} \quad \text{via} \quad h_m(\alpha, \beta, \vec{y}) = (\alpha + \sum_i y_i, \beta + \sum_i (1 - y_i)) \quad (22)$$

Equation (21) then holds: using a reformulation in the style of (18), and with pdf's u, v , as in Example 5.3, it suffices to prove:

$$\begin{aligned} & \left(\int u(\alpha, \beta, x) \cdot v(x, y_1) \cdot \dots \cdot v(x, y_m) \, dx \right) \cdot \left(\int_M u(h_m(\alpha, \beta, \vec{y}), x) \, dx \right) \\ &= \int_M u(\alpha, \beta, x) \cdot v(x, y_1) \cdot \dots \cdot v(x, y_m) \, dx. \end{aligned}$$

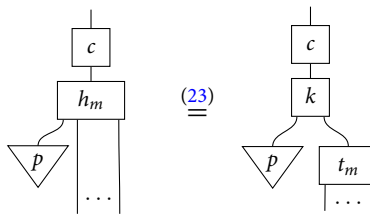
We see that for this multiple updating, we do not need the individual observations, $y_1, \dots, y_m \in 2$, but only the sums $\sum_i y_i$ and $\sum_i (1 - y_i) = m - \sum_i y_i$. In many such cases the situation can be simplified via factorisation, using the notion of sufficient statistic, see e.g. Koopman (1936) and Bishop (2006). We adapt it to the current setting.

Definition 7.1. A sufficient statistic for a collection of parameter update functions $(h_m : P \times O^m \rightarrow P)_{m \in \mathbb{N}}$ is a collection of functions $(t_m : O^m \rightarrow Z)_{m \in \mathbb{N}}$ with $k : P \times Z \rightarrow P$ such that the following triangle commutes for each $m \in \mathbb{N}$.

$$\begin{array}{ccc} P \times O^m & \xrightarrow{h_m} & P \\ \text{id} \times t_m \downarrow & \nearrow k & \\ P \times Z & & \end{array} \quad (23)$$

Crucially, the set Z does not depend on m and t_m does not act on the parameters P . Moreover, in examples, Z is often simpler than these O^m .

Let's return to the situation in (21) where a channel c is conjugate prior to an m -tuple $\langle d, \dots, d \rangle : X \rightarrow O^m$. In line with (the proof of) Theorem 6.3, for an element $p \in P$ we obtain a Bayesian inversion $O^m \rightarrow X$ of the tuple $\langle d, \dots, d \rangle$ wrt. state $c \ggg p$ on the left:



In the presence of a sufficient statistic, this Bayesian inversion takes the form as on the right.

The factorisation in the above triangle (23) may satisfy an obvious universal property once we add the requirement that the functions t_m satisfy a certain invariance property. The latter property does not seem to be uniform across examples, as will be illustrated below.

Example 7.2. We mention sufficient statistics for some of the examples from Section 5.

- (1) As discussed in the beginning of this section, for Beta-Flip in Example 5.3 we have $P = \mathbb{N}_{>0} \times \mathbb{N}_{>0}$ and $O = 2 = \{0, 1\}$. The sufficient statistic has $Z = \mathbb{N} \times \mathbb{N}$ with:

$$\begin{array}{ccc} 2^m & \xrightarrow{t_m} & \mathbb{N} \times \mathbb{N} & \quad & \mathbb{N}_{>0} \times \mathbb{N}_{>0} \times \mathbb{N} \times \mathbb{N} & \xrightarrow{k} & \mathbb{N}_{>0} \times \mathbb{N}_{>0} \\ \vec{y} \mapsto & & (\sum_i y_i, \sum_i (1 - y_i)) & & (\alpha, \beta, n_1, n_0) \mapsto & & (\alpha + n_1, \beta + n_0). \end{array}$$

Clearly, $k \circ (\text{id} \times t_m) = h_m$, for the functions h_m in (22).

We notice that the functions t_m are ‘invariant under permutation’ in the sense that:

$$t_m(y_1, \dots, y_i, \dots, y_j, \dots, y_m) = t_m(y_1, \dots, y_j, \dots, y_i, \dots, y_m),$$

for each $1 \leq i, j \leq m$. This invariance allows us to state that the above (t_m) and k satisfy the following universal property: for each collection of functions $s_m: O^m \rightarrow W$ invariant under permutation, together with $\ell: P \times W \rightarrow P$ such that $\ell \circ (\text{id} \times s_m) = h_m$, there is a unique $f: Z \rightarrow W$ with $f \circ t_m = s_m$, for each $m \in \mathbb{N}$, and $\ell \circ (\text{id} \times f) = k$.

Indeed, we can define $f: \mathbb{N} \times \mathbb{N} \rightarrow W$ as:

$$f(n_1, n_0) = s_{n_1+n_0}(\underbrace{1, \dots, 1}_{n_1 \text{ times}}, \underbrace{0, \dots, 0}_{n_0 \text{ times}}).$$

Since these s functions are invariant under permutation, the precise order of the 1’s and 0’s does not matter. Then:

$$\begin{aligned} (f \circ t_m)(y_1, \dots, y_m) &= f(\sum_i y_i, \sum_i 1 - y_i) \\ &= s_m(1, \dots, 1, 0, \dots, 0) = s_m(y_1, \dots, y_m) \\ (\ell \circ (\text{id} \times f))(\alpha, \beta, n_1, n_0) &= \ell(\alpha, \beta, s_{n_1+n_0}(1, \dots, 1, 0, \dots, 0)) \\ &= h_{n_1+n_0}(\alpha, \beta, 1, \dots, 1, 0, \dots, 0) \\ &\stackrel{(22)}{=} (\alpha + n_1, \beta + n_0) = k(\alpha, \beta, n_1, n_0). \end{aligned}$$

We further have to check uniqueness: if $g: \mathbb{N} \times \mathbb{N} \rightarrow W$ also satisfies $g \circ t_m = s_m$ and $\ell \circ (\text{id} \times g) = k$, then $f = g$ since:

$$f(n_1, n_0) = s_{n_1+n_0}(1, \dots, 1, 0, \dots, 0) = g(t_m(1, \dots, 1, 0, \dots, 0)) = g(n_1, n_0).$$

- (2) For the Gam–Pois conjugate priorship in Example 5.6, we have $P = \mathbb{R}_{>0} \times \mathbb{R}_{>0}$ and $O = \mathbb{N}$ as spaces of parameters and observations, and $h_m: P \times O^m \rightarrow P$ is defined by $h_m(\alpha, \beta, \vec{y}) = (\alpha + \sum_i y_i, \beta + m)$. The sufficient statistic has the same $Z = \mathbb{N} \times \mathbb{N}$ and $k(\alpha, \beta, z, n) = (\alpha + z, \beta + n)$, much as in the previous point, but with:

$$\mathbb{N}^m \xrightarrow{t_m} \mathbb{N} \times \mathbb{N} \quad \text{given by} \quad t_m(\vec{y}) = (\sum_i y_i, m).$$

The invariance property in this case for these t_m is:

$$t_m(y_1, \dots, y_i + z, \dots, y_j, \dots, y_m) = t_m(y_1, \dots, y_i, \dots, y_j + z, \dots, y_m),$$

With this invariance, we can prove a universal property of the factorisation (t_m) and k .

- (3) For the Norm–Norm conjugate priorship from Example 5.7, we have $P = \mathbb{R} \times \mathbb{R}_{>0}$ and $O = \mathbb{R}_{>0}$. We can take $Z = \mathbb{R}_{>0} \times \mathbb{N}$ with:

$$\begin{aligned} (\mathbb{R}_{>0})^m &\xrightarrow{t_m} \mathbb{R}_{>0} \times \mathbb{N} & \mathbb{R} \times \mathbb{R}_{>0} \times \mathbb{R}_{>0} \times \mathbb{N} &\xrightarrow{k} \mathbb{R} \times \mathbb{R}_{>0} \\ \vec{y} &\longmapsto (\sum_i y_i, m) & (\mu, \sigma, s, n) &\longmapsto \left(\frac{\mu v^2 + s \sigma^2}{v^2 + n \sigma^2}, \frac{v \sigma}{\sqrt{v^2 + n \sigma^2}}\right) \end{aligned}$$

The proof that this works is laborious but essentially straightforward. Universality is obtained if we use the same invariance property as in the previous point.

8. Conclusions

This paper contains a novel view on conjugate priors, using the concept of channel in a systematic manner. It introduces a precise definition for conjugate priorship, using a pair of composable channels $P \rightarrow X \rightarrow O$ and a parameter translation function $P \times O \rightarrow P$, satisfying a non-trivial

equation, see Definition 5.1. It has been shown that this equation holds for several standard conjugate prior examples. There are many more examples that have not been checked here. One can be confident that the same equation holds for those unchecked examples too, since it has been shown here that conjugate priors amount to Bayesian inversions. This inversion property is the essential characteristic for conjugate priors. It has been extended to multiple updates, using the concept of sufficient statistic.

Acknowledgements. Thanks are due to the two anonymous reviewers, one of whose comments significantly improved Section 7.

Notes

1 See https://en.wikipedia.org/wiki/Conjugate_prior or online lists, such as <https://www.johndcook.com/CompendiumOfConjugatePriors.pdf>, consulted on Sept. 10, 2018

2 Sometimes these distributions $\sum_i r_i |x_i\rangle$ are called ‘multinomial’ or ‘categorical’; the latter terminology is confusing in the present context.

3 In Jacobs (2013), integration $\int f d\omega$ is defined only for $[0, 1]$ -valued functions f , but that does not matter for the relevant equations, except that integrals may not exist for $\mathbb{R}_{\geq 0}$ -valued functions (or have value ∞). These integrals are determined by their valued $\int \mathbf{1}_M d\omega = \omega(M)$ on indicator functions $\mathbf{1}_M$ for measurable subsets, via continuous and linear extensions, see also Jacobs and Westerbaan (2015).

References

- Abramsky, S. and Coecke, B. (2009). A categorical semantics of quantum protocols. In: Engesser, K., Gabbay, D. M. and Lehmann, D. (eds.) *Handbook of Quantum Logic and Quantum Structures: Quantum Logic*, North-Holland, Elsevier, Computer Science Press, 261–323.
- Ackerman, N., Freer, C. and Roy, D. (2011). Noncomputable conditional distributions. In: Dawar, A. and Grädel, E. (eds.) *Logic in Computer Science*, IEEE, Computer Science Press.
- Alpaydin, E. (2010). *Introduction to Machine Learning*, 2nd edn., Cambridge, MA, MIT Press.
- Bernardo, J. and Smith, A. (2000). *Bayesian Theory*. Chichester, John Wiley & Sons.
- Bishop, C. (2006). *Pattern Recognition and Machine Learning*. Information Science and Statistics. Chichester, Springer.
- Blackwell, D. (1951). Comparison of experiments. In: *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability*, Springer/British Computer Society, 93–102.
- Cho, K. and Jacobs, B. (2019). Disintegration and Bayesian inversion, both abstractly and concretely. In: *Mathematical Structures in Computer Science*. See also arxiv.org/abs/1709.00322.
- Cho, K., Jacobs, B., Westerbaan, A. and Westerbaan, B. (2015). An introduction to effectus theory. See arxiv.org/abs/1512.05813.
- Clerc, F., Dahlqvist, F., Danos, V. and Garnier, I. (2017). Pointless learning. In: Esparza, J. and Murawski, A. (eds.) *Foundations of Software Science and Computation Structures*, Lecture Notes in Computer Science, vol. 10203, Berlin, Springer, 355–369.
- Culbertson, J. and Sturtz, K. (2014). A categorical foundation for Bayesian probability. *Applied Categorical Structures* 22(4) 647–662.
- Diaconis, P. and Ylvisaker, D. (1979). Conjugate priors for exponential families. *Annals of Statistics* 7(2) 269–281.
- Faden, A. (1985). The existence of regular conditional probabilities: Necessary and sufficient conditions. *The Annals of Probability* 13(1) 288–298.
- Giry, M. (1982). A categorical approach to probability theory. In: Banaschewski, B. (ed.) *Categorical Aspects of Topology and Analysis*, Lecture Notes in Computer Science, vol. 915, Berlin, Springer, 68–85.
- Jacobs, B. (2013). Measurable spaces and their effect logic. In: *Logic in Computer Science*, IEEE, Computer Science Press.
- Jacobs, B. (2015). New directions in categorical logic, for classical, probabilistic and quantum logic. *Logical Methods in Computer Science* 11(3). See <https://lmcs.episciences.org/1600>.
- Jacobs, B. (2018). From probability monads to commutative effectuses. *Journal of Logical and Algebraic Methods in Programming* 94 200–237.
- Jacobs, B. and Westerbaan, A. (2015). An effect-theoretic account of Lebesgue integration. In: Ghica, D. (ed.) *Mathematical Foundations of Programming Semantics*, Electronic Notes in Theoretical Computer Science, vol. 319, Amsterdam, Elsevier, 239–253.
- Jacobs, B. and Zanasi, F. (2019). The logical essentials of Bayesian reasoning. In: Barthe, G. and Katoen, J.-P. and Silva, A. *Probabilistic Programming*, Cambridge, Cambridge University Press. See arxiv.org/abs/1804.01193.

Koopman, B. (1936). On distributions admitting a sufficient statistic. *Transactions of the American Mathematical Society* **39** 399–409.

Panangaden, P. (2009). *Labelled Markov Processes*. London, Imperial College Press.

Russell, S. and Norvig, P. (2003). *Artificial Intelligence. A Modern Approach*. Englewood Cliffs, NJ, Prentice Hall.

Selinger, P. (2007). Dagger compact closed categories and completely positive maps (extended abstract). In: Selinger, P. (ed.) *Proceedings of the 3rd International Workshop on Quantum Programming Languages (QPL 2005)*, Electronic Notes in Theoretical Computer Science, vol. 170, Amsterdam, Elsevier, 139–163. doi: <http://dx.doi.org/10.1016/j.entcs.2006.12.018>.

Selinger, P. (2011). A survey of graphical languages for monoidal categories. In: Coecke, B. (ed.) *New Structures in Physics*, Lecture Notes in Physics, vol. 813, Berlin, Springer, 289–355.

Stoyanov, J. (2014). *Counterexamples in Probability*, 2nd rev. edn. Chichester, Wiley.

Appendix A. Calculation laws for Giry–Kleisli Maps with pdf’s

We assume that for a probability distribution (state) $\omega \in \mathcal{G}(X)$ and a measurable function $f: X \rightarrow \mathbb{R}_{\geq 0}$, the integral $\int f d\omega \in [0, \infty]$ can be defined as a limit of integrals over simple functions that approximate f . We shall follow the description of Jacobs (2013), to which we refer for details.³ This integration satisfies the Fubini property, which can be formulated, for states $\omega \in \mathcal{G}(X)$, $\rho \in \mathcal{G}(Y)$ and measurable function $h: X \times Y \rightarrow \mathbb{R}_{\geq 0}$, as:

$$\int h d(\omega \otimes \rho) = \int \int h d\omega d\rho. \tag{A.1}$$

The product state $\omega \otimes \rho \in \mathcal{G}(X \times Y)$ is defined by $(\omega \otimes \rho)(M \times N) = \omega(M) \cdot \rho(N)$.

A.1 States via pdf’s

For a subset $X \subseteq \mathbb{R}$, a measurable function $f: X \rightarrow \mathbb{R}_{\geq 0}$ is called a probability density function (pdf) for a state $\omega \in \mathcal{G}(X)$ if $\omega(M) = \int_M f(x) dx$ for each measurable subset $M \subseteq X$. In that case we simply write $\omega = \int f(x) dx$, or even $\omega = \int f$. If ω is given by such a pdf f , integration with state ω can be described as:

$$\int g d\omega = \int f(x) \cdot g(x) dx. \tag{A.2}$$

A.2 Channels via pdf’s

Let channel $c: X \rightarrow \mathcal{G}(Y)$ be given as $c = \int u$ by pdf $u: X \times Y \rightarrow \mathbb{R}_{\geq 0}$ as $c(x)(N) = \int_N u(x, y) dy$, for each $x \in X$ and measurable $N \subseteq Y$, like in (9). If $\omega = \int f$ is a state on X , then state transformation $c \gg \omega \in \mathcal{G}(Y)$ is given by:

$$\begin{aligned} (c \gg \omega)(N) &\stackrel{(7)}{=} \int c(-)(N) d\omega \stackrel{(A.2)}{=} \int f(x) \cdot c(x)(N) dx \\ &= \int_N \int f(x) \cdot u(x, y) dx dy. \end{aligned} \tag{A.3}$$

Hence the pdf of the transformed state $c \gg \omega$ is $y \mapsto \int f(x) \cdot u(x, y) dx$.

Given a channel $d: Y \rightarrow \mathcal{G}(Z)$, say with $d = \int v$, then sequential channel composition $d \circ c$ is given, for $x \in X$ and $K \subseteq Z$, by:

$$\begin{aligned} (d \circ c)(x)(K) &\stackrel{(8)}{=} \int d(-)(K) dc(x) \stackrel{(A.2)}{=} \int u(x, y) \cdot d(y)(K) dy \\ &= \int_K \int u(x, y) \cdot v(y, z) dy dz \end{aligned} \tag{A.4}$$

We see that the pdf of the channel $d \circ c$ is $(x, z) \mapsto \int u(x, y) \cdot v(y, z) dy$.

For a channel $e = \int w: A \rightarrow \mathcal{G}(B)$, we get a parallel composition channel $c \otimes e: X \times A \rightarrow \mathcal{G}(Y \times B)$ given by:

$$\begin{aligned} (c \otimes e)(x, a)(M \times N) &= c(x)(M) \otimes e(a)(N) \\ &= \left(\int_M u(x, y) \, dy \right) \cdot \left(\int_N w(a, b) \, db \right) \\ &= \int_{M \times N} u(x, y) \cdot w(a, b) \, d(y, b). \end{aligned} \tag{A.5}$$

Hence the pdf of the channel $c \otimes d$ is $(x, a, y, b) \mapsto u(x, y) \cdot w(a, b)$.

A.3 Graph channels and pdf's

For a channel $c: X \rightarrow \mathcal{G}(Y)$, we can form ‘graph’ channels $\langle \text{id}, c \rangle = (\text{id} \otimes c) \circ \Upsilon: X \rightarrow \mathcal{G}(X \times Y)$ and $\langle c, \text{id} \rangle = (c \otimes \text{id}) \circ \Upsilon: X \rightarrow \mathcal{G}(Y \times X)$. For $x \in X$ we have:

$$\langle \text{id}, c \rangle(x) = \eta(x) \otimes c(x) \quad \text{and} \quad \langle c, \text{id} \rangle(x) = c(x) \otimes \eta(x). \tag{A.6}$$

If $c = \int u$ and $\omega = \int f$ is a state on X , then:

$$\begin{aligned} (\langle \text{id}, c \rangle \gg \omega)(M \times N) &\stackrel{(A.2)}{=} \int f(x) \cdot \langle \text{id}, c \rangle(x)(M \times N) \, dx \\ &\stackrel{(A.6)}{=} \int f(x) \cdot \eta(x)(M) \cdot c(x)(N) \, dx \\ &= \int_N \int_M f(x) \cdot u(x, y) \, dx \, dy. \end{aligned} \tag{A.7}$$

We also consider the situation where $d: X \times Y \rightarrow \mathcal{G}(Z)$ is of the form $d = \int v$, with composition:

$$\begin{aligned} (d \circ \langle \text{id}, c \rangle)(x)(K) &\stackrel{(A.6)}{=} \int d(-)(K) \, d(\eta(x) \otimes c(x)) \\ &\stackrel{(A.1)}{=} \int d(-)(K) \, d\eta(x) \, dc(x) \\ &= \int d(x, -)(K) \, dc(x) \\ &\stackrel{(A.2)}{=} \int u(x, y) \cdot d(x, y)(K) \, dy \\ &= \int_K \int u(x, y) \cdot v(x, y, z) \, dy \, dz. \end{aligned} \tag{A.8}$$

Hence the pdf of the channel $d \circ \langle \text{id}, c \rangle$ is $(x, z) \mapsto \int u(x, y) \cdot v(x, y, z) \, dy$.

Cite this article: Jacobs B (2020). A channel-based perspective on conjugate priors. *Mathematical Structures in Computer Science* 30, 44–61. <https://doi.org/10.1017/S0960129519000082>