# Diversity, dependence and independence

**Pietro Galliani**[1] [iD] · **Jouko Väänänen**[2,3]

## Abstract

We propose a very general, unifying framework for the concepts of dependence and independence. For this purpose, we introduce the notion of *diversity rank*. By means of this diversity rank we identify *total determination* with the inability to create more diversity, and *independence* with the presence of maximum diversity. We show that our theory of dependence and independence covers a variety of dependence concepts, for example the seemingly unrelated concepts of linear dependence in algebra and dependence of variables in logic.

## 1 Introduction

The concepts of dependence and independence occur widely in science. The exact study of these concepts has taken place at least in four different contexts:

- **Mathematics:** Dependence and independence are fundamental concepts in algebra: linear dependence in linear algebra and algebraic dependence in field theory. In both cases independence is defined as the lack of dependence: elements $\{x_1, \ldots, x_n\}$ are *independent* if no $x_i$ is dependent on the rest. Whitney [1] and van der Waerden [2] pointed out

These authors contributed equally to this work.

✉ Pietro Galliani
  Pietro.Galliani@unibz.it

  Jouko Väänänen
  jouko.vaananen@helsinki.fi

1 Faculty of Computer Science, Free University of Bozen-Bolzano, piazza Domenicani, 3, Bozen-Bolzano, 39100, Italy

2 Department of Mathematics and Statistics, University of Helsinki, Pietari Kalmin katu 5, Helsinki, PL 68 FIN-00014, Finland

3 FNWI, ILLC, Universiteit van Amsterdam, P.O. Box 94242, Amsterdam, 1090 GE, The Netherlands

the similarity between these two notions of dependence and proposed axioms that cover both cases. Whitney suggested the name *matroid* for the general dependence structure inherent in algebra, giving rise to *matroid theory*, nowadays a branch of discrete mathematics.

- **Computer science:** *Functional* dependence [3] is a fundamental concept of database theory. The design and analysis of so called relational databases is often based on a careful study of the functional dependencies between attributes of various parts of the database. The more general *multivalued* dependencies are analogous to what we call independence relations between attributes.
- **Statistics and probability theory:** Dependence and independence of events (or random variables) is the basis of probability theory and statistical analysis of data.
- **Logic:** Dependence of a variable on another is the basic concept in quantification theory. In *Dependence Logic* [4] this concept is separated from quantification, making it possible, as in *Independence-Friendly Logic* [5], to write formulas with more complicated dependence relations between variables than what first order logic allows. Likewise *Independence Logic* [6] extends First Order Logic by an atom $\vec{x} \perp \vec{y}$ that states that the tuples of quantified variables $\vec{x}$ and $\vec{y}$ are chosen independently, in the sense that every possible choice of $\vec{x}$ and of $\vec{y}$ may occur together. These logics – and the generalization of Tarski's Semantics used for their analysis, commonly called *Team Semantics* – have lead in the last decade to a considerable amount of research regarding logics augmented by various *notions of dependence and independence*, in the first order case but also in the propositional case [7], in the modal case [8, 9], in the temporal case [10], and recently even in probabilistic cases [11–13].

In this paper, we show that these seemingly very diverse notions of dependence and independence can be captured by a simple, very general, unifying framework.

Our starting point is very general. Suppose we have a set $M$ of objects. We want to make sense of the concept that a finite subset[1] $x \subseteq_f M$ **depends** on another subset $y \subseteq_f M$, or that a subset $x \subseteq_f M$ is **independent** of another subset $y \subseteq_f M$. To accomplish this in the most general sense, we define the concept of the **diversity** of a set $x \subseteq_f M$. A small set has less diversity than a bigger set, hence our diversity function is monotone. Also, the diversity of $x$ arises from properties of the individual elements, hence our diversity function satisfies certain further conditions. The connection between diversity and dependence arises from the idea that dependence reduces diversity, and respectively independence preserves diversity. If $y$ is totally determined by $x$, then adding $y$ to $x$ does not increase the diversity of $x$ at all. On the other hand, if $x$ and $y$ are independent, then putting them together means simply adding the diversities together: nothing is lost, because there is no interaction between $x$ and $y$.

Because of the generality of our approach, according to which $M$ is just a set of objects about which we a priori know nothing, we do not *define* the diversity function explicitly, but rather give a few conditions it ought to satisfy. The point is that on the basis of these conditions we can introduce natural notions of **dependence** and **independence** with a variety of applications, in particular the ones mentioned above.

We will now give an overview of our results. In Section 2 we introduce the concept of diversity rank $\|x\|$, a non-negative real number, that the whole paper is about. It is defined in a completely general setting by means of four axioms for $\|x\|$. We show that another general and closely related concept, that of a matroid, is a special case. Thus our approach generalizes the matroid approach. We use our diversity rank to define two binary relations, namely

---

[1]In this work, we will write $x \subseteq_f M$ for "$x$ is a finite subset of $M$".

$= (x, y)$ (the dependence relation "$y$ is totally determined on $x$") and $x \perp y$ (the independence relation "$x$ and $y$ are independent"). Previously dependence and independence relations were studied under various guises in particular contexts such as team semantics, databases, algebraic structures, and statistics. In our approach these concepts are defined in a general setting which covers the special cases mentioned.

In Section 3 we give examples of diversity rank functions. The most interesting examples are relational diversity, algebraic diversity and entropy. We show that relational diversity does not satisfy submodularity, a property that instead matroids (such as the ones corresponding to algebraic diversity) necessarily satisfy. The entropy of a set of random variables, as we will see, also satisfies the submodularity property; however, it is not necessarily integer, and thus it does not correspond to a matroid either.

The main benefit of our approach as compared to the matroid approach is that we can develop a theory of diversity which covers matroids, the relational case, and the probabilistic (entropy) case. The relational case is important because it arises naturally in database theory, where our dependence is called functional dependence and our independence is called embedded multivalued dependence. Likewise, the entropy diversity rank function gives rise to the known probability-theoretic notions of independence and functional dependence between (sets of) random variables.

As was observed above, our diversity rank $\|x\|$ makes it possible to define a dependence relation $=(x, y)$. In section 4 we use the diversity rank axioms to prove the so-called Armstrong Axioms for this relation. In Theorem 1 we prove the completeness of Armstrong Axioms in our more general context. Thus the well-known completeness of Armstrong Axioms in team semantics or in the theory of database dependencies is a more general phenomenon.

In Section 5 we use diversity rank axioms to prove (a simple extension of) the Geiger-Paz-Pearl axioms for the independence relation $x \perp y$ derived from a diversity rank $\|x\|$. In Theorem 2 we prove the completeness of these axioms in our general context.

In Section 6 we address the "reverse" question, whether every binary relation $=(x, y)$ satisfying the Armstrong Axioms of dependence arises from a diversity rank function $\|x\|$. In Theorem 3 we give a positive answer in the case of countable domains. The corresponding question for the Geiger-Paz-Pearl axioms for the independence relation remains open.

## 2 Diversity rank in a general setting

We now define the concept of *diversity rank* in an entirely general setting. We use the notation $xy$ to denote the union $x \cup y$ of subsets $x$ and $y$ of a fixed set $M$. The following is the key definition of this paper:

**Definition 1** Suppose $M$ is an arbitrary set. A function $x \mapsto \|x\|$ from the finite subsets of $M$ to $\mathbb{R}^+ \cup \{0\}$ is called a *diversity rank function on $M$* if it satisfies the following conditions for all $x, y, z \subseteq_f M$:

**R1:** $\|\emptyset\| = 0$;
**R2:** $\|x\| \leq \|xy\| \leq \|x\| + \|y\|$;
**R3:** If $\|xy\| = \|x\|$ then $\|xyz\| = \|xz\|$;
**R4:** If $\|xyz\| = \|x\| + \|yz\|$ then $\|xy\| = \|x\| + \|y\|$.[2]

---

[2]Since set union is commutative, it also follows that if $\|xyz\| = \|x\| + \|yz\|$ then $\|xz\| = \|x\| + \|z\|$.

Note that the set $M$ may be infinite, but the subsets $x$, $y$, ... in the domain of the diversity rank function must be finite. This is intentional: in the case of relational diversity (Section 3.7), for example, we can have potentially infinitely many variables, but we are only talking about dependence/independence between finite sets of variables.

Intuitively, the diversity rank of $x$ is the amount of "diversity" or "variation" that $x$ contains. For example, if $x$ is a sequence of vectors in a vector space, the amount of diversity in $x$ is revealed by the dimension of the subspace spanned by $x$. If $x$ is the set of attributes in a relation schema in a given database schema, the amount of diversity in $x$ is revealed by the maximum number of different tuples (records) that may exist in the corresponding relation in a given database instance that can be considered as valid for that schema. Finally, if $x$ is simply a word in a finite alphabet, a possible measure of the amount of diversity in $x$ is the number of different letters in $x$, so that for example the diversity of "abbab" is 2 and the diversity of "abcdda" is 4.

It is obvious that we have to require **R1** and **R2**. The empty set cannot manifest any diversity, more elements means more diversity, and the amount of diversity manifested by two sets taken together is at most the sum of the amounts of diversity occurring in each of them separately.

The axioms **R3** and **R4** are less intuitive. In brief, Axiom **R3** states that if adding $y$ to $x$ does not increase the amount of diversity of $x$ (that is, $y$ is "trivial" given $x$), then adding it to $xz$ does not increase the amount of diversity of $xz$ (that is, $y$ is also "trivial" given $xz$) either; and Axiom **R4** states that if adding $yz$ to $x$ increases the amount of diversity of the maximum amount possible (that is, $yz$ is "maximally non-trivial" given $x$) then adding $y$ to $x$ also increases the amount of diversity of the maximum amount possible (that is, $y$ is also "maximally non-trivial" given $x$).

To better understand the roles of these axioms, let us briefly consider toy examples that would violate them.

For **R3** let us consider the function, defined over the two-element set $M = \{a, b\}$, such that $\|\emptyset\| = \|a\| = \|b\| = 0$ but $\|ab\| = 1$. According to this candidate diversity rank function, $\{a\}$ is "trivial" with respect to the empty set (i.e., it adds no further diversity if added to it), but it is not so with respect to the bigger set $\{b\}$: indeed, $\|\emptyset a\| = \|\emptyset\| = 0$ but $\|\emptyset ab\| = 1 > \|\emptyset b\| = 0$. The purpose of **R3** is to prevent this kind of scenario, in which $a$ contributes more diversity in the presence of some other element $b$ than in its absence.

For **R4**, instead, let us consider a function over three elements $a$, $b$ and $c$ such that $\|a\| = \|b\| = \|c\| = 1$, $\|ab\| = 1.5$ , $\|ac\| = 2$ and $\|abc\| = 3$. Here $b$ is not maximally diverse with respect to $a$, in the sense that $\|ab\| < \|a\| + \|b\|$; however, $b$ is maximally diverse with respect to $ac$, in the sense that $\|abc\| = \|ac\| + \|b\|$. This type of scenario, in some of the diversity in $b$ would appear to be subsumed by the diversity in $a$ but not by the diverstiy in $ac$, is what we aim to prevent by this rule.

**R3** and **R4** (as well as the right part of **R2**) would follow immediately if we assumed that our diversity rank function is *submodular*, in the sense that it satisfies the condition

**SUBM:**    $\|xyz\| + \|z\| \leq \|xz\| + \|yz\|$ :

**Proposition 1** *Every function from finite subsets of some set $M$ to non-negative real numbers satisfying **R1**, the left part of **R2** and **SUBM** is a diversity rank function in the sense of Definition 1.*

*Proof* Let $\|x\|$ satisfy **R1**, the left part of **R2** and **SUBM**. We need to show that $\|x\|$ also satisfies the right part of **R2** as well as **R3** and **R4**. Choosing $z = \emptyset$, we obtain immediately from submodularity that $\|xy\| + \|\emptyset\| \leq \|x\| + \|y\|$. But by **R1** we know $\|\emptyset\| = 0$, and subadditivity (that is, the right part of **R2**) follows. As for **R3**, suppose $\|xy\| = \|x\|$. Then by **SUBM**, $\|xyz\| \leq \|xy\| + \|xz\| - \|x\| = \|x\| + \|xz\| - \|x\| = \|xz\|$; but on the other hand $\|xz\| \leq \|xyz\|$ by the left part of **R2** and so $\|xyz\| = \|xz\|$ as required.

Finally, **R4** holds. Indeed, by **SUBM** we know $\|xyz\| + \|y\| \leq \|xy\| + \|yz\|$. Thus, if $\|xyz\| = \|x\| + \|yz\|$ we have immediately that $\|x\| + \|yz\| + \|y\| \leq \|xy\| + \|yz\|$, that is, $\|x\| + \|y\| \leq \|xy\|$. But $\|xy\| \leq \|x\| + \|y\|$ by the right part of **R2**, which we already proved, and hence $\|xy\| = \|x\| + \|y\|$ as required. □

The converse of the above result, however, is not true: there exist diversity rank functions that do not satisfy submodularity. As a toy example, consider the diversity rank function over the set $\{a, b, c\}$ defined as

$$\begin{aligned}
&\|\emptyset\| = 0; &\quad &\|ab\| = 2.1; \\
&\|a\| = \|b\| = 1.5; &\quad &\|ac\| = \|bc\| = 1.6; \\
&\|c\| = 1; &\quad &\|abc\| = 3.
\end{aligned}$$

Here $\|abc\| + \|c\| = 3 + 1 = 4$, but $\|ac\| + \|bc\| = 1.6 + 1.6 = 3.2$; therefore, **SUBM** fails.

Of course, one needs to verify that this is indeed a diversity rank function. **R1** holds, because $\|\emptyset\|$ is indeed 0; **R3** holds, because its premise is never satisfied non-trivially (it is never the case that $\|xy\| = \|x\|$ for $y \neq \emptyset$); **R4** likewise holds because its premise is never satisfied non-trivially; the left side of **R2** is easily verified by inspection; and as for the right part of **R2** we need to verify the following six cases:

1.  $\|ab\| = 2.1 \leq \|a\| + \|b\| = 1.5 + 1.5 = 3$;
2.  $\|ac\| = 1.6 \leq \|a\| + \|c\| = 1.5 + 1 = 2.5$;
3.  $\|bc\| = 1.6 \leq \|b\| + \|c\| = 1.5 + 1 = 2.5$;
4.  $\|abc\| = 3 \leq \|a\| + \|bc\| = 1.5 + 1.6 = 3.1$;
5.  $\|abc\| = 3 \leq \|b\| + \|ac\| = 1.5 + 1.6 = 3.1$;
6.  $\|abc\| = 3 \leq \|c\| + \|ab\| = 1 + 2.1 = 3.1$.

(The other possible cases follow from these because of the left side of rule **R2**, e.g. $\|ab\| + \|ac\| \geq \|ab\| + \|c\| \geq \|abc\|$).

It could appear tempting at this point of the work to exclude counterexamples such as this one from consideration by adding submodularity to the definition of diversity rank function. But as Proposition 3 will show, there exist diversity rank functions of practical interest that fail to satisfy **SUBM**.

A direct consequence of Proposition 1 is that our diversity rank functions generalize matroids:

**Definition 2** (Matroid) A matroid $r$ over some finite set $E$ is a function from subsets of $E$ to non-negative integers satisfying the following conditions for all subsets $x$ and $y$ of $E$.

**M1**   $r(x) \leq |x|;$[3]
**M2**   $r(xy) + r(x \cap y) \leq r(x) + r(y);$
**M3**   If $|y| = 1$ then $r(x) \leq r(xy) \leq r(x) + 1$.

---

[3] As is customary, $|x|$ here defines the cardinality of the set $x$. This is not necessarily the same as $\|x\|$, which in this work is the diversity rank of $x$ with respect to some diversity rank function $\| \cdot \|$.

The above is just one of several equivalent definitions used in the literature. It is equally possible to define a matroid in terms of its *independent sets* (that is, the $x$ such that $r(x) = |x|$), in terms of *bases* (maximal independent sets), in terms of *circuits* (minimal non-independent sets), or in terms of *closure operations*. All these definitions can be shown to give rise to the same class of mathematical objects: we refer the reader to [14] for more details.

**Corollary 2** *Let $r$ be a matroid over some finite set $E$. Then $r$ is a diversity rank function over $E$.*

*Proof* **R1** holds for $r$. Indeed, by **M1**, $r(\emptyset) \leq |\emptyset| = 0$, and therefore the only possibility is that $r(\emptyset) = 0$. The left part of **R2** holds immediately because of **M3** (this can be shown by easy induction on the the number of elements in $y \backslash x$). Moreover, **SUBM** also holds of $r$. Indeed, since $xyz = xzyz$, we have by **M2**, $r(xyz) + r(xz \cap yz) \leq r(xz) + r(yz)$. But $z \subseteq xz \cap yz$, and so by the left part of **R2** (which we already proved), $r(z) \leq r(xz \cap yz)$ and hence $r(xyz) + r(z) \leq r(xz) + r(yz)$, as required. The conclusion then follows, since because of Proposition 1 the right part of **R2** as well as **R3** and **R4** are also true of $r$. $\square$

Given a notion of diversity, it is easy to define dependence and independence in terms of *minimal* and *maximal* diversity contributions:

**Definition 3** Suppose $M$ is a set and $\| \cdot \|$ a diversity rank function on $M$. We can now define *dependence* relations between finite subsets of $M$ (with respect to $\| \cdot \|$) as follows:

- **Dependence:** $y$ *is totally determined by* (or *depends on*) $x$, in symbols $=(x, y)$, if and only if $\|xy\| = \|x\|$.
- **Constancy:** $x$ *is constant*, in symbols $=(x)$, if and only if $\|x\| = 0$.
- **Independence:** $x$ and $y$ are *independent*, in symbols $x \perp y$, if and only if $\|x\| + \|y\| = \|xy\|$.

The idea is that $=(x, y)$ holds under a diversity rank function if the amount of diversity inherent in $x$ in terms of the rank function does not increase when $y$ is added. Simply put, $x$ determines $y$, so no new diversity occurs. $=(x)$, on the other hand, holds if $x$ has no diversity at all; and $x \perp y$ holds if the diversity inherent in $x$ is so unrelated to the diversity inherent in $y$ that when the two are put together into $xy$, the diversity is the sum of the diversity of $x$ and the diversity of $y$: no loss of diversity occurs because there is—intuitively—no connection between $x$ and $y$.

## 3 Examples

Let us now consider some examples of our definitions, in order to get a better feel of their applicability and consequences. As we will see, aside from the example of Section 3.7 (relational diversity), all our examples also satisfy the submodularity condition. Nonetheless, it is worth emphasizing that, when the ranks are not necessarily integer, these diversity rank functions are not matroids.

There are various diversity ranks, also called diversity measures or diversity indexes, in biology, such as Shannon index, Simpson index, Renyi index, $\alpha$-diversity, $\beta$-diversity or

$\gamma$-diversity, each attempting to describe in terms of a single real number the number of different species in an area, the number of individuals of each species, and other parameters relevant in characterizing diversity of a biological system [15–17]. Similar diversity ranks have been introduced in other areas, as well, and usually take the concept of entropy (see subsection 3.6 below) as a starting point.

## 3.1 Constant diversity

One extreme case is the constant rank for which $\|\emptyset\| = 0$ and, for some number $c$, $\|x\| = c$ for all $x \subseteq_f M$ with $x \neq \emptyset$. If $c = 0$, there is no diversity: every set depends on every other set and is also independent of every other set. If $c \neq 0$, then every set $y$ is still dependent on any non-empty set $x$, because $\|xy\| = c = \|x\|$, and every set $x$ is still independent from the empty set $\emptyset$, because $\|x\emptyset\| = \|x\| = \|x\| + \|\emptyset\|$; but two non-empty sets $x$ and $y$ are not independent, because $\|xy\| = c \neq c + c = \|x\| + \|y\|$.

Constant diversity is trivially submodular: if $z$ is empty then $\|xyz\| + \|z\| = 2c = \|xz\| + \|yz\|$, and if instead $z$ is empty then submodularity reduces to $\|xy\| \leq \|x\| + \|y\|$ (which is part of rule **R2**, and which is also easily verified for this diversity function).

## 3.2 Singular diversity

Let $a_0 \in M$ be fixed. Let
$$\|x\| = \begin{cases} 1, & \text{if } a_0 \in x, \\ 0, & \text{otherwise.} \end{cases}$$
In this case a selected element $a_0$ is the only source of "diversity" there is. A set has diversity 1 if and only if it contains $a_0$. In this case $y$ depends on $x$ if
$$a_0 \in y \rightarrow a_0 \in x$$
and two sets $x$ and $y$ are independent if at most one of them contains $a_0$. So dependence reduces in this case to implication and independence to the Sheffer stroke (also known as NAND).

Again, this diversity rank function is submodular: if $a_0 \in z$ then $\|xyz\| + \|z\| = 2 = \|xz\| + \|yz\|$, and if $a_0 \notin z$ then $\|xyz\| + \|z\| = \|xy\| \leq \|x\| + \|y\| = \|xz\| + \|yz\|$.

Note that the notion of dependence arising from singular diversity is not symmetric: in particular, the empty set depends on $\{a_0\}$ but $\{a_0\}$ does not depend on the empty set.

## 3.3 Two-valued diversity

Suppose $\|\{a\}\|$ is either 0 or 1 for all $a \in M$. Then, by Rule **R2**, it follows that
$$\|x\| = \max\{\|\{a\}\| : a \in x\}$$
for all $x \subseteq_f M$. Hence, it suffices in this case to declare which elements have diversity 1; a set has diversity 1 if and only one if it contains one of those elements.[4]

In this case diversity is an on/off phenomenon, either it exists (1) or it does not (0), and a set has diversity if it includes some singleton that has it. In an extreme case $\|\{a\}\| = 0$ for *all* singletons $\{a\}$, $a \in M$, and we have zero constant diversity: every set has diversity 0. In

---

[4]Hence, a singular diversity is a special case of a two-valued diversity where only one element is declared to have diversity 1.

another extreme case $\|\{a\}\| = 1$ for *all* singletons $\{a\}$, $a \in M$, and we are again in constant diversity: every non-empty set has diversity 1.

According to this diversity notion, $=(x, y)$ if and only if $\exists a \in x$ such that $\|\{a\}\| = 1 \Rightarrow \exists b \in y$ such that $\|\{b\}\| = 1$; and $x$ is independent from $y$ if and only if at most one of $x$ and $y$ contain an element $c$ with $\|\{c\}\| = 1$.

Submodularity is, once more, straightforwardly verified. $\|xyz\| + \|z\|$ is 2 if and only if $\|z\| = 1$, and in that case it is easy to see that $\|xz\| + \|yz\|$ is also 2. $\|xyz\| + \|z\|$ is 1 if $\|xy\| = 1$, in which case at least one of $\|x\|$ and $\|y\|$ is 1 and hence $\|xz\| + \|yz\| \geq 1$; and finally, if $\|xyz\| + \|z\| = 0$ then it is trivially true that $\|xyz\| + \|z\| \geq \|xz\| + \|yz\|$.

## 3.4 Uniform diversity

Suppose

$$\|x\| = |x|.$$

This is the choice of taking the cardinality of the (finite) set $x \subseteq_f M$ as the measure of its diversity. Dependence means inclusion: $y$ is totally determined by $x$ if and only if $|xy| = |x|$, that is, if and only if $y \subseteq x$. Independence is disjointness: $x$ and $y$ are independent if and only if $|xy| = |x| + |y|$, that is, if and only if $x \cap y = \emptyset$. Once more, submodularity is easily verified by observing that $|xyz| = |xz| + |yz| - |(xz) \cap (yz)|$ and that $|(xz) \cap (yz)| > |z|$.

**Remark** Uniform diversity shows that if $M$ has at least three elements $a, b, c$ then independence is not necessarily equivalent to the failure of dependence both ways. Indeed, $\{a, b\}$ is not dependent on $\{b, c\}$ or vice versa, since neither set is contained in the other, but $\{a, b\}$ and $\{b, c\}$ are not independent either since they are not disjoint. This is not surprising: in our framework, $y$ depends on $x$ if adding $y$ to $x$ contributes no diversity whatsoever to it, while $x$ and $y$ are independent if adding $y$ to $x$ contributes the maximal amount $\|y\|$ of diversity to it. But it is certainly possible for neither extreme to be the case.

## 3.5 Coverage diversity

Suppose $U$ is a finite set and choose $A_a \subseteq U$ for each $a \in M$. For $a_1 \ldots a_n \in M$, let

$$\|\{a_1, \ldots, a_n\}\| = |A_{a_1} \cup \ldots \cup A_{a_n}|.$$

We can think of each $A_a$ as "data", about the element $a$ of $M$. The more data we have the more diversity we give to the element, and the diversity of a set is obtained by simply putting together all the data we have. In this simple example the data is not thought to be specific to the elements of $M$, so the data about different elements is just lumped together. For example, if $a$ and $b$ are two botanic genera in the bean family, the diversity of $\{a, b\}$ in a set $U$ of data about species (e.g. in some location) is obtained by counting how many different species of the two genera there are in $U$.

According to this diversity notion, $y$ is dependent on $x$ if and only if $\bigcup\{A_a : a \in y\} \subseteq \bigcup\{A_b : b \in x\}$, that is, every data point corresponding to some element of $y$ also corresponds to some element of $x$; and $y$ is independent from $x$ if and only if $\bigcup\{A_a : a \in y\} \cap \bigcup\{A_b : b \in x\} = \emptyset$, that is, if no data point corresponds to some element of $x$ *and* to some element of $y$.

Submodularity can be verified much as in the previous example: indeed, $\|xyz\| = |\bigcup A_a : a \in x \cup y \cup z| = \|xy\| + \|xz\| - |C|$ for $C = \{u : u \in A_b \cap A_c$ for some $b \in x \cup z, c \in y \cup z\}$, and it is easy to see that $|C| \geq \|z\|$.

## 3.6 Entropy

Let us think of the elements of $M$ as *discrete random variables* $v_1, v_2, \ldots$ over some probability space and with outcomes in some finite set $A$.[5] Then for any $x = \{v_1 \ldots v_k\} \subseteq_f M$ we can define $\|x\|$ as the joint entropy [18] $H(x)$ of $v_1 \ldots v_k$, that is, as[6]

$$- \sum_{(m_1 \ldots m_k) \in A^k} P(v_1 \ldots v_k = m_1 \ldots m_k) \log P(v_1 \ldots v_k = m_1 \ldots m_k).$$

This definition clearly satisfies rule **R0**, since the entropy of the only possible distribution over the empty space is zero; moreover, it is not hard to convince oneself that it is monotone and submodular. In brief, this can be shown by considering the *conditional entropy* $H(y \mid x) = H(xy) - H(x)$.

Indeed, it can be proved (see any Information Theory textbook, for instance Theorem 2.2.1 of [19]) that the conditional entropy $H(y \mid x)$ is always non-negative[7], from which we have the left part of **R2**; furthermore (see e.g. Theorem 2.6.5 of [19])[8] $H(x \mid yz) \leq H(x \mid z)$, from which we obtain immediately $H(xyz) - H(yz) \leq H(xz) - H(z)$, that is, Axiom **SUBM**.

From Proposition 1, we can immediately conclude that entropy is an example of a (submodular) diversity rank function. Here, $y$ depends on $x$ according to the entropy diversity rank if and only if $H(xy) = H(x)$, that is, if and only if the relative entropy of $y$ given $x$ is 0, or in other words if the value of $y$ is completely determined by the value of $x$; and $x$ and $y$ are independent according to this rank if and only if they are independent sets of random variables, that is, $P(x = m, y = m') = P(x = m)P(y = m')$ for all possible choices of values $a$ and $b$ for $x$ and $y$. We observe that this notion of probabilistic independence of random variables is exactly the one axiomatized by Geiger, Paz and Pearl in [20]. Their axiomatization will be the base for our axiomatization of independence in our more general setting in Section 5.

## 3.7 Relational diversity

Suppose $X$ is a nonempty, finite set of *variable assignments $s$* from a set $V$ of variables to a set $A$ of elements (in the language of Dependence and Independence Logic, $X$ is said to be a *team* over $A$ with domain $V$; and it is not hard to see that it is equivalent to a relational table in which every variable indicates a different column).[9] Given some $x = \{v_1 \ldots v_n\} \subseteq_f V$, let

$$\|x\| = \log(\#\mathrm{rows}_X(v_1 \ldots v_n)).$$

---

[5]Nothing in this example hinges on $A$ being the same for all $v \in M$, but we will assume so for simplicity.

[6]In this work, log will always represent the base-2 logarithm.

[7]More precisely, this theorem shows that $H(xy) - H(x) = -\sum_m P(x = m) \sum_{m'} P(y = m' \mid x = m) \log P(y = m' \mid x = m)$, and the right hand side is straightforwardly seen to be non-negative.

[8]Strictly speaking, this theorem states that $H(x) - H(x \mid y) \geq 0$, but if we consider the above inequality with respect to distributions already conditioned on $z$ the result follows.

[9]In general, in Dependence and Independence Logic teams do not necessarily have to be finite, but we will focus on the finite case in this example.

where

$$\#\mathrm{rows}_X(v_1 \ldots v_n) = |\{(s(v_1), \ldots, s(v_n)) : s \in X\}|$$

is the number of different values that $x = v_1 \ldots v_n$ takes in $X$.[10]

We can think of each $s \in X$ as an "observation", or "data", about the possible values that the variables in $V$ can take. The more different observations we have the more diversity we give to the element. Note the difference with coverage diversity, where the data was not specific to the element of $A$. Here what matters is the relationships of the different observations to each other. Thus

$$\|\{v\}\| = \log |\{s(v) : s \in X\}|,$$

that is, the diversity rank of a single element $v$ of $V$ is the (logarithm of the) number of different observations about $v$. The diversity of a pair $\{v, w\}$ is the (logarithm of the) number of different combinations of observations of $v$ and $w$. For example, if $v$ and $w$ are two genera, the diversity of $\{v, w\}$ in a set $X$ of observations is calculated by counting how many different pairs of observations of a specimen of $v$ and a specimen of $w$ there are in $X$.

The presence of the logarithm operator in this definition may appear at first sight somewhat outlandish, but it is in fact quite natural. Indeed, one may recall that, in information theory, the *information content* of an event is defined as the negative logarithm of its probability, and the information content of a random variable (i.e. its entropy) is the expectation of the information content of it taking all possible values [18, 21]. Therefore, if we associate a relation $R$ with the probability distribution selecting any tuple in $R$ with equal probability $1/|R|$, the information content of this distribution (and, hence, of the relation) is precisely $-\log(1/|R|) = \log(|R|)$.

The dependence relation arising from the relational diversity rank is the usual functional dependence relation of database theory and dependence logic [3, 4, 22]. Why? By definition, $=(x, y)$ if and only if $\log(\#\mathrm{rows}_X(xy)) = \log(\#\mathrm{rows}_X(x))$, that is, if and only if $\#\mathrm{rows}_X(xy) = \#\mathrm{rows}_X(x)$. This can be the case if and only if any two $s, s' \in X$ which differ with respect to $xy$ differ already on $x$ alone, or, by contraposition, if and only if any two $s, s' \in X$ which are the same with respect to $x$ are also the same with respect to $y$. This is precisely the usual notion of database-theoretic functional dependence [3], which is arguably the most important (although by no means the only) notion of dependence studied in the context of database theory.

The independence relation arising from the relational diversity rank is also the independence relation of Independence Logic [6]. Indeed, by definition, $x \perp y$ if and only if $\#\mathrm{rows}_X(xy) = \#\mathrm{rows}_X(x) \cdot \#\mathrm{rows}_X(y)$. Then, enumerating the elements of $x$ as $(v_1 \ldots v_k)$ and the elements of $y$ as $(w_1 \ldots w_t)$, it is always the case that $\#\mathrm{rows}_X(xy) = |\{(s(v_1) \ldots s(v_k), s(w_1) \ldots s(w_t)) : s \in X\}| \leq |\{(s(v_1) \ldots s(v_k)) : s \in X\}| \cdot |\{(s'(w_1) \ldots s'(w_t)) : s' \in X\}| = \#\mathrm{rows}_X(x) \cdot \#\mathrm{rows}_X(y)$, since every possible value for $(v_1 \ldots v_k, w_1 \ldots w_t)$ in $X$ corresponds to one possible value for $(v_1 \ldots v_k)$ and one possible value for $(w_1 \ldots w_k)$ in $X$. Equality holds if and only if the converse also holds, i.e., if and only if for any two $s, s' \in X$ there exists some $s'' \in X$ such that $s''(v_1 \ldots v_k) = s(v_1 \ldots v_k)$ and $s''(w_1 \ldots w_t) = s'(w_1 \ldots w_t)$. In particular, this implies at once that if $x \perp y$ the variables occurring in both $x$ and $y$ must take only one value in $X$, and that if $x$ and $y$ are disjoint

---

[10]It is easy to check that this value does not depend on the ordering of $v_1 \ldots v_n$.

and $x \perp y$ then the projection of $X$ onto $xy$ must be the Cartesian product of its projections onto $x$ and $y$.[11]

It may be instructive to verify that the relational diversity notion of rank satisfies our axioms:

**R1:**   Since #rows($\emptyset$) = $|\{()\}|$ = 1 for any choice of $X$, where () represents the empty tuple, we have $\|\emptyset\| = 0$, as required.

**R2:**   Since #rows($x$) $\leq$ #rows($xy$) and the logarithm is a monotone function, we have immediately that $\|x\| \leq \|xy\|$; and since #rows($xy$) $\leq$ #rows($x$) · #rows($y$), we have immediately that $\|xy\| \leq \|x\| + \|y\|$.

**R3:**   If $\|xy\| = \|x\|$, #rows($xy$) = #rows($x$) and hence every possible value of $x$ occurs together with only one possible value of $y$. But then every possible value of $xz$ occurs together with only one possible value of $y$, and hence #rows($xyz$) = #rows($xz$) and $\|xyz\| = \|xz\|$;

**R4:**   If $\|xyz\| = \|x\| + \|yz\|$, it must be the case that #rows($xyz$) = #rows($x$)·#rows($yz$), and hence that every possible value of $x$ occurs together with every possible value for $yz$. But then in particular every possible value for $x$ occurs together with every possible value for $y$, and so #rows($xy$) = #rows($x$) · #rows($y$) and $\|xy\| = \|x\| + \|y\|$.

In contrast to our other examples, this notion of relational diversity is *not* submodular, as the following counterexample, which we owe to Tong Wang[12] shows:

**Proposition 3** *Relational diversity fails to satisfy* ***SUBM****.*

*Proof* Consider the relation $X$

$$
\begin{array}{ccc}
v_1 & v_2 & v_3 \\
1 & 1 & 1 \\
1 & 1 & 2 \\
2 & 1 & 1 \\
1 & 2 & 1 \\
2 & 1 & 2 \\
\end{array}
$$

Then #rows$_X(v_1 v_2 v_3)$ = 5, #rows$_X(v_2)$ = 2, and #rows$_X(v_1 v_2)$ = #rows$_X(v_2 v_3)$ = 3. Thus, #rows$_X(v_1 v_2 v_3)$ · #rows$_X(v_2)$ = 10 > 9 = #rows$_X(v_1 v_2)$ · #rows$_X(v_2 v_3)$, and hence $\|v_1 v_2 v_3\| + \|v_2\| > \|v_1 v_2\| + \|v_2 v_3\|$. $\square$

## 3.8 Algebraic diversity

Suppose that $V$ is a vector space and that $h$ maps $M$ into $V$. We obtain a diversity rank function by letting for finite $x \subseteq M$:[13]

$$\|x\| = \text{ the dimension of the subspace generated by } \{h(a) : a \in x\}.$$

---

[11] Note that this is not the same as saying that the projection of $X$ onto $xy$ is the *Natural Join* of the projections of $X$ onto $x$ and onto $y$. For example, consider the relation $X = \{(0, 1, 0), (1, 2, 1)\}$ over three variables named $v_1, v_2, v_3$ respectively.

Then the natural join of the projections of $X$ onto $v_1 v_3$ and $v_2 v_3$ is its projection onto $v_1 v_2 v_3$; however, $v_1 v_3 \perp v_2 v_3$ is not the case (indeed, #rows($v_1 v_3$) · #rows($v_2 v_3$) = 2 · 2 = 4 but #rows($v_1 v_2 v_3$) = 2).

[12] Personal communication.

[13] We refer the reader to any algebra textbook, for example to [23], for the relevant algebra background.

Submodularity **SUBM** follows from the well known fact that if $U$ and $V$ are vector subspaces,

$$\dim(U \cup V) = \dim(U) + \dim(V) - \dim(U \cap V).$$

In this context it is not hard to verify that $V$ is dependent on $U$ if and only if $\dim(U \cup V) = \dim(U)$, that is, if and only if every vector of $V$ is a linear combination of vectors in $U$; and that, on the other hand, $U$ and $V$ are independent if and only if $\dim(U \cup V) = \dim(U) + \dim(V)$, that is, if and only if the subspace generated by $U$ and the subspace generated by $V$ do not share a nonzero vector.

Likewise if $F$ is a field, we get a diversity rank function by letting for finite $x \subseteq M$ and letting $h$ map $M$ into $F$ instead:

$$\|x\| = \text{ the transcendence degree of the subfield generated by } \{h(a) : a \in x\}.$$

This gives rise to the concepts of algebraic dependence and independence.

As mentioned in the Introduction, this notion of rank defines a matroid (in fact, it was one of the original motivations for the development of Matroid Theory); and thus, by Corollary 2, it is also a diversity rank function according to our definition.

## 4 From diversity to dependence

Given a diversity function $\| \cdot \|$, we have defined *dependence* $=(x, y)$ of $y$ on $x$ by $\|xy\| = \|x\|$.

It is easy to verify that dependence satisfies the following axioms:

**Proposition 4** *Dependence satisfies the following properties:*

*1. Reflexivity:*    $=(xy, x)$.
*2. Augmentation:*    $=(x, y)$ *implies* $=(xz, yz)$.
*3. Transitivity:*    *If* $=(x, y)$ *and* $=(y, z)$, *then* $=(x, z)$.

*Proof* Reflexivity:    Clearly $\|xyx\| = \|xy\|$. Therefore, $=(xy, x)$.
Augmentation:    Suppose $\|xy\| = \|x\|$. Then, by **R3**, $\|xyz\| = \|xz\|$; and therefore, $\|xzyz\| = \|xz\|$, or, in other words, $=(xz, yz)$.
Transitivity:    Suppose $\|x\| = \|xy\|$ and $\|y\| = \|yz\|$. Again, by **R3**, from $\|x\| = \|xy\|$ we get $\|xz\| = \|xyz\|$; and similarly, from $\|y\| = \|yz\|$ we get $\|xy\| = \|xyz\|$. By the transitivity of equality, we can conclude $\|xz\| = \|xy\|$. But we have as an hypothesis $\|xy\| = \|x\|$, and therefore we can conclude $\|x\| = \|xz\|$, or, in other words, $=(x, z)$. □

We can use the above rules as the axioms of a proof system for inferring the consequences of a set of dependence assertions. More precisely, given a set $\Sigma$ of dependence assertions and a dependence assertion $=(x, y)$ for $x, y \subseteq_f M$, we will write that $\Sigma \vdash_D =(x, y)$ if it is possible to derive $=(x, y)$ from $\Sigma$ through applications of the rules of Reflexivity, Augmentation and Transitivity.

It follows from these axioms that a dependency notion is entirely defined even if we only consider singletons on the right-hand side of it:

**Proposition 5** *Let $\Sigma$ be a set of assertions of the form $=(z, w)$ for $z, w \subseteq_f M$, and let also $x, y \subseteq_f M$. Then $\Sigma \vdash_D =(x, y)$ if and only if $\Sigma \vdash_D =(x, \{a\})$ for all $a \in y$.*

*Proof* By Reflexivity, if $a \in y$ then it is always the case that $\Sigma \vdash_D =(y, \{a\})$. If $\Sigma \vdash_D =(x, y)$, by Transitivity it is thus the case that $\Sigma \vdash_D =(x, \{a\})$ for all such $a$. Conversely, suppose $\Sigma \vdash_D =(x, \{a\})$ for all $a \in y$. Then, in order to reach our conclusion that $\Sigma \vdash_D =(x, y)$, it suffices to verify that whenever $\Sigma \vdash_D =(x, y_1)$ and $\Sigma \vdash_D =(x, y_2)$ it is also the case that $\Sigma \vdash_D =(x, y_1 y_2)$. This is easily shown: if $\Sigma \vdash_D =(x, y_1)$, by Augmentation we have $\Sigma \vdash_D =(x, xy_1)$ (remember that in our notation $xx = x \cup x = x$), and if $\Sigma \vdash_D =(x, y_2)$ again by Augmentation we have $\Sigma \vdash_D =(xy_1, y_1 y_2)$, and an application of Transitivity gives us $\Sigma \vdash_D =(x, y_1 y_2)$. The conclusion follows immediately. $\square$

The following is essentially proved in [3], albeit in the special case of relations and functional dependencies:

**Theorem 6** (Completeness of the Dependence Axioms) *Let $\Sigma$ be a set of assertions of the form $=(z, w)$, where all $z$ and $w$ are finite subsets of some set $M$ and let also $x, y \subseteq_f M$. The following properties are equivalent:*

1.  *$=(x, y)$ holds under every diversity rank function on $M$ under which $\Sigma$ holds.*
2.  *$=(x, y)$ holds under every two-valued diversity rank function $\mathbf{P}(M) \rightarrow \{0, 1\}$ under which $\Sigma$ holds.*
3.  *$=(x, y)$ holds under every relational diversity rank function under which $\Sigma$ holds.*
4.  *$=(x, y)$ follows from $\Sigma$ by the rules of Proposition 4.*

*Proof* Trivially, (1) implies (2) and (3). Furthermore, (4) implies (1) by Proposition 4. We demonstrate that (2) implies (4) and (3) implies (4). Let us first assume (2). Suppose $=(x, y)$ does not follow from $\Sigma$ by the rules of Proposition 4. Let $V$ be the set of $a \in M$ such that $=(x, \{a\})$ follows from $\Sigma$ by these rules. By Proposition 5, for all $w \subseteq M$, we have that $\Sigma \vdash_D =(x, w)$ if and only if $w \subseteq V$.

Let $W$ be all the remaining elements of $M$. Since $y \nsubseteq V$, $W \neq \emptyset$. Let us define a diversity rank function on $M$ by letting for $a \in M$:

$$\|\{a\}\| = \begin{cases} 0, & \text{if } a \in V \\ 1, & \text{if } a \in W, \end{cases}$$

and otherwise

$$\|\{a_1, \ldots, a_n\}\| = \max\{\|\{a_1\}\|, \ldots, \|\{a_n\}\|\}.$$

Note that $\|xy\| = 1$, while $\|x\| = 0$. Thus the relation $=(x, y)$ does not hold under this diversity rank function. Suppose then $=(z, w) \in \Sigma$. If $z \subseteq V$, this means $\Sigma \vdash_D =(x, z)$; and then, by Transitivity, $\Sigma \vdash_D =(x, w)$ and so $w \subseteq V$ as well. So $\|zw\| = \|z\| = 0$ and $=(z, w)$ holds. On the other hand, if $z \nsubseteq V$, then $\|z\| = 1$. So $\|zw\| = \|z\| = 1$, whence $=(z, w)$ holds again.

Let us then assume (3). We proceed as above. Let $X$ consist of the two functions $\{s_1, s_2\}$, where $s_1(a) = 0$ for all $a \in M$, $s_2(a) = 0$ for $a \in V$ and $s_2(a) = 1$ for $a \in W$. We obtain the same rank as above, so we are done. $\square$

Since – as we saw – functional dependence is exactly the dependency notion generated by the relational diversity rank function, we obtain:

**Corollary 7** (Armstrong) *A functional dependence follows semantically, for all relations, from a given set of functional dependencies if and only if it follows by the rules of Proposition 4.*

Theorem 6 shows that Armstrong's completeness theorem for functional dependence is actually a more general completeness theorem of dependence relations arising from diversity ranks.

## 5 From diversity to independence

We shall now study the properties of the notions of independence arising from our diversity ranks. Let us recall that, according to our definition, $x$ and $y$ are independent ($x \perp y$) if and only if $\|xy\| = \|x\| + \|y\|$.

Let us begin by observing that, by our definition, $x \perp x$ if and only if $\|x\| = \|x\| + \|x\|$, that is, if and only if $\|x\| = 0$: $x$ is independent of itself if and only if $x$ contains no diversity whatsoever according to our diversity rank function, that is, if and only if $x$ is *constant* in the sense of Definition 3.

Note that in the probabilistic case, a random variable $v$ is independent of itself, in the sense that $P(v = a \text{ and } v = b) = P(v = a)P(v = b)$ for all possible values $a, b$ of $X$, if and only if $v$ may take only one value with probability 1, and hence it is constant in the ordinary sense of the word; and similarly, in the logical/relational case a variable $v$ is independent on itself if and only if it takes only possible value for all variable assignments in the set being considered. Thus, our use of the term "constancy" is not unmotivated.

**Proposition 8** *Independence satisfies the following properties:*

*1. Empty Set:*    $x \perp \emptyset$.
*2. Symmetry:*    *If $x \perp y$, then $y \perp x$.*
*3. Decomposition:*    *If $x \perp yz$, then $x \perp y$.*[14]
*4. Mixing:*    *If $x \perp y$ and $xy \perp z$, then $x \perp yz$.*
*5. Constancy:*    *If $z \perp z$ then $z \perp x$.*[15]

*Proof* Let us prove that these axioms follow from our notion of independence:

Empty Set:    Since $\|\emptyset\| = 0$, $\|x\| + \|\emptyset\| = \|x\| + 0 = \|x\| = \|x\emptyset\|$.
Symmetry:    This follows easily from the commutativity of sum and union. If $\|x\| + \|y\| = \|xy\|$ then $\|y\| + \|x\| = \|xy\| = \|yx\|$.
Decomposition:    Suppose $x \perp yz$, that is, $\|x\| + \|yz\| = \|xyz\|$.
    By **R4**, we then have $\|xy\| = \|x\| + \|y\|$ and $x \perp y$.
Mixing:    Suppose $\|xy\| = \|x\| + \|y\|$ and $\|xyz\| = \|xy\| + \|z\|$. We need to prove $\|xyz\| = \|x\| + \|yz\|$.
    We begin by observing that $\|x\| + \|y\| + \|z\| = \|xy\| + \|z\| = \|xyz\|$. But by **R2** $\|yz\| \leq \|y\| + \|z\|$, and therefore $\|x\| + \|yz\| \leq \|x\| + \|y\| + \|z\| = \|xyz\|$.

---

[14]By the symmetry of union, the Decomposition rule implies that if $x \perp yz$ then $x \perp z$ as well.
[15]If one is uninterested in independence assertions $x \perp y$ in which $x$ and $y$ overlap, this axiom can be removed. Our proof of Theorem 11 then reduces essentially to the proof in [20].

On the other hand, again by **R2**, $\|xyz\| \leq \|x\| + \|yz\|$, and so in conclusion $\|xyz\| = \|x\| + \|yz\|$, as required.

Constancy:  If $z \perp z$ then $\|z\| = \|z\| + \|z\|$, and hence $\|z\| = 0$. But then by **R2** $\|x\| \leq \|xz\| \leq \|x\| + \|z\| = \|x\|$, and thus $\|xz\| = \|x\| + \|z\|$ and $z \perp x$.

$\square$

Given a set $\Sigma$ of independence assertions and an independence assertion $x \perp y$ for $x, y \subseteq_f M$, we will write that $\Sigma \vdash_I x \perp y$ if it is possible to derive $x \perp y$ from $\Sigma$ through applications of the rules of Empty Set, Symmetry, Decomposition, Mixing and Constancy.

The following derived rule will be useful:

**Proposition 9** (Constancy Augmentation) *Given a set $M$, let $\Sigma$ be a set of independence assertions over $M$ of the form $z \perp w$ for $z, w \subseteq_f M$, and suppose $\Sigma \vdash_I u \perp u$ and $\Sigma \vdash_I x \perp y$. Then $\Sigma \vdash_I xu \perp y$*

*Proof* By Constancy, if $\Sigma \vdash_I u \perp u$ then $\Sigma \vdash_I u \perp xy$, and so by Symmetry $\Sigma \vdash_I xy \perp u$. If furthermore $\Sigma \vdash_I x \perp y$, by Symmetry $\Sigma \vdash_I y \perp x$; and thus, by Mixing, $\Sigma \vdash_I y \perp xu$, and by Symmetry once more $\Sigma \vdash_I xu \perp y$ as required. $\square$

In [20], a sound and complete axiomatization for independence of tuples of *random variables* (as derived from the definition of entropy in Section 3.6) was found, with the additional requirement that the left- and right-hand sides of the independence assertion are disjoint.

Theorem 11 below is a generalization of that result to the case of general diversity rank functions, and without that additional requirement.

First, we will show that that the axioms given above are complete for assertions of the form $\{a\} \perp \{a\}$:

**Lemma 10** (Completeness of Independence Axioms wrt Constancy Assertions) *Given a set $M$, let $\Sigma$ be a set of independence assertions over $M$, and let $a \in M$. Then the following properties are equivalent:*

1. *$\{a\} \perp \{a\}$ holds under every diversity rank function on $M$ under which $\Sigma$ holds.*
2. *$\{a\} \perp \{a\}$ holds under every relational diversity rank function under which $\Sigma$ holds.*
3. *$\{a\} \perp \{a\}$ follows from $\Sigma$ by the rules of Proposition 8.*

*Proof* Trivially (1) implies (2) and (3) implies (1). Let us verify that (2) implies (3). Suppose that $\{a\} \perp \{a\}$ does not follow from $\Sigma$ by the above rules. Then let $V$ contain all $b \in M$ such that $\Sigma \vdash_I \{b\} \perp \{b\}$ and let $S$ be a team with domain $M$ over $\{0, 1\}$ (that is, a set of functions from $M$ to $\{0, 1\}$) that contains all $s : M \to \{0, 1\}$ such that $s(b) = 0$ for all $b \in V$.

Now let $\| \cdot \| = \log(\#\mathrm{rows}_S(\cdot))$ be the relational diversity rank function induced by $S$: as already discussed, such a diversity rank function satisfies an independence assertion $z \perp w$ if and only if any possible values of $z$ and $w$ in $S$ may occur together, or, in other words, if and only if for all $s, s' \in S$ there exists some $s'' \in S$ that agrees with $s$ on $z$ and with $s'$ on $w$. In particular, for the $S$ given, this means that $z \perp w$ is satisfied if and only if $z \cap w \subseteq V$. Thus, $S$ does not satisfy $\{a\} \perp \{a\}$, since by assumption $a \notin V$.

On the other hand, $S$ satisfies all assertions of $\Sigma$. Indeed, consider any $z \perp w \in \Sigma$. By Decomposition and Symmetry, every element of $c \in M$ which is in both $z$ and $w$ is such

that $\Sigma \vdash_I \{c\} \perp \{c\}$, that is, such that $c \in V$. Thus, $z \cap w \subseteq V$ and therefore $z \perp w$ is satisfied by (the relational diversity rank function corresponding to) $S$, as required.

In conclusion, from the assumption that $\{a\} \perp \{a\}$ does not follow from $\Sigma$ according to the rules we were able to find a relational diversity rank function that satisfies $\Sigma$ but not $\{a\} \perp \{a\}$. Thus (2) implies (3), and this concludes the proof. □

Now we can generalize the completeness result of Lemma 10 to arbitrary independence assertions. The proof is an adaptation of the proof of [20] to our framework.

**Theorem 11** (Completeness of the Independence Axioms) *Let M be a set and let $\Sigma$ be a set of independence assertions $z \perp w$ for $z, w \subseteq_f M$. Then the following conditions are equivalent for any $x, y \subseteq_f M$:*

1.  *$x \perp y$ holds under every diversity rank function on M under which $\Sigma$ holds.*
2.  *$x \perp y$ holds under every relational diversity rank function under which $\Sigma$ holds.*
3.  *$x \perp y$ follows from $\Sigma$ by the rules of Proposition 8.*

*Proof* We adapt the proof of [20] to our framework. Trivially (1) implies (2). Furthermore, (3) implies (1) by Proposition 8. We shall prove that (2) implies (3) by contraposition: supposing that $\Sigma \nvdash_I x \perp y$, we shall construct a relational diversity rank function that does not satisfy $x \perp y$ but satisfies $\Sigma$.

Thus, suppose that $\Sigma \nvdash_I x \perp y$. Without loss of generality, we can assume that $\Sigma$ is closed under the rules. Note that $x \neq \emptyset$ and $y \neq \emptyset$, since otherwise $x \perp y$ would follow from $\Sigma$ by the Empty Set and Symmetry rules. We can also assume that $x$ and $y$ are *minimal*, in the sense that if $x' \subseteq x$ and $y' \subseteq y$ and at least one containment is proper then $\Sigma \vdash_I x' \perp y'$ (if $x$ and $y$ are not minimal in this sense, we can replace them with minimal subsets $x_{\min} \subseteq x$ and $y_{\min} \subseteq y$ such that $\Sigma \nvdash_I x_{\min} \perp y_{\min}$: then any diversity rank function that does not satisfy $x_{\min} \perp y_{\min}$ will not satisfy $x \perp y$ either, by the soundness of the Decomposition rule).

If $x = y = c$ for some $c \in M$, the existence of a relational diversity rank function that satisfies $\Sigma$ but not $c \perp c$ follows immediately from Lemma 10 and the proof is concluded.

Let us suppose instead that this is not the case, and let $V = \{c \in M : \Sigma \vdash_I c \perp c\}$. Notice that $x \cap V = \emptyset$, because if $x = x'c$ for $c \in V$ then by Constancy Augmentation we could derive $x \perp y$ from $x' \perp y$ (which follows from $\Sigma$ because of our minimality assumption) and $c \perp c$. Similarly, $y \cap V = \emptyset$.

Next, we show that $x \cap y \subseteq V$: since we already proved that $x \cap V = y \cap V = \emptyset$, this will imply at once that $x \cap y = \emptyset$. Let $c \in M$ is such that $c \in x$ and $c \in y$. Since as we said we can exclude the case that $x = y = c$, at least one of them contains another element; and therefore, by minimality, $\Sigma \vdash_I c \perp c$, i.e., $c \in V$. Therefore $x \cap y = \emptyset$, as stated.

Now, let $S$ be the team with domain $M$ over $\{0, 1\}$ consisting of all functions $s : M \to \{0, 1\}$ satisfying the two conditions

-   $\sum_{a \in x} s(a) \equiv \sum_{b \in y} s(b) \mod 2$;
-   $s(c) = 0$ for all $c \notin xy$

and consider the diversity rank function $\| \cdot \| = \log(\#\mathrm{rows}_S(\cdot))$ induced by $S$. Recall that, for any two disjoint $z$ and $w$, this diversity rank function satisfies $z \perp w$ if and only $S(zw) = S(z) \times S(w)$, i.e. if and only if all possible values in $S$ for $z$ and for $w$ appear jointly in some row of $S$.

This diversity rank function does not satisfy $x \perp y$. Indeed, let $a \in x$ and $b \in y$ and consider the two assignments $s_1, s_2 \in S$ defined by

- $s_1(a) = s_1(b) = 1, s_1(c) = 0$ for $c \notin \{a, b\}$;
- $s_2(c) = 0$ for all $c \in M$.

There exists no $s \in S$ that agrees with $s_1$ over $x$ and with $s_2$ over $y$, because this would violate the parity condition; therefore, it is indeed the case that the induced diversity rank function does not satisfy $x \perp y$.

It remains to show that this induced diversity rank function satisfies all independence assertions in $z \perp w \in \Sigma$ (remember that we assumed that $\Sigma$ is closed under our rules). Let us consider the following cases:

1. $z = \emptyset$ or $w = \emptyset$:

    By the soundness of Empty Set and Symmetry rules, every diversity rank function (and thus, in particular, our diversity rank function $\| \cdot \|$ induced by $S$) satisfies $z \perp w$.
2. $z \cap w \neq \emptyset$:

    By Decomposition and Symmetry, for all $c \in z \cap w$ we have that $\Sigma \vdash_I c \perp c$. Thus, $c \in V$ and therefore (since as we saw $xy \cap V = \emptyset$) $s(c) = 0$ for all $s \in S$. Thus, for $z' = z \backslash w$ and $w' = w \backslash z$, $\#\text{rows}_S(z) = \#\text{rows}_S(z')$, $\#\text{rows}_S(w) = \#\text{rows}_S(w')$ and $zw = z'w'$. Therefore, $\|zw\| = \|z\| + \|w\|$ if and only if $\|z'w'\| = \|z'\| + \|w'\|$, and therefore it suffices to check whether our diversity rank function satisfies $z' \perp w'$.
3. $z \cap w = \emptyset, z \neq \emptyset$ and $w \neq \emptyset$:

    (a) $z \backslash xy \neq \emptyset$ or $w \backslash xy \neq \emptyset$:

    Since all $s \in S$ are such that $s(c) = 0$ for all $c \notin xy$, the induced diversity rank function $\| \cdot \|$ satisfies $z \perp w$ if and only if it satisfies $(z \cap xy) \perp (w \cap xy)$. Hence, we need not consider this case further.

    (b) $zw = xy$:

    As we will now see, in this case we could conclude that $\Sigma \vdash_I x \perp y$, which contradicts our hypothesis. Thus, this case is not possible.

    Let $x_z = x \cap z$, $x_w = x \cap w$. $y_z = y \cap z$ and $y_w = y \cap w$. Then $z = x_z y_z$, $w = x_w y_w$, $x = x_z x_w$ and $y = y_z y_w$. Suppose $x_z = \emptyset$. Then $x_w \neq \emptyset$ since $x \neq \emptyset$ and $y_z \neq \emptyset$ since $z \neq \emptyset$. Moreover $y_w \neq \emptyset$ for, otherwise, $x = w$ and $y = z$, which is impossible because $\Sigma \nvdash_I x \perp y$. By symmetry, we conclude that at most one of $x_z$, $x_w$, $y_z$ and $y_w$ is empty. Without loss of generality, assume that both $x_z \neq \emptyset$ and $x_w \neq \emptyset$. By the minimality of $x \perp y$, it follows that $\Sigma \vdash_I x_z \perp y_z$. From this and $\Sigma \vdash_I x_z y_z \perp x_w y_w$, it follows by Mixing that $\Sigma \vdash_I x_z \perp x_w y_z y_w$. Again, by minimality of $x \perp y$. we have that $\Sigma \vdash_I x_w \perp y_z y_w$. Applying Mixing again and Symmetry, it follows that $\Sigma \vdash_I x \perp y$, i.e. $\Sigma \vdash_I x \perp y$. But this contradicts our assumption, and therefore it cannot be the case that $zw = xy$.

    (c) $zw \subsetneq xy$:

    It follows immediately that the relational diversity rank induced by $S$ satisfies $z \perp w$, since we can use the variable(s) in $xy$ but not in $zw$ to tweak the parity condition.

Thus, we saw that the diversity rank function induced by $S$ does not satisfy $x \perp y$ but satisfies all independence statements of $\Sigma$, as required. $\square$

**Corollary 12** [20] *An independence assertion follows semantically, in all databases, from a given set of independence assertions if and only if it follows by the rules of Proposition 8.*

At this point it would be natural to ask the following

**Open Problem** What are the rules that govern the interaction between dependence and independence in our framework?

In general, inference problems for dependence/independence assertions are not necessarily decidable in a relational setting [24–26]; but in the setting of general diversity rank functions, the decision problem for sets of dependence and independence assertions is necessarily decidable. Indeed, the axioms of diversity rank functions as well as dependence and independence assertions may be translated into the first-order arithmetic of the reals, replacing each $\|x\|$ expression with a variable $r_x$ that ranges over non-negative real numbers and – for example – writing a dependence atom $=(x, y)$ as $r_{xy} = r_x$, an independence atom $x \perp y$ as $r_{xy} = r_x + r_y$, and Axiom **R3** as an axiom schema of the form $(r_{xy} = r_x) \rightarrow (r_{xyz} = r_{xz})$. But the arithmetic of the reals is decidable [27], and therefore it is decidable whether some dependence or independence assertion follows from a set $\Sigma$ of dependence and independence assertions in our setting.

We leave the problem of searching for an axiom system for dependence and independence assertions combined to future work. Here we only point out two simple axioms that govern the interaction between dependence and independence in our setting:

- *Constancy Equivalence:* $x \perp x$ if and only if $=(\emptyset, x)$;
- *Propagation:* If $x \perp y$ and $=(y, z)$ then $x \perp yz$.

Both of these can be shown to follow from our notion of rank:

- *Constancy Equivalence:* Suppose $x \perp x$. Then, by definition, $\|x\| + \|x\| = \|xx\|$. But on the other hand, $xx = x$ and therefore, $\|x\| = 0$ and $\|x\| = \|\emptyset x\| = \|\emptyset\| = 0$. Conversely, suppose $=(\emptyset, x)$. Then $\|x\| = \|\emptyset x\| = \|\emptyset\| = 0$, and therefore $\|x\| + \|x\| = 0 = \|xx\|$.
- *Propagation:* Suppose that $\|x\| + \|y\| = \|xy\|$ and that $\|yz\| = \|y\|$. From the second hypothesis, by **R3**, we can show that $\|xy\| = \|xyz\|$ and therefore in the first hypothesis we can replace $\|y\|$ with $\|yz\|$ and $\|xy\|$ with $\|xyz\|$, thus obtaining $\|x\| + \|yz\| = \|xyz\|$. Therefore $x \perp yz$, as required.

## 6 A representation theorem for dependence atoms

As we saw, every dependence notion induced by a diversity rank function satisfies Armstrong's Axioms, which furthermore are complete for diversity rank functions. However, a question that remains open is whether every dependency notion that satisfies Armstrong's Axioms is induced by a diversity rank function.

More formally, let $M$ be a set, and let $\Sigma$ be a set of dependence assertions over $M$ which is closed under Armstrong's Axioms. Is there a diversity rank function $\| \cdot \|$ for which $\Sigma = \{=(x, y) : x, y \subseteq_f M, \|xy\| = \|x\|\}$?[16]

An example may be helpful here. Suppose that $M = \{a, b, c\}$ and $\Sigma$ is the closure under Armstrong's Axioms of $\{=(ab, c)\}$, i.e.

$$\Sigma = \{=(x, y) : x, y \subseteq M, y \subseteq x\} \cup \{=(x, y) : ab \subseteq x \text{ and } y \subseteq M\}.$$

---

[16]We thank Samson Abramsky for asking this question in a personal communication.

It is not hard to verify that this set is indeed closed under Reflexivity, Augmentation and Transitivity; but can we find a diversity rank function over $\{a, b, c\}$ such that, for $x, y \subseteq \{a, b, c\}$, $y$ depends on $x$ according to it if and only if $=(x, y) \in \Sigma$?

By Theorem 13 below, such a diversity function does exist.

**Theorem 13** *Let $M$ be a countable set and let $\Sigma$ be a set of dependency assertions of the form $=(x, y)$ (for $x, y \subseteq_f M$) that is closed under Armstrong's Axioms. Then there exists a diversity rank function $\| \cdot \|$ over $M$ such that $\Sigma_{\| \cdot \|} = \Sigma$.*

*Proof* We first define the relation $\equiv$ on the finite subsets of $M$, as follows. For $x, y \subseteq_f M$, $x \equiv y$ if both $=(x, y) \in \Sigma$ and $=(y, x) \in \Sigma$. Using that $\Sigma$ satisfies Armstrong's axioms, it is straightforwardly seen that $\equiv$ is an equivalence relation. For $x \subseteq_f M$, let $E_x$ be the equivalence class of $x$, and let $\mathcal{E} = \{E_x : x \subseteq_f M\}$. We first observe that if $E_{x_1} = E_{x_2}$ and $E_{y_1} = E_{y_2}$, then $=(x_1, y_1) \in \Sigma$ if and only if $=(x_2, y_2) \in \Sigma$. To see this, it suffices to prove one direction, by symmetry. Thus suppose $=(x_1, y_1) \in \Sigma$. Since $E_{x_1} = E_{x_2}, =(x_2, x_1) \in \Sigma$. Since $E_{y_1} = E_{y_2}, =(y_1, y_2) \in \Sigma$. By two applications of Transitivity, it now follows that $=(x_2, y_2) \in \Sigma$. This observation justifies the following definition of the relation $\leq$ on $\mathcal{E}$: for $x, y \subseteq_f M$, $E_y \leq E_x$ if $=(x, y) \in \Sigma$. It is straightforwardly seen that $\leq$ is a partial order on $\mathcal{E}$: reflexivity and transitivity follow from $\Sigma$ satisfying Reflexivity and Transitivity, and antisymmetry follows from the definition of the equivalence relation $\equiv$.

Next, we observe that the set of all finite subsets of $M$ is countable, because $M$ is. Hence $\mathcal{E}$ is also countable. Let $\mathcal{E} = \{E_0, E_1, E_2, E_3, ...\}$ be an enumeration of $\mathcal{E}$ with $E_0 = E_\emptyset$. We now construct an order-preserving injection $f : \mathcal{E} \to \{0\} \cup ]1, 2[$ recursively as follows:

1. $f(E_0) = 0$.
2. Assume that, for $i = 0, ..., n$, $f(E_i)$ has been defined. Let

$$l = \max(\{f(E_i) : 1 \leq i \leq n, E_i < E_{n+1}\} \cup \{1\}),$$

$$r = \min(\{f(E_i) : 1 \leq i \leq n, E_i > E_{n+1}\} \cup \{2\}).$$

Then, let $f(E_{n+1})$ be any value in $]l, r[ \setminus \{f(E_1), ..., f(E_n)\}$.

Notice that, for all $x \subseteq_f M$, $E_0 = E_\emptyset \leq E_x$, since $=(x, \emptyset) \in \Sigma$ by Reflexivity. Hence, the construction of $f$ guarantees that it is indeed order-preserving.

Finally, we define a diversity rank function $\| \cdot \|$ on M as follows: for $x \subseteq_f M$, $\|x\| = f(E_x)$. We prove that $\| \cdot \|$ satisfies R1-R4.

R1: We have that $\|\emptyset\| = f(E_\emptyset) = 0$.

R2: Let $x, y \subseteq_f M$. Since $x \subseteq xy$, $=(xy, x) \in \Sigma$ by Reflexivity. Hence $E_x \leq E_{xy}$. Since $f$ is order-preserving, $\|x\| = f(E_x) \leq f(E_{xy}) = \|xy\|$. If $x = \emptyset$, or $y = \emptyset$, then $\|xy\| \leq \|x\| + \|y\|$ is satisfied by R1; if $x \neq \emptyset$, and $y \neq \emptyset$, then $\|xy\| \leq \|x\| + \|y\|$ is satisfied since, by construction of $f$, $1 < \|x\|, \|y\|, \|xy\| < 2$.

R3: Let $x, y \subseteq_f M$ and assume that $\|x\| = \|xy\|$. Then, $f(E_x) = f(E_{xy})$, hence $E_x = E_{xy}$ since $f$ is injective. If follows that both $=(xy, x) \in \Sigma$ and $=(x, xy) \in \Sigma$.[17] Now, let $z \subseteq_f M$. By Augmentation, we have that both $=(xyz, xz) \in \Sigma$ and $=(xy, xyz) \in \Sigma$. Hence $E_{xy} = E_{xyz}$. It follows that $\|xy\| = f(E_{xy}) = f(E_{xyz}) = \|xyz\|$.

---

[17]The former is of course always the case, because of Reflexivity.

R4:    Notice that this property is satisfied by R1 if $x = \emptyset$ or $y = \emptyset$, and is trivially satisfied
if $z = \emptyset$. If $x \neq \emptyset$, $y \neq \emptyset$, and $z \neq \emptyset$, the property is satisfied because the antecedent is
false, as, by construction of $f$, $1 < \|xyz\|$, $\|x\|$, $\|yz\| < 2$.

It remains to show that $\Sigma = \{=(x, y) : x, y \in M, \|xy\| = \|x\|\}$. If $=(x, y) \in \Sigma$, then
by Augmentation, $=(x, xy) \in \Sigma$. By Reflexivity, $=(xy, x) \in \Sigma$. Hence, $E_x = E_{xy}$ and
$\|x\| = f(E_x) = f(E_{xy}) = \|xy\|$. Conversely, if $\|x\| = \|xy\|$, then $f(E_x) = f(E_{xy})$.
Hence, $E_x = E_{xy}$, since $f$ is injective. Consequently, $=(x, xy) \in \Sigma$, and, by Reflexivity
and Transitivity, $=(x, y) \in \Sigma$. □

As an example, let us apply the procedure described here to $M = \{a, b, c\}$ and to the set
$\Sigma$ we mentioned before, i.e.

$$\Sigma = \{=(x, y) : x, y \subseteq M, y \subseteq x\} \cup \{=(x, y) : ab \subseteq x \text{ and } y \subseteq M\}.$$

It is readily seen that $E_{ab} = E_{abc} = \{ab, abc\}$. For all other $x \subseteq_f M$, $E_x = \{x\}$. Let us
now consider the following enumeration of the finite subsets of $\mathcal{E} = \{E_x : x \subseteq_f M\}$:

$$E_\emptyset, E_a, E_b, E_c, E_{ab} = E_{abc}, E_{ac}, E_{bc}.$$

We now construct the function $f$ in the proof of Theorem 3 recursively according to this
enumeration:

1.  $f(\emptyset) = 0$, by construction.
2.  Since $E_a$, $E_b$, and $E_c$ are mutually incomparable, we may assign to them arbitrary
    mutually different values between 1 and 2, e.g., $f(E_a) = 1.5$, $f(E_b) = 1.6$, and
    $f(E_c) = 1.1$.
3.  Since, for every $x \subseteq_f M$, $=(ab, x) \in \Sigma$, we have in particular that $E_a < E_{ab}$, $E_b <$
    $E_{ab}$, and $E_c < E_{ab}$. Therefore, we must assign to $E_{ab}$ a value between 1.6 and 2, say,
    $f(E_{ab}) = 1.8$.
4.  Since $=(ac, a) \in \Sigma$, $=(ac, c) \in \Sigma$, and $=(ab, ac) \in \Sigma$, we have $E_a < E_{ac}$, $E_c < E_{ac}$,
    and $E_{ac} < E_{ab}$. Since $=(ac, b) \notin \Sigma$ and $=(b, ac) \notin \Sigma$, $E_b$ and $E_{ac}$ are incomparable.
    Therefore, we must assign to $E_{ac}$ an unused value between 1.5 and 1.8, say, $f(E_{ac}) =$
    1.7.
5.  Since $=(bc, b) \in \Sigma$, $=(bc, c) \in \Sigma$, and $=(ab, bc) \in \Sigma$, $E_b < E_{bc}$, $E_c < E_{bc}$, and
    $E_{bc} < E_{ab}$. Since $=(bc, a) \notin \Sigma$ and $=(a, bc) \notin \Sigma$, $E_a$ and $E_{bc}$ are incomparable.
    Since $=(bc, ac) \notin \Sigma$ and $=(ac, bc) \notin \Sigma$, $E_{ac}$ and $E_{bc}$ are incomparable. Therefore,
    we must assign to $E_{bc}$ an unused value between 1.6 and 1.8, say, $f(E_{bc}) = 1.65$.

According to the proof of Theorem 13, $\Sigma = \{=(x, y) : x, y \subseteq_f M, \|xy\| = \|x\|\}$, as
can be easily verified. It is also easily verified that this is indeed a diversity rank function.

We leave for future work the question whether this representation theorem may be
generalized to independence atoms, that is to say, whether every set of dependence *and
independence* assertions satisfying the axioms for dependence, the ones for independence,
and some additional axioms for dependence/independence interactions arises from some
diversity rank function.

# 7 Conclusions

In this work, we showed how many distinct notions of dependence and independence,
originating in different branches of mathematics and computer science, may be treated as
instances of the same framework: one which can be seen as a generalization of matroid

theory allowing for non-integer ranks and weakening the submodularity condition. In this framework, $y$ is said to be dependent from $x$ if adding it to $x$ does not increase the amount of diversity, while $y$ is said to be independent from $x$ if adding it to $x$ increases *maximally* the amount of diversity.

Despite its generality, this framework is nonetheless powerful enough to prove non-trivial results - including, in particular, completeness theorems for the corresponding dependence and independence notions. These results generalize to our entire setting the completeness theorems by Armstrong and by Geiger-Paz-Pearl for database-theoretic functional dependence and for probabilistic independence respectively. In addition, we have obtained a representation theorem showing that every set of dependence assertions that satisfies Armstrong's Axioms arises from some diversity rank function.

One natural next step would be to investigate further the properties of this formalism, in particular with respect to the interaction between independence and dependence assertions. Combinatorial properties of this system would also be worth investigating, as would the study of possible operations that *combine* different diversity rank functions. This could also contribute to the logical study of notions of dependence and independence in the context of Team Semantics, in particular providing a unifying approach for its different variants (e.g. probabilistic, modal, propositional, . . . ).

Another important next step would be to extend our representation result to dependence *and independence* assertions combined, as well as investigate the connections between our approach and other approaches, such as the study of measure-based constraints in the context of Database Theory [28], and the lattice-theoretic study of conditional independence of [29].

Finally, it would be interesting to investigate potential applications of our framework. Our approach in this work has been one of synthesis, motivated by the search of a framework that captures a wide array of notions of dependence and independence studied in different areas. Such a framework had to be more general than matroids, which – despite their success and undeniable importance – are limited by their commitment to submodularity and to integer values; had to be formally simple and widely applicable; and it had to be specific enough to highlight true commonalities and analogies between notions of dependence and independence developed in different areas. We think that our notion of diversity rank function meets these objectives, and that this has the potential to lead to interesting connections and fruitful cross-topic fertilization between the study of notions of dependence and independence in these different areas.

# References

1. Whitney, H.: On the Abstract Properties of Linear Dependence. Amer. J. Math. **57**(3), 509–533 (1935). https://doi.org/10.2307/2371182
2. van der Waerden, B.L.: Moderne Algebra. J. Springer, Berlin (1940)
3. Armstrong, WW.: Dependency structures of data base relationships. Inf Process. **74** (1974)
4. Väänänen, J.: Dependence logic. London Mathematical Society Student Texts, vol. 70. Cambridge University Press, Cambridge (2007)
5. Mann, AL., Sandu, G., Sevenster, M.: Independence-friendly logic. London Mathematical Society Lecture Note Series, vol. 386. Cambridge University Press, Cambridge (2011). https://doi.org/10.1017/CBO9780511981418. A game-theoretic approach
6. Grädel, E., Väänänen, J.: Dependence and independence. Studia Log. **101**(2), 399–410 (2013). https://doi.org/10.1007/s11225-013-9479-2
7. Yang, F., Väänänen, J.: Propositional logics of dependence. Ann. Pure Appl. Log. **167**(7), 557–589 (2016)
8. Väänänen, J.: Modal dependence logic. In: New perspectives on games and interaction, Texts in Logic and Games, vol. 4, pp. 237–254. Amsterdam Univ. Press, Amsterdam (2008)
9. Hella, L., Luosto, K., Sano, K., Virtema, J.: The expressive power of modal dependence logic. In: Advances in modal logic. Vol. 10, pp. 294–312. Coll. Publ., London (2014)
10. Krebs, A., Meier, A., Virtema, J., Zimmermann, M.: Team semantics for the specification and verification of hyperproperties. In: 43rd International Symposium on Mathematical Foundations of Computer Science, LIPICs. Leibniz Int. Proc. Inform., vol. 117, pp. 10–16. Schloss Dagstuhl. Leibniz-Zent. Inform., Wadern (2018)
11. Hyttinen, T., Paolini, G., Väänänen, J.: A logic for arguing about probabilities in measure teams. Arch. Math. Log. **56**(5-6), 475–489 (2017)
12. Durand, A., Hannula, M., Kontinen, J., Meier, A., Virtema, J.: Probabilistic team semantics. In: Foundations of information and knowledge systems, Lecture Notes in Comput. Sci., vol. 10833, pp. 186–206. Springer, Cham (2018). https://doi.org/10.1007/978-3-319-90050-6_1
13. Hannula, M., Hirvonen, Å., Kontinen, J., Kulikov, V., Virtema, J.: Facets of distribution identities in probabilistic team semantics. In: Logics in artificial intelligence, Lecture Notes in Comput. Sci., vol. 11468, pp. 304–320. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-19570-0_2
14. Oxley, J.G.: Matroid theory, vol. 3. Oxford University Press, USA (2006)
15. Good, I.J.: The population frequencies of species and the estimation of population parameters. Biometrika **40**(3/4), 237–264 (1953)
16. Hill, M.: Diversity and evenness: A unifying notation and its consequences. Ecology **54**, 427–432 (1973). https://doi.org/10.2307/1934352
17. Magurran, A.E.: Measuring biological diversity. Wiley, New Jersey (2004). https://books.google.fi/books?id=tUqzLSUzXxcC
18. Borda, M.: Fundamentals in information theory and coding. Springer Science & Business Media, Berlin (2011)
19. Cover, T.M., Thomas, J.A.: Entropy, relative entropy and mutual information. Elem. Inf. Theory **2**, 1–55 (1991)
20. Geiger, D., Paz, A., Pearl, J.: Axioms and algorithms for inferences involving probabilistic independence. Inform. and Comput. **91**(1), 128–141 (1991). https://doi.org/10.1016/0890-5401(91)90077-F
21. Shannon, C.E.: A mathematical theory of communication. Bell syst. Techn. J. **27**(3), 379–423 (1948)
22. Codd, E.F.: A relational model of data for large shared data banks communications. Commun. ACM **26**(1) (1970)
23. Lang, S., et al.: Algebra. Springer, Berlin (2002)
24. Herrmann, C.: On the undecidability of implications between embedded multivalued database dependencies. Inf. Comput. **122**(2), 221–235 (1995)
25. Herrmann, C.: Corrigendum to: On the undecidability of implications between embedded multivalued database dependencies. Inform. and Comput. **204**(12), 1847–1851 (2006)
26. Hannula, M., Link, S.: On the interaction of functional and inclusion dependencies with independence atoms. In: International Conference on Database Systems for Advanced Applications, pp. 353–369, Springer (2018)

27. Tarski, A.: A Decision Method for Elementary Algebra and Geometry. RAND Corporation, Santa Monica (1948)
28. Sayrafi, B., Van Gucht, D., Gyssens, M.: The implication problem for measure-based constraints. Inf. Syst. **33**(2), 221–239 (2008)
29. Niepert, M., Gyssens, M., Sayrafi, B., Van Gucht, D.: On the conditional independence implication problem: A lattice-theoretic approach. Artif. Intell. **202**, 29–51 (2013)

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.