
Reasoning with Conditional Probabilities and Joint Distributions in Coq

Reynald Affeldt Jacques Garrigue Takafumi Saikawa

Probabilities occur in many applications of computer science, such as communication theory and artificial intelligence. These are critical applications that require some form of verification to guarantee the quality of their implementations. Unfortunately, probabilities are also the typical example of a mathematical theory whose abuses of notations make pencil-and-paper proofs difficult to formalize. In this paper, we experiment a new formalization of conditional probabilities that we validate with two applications. First, we formalize the foundational definitions and theorems of information theory, extending previous work with new lemmas. Second, we formalize the notion of conditional independence and its properties, paving the road for a formalization of probabilistic graphical models.

1 Introduction

Probabilities occur in many applications of computer science. One finds them in information theory to reason about the size of compressed data, in quantitative information flow analysis, in the analysis of side-channel attacks, etc. They are also omnipresent in artificial intelligence. These are critical applications that require some form of verification to guarantee the quality of their implementations.

Formal verification using proof-assistants has emerged as a trustful technique to guarantee the correctness of important mathematical theorems (e.g., [14]), as well as the correctness of critical computer software (e.g., [19]). Regarding more

specifically mathematics, formal verification (hereafter, we also say *formalization*) consists in turning a proof from a textbook or a scientific paper (hereafter, we refer to such proofs as “pencil-and-paper proofs”) into the language of the proof-assistant. The latter is a software mechanization of some mathematical foundations, such as set theory, higher-order logic, or type theory (as it is the case for the COQ proof-assistant). The success of this transcription guarantees that the original proof was indeed correct but it requires a lot of *proof-engineering*: make explicit all assumptions (particularly hidden ones), build libraries of intermediate definitions and lemmas, find abstractions to reduce complexity, etc.

In particular, pencil-and-paper proofs about probabilities are not easily amenable to formal verification. A **first reason** is that probabilities deal with a wide range of mathematics (from combinatorics to real analysis), but, as of today, no proof-assistant has enough libraries to bring the practitioner the same flexibility as pencil-and-paper proofs. The “ideal” approach (in the sense of be-

Coq による条件付き確率と同時分布の形式化とその応用
Reynald Affeldt, 産業技術総合研究所, National Institute of Advanced Industrial Science and Technology.
Jacques Garrigue, 名古屋大学, Graduate School of Mathematics, Nagoya University.
Takafumi Saikawa, 名古屋大学, Graduate School of Mathematics, Nagoya University.
コンピュータソフトウェア, Vol.0, No.0 (0), pp.0-0.
[研究論文] 2019 年???月???日受付.

ing the most faithful to the original pencil-and-paper proof) seems to be to derive the formalization of probabilities from a formalization of measure theory and Lebesgue integration. Unfortunately, not all proof assistants have such formalizations available (at the time of this writing, this is work in progress for COQ). There is one in HOL [21] and one in Isabelle/HOL [16], but these proof-assistants are lacking the theories for combinatorics and manipulation of iterated operators provided by the MATHCOMP [20][1] library of the COQ proof-assistant. Actually, one could even say that what is a good formalization of real analysis is still a topic of research (see for example [2]). A **second reason** why formalization of probabilities is difficult is that probabilities are the typical example of a mathematical theory where abuses of notations are omnipresent. An obvious example is the omission of the underlying distribution when talking about probabilities or expected values, often ruled out as “obvious from context”. The notations for random variables are also often abused. For example, the expression $P(X = a)$ is often abbreviated as $P(X)$ when a is universally quantified or $P(a)$ when a is bound by a summation operator (the proofs of Sect. 4.3.2 and Sect. 4.3.5 provide concrete examples, the definition of entropy in Sect. 3.1.1 is another illustration). In the formalization, we have to fill in all implicit elements. Though this could be done explicitly by the user, we should program the proof-assistant to help automatically.

Regarding the first reason explained above, the formalization of probabilities can be a bit simplified, thanks to our focus on computer science applications. We explained above that the “ideal” approach, which consists in deriving a formal theory of probability from a full-fledged formalization of measure theory, seems out of reach for now. Yet, many applications of information theory target the manipulation of data in computers (like in artificial

intelligence) which are always finite. Such applications can be verified by formalizing the finite part of probability and information theory, so that, in our case, a full-fledged formalization of measure theory and Lebesgue integration would look like using a sledgehammer to crack a nut. Moreover, a direct formalization of discrete probabilities makes it unnecessary to verify numerical integration and easier to define distributions and entropy as computable functions. We also expect a finite setting to lend itself more easily to the development of verified implementations when it comes to the verification of advanced topics of information theory such as coding methods (a topic to which we already hinted at in [3]). There therefore seems to be value in exploring a mechanization of discrete probabilities that directly allows for formal reasoning; this is why we restrict ourselves to a finite setting.

In this paper, we focus on the proposal of formal definitions and proofs that mimic the ones of information theory or artificial intelligence textbooks. However, it is not obvious to recover the terseness of the pencil-and-paper presentation. The main issue is the one of implicit assumptions. We already mentioned the problem of abuses of notations (see the second reason explained above), in which implicit elements can be seen as a form of implicit assumptions. To manage abuses of notations, we can leverage the user-defined notations provided by COQ and its dependent-type system to infer automatically implicit elements (see the definition of conditional independence for example—Sect. 4.2). Implicit assumptions can also take the form of definitions or rules that are not clearly stated or recalled when used. One example is the calculation rules for $\log P(x)$ when $P(x) = 0$ as they appear in the definition of Shannon’s entropy. A judicious choice of formalization turns out to provide a simple justification (see Sect. 3.1.1). Another example is the convexity statements for mutual informa-

tion. The standard textbook on information theory is not clear about the definition of convexity when the domain is the set of distributions [Thm 2.7.4 of [11]]. Our development on convex spaces is a useful abstraction to clarify this statement and its proof (see Sect. 3.3). See also our previous work for the formalization of more implicit assumptions in information theory [7]. On the bright side, and since we restrict ourselves to a finite setting, we can reuse existing libraries and in particular the MATHCOMP library that already provides many definitions and lemmas needed (most notably summations) and also the previous version of the INFOTHEO library [9]. We can also count on the expressiveness of the SSREFLECT language to provide short proofs despite the clutter of explicit assumptions.

More precisely, in this paper, we focus on the problem of formalizing and using conditional probabilities and joint distributions. Concretely, we extend an existing formalization of information theory [9] with an “explicit” formalization of conditional probability. We say “explicit” because the previous work [9] we extend was already dealing with conditional probabilities but given in the form of stochastic matrices, as used to model channels in information theory (see Sect. 5 for a more precise explanation). To show that our extension is indeed useful, we validate it with two substantial applications:

- First, we formalize the foundations of information theory. This is not a new topic but we are able to improve on previous work by formalizing lemmas there were not formalized before. In fact, we are on the verge of completing the formalization of [Chapter 2 of [11]], which provides the basic definitions for information theory.
- Second, we formalize the notion of conditional independence. We validate this definition by

proving the so-called *graphoid axioms*. These axioms have been introduced to reason symbolically about conditional independence using graphs [23]. The idea can be explained as follows. Conditional independence relations among several random variables can be represented by an undirected graph. In such a graph, a set A of variables is conditionally independent of a set B given C if the vertices corresponding to C separate the vertices of A from the vertices of B . Instead of reasoning directly using the probabilistic graphical model, one can alternatively use the graphoid axioms. We are not aware of a previous attempt to formalize the graphoid axioms from the ground up.

Paper Outline

In Sect. 2, we explain how we formalize distributions, joint distributions, and conditional probabilities. In Sect. 3, we explain our formalization of information theory. In Sect. 4, we explain our formalization of conditional independence. We review related work in Sect. 5 and conclude in Sect. 6.

About Notations

This paper displays COQ code verbatim. We rely on the MATHCOMP library whose notations are explained when used for the first time. We enforced a uniform naming convention for variables. We use A, B, C, D , etc. for finite types. We use a, b , etc. for elements of resp. A, B , etc. We use $E : \{\text{set } A\}, F : \{\text{set } B\}, G : \{\text{set } C\}, H : \{\text{set } D\}$, etc. for events. We use P, Q , etc. to range over distributions. We use X, Y, Z, W , etc. for random variables with values in A, B, C, D , resp.

2 Formalization of Conditional Proba-

bility

2.1 Background: Formalization of Finite Distributions

We define a finite distribution as a function f from a finite set A to the set \mathbb{R}^+ of nonnegative real numbers, such that the sum of the values of f equals 1. The formal definition in COQ is as follows (it comes from previous work [8]):

```
(* Module FDist *)
Record t := mk {
  f :> A → R+ ;
  _ : \sum_(a in A) f a == 1}.

```

The type A has type `finType`, it is a type with a finite number of elements [13]. The type $A \rightarrow \mathbb{R}^+$ is for functions with domain A and real, nonnegative outputs. The notation `\sum_(a in A)` (where a is a binder) is for the sum over A . Hereafter, the notation `{fdist A}` causes the type inference to correctly identify A as a `finType` *even when* A is a composite `finType` (e.g., the product of two finite types $A1 * A2$).

Given a distribution $P : \{fdist A\}$, we can formalize an event as a finite set E over A (i.e., a Boolean predicate over A represented as a list, whose type is denoted `{set A}` [20]) and compute its probability by summing the individual probabilities of its elements:

```
Variables (A : finType) (P : {fdist A}).
Definition Pr (E : {set A}) :=
  \sum_(a in E) P a.

```

2.2 Joint Distributions

Using the definition of Sect. 2.1, a joint distribution can be represented by a distribution over a product type such as `{fdist A * B}` or `{fdist 'rV[A]_n}` (where `'rV[A]_n` represents the type of row vectors of size n [20]).

Given a joint distribution, we often need to consider marginal distributions. For that purpose, it is useful to consider *probability monad*.

2.2.1 The Map Distribution from the

Probability Monad

We compute most marginal distributions using the *map distribution*. Given a distribution P of type `{fdist A}` and a function g of type $A \rightarrow B$, the map distribution is the distribution of type `{fdist B}` with probability mass function $b \mapsto \sum_{\substack{a \in A \\ g a = b}} P a$. In practice, this distribution is defined using the probability monad with unit `FDist1.d` and bind operator `FDistBind.d` [5]:

```
(* Module FDistMap. *)
Variables (A B : finType)
  (p : {fdist A}) (g : A → B).

```

```
Definition d : {fdist B} :=
  FDistBind.d p (fun a => FDist1.d (g a)).

```

As hinted at by the comment, this definition occurs inside the module `FDistMap`, which we use as a namespace; hence, the distribution d hereafter appears as `FDistMap.d`.

2.2.2 Joint Distributions over a Product Type

Given a joint distribution over a product type, `Bivar.fst` (`Bivar` for “bivariate”) builds the left marginal. It is defined as follows:

```
(* Module Bivar *)
Variables (A B : finType) (P : {fdist A * B}).
Definition fst : {fdist A} := FDistMap.d fst P.

```

Regarding joint distributions over a product type, one might also want to build the right marginal `Bivar.snd` (which is of course defined similarly) and the distribution `Swap.d` that swaps the left and the right elements.

Table 1 summarizes the distributions that one might want to build from of a joint distribution over a product type.

Similarly to joint distributions over a product type, one might want to build distributions from joint distributions over a triple of distributions. For example, in Table 2, `TripC12.d` permutes the first and the second element, while `Proj13.d` build the marginal w.r.t. the second element. In Sect. 4, we also deal with distributions made of quadruples of distributions but rarely enough so that we do not

Table 1 Distributions using the product type

<i>Original distribution</i>	<i>Built distribution</i>	<i>and its type</i>
$P : \{\text{fdist } A * B\}$	<code>Swap.d P</code>	$\{\text{fdist } B * A\}$
	<code>Bivar.fst P</code>	$\{\text{fdist } A\}$
	<code>Bivar.snd P</code>	$\{\text{fdist } B\}$

Table 2 Distribution using triples

<i>Original distribution</i>	<i>Built distribution</i>	<i>and its type</i>
$P : \{\text{fdist } A * B * C\}$	<code>TripA.d P</code>	$\{\text{fdist } A * (B * C)\}$
	<code>TripC12.d P</code>	$\{\text{fdist } (B * A) * C\}$
	<code>TripAC.d P</code>	$\{\text{fdist } (A * C) * B\}$
	<code>TripC13.d P</code>	$\{\text{fdist } (C * B) * A\}$
	<code>Proj13.d P</code>	$\{\text{fdist } A * C\}$
	<code>Proj23.d P</code>	$\{\text{fdist } B * C\}$

introduce specific definitions for them.

In the case of multivariate distributions (i.e., distributions of type $\{\text{fdist } \text{rv}[A]_n\}$), one might want to build the marginals made of the head or the tail, like `Multivar.head_of` or `Multivar.tail_of` in Table 3. In the same table, `Multivar.belast_last` is a distribution over the pairs formed by, on one side, all the elements but the last, and on the other side, the last element. The function `Nth.d` builds a marginal from any index and `Take.d` builds a marginal with a prefix. `PairNth.d` builds a marginal from a distribution over a product type and any index. `PairTake.d` operates a more complicated rearrangement (intuitively, a combination of `Take.d`, `Nth.d`, and `Bivar.snd`). Finally, `PairNth.d` and `PairTake.d` will find their use when stating the chain rule for information in Sect. 3.2.3.

In the following, we will need to make explicit mention of the distributions above, whereas they are often implicit in pencil-and-paper proofs.

2.3 Conditional Probability

We use joint distributions to define conditional probabilities.

Given a joint distribution $P : \{\text{fdist } A * B\}$ and

two events $E : \{\text{set } A\}$ and $F : \{\text{set } B\}$, the probability of E knowing F (i.e., $\text{Pr}[E|F]$) is formalized as follows:

Variables $(A B : \text{finType}) (P : \{\text{fdist } A * B\})$.

Definition $\text{cPr } E F :=$

$$\text{Pr } P (\text{setX } E F) / \text{Pr } (Bivar.snd P) F.$$

The set $\text{setX } E F$ is the Cartesian product of the sets E and F . Using a pencil-and-paper prose that mentions explicitly the underlying distributions (what many textbooks do not do), the definition of $\text{cPr } P E F$ would read as $\text{Pr}_P[E|F] \stackrel{\text{def}}{=} \frac{\text{Pr}_P(E \times F)}{\text{Pr}_{P_2}(F)}$, where P_2 represents the right marginal of P . This is why we use the notation $\backslash\text{Pr}_P[E | F]$ for $\text{cPr } P E F$ in the COQ scripts.

We formalized a number of lemmas about conditional probabilities, such as Bayes' theorem, its general form, etc. We do not provide a complete presentation because their proofs follow known formal proofs [15]. Let us just provide two examples that are used to prove the graphoid axioms. The proofs of most of the graphoid axioms (see Sect. 4) use the product rule (i.e., $\text{Pr}[E \times F|G] = \text{Pr}[E|F \times G] \text{Pr}[F|G]$), which we state formally as follows:

Lemma `product_rule E F G :`

$$\backslash\text{Pr}_P [\text{setX } E F | G] = \backslash\text{Pr}_P (\text{TripA.d } P) [E | \text{setX } F G] *$$

Table 3 Distributions using vectors

Original distribution	Built distribution	and its type
$P : \{\text{fdist 'rV[A]_n.+1}\}$	<code>Multivar.head_of P</code>	$\{\text{fdist A}\}$
	<code>Multivar.tail_of P</code>	$\{\text{fdist 'rV[A]_n}\}$
	<code>Multivar.belast_last P</code>	$\{\text{fdist 'rV[A]_n * A}\}$
$P : \{\text{fdist 'rV[A]_n}\}$	<code>Nth.d P i</code>	$\{\text{fdist A}\}$
	<code>Take.d P i</code>	$\{\text{fdist 'rV[A]_i}\}$
$P : \{\text{fdist 'rV[A]_n * B}\}$	<code>PairNth.d P i</code>	$\{\text{fdist A * B}\}$
$P : \{\text{fdist 'rV[A]_n.+1 * B}\}$	<code>PairTake.d P</code>	$\{\text{fdist ('rV[A]_i * A) * B}\}$

`\Pr_(Proj23.d P) [F | G].`

The proof of intersection uses the conditional property of the universal set, i.e., the fact that $\Pr[U|A] = 1$ where U is the universal set and A is an event with non-zero probability:

Variables $(A B : \text{finType}) (P : \{\text{fdist A * B}\})$.

Lemma `cPr_1 a : Bivar.snd P a != 0 →`

`\sum_(b in B)`

`\Pr_(Swap.d P)[[set b] | [set a]] = 1.`

In the formal definition, `[set x]` is a MATHCOMP notation for the singleton set $\{x\}$.

3 Application 1: Formalization of Information Theory

As a first application of our formalization of conditional probability and joint distributions, we formalize the basic elements of information theory. This includes lemmas that (to the best of our knowledge) were not formalized before (e.g., the chain rule for information).

3.1 Entropy and Conditional Entropy

3.1.1 Background: Entropy

Entropy is a measure of the uncertainty of a random variable. In [Sect. 2.1 of [11]], it is defined as follows (for a random variable X with alphabet \mathcal{X} and probability mass function $p(x) \stackrel{\text{def}}{=} \Pr(X = x)$):

$$\mathcal{H}(X) = - \sum_{x \in \mathcal{X}} p(x) \log p(x)$$

Since the body of the definition uses only distributions, we formalized it as follows in [8] (CoQ notation: `^H`):

Variables $(A : \text{finType}) (P : \{\text{fdist A}\})$.

Definition `entropy :=`

`- \sum_(a in A) P a * log (P a).`

One may wonder what this definition means if $Pa = 0$ for some a . Here we rely on CoQ's choice that the values of `log` be 0 for non-positive values^{†1}. As a result, the product of 0 with `log 0` is equal to 0, which happens to be what we want here. While CoQ uses completed versions of logarithm and division for instance, assuming both of them to be 0 outside of their usual definition domains, theorems on specific functions are restricted to their usual definition domains. For instance $x/x = 1$ is true only if $x \neq 0$, cf. lemma `cPr_1` (Sect. 2.3).

3.1.2 Conditional Entropy

When it comes to conditional entropy, the standard textbook [11] gives a choice of several definitions (of course all equivalent) using slightly different notions: a more primitive notion of conditional entropy (which is defined simultaneously), using conditional probabilities, using joint distributions, or using the expectation of a random variable [Equations 2.10–2.13 of [11]]. For example, using conditional probabilities, the definition of conditional entropy reads as follows:

$$\mathcal{H}(Y | X) = - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x) \Pr[y|x] \log \Pr[y|x]$$

Let us use the definition of conditional probability from Sect. 2.3 to write this definition in

^{†1} CoQ assumes all functions to be total. Hence `log` is defined for every real number.

COQ. Like for entropy (Sect. 3.1.1), we only care about distributions. We are given a joint distribution $QP : \{\text{fdist } B * A\}$ (to be able to write down the conditional probability $\Pr[y|x]$) from which we compute the marginal $P := \text{Bivar.snd } QP$ (so as to be able to write $p(x)$):

```
(* Module CondEntropy *)
Variables (A B : finType) (QP : {fdist B * A}).
Definition h1 a := - \sum_(b in B)
  \Pr_QP[[set b] | [set a]] *
  log (\Pr_QP[[set b] | [set a]]).
Let P := Bivar.snd QP.
Definition h := \sum_(a in A) P a * h1 a.
```

Let us check the validity of our formal definition with an example of computation and a (non-trivial) property of conditional entropy.

3.1.2.1 Computation of Conditional Entropy

Let us consider the joint distribution with the following probability mass function [Example 2.2.1 of [11]]:

```
(* Module conditional_entropy_example. *)
Definition f : 'I_4 * 'I_4 → ℝ :=
  [eta (fun => 0) with
    (zero,zero) ↦ 1/8,      (zero,one) ↦ 1/16,
    (zero,two) ↦ 1/16,     (zero,three) ↦ 1/4,
    (one,zero) ↦ 1/16,     (one,one) ↦ 1/8,
    (one,two) ↦ 1/16,     (one,three) ↦ 0,
    (two,zero) ↦ 1/32,     (two,one) ↦ 1/32,
    (two,two) ↦ 1/16,     (two,three) ↦ 0,
    (three,zero) ↦ 1/32,   (three,one) ↦ 1/32,
    (three,two) ↦ 1/16,   (three,three) ↦ 0].
```

Since it is nonnegative and sums to 1, it is a distribution:

```
Lemma f0 : forall x, 0 ≤ f x.
Proof. ... Qed.
Lemma f1 :
  \sum_(x in {:'I_4 * 'I_4}) f x = 1.
Proof. ... Qed.
Definition d : {fdist 'I_4 * 'I_4} :=
  FDist.make f0 f1.
```

The constructor `FDist.make` is a variant of the constructor `FDist.mk` of Sect. 2.1. We compute $H(X|Y)$ where X and Y are random variables whose joint distribution is d . Using the symbolic computations capabilities of COQ as implemented by tactics such as `lra`, we recover the expected re-

sult:

```
Lemma conditional_entropyE :
  CondEntropy.h d = 11/8.
```

3.1.2.2 “Information Can’t Hurt”

One important property of information theory is that conditioning reduces entropy, i.e., $H(X|Y) \leq H(X)$ [Thm 2.6.5 of [11]]. Using the formal definitions explained so far, this can be stated in COQ as follows:

```
Variables (A B : finType) (PQ : {fdist A * B}).
Let P := Bivar.fst PQ.
Lemma information_cant_hurt :
  CondEntropy.h PQ ≤ `H P.
```

We provide a succinct COQ proof in Sect. 3.2.1, using the properties of mutual information.

3.1.3 The Chain Rule for Entropy

The chain rule for entropy is stated for n random variables X_1, X_2, \dots, X_n drawn according to $p(x_1, x_2, \dots, x_n)$ [Thm 2.5.1 of [11]]:

$$\mathcal{H}(X_1, \dots, X_n) = \sum_{i=1}^n \mathcal{H}(X_i | X_{i-1}, \dots, X_1)$$

First, we observe that $0 < n$. We reflect this in the formal statement below by considering $n.+1$ instead of n . Second, we observe that the pencil-and-paper notation is such that $\mathcal{H}(X_i | X_{i-1}, \dots, X_1) = \mathcal{H}(X_1)$ when $i = 1$. This is why, in the formal statement below, we distinguish the case where $i = 1$ from the case where $1 < i$ (which become $i = 0$ and $0 < i$ in COQ where we count from 0):

```
Lemma chain_rule_rV
  A n (P : {fdist 'rV[A]_n.+1}) :
  `H P = \sum_(i < n.+1)
    if i == 0 then
      `H (Multivar.head_of P)
    else
      CondEntropy.h (Swap.d (Multivar.belast_last
        (Take.d P (lift ord0 i)))).
```

In the case where $i = 0$, the chain rule returns $\text{`H (Multivar.head_of P)}$, which is equal to `H P since P is a vector with only one element. The case with $0 < i$ is a little bit involved. The expression `lift ord0 i` means $i + 1$; it is written with `lift` so that it has the appropriate type `'I_n.+2`, the type of natural number strictly less than $n.+2$ [20]. The calls to `Take.d`, `Multivar.belast_last`, and `Swap.d`

(see Tables 1 and 3) build the product distribution corresponding to $(X_{i+1}, (X_i, \dots, X_1))$. This may look a little bit more complicated than what one could expect with lists (instead of (finite-size) vectors), but since distributions have finite support we cannot easily recast the problem in terms of “distributions of lists”.

The proof is by induction on n and uses as an intermediate lemma the “chain rule” [Thm 2.2.1 of [11]], which is actually a version of the chain rule for entropy but restricted to three random variables. One can find this restricted version in related work (e.g., [22]), but we are not aware of an existing generalization to n distributions. The formal statement above shows that the need to build a variety of joint distributions makes the generalization not immediate. It is nevertheless useful, for example to prove the chain rule for information (see Sect. 3.2.3) or Han’s inequality (see [9]).

3.2 Mutual Information and Conditional Mutual Information

In this section, we give an overview of our formalization of mutual information, conditional mutual information, and the chain rule for information. It is technically a bit more involved than the formalization of conditional entropy in Sect. 3.1, but it uses similar ideas, hence the shorter presentation.

3.2.1 Mutual Information

Mutual information quantifies the information shared by two random variables. Given two random variables X and Y with joint probability mass function $p(x, y)$ and marginals $p(x)$ and $p(y)$, mutual information is defined as $\mathcal{I}(X; Y) = \mathcal{D}(p(x, y) \parallel p(x)p(y))$ where $\mathcal{D}(\cdot \parallel \cdot)$ is the Kullback-Leibler distance [Sect. 2.3 of [11]] (also called *relative entropy*). The formal definition is immediate given that we have already defined in previous work [8] the relative entropy $\mathcal{D}(\dots \parallel \dots)$ and the product distribution `x :

```
Variables (A B : finType) (PQ : {fdist A * B}).
Let P := Bivar.fst PQ.
Let Q := Bivar.snd PQ.
Definition mi := D(PQ || P `x Q).
```

The notation $P \text{ `x } Q$ is for the distribution with probability mass function $(x, y) \mapsto P(x) \times Q(x)$. It is defined using a more generic construct (where the second projection may depend on the first one) that will be better explained in the context of Sect. 3.3.

This formal definition is equipped with the mandatory lemmas that, among others, provide a short proof for the “information can’t hurt” property of Sect. 3.1.2:

```
Lemma information_cant_hurt :
  CondEntropy.h PQ ≤ `H P.
```

Proof.

```
rewrite -subR_ge0 -MutualInfo.miE;
exact: MutualInfo.mi_ge0.
```

Qed.

`MutualInfo.miE` provides an alternative definition of mutual information in terms of entropy (i.e., $\mathcal{I}(X; Y) = \mathcal{H}(X) - \mathcal{H}(X|Y)$ [Thm 2.4.1 of [11]]) and `MutualInfo.mi_ge0` establishes that mutual information is nonnegative [p. 28 of [11]].

3.2.2 Conditional Mutual Information

Conditional mutual information quantifies the information shared by two random variables *under the assumption that the result of a third random variable is already known*: $\mathcal{I}(X; Y|Z) = \mathcal{H}(X|Z) - \mathcal{H}(X|Y, Z)$ [Sect. 2.5 of [11]]. The formalization is immediate given the definition of conditional entropy from Sect. 3.1.2:

```
Variables (A B C : finType)
  (PQR : {fdist A * B * C}).
```

Definition cmi :=

```
  CondEntropy.h (Proj13.d PQR) -
  CondEntropy.h (PairA.d PQR).
```

3.2.3 The Chain Rule for Information

Mutual information and conditional mutual information satisfy a chain rule [Thm 2.5.2 of [11]]:

$$\mathcal{I}(X_1, X_2, \dots, X_n; Y) = \sum_{i=1}^n \mathcal{I}(X_i; Y|X_{i-1}, \dots, X_1)$$

It looks similar to the chain rule for entropy and conditional entropy of Sect. 3.1.3.

For the formal statement, we first prepare the

joint distribution X_1, X_2, \dots, X_n, Y as the distribution `PY` of type `{fdist 'rV[A]_n.+1 * B}`:

```
Variables (A : finType). Let B := A.
Variables (n : nat)
(PY : {fdist 'rV[A]_n.+1 * B}).
```

For the case where $i = 1$, we take the mutual information of the distribution X_1, Y , i.e., `(PairNth.d PY ord0)` (`ord0` is the natural number 0 but with the type `'I_n.+1`). For the case where $1 < i$, we take the conditional mutual information of the distribution $X_i, Y, (X_{i-1}, \dots, X_1)$, i.e., the distribution `TripAC.d (TripC12.d (PairTake.d PY i))`

that we build as the function `fAC`:

```
Let f (i : 'I_n.+1) : {fdist A * 'rV[A]_i * B} :=
  TripC12.d (PairTake.d PY i).
Let fAC (i : 'I_n.+1) : {fdist A * B * 'rV[A]_i} :=
  TripAC.d (f i).
```

Lemma `chain_rule_information` :

```
MutualInfo.mi PY = \sum_(i < n.+1)
  if i == 0 then
    MutualInfo.mi (PairNth.d PY ord0)
  else
    cmi (fAC i).
```

The proof uses the chain rule for entropy of Sect. 3.1.3.

3.3 Convexity of Entropy and Mutual Information

The goals of this section are formal proofs of convexity for entropy and mutual information. We start by defining convexity using the notion of *convex space* [17].

Let A and B be convex spaces. A function f of type $A \rightarrow B$ is convex when it satisfies the predicate `convex_function`:

```
Definition convex_function_at a b (p : prob) :=
  f (a <| p |> b) ≤ f a <| p |> f b.
```

```
Definition convex_function :=
  forall a b (p : prob),
    convex_function_at f a b p.
```

The type `prob` is for reals between 0 and 1; `a <| p |> b` is a notation for $p \times a + (1-p) \times b$. Similarly, the predicate `concave_function` is for functions f such that $-f$ is convex. The above defini-

tion applies to the type of finite distributions and to the type of CoQ reals \mathbb{R} which both happen to be convex spaces; in particular, to indicate to the type system of CoQ that finite distributions over A are to be seen as forming a convex space, we use the type `fdist_convType A` instead of `{fdist A}`. We do not explain here in details the formalization of convex spaces, it is the object of a forthcoming publication, in the meantime all the details are in the CoQ formalization [file `convex_choice.v` of [9]].

Given the above definitions, we can state the convexity of entropy as follows [Thm 2.7.3 of [11]]:

```
Variable (A : finType).
Hypothesis A_not_empty : 0 < #|A|.
Lemma entropy_concave : concave_function
  (fun P : fdist_convType A => `H P).
```

The proof relies on a convexity property of relative entropy.

The convexity of mutual information requires an additional construct. Let us first recall the pencil-and-paper statement. Let (X, Y) be two random variables drawn according to $p(x, y) = p(x) \Pr[y|x]$. For fixed $\Pr[y|x]$, $\mathcal{I}(X; Y)$ is a concave function of $p(x)$ [first part of Thm 2.7.4 of [11]]. In the formal statement, we represent $\Pr[y|x]$ as $W : A \rightarrow \{\text{fdist } B\}$ (i.e., a stochastic matrix), and $p(x)$ is given as the parameter $P : \{\text{fdist } A\}$. We build the joint distribution $p(x, y)$ using the function `make_joint` defined as follows:

```
(* Module CJFDist *)
Variables (A B : finType).
Record t :=
  mkt {P : {fdist A} ; W :> A → {fdist B}}.
Definition joint_of (x : t) : {fdist A * B} :=
  ProdDist.d (P x) (W x).
Definition make_joint (P : dist A)
  (W : A → dist B) : {fdist A * B} :=
  joint_of (mkt P W).
```

Given a distribution P and a stochastic matrix W , `ProdDist.d` builds the “product distribution” with probability mass function $(x, y) \mapsto P(x) \times Q(x, y)$. (This explains the product distribution of Sect. 3.2.1: $P \times Q$ is actually a notation for `ProdDist.d P (fun => Q)`.) By construction, the left

marginal of the joint distribution is P , which is one part of the original condition (specifically the condition $p(x, y) = p(x) \Pr[y|x]$) and lets us prove the concavity of mutual information (when its second argument is blocked):

Variables (A B : finType) (Q : A → {fdist B}).

Hypothesis B_not_empty : 0 < #|B|.

Lemma mutual_information_concave :

```
concave_function
(fun P : fdist_convType A =>
  MutualInfo.mi (CJFDist.make_joint P Q)).
```

The proof uses the convexity of entropy as an intermediate step.

The convexity of mutual information (as a function of $\Pr[y|x]$ for fixed $p(x)$) [second part of Thm 2.7.4 of [11]] follows along similar lines (see [9]).

3.4 More Information Theory

We were actually able to formalize more information theory than what we have explained so far. At the time of this writing, we have almost completed the formalization of [Chapter 2 of [11]]: only the sections about sufficient statistics and Fano’s inequality are left unexplored.

Among theorems that we have not explained are also the chain rule for relative entropy [Thm 2.5.3 of [11]], the independence bound on entropy [Thm 2.6.6 of [11]], the definition of Markov chains, and the data-processing inequality [Thm 2.8.1 of [11]].

This work uses and complements our previous work that already provides in COQ several theorems from [Chapter 2 of [11]]: Jensen’s inequality [Thm 2.6.2 of [11]] (see [6]), the log sum inequality [Thm 2.7.1 of [11]] (see [8]), and the proof that the nonnegativity of the second derivative implies convexity [Thm 2.6.1 of [11]].

4 Application 2: Formalization of Con-

ditional Independence

As a second application of our formalization of conditional probability and joint distributions, we formalize conditional independence and we validate this formal definition by establishing the graphoid axioms.

4.1 Random Variables

Given a distribution $P : \{\text{fdist } U\}$, a random variable over A is defined as a function of type $U \rightarrow A$, where A is expected to be some finite type^{†2} for the values that the random variable can take:

Definition RV {U : finType} (P : {fdist U})

(A : eqType) := U → A.

This is the same definition as Moreira [22]. In the following, we use $\{\text{RV } P \rightarrow A\}$ as a notation for $\text{RV } P \text{ A}$.

A random variable $X : \{\text{RV } P \rightarrow A\}$ induces a distribution over its codomain. This is the distribution whose probability mass function associates to each element a of A the probability of the event corresponding to the pre-image of a via X , i.e., $\text{FDistMap.d } X \text{ P}$ (see Sect. 2.2.1). We introduce the COQ notation $\backslash\text{Pr}[X = a]$ corresponding to the probability $\text{Pr}(\text{FDistMap.d } X \text{ P})[\text{set } a]$, and more generally the COQ notation $\backslash\text{Pr}[X = a \mid Y = b]$ for $\backslash\text{Pr}(\text{FDistMap.d } [\% X, Y] \text{ P})[[\text{set } a] \mid [\text{set } b]]$. Recall the notation for conditional probability from Sect. 2.3. The complete COQ definitions and proofs are available online [files `proba.v` and `cinde.v` of [9]].

In this setting, we observe in particular that we can create random variables by pairing existing random variables:

Variables (U : finType) (P : {fdist U}).

Variables (A : finType) (X : {RV P → A})

(B : finType) (Y : {RV P → B}).

Definition RV2 : {RV P → A * B} :=

`fun x => (X x, Y x).`

^{†2} This is not enforced by the definition that only requires a type with a decidable equality but this is required by most useful lemmas.

In the following, $[X, Y, \dots, Z]$ denotes the iterated pairing of the random variables X, Y, \dots, Z .

4.2 Conditional Independence

We are concerned with the formalization of the predicate $X \perp Y | Z$, which intuitively means that X is independent of Y given Z . We first recall the textbook definition.

Definition 1 (Conditional Independence [Definition 2.4 of [18]]). *Let X, Y , and Z be random variables and P be a distribution. X is conditionally independent of Y given Z in the distribution P if for all values a, b , and c (belonging resp. to the codomains of X, Y , and Z), we have:*

$$\begin{aligned} \Pr(X = a, Y = b | Z = c) \\ = \Pr(X = a | Z = c) \Pr(Y = b | Z = c). \end{aligned}$$

We now move on to the formal definition of conditional independence. First, we give ourselves a sample space U with a distribution $P : \{\text{fdist } U\}$:

Variables ($U : \text{finType}$) ($P : \{\text{fdist } U\}$).

Second, we declare three random variables X, Y , and Z :

Variables ($A B C : \text{finType}$).

Variables ($X : \{\text{RV } P \rightarrow A\}$) ($Y : \{\text{RV } P \rightarrow B\}$)

($Z : \{\text{RV } P \rightarrow C\}$).

Finally, we state the property of conditional independence (COQ notation: $X \perp Y | Z$):

Definition `cinde_drv` := `forall a b c,`

$$\begin{aligned} \backslash\Pr [[X, Y] = (a, b) \mid Z = c] = \\ \backslash\Pr [X = a \mid Z = c] * \backslash\Pr [Y = b \mid Z = c]. \end{aligned}$$

Comparison with the pencil-and-paper definition is immediate. The fact that each probability appearing in this definition relies on different distributions can still be guessed by the reader from the occurrences of the random variables X, Y, Z , but moreover it is now possible to unfold formal definitions to discover explicitly that, say, distributions in the right hand-side are marginals of the distribution in the left hand-side. One goal of our formalization is however to provide enough definitions and lemmas to avoid doing unfolding and keeping formal proofs as simple as possible.

4.3 Conditional Independence is a Graphoid

Conditional independence satisfies the graphoid axioms [23]:

- Symmetry: $X \perp Y | Z$ implies $Y \perp X | Z$.
- Decomposition: $X \perp Y, W | Z$ implies $X \perp Y | Z$.
- Weak union: $X \perp Y, W | Z$ implies $X \perp Y | Z, W$.
- Contraction: $X \perp W | Z, Y$ and $X \perp Y | Z$ imply $X \perp Y, W | Z$.
- Intersection: $X \perp Y | Z, W$ and $X \perp W | Z, Y$ imply $X \perp Y, W | Z$ for positive distributions.

These axioms have an intuitive interpretation. For example, symmetry means that in a state Z where X says nothing new about Y , then Y says nothing new about X ; decomposition means that if Y, W say nothing new about X , then Y alone does not say anything new about X ; etc.

In order to demonstrate the expressiveness of our approach, we have proved that the graphoid axioms for conditional independence are satisfied by the construction of finite discrete probability theory on top of finite types and real numbers seen in Sect. 2.1. The formal proofs, detailed below, are written following pencil-and-paper proofs that have themselves been reconstructed from [18][24], trying to keep them as close as possible. The full scripts are available online [file `cinde.v` of [9]]. The proof scripts for symmetry and decomposition are almost identical to the pencil-and-paper proofs, but the remaining ones require extra handling for corner cases.

4.3.1 The Proof of Symmetry in Coq

Let us recall the pencil-and-paper proof:

Lemma 1 (Symmetry). $X \perp Y | Z$ implies $Y \perp X | Z$.

Proof.

$$\begin{aligned}
P(Y, X|Z) &= P(X, Y|Z) \\
&\quad (\text{by set-theoretic reasoning}) \\
&= P(X|Z)P(Y|Z) \\
&\quad (\text{because } X \perp Y | Z) \\
&= P(Y|Z)P(X|Z) \\
&\quad (\text{by commutativity of multiplication}) \quad \square
\end{aligned}$$

In Coq, we write the statement “ $X \perp Y | Z$ implies $Y \perp X | Z$ ” as follows:

Variable (U : finType) (P : {fdist U}).

Variables (A B C : finType) (X : {RV P → A}).

(Y : {RV P → B}) (Z : {RV P → C}).

Lemma symmetry : X ⊥ Y | Z → Y ⊥ X | Z.

The proof script follows the pencil-and-paper proof:

```

1 Lemma symmetry : X ⊥ Y | Z → Y ⊥ X | Z.
2 Proof.
3 move=> H b a c.
4 rewrite RV_Pr_1C.
5 rewrite H.
6 by rewrite mulRC.
7 Qed.

```

The first line introduces the hypothesis, and moves the universally quantified variables to the context. We are then working on the equation defined in `cinde_drv`, and do rewritings until both sides are syntactically equal. At first, there is no syntactic match between the random variables in the hypothesis and in the goal: the hypothesis is about the random variables $[X, Y, Z]$ whereas the goal is about $[Y, X, Z]$. The purpose of line ?? is to swap the two leading random variables so as to be able to rewrite with the hypothesis at line ?. We apply the commutativity of multiplication at line ??, which completes the proof.

4.3.2 The Proof of Decomposition in Coq

The proof of decomposition is more involved than symmetry; it relies in particular on “reasoning by cases”, i.e., $P(X|Y) = \sum_z P(X, z|Y)$, which becomes in our formalization:

Lemma reasoning_by_cases E F :

$\Pr[X \in E | Y \in F] =$
 $\sum_{z \leftarrow \text{fin_img } Z}$

$\Pr[[X, Z] \in \text{setX } E [\text{set } z] | Y \in F].$

where E and F are two events and `fin_img` is the (finite) image of a random variable. The notation $\Pr[X \in E | Y \in F]$ stands for $\Pr_{\text{(FDistMap.d } [X, Y] P)}[E | F]$ and comes as a slight generalization of the notation introduced in Sect. 4.1.

Lemma 2 (Decomposition). $X \perp Y, W | Z$ implies $X \perp Y | Z$.

Proof. [Sect. 2.1.4.3 of [18]]

$$\begin{aligned}
P(X, Y|Z) &= \sum_w P(X, Y, w|Z) \\
&\quad (\text{reasoning by cases}) \\
&= \sum_w P(X|Z)P(Y, w|Z) \\
&\quad (\text{by } X \perp Y, W | Z) \\
&= P(X|Z) \sum_w P(Y, w|Z) \\
&\quad (\text{by distributivity}) \\
&= P(X|Z)P(Y|Z) \\
&\quad (\text{reasoning by cases}) \quad \square
\end{aligned}$$

Here follows the proof statement in Coq:

Variables (U : finType) (P : {fdist U}).

(A B C D : finType).

Variables (X : {RV P → A}) (Y : {RV P → B}).

(Z : {RV P → C}) (W : {RV P → D}).

Lemma decomposition :

X ⊥ [Y, W] | Z → X ⊥ Y | Z.

Figure 1 displays a proof script that proves decomposition. Although it is not necessary (and actually lengthen the proof artificially), we have added `transitivity` steps so that the flow of the proof script matches the pencil-and-paper proof. If the goal is $A = B$, then `transitivity C` will split it into two goals, $A = C$ and $C = B$. The first intermediate step of the pencil-and-paper proof appears explicitly as a `transitivity` subgoal at line ?. Solving it is essentially a matter of reasoning by cases (see line ?). The second intermediate step of the pencil-and-paper proof can be seen explicitly at line ?. It is essentially solved by using the hypothesis of conditional independence (which is

```

1 Lemma decomposition : X ⊥ [% Y, W] | Z → X ⊥ Y | Z.
2 Proof.
3 move=> H a b c.
4 transitivity (\sum_(d <- fin_img W) \Pr[ [% X, [% Y, W]] = (a, (b, d)) | Z = c]).
5   rewrite (reasoning_by_cases W); apply eq_bigr => /= d _.
6   by rewrite RV_Pr_1A setX1.
7 transitivity (\sum_(d <- fin_img W)
8   \Pr[ X = a | Z = c] * \Pr[ [% Y, W] = (b, d) | Z = c]).
9   by apply eq_bigr => d _; rewrite H.
10  rewrite -big_distr /=; congr (_ * _).
11  rewrite (reasoning_by_cases W); apply eq_bigr => d _.
12  by rewrite setX1.
13 Qed.

```

Figure 1 Proof script for the lemma decomposition

bound to the identifier H inside the proof script) at line ???. Line ??? is where we use the distributivity of multiplication w.r.t. addition (lemma `big_distr`). Line ??? is the conclusive reasoning by cases, like the first proof-step. The lines that we have not explained in details are just trivial manipulations of random variables and finites sets. For example, `RV_Pr_1A` is a lemma to rearrange random variables appearing in conditional probabilities by appealing to their associativity.

4.3.3 The Proof of Weak Union in Coq

The pencil-and-paper proof of weak union is shorter than decomposition, but uses it as an intermediate step.

Lemma 3 (Weak Union). $X \perp Y, W | Z$ implies $X \perp Y | Z, W$.

Proof.

$$\begin{aligned}
 P(X, Y | Z, W) &= P(X | Y, Z, W) P(Y | Z, W) \\
 &\quad \text{(by the product rule)} \\
 &= P(X | Z) P(Y | Z, W) \\
 &\quad \text{(because } X \perp Y, W | Z) \\
 &= P(X | Z, W) P(Y | Z, W)
 \end{aligned}$$

The last step is because $X \perp W | Z$ by decomposition (Lemma 2) from $X \perp Y, W | Z$. \square

Here follows the proof statement in COQ:

```

Variables (U : finType) (P : {fdist U})
(A B C D : finType).
Variables (X : {RV P → A}) (Y : {RV P → B})
(Z : {RV P → C}) (W : {RV P → D}).

```

Lemma weak_union :

$$X \perp [% Y, W] | Z \rightarrow X \perp Y | [% Z, W].$$

See Fig. 2 for the complete proof script. The first step appears at line ??? and is completed by application of the product rule at line ???. Note that this is a variant of the product rule presented in Sect. 2.3: it is specialized to random variables. The second step appears at line ??? and is completed by using the conditional independence hypothesis at line ???. The use of decomposition can be observed at line ??? (where the lemma `cinde_drv_2C` just operates a reordering of random variables) and is followed by another use of the conditional independence hypothesis at line ???.

The proof script is maybe longer than what one could expect given the previous examples. This is because of corner cases with null probabilities, that are often ignored in pencil-and-paper proofs. Here they come from the application of the lemma `cindeP`, which establishes

$$X \perp Y | Z \rightarrow P(X | Y, Z) = P(X | Z)$$

under the hypothesis that $P(Y, Z) \neq 0$.

4.3.4 The Proof of Contraction in Coq

The proof of contraction uses that same ideas as previous proofs:

Lemma 4 (Contraction). $X \perp W | Z, Y$ and $X \perp Y | Z$ imply $X \perp Y, W | Z$.

```

1 Lemma weak_union : X ⊥ [% Y, W] | Z → X ⊥ Y | [% Z, W].
2 Proof.
3 move=> H a b [c d].
4 transitivity (\Pr[ X = a | [% Y, Z, W] = (b, c, d)] *
5 \Pr[ Y = b | [% Z, W] = (c, d)]).
6 by rewrite RV_product_rule RV_Pr_rA.
7 transitivity (\Pr[ X = a | Z = c] * \Pr[ Y = b | [% Z, W] = (c, d)]).
8 rewrite RV_Pr_rAC.
9 case/boolP : (\Pr[ [% Y, W, Z] = (b, d, c)] == 0) => [/eqP] H0.
10 - by rewrite [X in _ * X = _ * X]RV_cPrE RV_Pr_A RV_Pr_AC H0 divOR !mulR0.
11 - by rewrite (cindeP _ H).
12 case/boolP : (\Pr[ [% Z, W] = (c, d) ] == 0) => [/eqP] ?.
13 - by rewrite [X in _ * X = _ * X]RV_cPrE RV_Pr_domin_snd ?(divOR,mulR0).
14 - have {H}H : X ⊥ W | Z by move/cinde_drv_2C : H; apply decomposition.
15 by rewrite [in X in _ = X * _]RV_Pr_rC (cindeP _ H) // RV_Pr_C.
16 Qed.

```

Figure 2 Proof script for the lemma weak_union

Proof.

$$\begin{aligned}
P(X, Y, W|Z) &= P(X|Y, W, Z)P(Y, W|Z) \\
&\quad \text{(by the product rule)} \\
&= P(X|Y, Z)P(Y, W|Z) \\
&\quad \text{(because } X \perp W | Z, Y) \\
&= P(X|Z)P(Y, W|Z) \\
&\quad \text{(because } X \perp Y | Z) \\
&\quad \square
\end{aligned}$$

We display the formal proofs of contraction in Fig. 3 for the sake of completeness but without further comments since it uses the same ingredients as Fig. 2.

4.3.5 The Proof of Intersection in Coq

The proof of intersection is the most involved of all graphoid axioms. In particular, intersection uses contraction, and reasoning by cases and the product rule as intermediate steps.

Lemma 5 (Intersection). $X \perp Y | Z, W$ (1) and $X \perp W | Z, Y$ (2) imply $X \perp Y, W | Z$ provided that the distribution Y, Z, W is positive (3) and the type of the codomain of W is not empty (4).

Proof. It suffices to show $X \perp Y | Z$ by contraction (Lemma 4) because we already have (2).

The goal is therefore $P(X, Y|Z) = P(X|Z)P(Y|Z)$

which is proved as follows:

$$\begin{aligned}
P(X|Z)P(Y|Z) &= \sum_w P(X, w, Z)P(Y|Z) \\
&\quad \text{(reasoning by cases)} \\
&= \sum_w P(X, Y|Z)P(w|Z) \\
&\quad \text{(see (*) below)} \\
&= P(X, Y|Z) \sum_w P(w|Z) \\
&\quad \text{(by distributivity)} \\
&= P(X, Y|Z)
\end{aligned}$$

The last step is by conditional probability of the universal set (Sect. 2.3) and (4).

Proof of step (*): By (2), (3), and the product rule, we have $P(X|Y, Z) = P(X|W, Z, Y)$.

By (1), (3), and the product rule, we have $P(X|W, Z) = P(X|Y, W, Z)$.

Therefore we have $P(X|Y, Z) = P(X|W, Z)$.

By (3) and the product rule, we derive $\frac{P(X, Y|Z)}{P(Y|Z)} = \frac{P(X, W|Z)}{P(W|Z)}$,

which implies

$$P(X, Y|Z)P(W|Z) = P(X, W|Z)P(Y|Z)$$

and

$$\sum_w P(X, Y|Z)P(w|Z) = \sum_w P(X, w|Z)P(Y|Z).$$

□

The longer proof of intersection naturally translated into a longer proof script that we partially display in Fig. 4. Its structure is however still read-

```

1 Lemma contraction : X ⊥ W | [% Z, Y] → X ⊥ Y | Z → X ⊥ [% Y, W] | Z.
2 Proof.
3 move=> H1 H2 a [b d] c.
4 rewrite RV_product_rule.
5 transitivity (\Pr[X = a | [% Y, Z] = (b, c)] * \Pr[[% Y, W] = (b, d) | Z = c]).
6 rewrite -RV_Pr_rA [in X in X * _ = _]RV_Pr_rC -RV_Pr_rA.
7 case/boolP : (\Pr[ [% W, [% Z, Y]] = (d, (c, b))] == 0) => [/eqP|] H0.
8   rewrite [in X in _ * X = _ * X]RV_cPrE.
9   by rewrite -RV_Pr_A RV_Pr_C -RV_Pr_A H0 divOR !mulRO.
10  by rewrite (cindeP _ H1) // RV_Pr_rC.
11 case/boolP : (\Pr[ [% Y, Z] = (b, c) ] == 0) => [/eqP|] H0.
12 - rewrite [X in _ * X = _ * X]RV_cPrE.
13   by rewrite RV_Pr_AC RV_Pr_domin_fst ?divOR ?mulRO.
14 - by rewrite (cindeP _ H2).
15 Qed.

```

Figure 3 Proof script for the lemma contraction

able. The first step using contraction appears at line ?? . It is a “it suffices to” step that is conveniently represented by the `suff` tactic of SSREFLECT. The other steps of the main proof appear respectively at line ?? (reasoning by cases), line ?? ((*)) step, lines ?? and ?? (distributivity), and after line ?? for the conclusion. The technicalities of the (*) step have been omitted but can be found online [file `cinde.v` of [9]].

4.4 Application of the Graphoid Axioms

Now that we have proved the graphoid axioms, we can derive other rules to reason formally about probabilistic models without much clutter. For example, we can prove the *chaining rule* without unfolding definitions. The proof script stays to-the-point and is short, even though it uses *all* the *semi-graphoid axioms* (i.e., the graphoid axioms except intersection):

```

Lemma chaining_rule :
  X ⊥ Z | Y ∧ [% X, Y] ⊥ W | Z →
  X ⊥ W | Y.
Proof.
case=> ? ?.
suff : X ⊥ [% W, Z] | Y by move/decomposition.
apply/cinde_drv_2C/contraction => //.
exact/cinde_drv_3C/symmetry/weak_union/symmetry.
Qed.

```

The lemmas `cinde_drv_xC` operate mere reorderings of random variables (we already saw the lemma

`cinde_drv_2C` in Sect. 4.3.3).

See also the *mixing rule* for another example of derived rule [file `cinde.v` of [9]].

5 Related Work

The theory of conditional probabilities (Bayes’ theorems, etc.) is developed in HOL and applied to the analysis of the binary asymmetric channel [15], but there does not seem to be any theory of joint distributions in this work. We were able to reproduce the same theory of conditional probabilities as [15] in the setting of Sect. 2.3, except for lemmas involving countable unions which are not handled by the theory of finite sets of MATHCOMP. Information-theoretic notions such as entropy, relative entropy, and mutual information are defined in HOL using a formalization of probability based on measure theory and Lebesgue integration [21], but there does not seem to be much lemmas about them. There are more information-theoretic notions (e.g., conditional mutual information) defined Isabelle/HOL using Lebesgue integration [16]. Compared to the work above, our work features much more lemmas (all the chain rules, lemmas about conditional independence, etc.).

Our work uses a COQ library called INFOTHEO [9]

```

1 Lemma intersection : X ⊥ Y | [% Z, W] → X ⊥ W | [% Z, Y] → X ⊥ [% Y, W] | Z.
2 Proof.
3 move=> H1 H2.
4 suff : X ⊥ Y | Z by apply contraction.
5 move=> a b c; apply/esym.
6 rewrite [in X in X * _ = _](reasoning_by_cases W).
7 under eq_bigr do rewrite setX1.
8 rewrite big_distr1 /=.
9 have <- : \sum_(d <- fin_img W)
10     \Pr[ [% X, Y] = (a, b) | Z = c] * \Pr[ W = d | Z = c] =
11     \sum_(d <- fin_img W)
12     \Pr[ [% X, W] = (a, d) | Z = c] * \Pr[ Y = b | Z = c].
13 suff H : forall d, \Pr[ [% X, Y] = (a, b) | Z = c] / \Pr[ Y = b | Z = c] =
14     \Pr[ [% X, W] = (a, d) | Z = c] / \Pr[ W = d | Z = c].
15 ... (* by using properties of positive distributions *)
16 suff H : forall d, \Pr[ X = a | [% Y, Z] = (b, c)] =
17     \Pr[ X = a | [% W, Z] = (d, c)].
18 ... (* by using the product rule *)
19 have {H2}H2 : forall d, \Pr[ X = a | [% Y, Z] = (b, c)] =
20     \Pr[ X = a | [% W, Z, Y] = (d, c, b)].
21 ... (* by using the product rule and the second conditional independence hypothesis *)
22 have {H1}H1 : forall d, \Pr[ X = a | [% W, Z] = (d, c)] =
23     \Pr[ X = a | [% Y, W, Z] = (b, d, c)].
24 ... (* by using the product rule and the first conditional independence hypothesis *)
25 by move=> d; rewrite {H2}(H2 d) {}H1 RV_Pr_rC RV_Pr_rA.
26 rewrite -big_distr // = cPr_1_RV ?mulR1 //.
27 move: (P0 b c D_not_empty); apply: contra.
28 by rewrite RV_Pr_AC => /eqP/(RV_Pr_domin_snd [% Y, W] (b, D_not_empty)) ->.
29 Qed.

```

Figure 4 Partial proof script for the lemma intersection

(see [file cinde.v of [9]] for the complete proof script)

that comes with a formalization of distributions and probabilities (including basic lemmas such as the weak law of large numbers) and applications to information theory (including proofs of Shannon’s theorems and linear error-correcting codes). This work already features definitions of conditional entropy, mutual information, and aposteriori probabilities, but specialized for a setting with one communication channel W (modeled as a stochastic matrix) together with an input distribution P . In this setting, the conditional entropy is denoted by $H(W|P)$ and the mutual information is denoted by $I(P, W)$, and both are defined using the joint distribution $J(P, W)$ (see [8] for details). These notations departs from [11] but are not uncommon (see, e.g., [12]). Still, we can express $H(W|P)$ using `CondEntropy.h` and $I(P, W)$ using `MutualInfo.mi`.

The first step is to show that each entry of the communication channel W corresponds to a conditional probability as defined in this paper using the joint distribution $J(P, W)$:

```

Variables (A B : finType) (W : `Ch_1(A, B))
(P : {fdist A}).
Let QP := Swap.d (`J(P, W)).
Lemma WcPr : forall a b, P a != 0 →
W a b = \Pr_QP[[set b] | [set a]].

```

Using this lemma, we show that $H(W|P)$ is equal to `CondEntropy.h QP` and $I(P, W)$ is equal to `MutualInfo.mi QP`. In other words, we do not need the specialized definitions of conditional entropy and of mutual information anymore in our previous work [8]. Similarly, in our previous work [3], we defined aposteriori probabilities using P and W as follows:

$$P^W(x|y) \stackrel{\text{def}}{=} \frac{P(x)W^n(y|x)}{\sum_{x' \in A^n} P(x')W^n(y|x')}.$$

We can now express a posteriori probabilities using the definition of conditional probability introduced in this paper. More precisely, $P^W(x|y)$ is equal to $\Pr_{\mathcal{J}(P, W \text{ } n)}([\text{set } x] | [\text{set } y])$ where $\mathcal{J}(P, W \text{ } n)$ is the joint distribution of the n^{th} extension of the channel W . See [file chap2.v of [9]] for details. In other words, the formal definitions we introduce in this paper are more generic because they can be used to factorize and thus simplify our previous formal developments.

There is another formalization of information theory that uses a similar setting [22]. Distributions are defined similarly, which is not surprising because it is natural to use the big operators of MATHCOMP in this way. The definition of random variable in Sect. 4.1 is the same definition as [22]. Yet, that work stops at an early stage with the chain rule for entropy restricted to three variables and the definition of mutual information. Another limitation of this formalization is that it axiomatizes the logarithm, whereas INFOTHEO is compatible with the logarithm defined in the standard library of COQ.

Our work connects concretely to the semantics of probabilistic programming languages. We are currently using the INFOTHEO library to provide a formal model for a monad combining probability and nondeterminism [4]. It happens that such a model requires the notion of convex spaces that we have developed in Sect. 3.3.

Regarding conditional independence, one can find a formal definition of the graphoid axioms in COQ/SSREFLECT [25] but this work does not seem to provide the underlying formalization of probabilities; instead, it focuses on a high-level algebraic presentation and its application.

The formalization of conditional independence can be tackled in a different way using *string diagrams* [Proposition 6.10 of [10]].

6 Conclusion

In this paper, we proposed a formalization of conditional probabilities and joint distributions validated by two original applications: a formalization of the foundations of information theory and a formalization of conditional independence. Our formalization of information theory extends previous work with lemmas that were not formalized before (e.g., the chain rule for information, convexity of entropy). Our second application is (to the best of our knowledge) the first formal account *from the ground up* of the graphoid axioms of conditional independence. We tried as much as possible to recover in the formalization the terseness of the presentation of information theory or artificial intelligence textbooks, resorting to COQ's dependent types and user-defined notations (to define conditional independence for example) and to the expressiveness of the SSREFLECT language of tactics (see the proof of the chaining rule for example). The complete formalization is available online as a conservative extension of the INFOTHEO library [9] (Table 4 summarizes the relevant files).

Using our work, it is now possible to address new challenges in formalization of information theory and artificial intelligence. Next, we plan to tackle the formalization of Fano's inequality, the formalization of Bayesian networks, and more generally formal reasoning about probabilistic graphical models.

Acknowledgments The authors are thankful to Kenta Cho, Shin-ya Katsumata, and Keisuke Nakano for fruitful discussions about the work presented in this paper, to Paul Fournier for his proofreading, and to the anonymous reviewers of PPL2019 for their helpful comments and for suggesting the application to Han's inequality. The authors acknowledge partial support from a Grant-in-Aid for Scientific Research (B) (project number

Table 4 Overview of the formalization in this paper (InfoTheo, version 0.0.6 [9])

File	Contents	Sections in this paper
<i>INFOtheo file that required extensions for this paper</i>		
<code>probability/fdist.v</code>	Probabilities over finite distributions	Sect. 2.1
	Joint distributions	Sect. 2.2
<i>New files that come with this paper</i>		
<code>probability/jfdist.v</code>	Conditional probabilities over joint distributions	Sect. 2.3
<code>information_theory/chap2.v</code>	Information theory	Sections 3.1, 3.2, and 3.4
<code>toy_examples/conditional_entropy.v</code>	Example of computation	Sect. 3.1.2.1
<code>probability/convex_choice.v</code>	Convexity properties	Sect. 3.3
<code>information_theory/convex_fdist.v</code>	Convexity properties	Sect. 3.3
<code>probability/cinde.v</code>	Conditional independence	Sect. 4

18H03204).

References

- [1] R. Affeldt. Introduction to Mathematical Components. *Computer Software*, 34:64–74, 2017.
- [2] R. Affeldt, C. Cohen, and D. Rouhling. Formalization techniques for asymptotic reasoning in classical analysis. *Journal of Formalized Reasoning*, 11(1):43–76, 2018.
- [3] R. Affeldt and J. Garrigue. Formalization of error-correcting codes: from Hamming to modern coding theory. In *6th Conference on Interactive Theorem Proving (ITP 2015), Nanjing, China, August 24–27, 2015*, volume 9236 of *Lecture Notes in Computer Science*, pages 17–33. Springer, Aug 2015. Superseded by [7].
- [4] R. Affeldt, J. Garrigue, D. Nowak, and T. Saikawa. Monadic equational reasoning in Coq. <https://github.com/affeldt-aist/monae> (last access: 2020/01/05). Main reference: [5].
- [5] R. Affeldt, D. Nowak, and T. Saikawa. A Hierarchy of Monadic Effects for Program Verification using Equational Reasoning. In *13th International Conference on Mathematics of Program Construction (MPC 2019), Porto, Portugal, October 7–9, 2019*, volume 11825 of *Lecture Notes in Computer Science*. Springer, Oct 2019.
- [6] R. Affeldt, J. Garrigue, and T. Saikawa. Examples of formal proofs about data compression. In *International Symposium on Information Theory and Its Applications (ISITA 2018), Singapore, October 28–31, 2018*, pages 665–669. IEICE. IEEE Xplore, Oct 2018.
- [7] R. Affeldt, J. Garrigue, and T. Saikawa. A Library for Formalization of Linear Error-correcting Codes. *Journal of Automated Reasoning*, 2020. To appear.
- [8] R. Affeldt, M. Hagiwara, and J. Sénizergues. Formalization of Shannon’s theorems. *Journal of Automated Reasoning*, 53(1):63–103, 2014.
- [9] R. Affeldt, M. Hagiwara, J. Sénizergues, J. Garrigue, K. Sakaguchi, T. Asai, T. Saikawa, and N. Obata. A Coq formalization of information theory and linear error-correcting codes. Version 0.0.6. <https://github.com/affeldt-aist/infotheo> (last access: 2020/01/05).
- [10] K. Cho and B. Jacobs. Disintegration and bayesian inversion via string diagrams. *Mathematical Structures in Computer Science*, 2019.
- [11] T. M. Cover and J. A. Thomas. *Elements of information theory*. Wiley, 2006. 2nd edition.
- [12] I. Csiszár and J. Körner. *Information Theory: Coding Theorems for Discrete Memoryless Systems*. Cambridge University Press, 2011. 2nd edition.
- [13] G. Gonthier, A. Mahboubi, and E. Tassi. A small scale reflection extension for the Coq system. Technical report, INRIA, 2008. Version 17 (Nov 2016).
- [14] G. Gonthier, A. Asperti, J. Avigad, Y. Bertot, C. Cohen, F. Garillot, S. Le Roux, A. Mahboubi, R. O’Connor, S. Ould Biha, I. Pasca, L. Rideau, A. Solovyev, E. Tassi, and L. Théry. A Machine-Checked Proof of the Odd Order Theorem. In *4th International Conference on Interactive Theorem Proving (ITP 2013), Rennes, France, July 22–26, 2013*, volume 7998 of *Lecture Notes in Computer Science*, pages 163–179. Springer, 2013.
- [15] O. Hasan and S. Tahar. Reasoning about conditional probabilities in a higher-order-logic theorem prover. *Journal of Applied Logic*, 9(1):23–40, 2011.
- [16] J. Hölzl. *Construction and Stochastic Applications of Measure Spaces in Higher-Order Logic*. PhD thesis, Technische Universität München Institut für Informatik, 2012.
- [17] B. Jacobs. Convexity, Duality and Effects. In *Theoretical Computer Science*, volume 323 of *IFIP Advances in Information and Communication Technology*, pages 1–19. Springer, 2010.
- [18] D. Koller and N. Friedman. *Probabilistic Graphical Models: Principles and Techniques*. MIT Press, 2009.
- [19] X. Leroy. Formal verification of a realistic compiler. *Communications of the ACM*, 52(7):107–115. ACM, 2009.
- [20] A. Mahboubi and E. Tassi. *Mathematical Com-*

ponents. Available at: <https://math-comp.github.io/mcb/> (last access: 2020/01/05), 2016. With contributions by Y. Bertot and G. Gonthier. Version of 2018/08/11.

[21] T. Mhamdi, O. Hasan, and S. Tahar. Formalization of entropy measures in HOL. In *2nd International Conference on Interactive Theorem Proving (ITP 2011), Berg en Dal, The Netherlands, August 22–25, 2011*, volume 6898 of *Lecture Notes in Computer Science*, pages 233–248. Springer, 2011.

[22] D. A. C. V. Moreira. Finite probability distributions in Coq. Master’s thesis, Universidade do Minho, Escola de Engenharia, 01 2012. http://mei.di.uminho.pt/sites/default/files/dissertacoes//eeum_di_dissertacao_pg16019.pdf (last access: 2020/01/05).

[23] J. Pearl and A. Paz. GRAPHOIDS: A graph-based logic for reasoning about relevance relations. Technical Report R-53 850038, UCLA Computer Science Department, Dec 1985.

[24] StackExchange. What does the decomposition, weak union and contraction rule mean for conditional probability and what are their proofs? <https://math.stackexchange.com/questions/855002/> (last access: 2020/01/05)

[25] R. Yamaguchi, K. Kin, S. Shimoyama, M. Hagiwara, M. Yamamoto, and J. Wang. Formalization of the conditional independence using Coq/SSReflect. Technical Reports of Mathematical Sciences, Chiba University, 29:1–40, Department of Mathematics and Informatics Graduate School of Science, Chiba University, Japan, 2016. In Japanese.

Reynald Affeldt

Reynald Affeldt received his M.S. and Ph.D. degrees in computer science from the University of Tokyo in 2001 and 2004 respectively. He

was a Research Scientist in the Graduate School of Information Science and Technology of the University of Tokyo, and joined AIST in 2005 as a Research Scientist specializing in software verification applied to secure computing. He is now a member of the Cyber Physical Security Research Center.

Jacques Garrigue

Alumnus of École Normale Supérieure in Paris. Received his D.S. degree from the University of Tokyo in 1995. He was a Research Associate at Kyoto University from 1995 to 2004, and is

now a Professor at Nagoya University. His interests are in the theory of programming languages, particularly type systems and proof of programs. He is a member of JSSST, IPSJ and ACM.

Takafumi Saikawa

Takafumi Saikawa is a Ph.D. candidate at Nagoya University with the advisory of Professor Jacques Garrigue. He is interested in applica-

tions of computer science in mathematics, and has worked on several formalization projects based on probability theory and category theory. He is also a collaborator of research projects on public health informatics in Okumura laboratory, Kitami Institute of Technology. He was a director of Peano System inc. from 2016 to 2020.