*I have a dream...*

# Ignorance in the Relational Model

Maarten M. Fokkinga
DB group, dept INF, fac EWI, University of Twente
PO Box 217, NL 7500 AE Enschede, Netherlands
m.m.fokkinga@utwente.nl

**Abstract**. We hypothesize that an extension-with-conditioning of Dempster-Shafer theory is suitable for encoding uncertainty and ignorance in the Relational Model. We present a formal and well-motivated definition of conditioning, and show the spirit of the required change in the Relational Model and some results that then follow. It remains to be investigated whether these results are satisfactory.

## Introduction

### 1 Ignorance

Ignorance is closely related to uncertainty. Commonly, we say that a property is *uncertain* if it is not considered true or false but, instead, it is assigned a probability of being true. Now consider a set of exhaustive and mutually disjoint properties. Probability theory requires that the probabilities assigned to these properties add up to 1. *Ignorance* is the phenomenon that the "probabilities" do not add up to 1. Formally, one axiom of probability theory is not fulfilled, and hence we speak of *belief* instead of probability. Dempster-Shafer theory gives a proper formalization (summarized in paragraph 8–10), and we shall build upon that theory (paragraph 11–15).

### 2 Setting

Our work is an attempt to improve upon Choenni *et al.* [1, 2] in the following aspects: a more fundamental approach and better motivated definition of *conditioning*, and a better formalization of an *extension* of the Relational Model in order to deal with ignorance. We borrow the following example from Choenni [1], and take it to be the *leading* example.

### 3 Example: CIA

The ship type department of the CIA has 0.6 evidence that the type of ship *Maria* is *Frigate*, and 0.3 evidence that it is *Tugboat*; for the remaining 0.1 there is ignorance. This is encoded in the one-row table *SHIP* below at the left. The type speed department of the CIA has evidence that 30% of the frigates has a max speed of 20 knots, and 70% has 30 knots, whereas all tugboats have a 15 knots max speed. This is encoded in the two-row table *DESC*ription at the right:

SHIP

| Name | Type | |
|---|---|---|
| Maria | Frigate | $\mapsto 0.6$ |
| | Tugboat | $\mapsto 0.3$ |
| | * | $\mapsto 0.1$ |

DESC

| Type | Speed | |
|---|---|---|
| Frigate | 20K | $\mapsto 0.3$ |
| | 30K | $\mapsto 0.7$ |
| Tugboat | 15K | $\mapsto 1.0$ |

For the purpose of decision making, the US government requests to join the information. Here are two candidate results that they might get offered:

*join candidate* 1

| Name | Type | Speed | |
|---|---|---|---|
| Maria | Frigate | 20K | $\mapsto 0.18$ |
| | Frigate | 30K | $\mapsto 0.42$ |
| | Tugboat | 15K | $\mapsto 0.3$ |
| | * | * | $\mapsto 0.1$ |

*join candidate* 2

| Name | Type | Speed | |
|---|---|---|---|
| Maria | Frigate | 20K | $\mapsto 0.18$ |
| | Frigate | 30K | $\mapsto 0.42$ |
| | Tugboat | * | $\mapsto 0.3$ |
| | * | 20K | $\mapsto 0.03$ |
| | * | 30K | $\mapsto 0.07$ |
| Maria | Frigate | * | $\mapsto 0.6$ |
| | Tugboat | 15K | $\mapsto 0.3$ |
| | * | 15K | $\mapsto 0.1$ |

The one-row table *join candidate* 1 is obtained by "intuitive combination". However, the information in this table is *too weak* in the sense that the probability of "the max speed of *Maria* is 20K" has an upperbound (when all ignorance goes to this case) of $0.18 + 0.1 = 0.28$, whereas that upperbound is $0.6 \times 0.3 + 0.1 \times 0.3 = 0.21$ according the original *SHIP* and *DESC*.

The two-row table *join candidate* 2 is proposed by Choenni *et al.* [1]. This information is *too strong*: The first row of the table expresses that the probability of "the max speed of *Maria* is 20K" has a lowerbound (when all ignorance about the speed is not in favor of this case) of $0.18 + 0.03 = 0.21$, whereas that lowerbound is only $0.6 \times 0.3 = 0.18$ according to *SHIP* informally joined with *DESC*.

## 4 Goal, plan

Our goal is to extend the Relational Model and relational operators (like projection, selection, and in particular the join) in such a way that we can offer the US government the right information. Moreover, we should also be able to relate in a formal way the above candidate joins to "the correct join" of *SHIP* and *DESC*. The next paragraph gives the outline of the theory that we want to develop, and paragraph 6 discusses the previous example in the theory that we envisage.

## 5 Hypothesis, focus

In order to deal with uncertainty and ignorance, Dempster has weakened probability theory to what currently is known as *Dempster-Shafer theory*. The primary notion is *bpa* (basic probability assignment), from which the notions of *belief*, *plausibility*, and *ignorance* can be defined; and conversely. Our *hypothesis* is that an extension of Dempster-Shafer theory provides a solution for the problem how to deal with uncertainty and ignorance in the Relational Model, and we want to investigate this hypothesis. The main line, then is as follows.

In the leading example of paragraph 3, we start out with bpa's as *attribute values* in the table, and observe that all our attempts for a join lead to a table with a "bpa covering several attributes", which we call *tupled-bpa*, or just *t-bpa* for short. It can be shown that this generalization (*our* generalization!) of bpa to t-bpa is not essential: t-bpa's can be expressed as bpa's (though at considerable loss of readability), and vice versa. Continuing with taking joins of the resulting table with other tables will lead to tables in which t-bpa's cover more and more attributes. Therefore we generalize the notion of relation right away to one where *each row is a t-bpa*.

Moreover, we see that in Choenni's attempt the failure is due to the omission of a *condition* in the bpa: the 0.03 evidence for max speed $20K$ cannot be given unconditionally, but is only valid if it is known that the type is *Frigate*. Therefore we generalize right away to one where each row is a "conditioned t-bpa", or *ct-bpa* for short; the definition of "conditioned bpa" (or conditioned t-bpa) is new and is the focus is this paper. This notion of conditioning differs entirely from the notion of conditioning briefly discussed by Shafer [3], from the notion defined by Choenni [1, 2], and from the notion of conditioning as known in probability theory.

## 6 Envisaged solution

Once the *conditioning* (and *tupling*) extension to Dempster-Shafer theory has been developed and, based on this, also a new Relational Model, we expect to be able to deal formally with the example in the following way.

To deal with uncertainty and ignorance is quite straightforward: let each row in each table be a ct-bpa. We call such relations: *ui-relations*, where the letters 'ui' derive from 'uncertainty and ignorance'. For example, we encode the one-row table *SHIP* of paragraph 3 as an ui-relation with one row (that is, one ct-bpa) with the following pretty-print:

| | Name | Type | Name | Type | |
|---|---|---|---|---|---|
| *SHIP'* | * | *Frigate* \| | *Maria* | * | $\mapsto 0.6$ |
| | * | *Tugboat* \| | *Maria* | * | $\mapsto 0.3$ |
| | * | * \| | *Maria* | * | $\mapsto 0.1$ |

Fully written out the relation reads as in Figure 1.

Each star, *, is pronounced "*unknown*" and stands for the entire domain of the corresponding attribute: the *Name*-star stands for {*Maria*, ...}, the *Type*-star stands for {*Frigate*, *Tugboat*, ...}, and so on. The meaning of the first line of the above one-row relation is, roughly: "there is evidence 0.6 for that the type is *Frigate* on the condition that the ship is *Maria*. More precisely, the line means:

> There is evidence 0.6 for the property
> > $(Name, Type) \subseteq (unknown, \{Frigate\})$
> on the condition that
> > $(Name, Type) \subseteq (\{Maria\}, unknown)$
> is true.

Here it is understood that $(U, V) \subseteq (X, Y)$ means: $U \subseteq X \wedge V \subseteq Y$.

So, the one-row *ui-relation SHIP'* given above encodes that the ship type department of the CIA has 0.6 evidence that the type of a ship is *Frigate if its name is Maria*, and 0.3 evidence that it is *Tugboat*; for the remainder there is ignorance.

Further, the type speed department of the CIA has evidence that 30% of the frigates has a max speed of 20 knots, and 70% has 30 knots, whereas all tugboats have a 15 knots max speed. This is encoded in the two-row ui-relation with the following pretty-print:

| *DESC'* | Type | Speed | Type | Speed | |
|---|---|---|---|---|---|
| | * | $20K$ \| | *Frigate* | * | $\mapsto 0.3$ |
| | * | $30K$ \| | *Frigate* | * | $\mapsto 0.7$ |
| | * | $15K$ \| | *Tugboat* | * | $\mapsto 1.0$ |

For the purpose of decision making, the US government requests to join the information. Here is what they get:

*SHIP' ⋈ DESC'*

| Name | Type | Speed | Name | Type | Speed | |
|---|---|---|---|---|---|---|
| * | *Frigate* | $20K$ \| | *Maria* | * | * | $\mapsto 0.18$ |
| * | *Frigate* | $30K$ \| | *Maria* | * | * | $\mapsto 0.42$ |
| * | *Tugboat* | * \| | *Maria* | * | * | $\mapsto 0.3$ |
| * | * | $20K$ \| | *Maria* | *Frigate* | * | $\mapsto 0.03$ |
| * | * | $30K$ \| | *Maria* | *Frigate* | * | $\mapsto 0.07$ |
| * | *Frigate* | * \| | *Maria* | * | * | $\mapsto 0.6$ |
| * | *Tugboat* | $15K$ \| | *Maria* | * | * | $\mapsto 0.3$ |
| * | * | $15K$ \| | *Maria* | *Tugboat* | * | $\mapsto 0.1$ |

Note the last two lines of the first row, and the last line of the second row: these say that the ship type is unknown but yet the speed is certain to some degree *if* the type happens to be frigate or tugboat, respectively. (And if

$$\begin{aligned}
\{ & \\
(\{Name \mapsto *, \quad Type \mapsto \{Frigate\}\} \quad &| \quad \{Name \mapsto \{Maria\}, \quad Type \mapsto *\}) \quad \mapsto 0.6, \\
(\{Name \mapsto *, \quad Type \mapsto \{Tugboat\}\} \quad &| \quad \{Name \mapsto \{Maria\}, \quad Type \mapsto *\}) \quad \mapsto 0.3, \\
(\{Name \mapsto *, \quad Type \mapsto \quad * \quad\} \quad &| \quad \{Name \mapsto \{Maria\}, \quad Type \mapsto *\}) \quad \mapsto 0.1 \\
\}&
\end{aligned}$$

**Figure 1:** *Fully written-out one-row ui-relation SHIP$'$ (see paragraph 6).*

the condition is not fulfilled, the evidence supports just *unknown* — see paragraph 12.) In order to eliminate this fine-grained conditioned information and obtain information that is somewhat easier to understand for the US government, we can weaken each row by "condition-restricting it to $\{Name\}$", that is, replacing all conditions except *Name* by *unknown*, which we will be able to denote formally by $(\{Name\}\blacktriangleleft_c) * (SHIP' \bowtie DESC')$:

| Name | Type | Speed | | Name | Type | Speed | |
|---|---|---|---|---|---|---|---|
| * | Frigate | 20K | \| | Maria | * | * | $\mapsto 0.18$ |
| * | Frigate | 30K | \| | Maria | * | * | $\mapsto 0.42$ |
| * | Tugboat | * | \| | Maria | * | * | $\mapsto 0.3$ |
| * | * | * | \| | Maria | * | * | $\mapsto 0.03{+}0.07$ |
| * | Frigate | * | \| | Maria | * | * | $\mapsto 0.6$ |
| * | Tugboat | 15K | \| | Maria | * | * | $\mapsto 0.3$ |
| * | * | * | \| | Maria | * | * | $\mapsto 0.1$ |

Taking the "least upper bound" of the two rows gives the following one-row ui-relation:

| Name | Type | Speed | | Name | Type | Speed | |
|---|---|---|---|---|---|---|---|
| * | Frigate | 20K | \| | Maria | * | * | $\mapsto 0.18$ |
| * | Frigate | 30K | \| | Maria | * | * | $\mapsto 0.42$ |
| * | Tugboat | 15K | \| | Maria | * | * | $\mapsto 0.3$ |
| * | * | * | \| | Maria | * | * | $\mapsto 0.1$ |

As observed in paragraph 3 this information is too weak.

Phrased in our concepts, the kind of join that Choenni [1] (and [2]?) proposes, yields for *SHIP$'$* and *DESC$'$* the following ui-relation (again, compare with paragraph 3):

| Name | Type | Speed | | Name | Type | Speed | |
|---|---|---|---|---|---|---|---|
| * | Frigate | 20K | \| | Maria | * | * | $\mapsto 0.18$ |
| * | Frigate | 30K | \| | Maria | * | * | $\mapsto 0.42$ |
| * | Tugboat | * | \| | Maria | * | * | $\mapsto 0.3$ |
| * | * | 20K | \| | Maria | * | * | $\mapsto 0.03$ |
| * | * | 30K | \| | Maria | * | * | $\mapsto 0.07$ |
| * | Frigate | * | \| | Maria | * | * | $\mapsto 0.6$ |
| * | Tugboat | 15K | \| | Maria | * | * | $\mapsto 0.3$ |
| * | * | 15K | \| | Maria | * | * | $\mapsto 0.1$ |

As observed in paragraph 3 this information is too strong.

## The well-known Dempster-Shafer theory

Although in the running example there are three domains (*Name*, *Type*, and *Speed*), Dempster-Shafer theory deals only with one domain. This restriction is without loss of generality, as discussed in paragraph 15.

## 7 Notational conventions

We consider functions to be *sets of argument-result pairs*, and use the notation $x \mapsto y$ as a suggestive synonym for the pair $(x, y)$. So, the set $\{a \mapsto 3, b \mapsto 2, c \mapsto 3\}$ is a function that maps $a$ to 3, maps $b$ to 2, and $c$ to 3. For summation we use a notation without subscripting:

$$\begin{aligned}
\Sigma\, x, y \bullet expr(x, y) \qquad & \Sigma\, x, y \mid cond(x, y) \bullet expr(x, y) \\
= \Sigma_{x,y}\, expr(x, y) \qquad & = \Sigma_{x, y\ s.\ th.\ cond(x,y)}\, expr(x, y)
\end{aligned}$$

We do so because in some cases the '$x, y \mid cond(x, y)$' part is just too large to be written as a subscript (see for example paragraph 13).

In the context of a set $D$ we let $P$, $Q$ vary over subsets of $D$, that is, $P, Q : \mathbb{P}\, D$, and we sometimes write $*$ for $D$.

## 8 Basic probability assignment

For the reader not familiar with Dempster-Shafer theory, we provide an intuition in the appendix paragraph 18. We build on this intuition later when we generalize the theory with conditioning.

Let $D$ be a finite set. A *basic probability assignment over* $D$, abbreviated *bpa*, is a total function $m : \mathbb{P}\, D \longrightarrow [0, 1]$ satisfying:

$$\begin{aligned}
m\, P \quad &= \quad 0 \qquad \text{whenever } P = \varnothing \\
\Sigma\, P \bullet m\, P \quad &= \quad 1
\end{aligned}$$

The latter equation means that the sum of values $m\, P$, for all subsets $P$ of $D$, equals 1.

Sometimes, a bpa is called a mass function, hence letter $m$ for bpa's. We omit the entries $\{\ldots\} \mapsto 0$ in the presentation of a bpa. A bpa $m$ induces a *belief Bel*, a *plausibility Pl*, and an *ignorance Ig* as total functions of type $\mathbb{P}\, D \longrightarrow [0, 1]$ as follows:

$$\begin{aligned}
Bel\, P \quad &= \quad \Sigma\, P' \mid P' \subseteq P \bullet m\, P' \\
Pl\, P \quad &= \quad 1 - Bel(D \setminus P) \\
Ig\, P \quad &= \quad Pl\, P - Bel\, P
\end{aligned}$$

A single value $d \in D$ may be considered as a bpa, namely the bpa $m_d$ that maps every $P \subseteq D$ to zero except $P = \{d\}$, that is, $m_d = \{\{d\} \mapsto 1\}$.

## 9 Example

Department $m$ of the CIA has 0.6 evidence that the type of a certain ship is *Frigate*, and 0.3 evidence that it is *Tugboat*; for the remainder there is ignorance. The department is characterized as follows, as a bpa over $\{Frigate, Tugboat, \ldots\}$:

$$m \quad = \quad \{\{Frigate\} \mapsto 0.6,\ \{Tugboat\} \mapsto 0.3,\ * \mapsto 0.1\}$$

Recall that $*$ stands for the full set $\{Frigate, Tugboat, \ldots\}$.

## 10 Combination of bpa's

Dempster defines a combination $\oplus$ of bpa's, now commonly known as *Dempster's combination rule*, or *orthogonal sum*. We give the formal definition here, and our intuition in appendix paragraph 19. We build on this intuition later when we generalize the theory with conditioning.

Let $m_1$ and $m_2$ be bpa's over $D$. If constant $\kappa$, defined below, equals 0, then the combination of $m_1$ and $m_2$ is said not to exist; if $\kappa$ differs from 0, then the combination $m_1 \oplus m_2$ is a bpa over $D$ defined as follows:

$$(m_1 \oplus m_2)\,P = 0 \qquad \text{if } P = \varnothing, \text{ else:}$$
$$(m_1 \oplus m_2)\,P = (\Sigma\,P_1, P_2 \mid P_1 \cap P_2 = P \bullet m_1\,P_1 \times m_2\,P_2)/\kappa$$
$$\text{where}$$
$$\kappa = (\Sigma\,P_1, P_2 \mid P_1 \cap P_2 \neq \varnothing \bullet m_1\,P_1 \times m_2\,P_2)$$

It is easily checked that whenever $m_1 \oplus m_2$ is defined, it is a bpa; notice that $\kappa$ equals "the sum of all $m_1 \oplus m_2$-results if normalization '$/\kappa$' were left out of $\oplus$'s definition". More precisely, $\kappa$ = "the above $(\Sigma \ldots)$ except that part '$P_1 \cap P_2 = P$' is extended with 'for some $P \neq \varnothing$'. "

# Generalization: Conditioning (and Tupling)

## 11 Conditioned bpa

Let $D$ be a finite set. A *conditioned bpa over $D$*, *c-bpa* for short, is a total function $m : D \times D \longrightarrow [0,1]$ such that:

$$m\,(P \mid Q) = 0 \qquad \text{whenever} \quad P \cap Q = \varnothing$$
$$\Sigma\,P, Q \bullet m\,(P \mid Q) = 1$$

Consistent with common practice in probability theory, we separate the two arguments of a c-bpa with a '|' rather than a comma, and interpret the second one as the condition and the first one as the conclusion.

The belief, plausibility, and ignorance induced by $m$ are defined as follows:

$$Bel\,(P \mid Q) = \Sigma\,P', Q' \mid P' \subseteq P \wedge Q' \supseteq Q \bullet m\,(P' \mid Q')$$
$$Pl\,(P \mid Q) = 1 - Bel\,(D \setminus P \mid Q)$$
$$Ig\,(P \mid Q) = Pl\,(P \mid Q) - Bel\,(P \mid Q)$$

## 12 Interpretation

A c-bpa $m$ is interpreted as an agent, having *conditioned* evidences. Specifically, the statement $m\,(P \mid Q) = x$ is interpreted as follows:

> Agent $m$ has evidence $x$ supporting *just $P$* provided that a proposition from $Q$ is true; if the condition is not fulfilled, the evidence supports just $D$ unconditionally.

Due to the exhaustiveness of the set $D$ of propositions, there cannot be any evidence for $P \mid Q$ in case $P \cap Q$ is empty. The interpretation of the condition plays also an important role in the combination of two c-bpa's.

## 13 Combination of c-bpa's

Let $m_1$ and $m_2$ be c-bpa's over $D$. If constant $\kappa$ defined below equals 0, then the combination of $m_1$ and $m_2$ is said not to exist; if $\kappa$ differs from 0, then the combination $m_1 \oplus m_2$ is

a c-bpa over $D$ defined as follows — *with an indispensable(!) explanation following the definition*:

$$(m_1 \oplus m_2)(P \mid Q) = 0 \qquad \text{if } P \cap Q = \varnothing \quad \text{else:}$$
$$(m_1 \oplus m_2)(P \mid Q) =$$
$$(\Sigma\,P_1, Q_1, P_2, Q_2$$
$$\mid \quad P_1' \cap P_2' = P \quad \wedge \quad Q_1' \cap Q_2' = Q$$
$$\text{where}$$
$$P_1', Q_1' = (\textbf{if} \qquad Q_1 \subseteq D \setminus P_2 \textbf{ then } *, *$$
$$\textbf{if } P_2 \subseteq Q_1 \not\subseteq D \setminus P_2 \textbf{ then } P_1, *$$
$$\textbf{if } P_2 \not\subseteq Q_1 \not\subseteq D \setminus P_2 \textbf{ then } P_1, Q_1\,),$$
$$P_2', Q_2' = (\textbf{if} \qquad Q_2 \subseteq D \setminus P_1 \textbf{ then } *, *$$
$$\textbf{if } P_1 \subseteq Q_2 \not\subseteq D \setminus P_1 \textbf{ then } P_2, *$$
$$\textbf{if } P_1 \not\subseteq Q_2 \not\subseteq D \setminus P_1 \textbf{ then } P_2, Q_2\,)$$
$$\bullet\, m_1\,(P_1 \mid Q_1) \times m_2\,(P_2 \mid Q_2)$$
$$)\,/\,\kappa$$

The first clause is an immediate consequence of the observation in paragraph 12. We shall now explain the second clause. Let $P, Q$ be arbitrary with non-empty intersection. The summation constraint in between '|' and '$\bullet$' characterizes the possible $P_1 \mid Q_1$ and $P_2 \mid Q_2$ that *in combination* support just $P \mid Q$; precisely for these $P_1, Q_1, P_2, Q_2$ the product $m_1(P_1 \mid Q_1) \times m_2(P_2 \mid Q_2)$ is taken into the summation for $P \mid Q$. We explain the characterization of $P_1, Q_1$ only; the characterization of $P_2, Q_2$ is similar.

> In the general case, the evidence that agents 1 holds in support for $P_1 \mid Q_1$ will, *in combination with the other agent*, support just $P_1 \cap ... \mid Q_1 \cap ...$ (where the dots ... stand for the contribution of the other agent): the intersection in the conclusion is the same one as for normal bpa's, and the intersection in the condition expresses a conjunction of the conditions. However, for agent 1 there are two circumstances that lead to a change of its contribution in the combination.

- First, agent 1's condition $Q_1$ might be *in*consistent with the other agents conclusion $P_2$ ($Q_1 \cap P_2 = \varnothing$ or, equivalently, $Q_1 \subseteq D \setminus P_2$). The interpretation of paragraph 12 says: "if the condition is not fulfilled, the evidence supports just $D$ unconditionally." So, agent 1's evidence supports, *in the combination*, just $P_1' \cap ... \mid Q_1' \cap ...$ where $P_1', Q_1' = *, *$. This is covered by the first branch for $P_1', Q_1'$.

- Second, agent 1's condition $Q_1$ might be *implied* by the other agents conclusion $P_2$ ($P_2 \subseteq Q_1$ or, equivalently, $P_2 \subseteq Q_1 \not\subseteq D \setminus P_2$). In the combination, then, condition $Q_1$ is fulfilled and may be weakened to $*$. So, agent 1's evidence supports, *in the combination*, just $P_1' \cap ... \mid Q_1' \cap ...$ where $P_1', Q_1' = P_1, *$. This is covered by the second branch for $P_1', Q_1'$.

  > *Note.* Recall that the condition of $m_1$ requests the *truth* of a member in $Q_1$ whereas $P_2 \subseteq Q_1$ only asserts some *evidence* supporting $Q_1$. Yet, the condition is discarded, by putting $Q_1' = *$. This is justified since the combination $m_1 \oplus m_2$ doesn't assert any truth but only some evidence.

- In the remaining ("general") case, as we have said above, agent 1's evidence supports, *in the combination*, just $P_1' \cap ... \mid Q_1' \cap ...$ where $P_1', Q_1' = P_1, Q_1$. This is covered by the third branch for $P_1', Q_1'$.

It remains to define $\kappa$; it must make the total sum over $m_1 \oplus m_2$ equal to one. Hence, $\kappa$ is defined to be "the sum of all $m_1 \oplus m_2$-results if normalization '$/\kappa$' were left out of $\oplus$'s definition". Equivalently, take $\kappa =$ "the above $(\Sigma \ldots)$ except that part '$P'_1 \cap P'_2 = P \wedge Q'_1 \cap Q'_2 = Q$' is extended with 'for some $P, Q$ with $P \cap Q \neq \varnothing$'. "

## 14 Example

Take $D$ to be a set of numbers, partitioned into *small* and *large*, with *tiny* $\subset$ *small* and *huge* $\subset$ *large*, and let *even*, *odd*, *prime* have their conventional meaning. Consider the following c-bpa's:

$$m_1 = \{(small \mid *) \mapsto 0.4, \quad (huge \mid prime) \mapsto 0.6\}$$
$$m_2 = \{(even \mid large) \mapsto 0.3, \quad (odd \mid tiny) \mapsto 0.7\}$$

Then the combination $m_1 \oplus m_2$ is computed as follows:

| $small \mid *$ | 0.4 | $small \cap * \mid * \cap *$ | $small \cap odd \mid * \cap tiny$ |
|---|---|---|---|
| $huge \mid prime$ | 0.6 | — (see note †) | $huge \cap * \mid prime \cap *$ |
| | | 0.3 | 0.7 |
| $m_1 \Uparrow \quad m_2 \Longrightarrow$ | | $even \mid large$ | $odd \mid tiny$ |

Note †: *huge* $\cap$ *even* $\mid$ *prime* $\cap$ * is discarded since *huge* $\cap$ *even* $\cap$ *prime* $\cap$ * $= \varnothing$

For the upper-left rectangle of the "square", note that $m_1$'s condition * is implied by $m_2$'s conclusion *even* (that is, * $\supseteq$ *even*), so $m_1$'s evidence for *small* $\mid$ * is dealt with as *small* $\mid$ * in the combination (the * is discarded and replaced by *); further, $m_2$'s condition *large* is inconsistent with $m_1$'s conclusion *small* (they have an empty intersection), hence $m_2$'s evidence for *even* $\mid$ *large* is dealt with in the combination as * $\mid$ *. Similarly for the lower-right rectangle. For the upper-right rectangle, note that $m_1$'s condition * is implied by $m_2$'s conclusion *odd* (that is, * $\supseteq$ *odd*), so that $m_1$'s evidence for *small* $\mid$ * is dealt with as *small* $\mid$ * in the combination (the * is discarded and replaced by *); and further, $m_2$'s condition *tiny* is not implied by $m_1$'s conclusion *small* and these two are not inconsistent ($tiny \cap small \neq \varnothing$), so $m_2$'s evidence for *odd* $\mid$ *tiny* is dealt with unchanged in the combination. For the lower-left rectangle, we have the same situation as for the upper-right rectangle. However, since $huge \cap even \cap prime \cap * = \varnothing$ (that is, "$P \cap Q = \varnothing$" — there are no huge even primes), the evidence for the combined case *huge* $\cap$ *even* $\mid$ *prime* $\cap$ * must be zero by definition of the notion of c-bpa. All together:

$$m_1 \oplus m_2 = \{ (small \mid *) \mapsto 0.12/\kappa,$$
$$(small \cap odd \mid tiny) \mapsto 0.28/\kappa,$$
$$(huge \mid prime) \mapsto 0.42/\kappa \}$$
$$\text{where}$$
$$\kappa = 0.12 + 0.28 + 0.42$$

## 15 Tupling

So far, the formal definitions assume that there is a single domain of discourse, $D$. However, in the CIA-example there are clearly several distinct domains: *Name*, *Type*, and *Speed*. In order to deal with such a situation, we need to extend Dempster-Shafer theory, and our generalization with

conditioning, in such a way that several domains can be dealt with simultaneously. This is achieved by the notion of "tupled-bpa" (*t-bpa*, for short). It is not hard to do so, but space limitations do not permit us to give the details.

In fact, the notion of t-bpa is superfluous in the sense that normal bpa's can already express (although in a rather complicated and unpractical way) what we wish to express with t-bpa's. This came for us as a little surprise, because in general $\mathbb{P} D_1 \times \cdots \times \mathbb{P} D_n$ and $\mathbb{P}(D_1 \times \cdots \times D_n)$ are quite different; the explanation, however, is that the former can be *embedded* in the latter. For example, take t-bpa $m$ over $D = (D_1, D_2)$ as follows:

$$m = \{\ldots, \quad (\{a, b, c\}, \{p, q\}) \mapsto x, \quad \ldots\}$$

This $m$ can be viewed as denoting the following normal bpa $m'$ over $D' = D_1 \times D_2$:

$$m'$$
$$= \{\ldots, \{a, b, c\} \times \{p, q\} \mapsto x, \ldots\}$$
$$= \{\ldots, \{(a, p), (a, q), (b, p), (b, q), (c, p), (c, q)\} \mapsto x, \ldots\}$$

So, $m'$ maps a subset $P$ of $D_1 \times D_2$ to 0 except when $P$ happens to be equal to $\pi_1 P \times \pi_2 P$, in which case $m' P = x$ iff $m(\pi_1 P, \pi_2 P) = x$ — where $\pi_i$ is the projection of a set of tuples to coordinate $i$, that is, $\pi_i P = \{(x_1, x_2) : P \bullet x_i\}$. Formally, the normal bpa $m'$ that represents the t-bpa $m$, is defined as follows:

$$m' P = \textbf{if } P = \pi_1 P \times \pi_2 P \textbf{ then } m(\pi_1 P, \pi_2 P) \textbf{ else } 0$$

The construction is fully general, as can be proved formally.

Similarly for the combination of conditioning and tupling: ct-bpa. Again, space restrictions do not permit us to give the details.

## Extending the Relational Model

Having generalized Dempster-Shafer theory with conditioning and extended it with tupling as well, the formal definitions of the well-known classical Relational Model and our new one look very similar — even for the join operation. This is one half of our goal (the other half being the condition that conditioning expresses indeed what we intuitively wish to express). We want to show the *similarity* here, in particular for the join operation, without intending or attempting to explain the formulas. For the die-hards that nevertheless *do* want to understand every detail (which is not necessary to observe the similarity!), we provide some missing definitions in the appendix: paragraph 21–23.

## 16 The classical Relational Model

Let $(A, D)$ be a *schema* (definition omitted). A *relation over* $(A, D)$ is a subset of $\Pi_A D$ (definition omitted). We let $R$ vary over relations, and $r$ over members of $R$. Then we define:

| **Projection** | $\pi_B R$ | $= \{r : R \bullet B \triangleleft r\}$ |
|---|---|---|
| **Transformation** | $f * R$ | $= \{r : R \bullet f\, r\}$ |
| **Selection** | $\sigma_P R$ | $= \{r : R \mid r \in P\} = R \cap P$ |

**Join**. For $i = 1, 2$, let $(A_i, D_i)$ be a schema, and $R_i$ a relation over $(A_i, D_i)$ such that the domain assignments agree on the common attributes: $D_1 a = D_2 a$ for all $a$ in $A_1 \cap A_2$. Then:

$$R_1 \bowtie R_2 =$$

$$\{r_1 : R_1; \; r_2 : R_2 \mid \text{``function } r_1 \cup r_2 \text{ exists''} \bullet r_1 \cup r_2\}$$

Recall that functions are sets of argument-result pairs, so that for functions $f$ and $g$ the union $f \cup g$ is a well-defined set; it denotes a *function* again if $f$ and $g$ agree on their common arguments. So, since $D_1$ and $D_2$ are assumed to yield the same results on $A_1 \cap A_2$, the expression $D_1 \cup D_2$ denotes a proper function. If some $D_1$ and $D_2$ have a different domain for some common attribute $a$, then $R_1 \bowtie R_2$ is not defined. Note that, here, the condition "function $r_1 \cup r_2$ exists" formally means: $(\forall a : A_1 \cap A_2 \bullet r_1 \, a = r_2 \, a)$. The join is a relation over $(A_1 \cup A_2, D_1 \cup D_2)$.

## 17 Relations with uncertainty and ignorance

Let $(A, D)$ be a finite schema. A *relation-with-uncertainty-and-ignorance over* $(A, D)$, *ui-relation* for short, is a set of ct-bpa's over $(A, D)$. We let letter $R$ range over ui-relations, and letter $m$ range over members (being ct-bpa's) of $R$. We define:

| | | | |
|---|---|---|---|
| **Projection** | $\pi_B R$ | $=$ | $\{m : R \bullet B \triangleleft m\}$ |
| **Transformation** | $f * R$ | $=$ | $\{m : R \bullet f \, m\}$ |
| **Selection** | $\sigma_P R$ | $=$ | $\{m : R \mid m \in P\} = R \cap P$ |

**Join**. For $i = 1, 2$, let $(A_i, D_i)$ be a finite schema, and $R_i$ be an ui-relation over $(A_i, D_i)$ such that the domain assignments agree on the common attributes: $D_1 \, a = D_2 \, a$ for all $a$ in $A_1 \cap A_2$. Then the *join of $R_1$ and $R_2$*, denoted $R_1 \bowtie R_2$, is the ui-relation over $(A_1 \cup A_2, \; D_1 \cup D_2)$ that contains for each pair $(m_1, m_2)$ in $R_1 \times R_2$ the combination $m_1' \oplus m_2'$ provided it exists (where $m_1'$ is the "obvious extension" of $m_1$ to $A_1 \cup A_2$, and similarly for $m_2'$):

$$R_1 \bowtie R_2 \;= \\ \{m_1 : R_1; \; m_2 : R_2 \mid \text{``} m_1' \oplus m_2' \text{ exists''} \bullet m_1' \oplus m_2'\}$$

Within the above right-hand side, ct-bpa $m_1'$ over schema $(A_1 \cup A_2, D_1 \cup D_2)$ is constructed out of ct-bpa $m_1$ over $(A_1, D_1)$ by "extending $m_1$ with $*$ in the entries for all $a \notin A_1$". Formally, $m_1'$ maps $(P \mid Q)$ to positive $x$ if and only if $m_1$ maps $(A_1 \triangleleft P \mid A_1 \triangleleft Q)$ to positive $x$ and for all $a \in A \backslash A_1$ we have $P \, a = Q \, a = * = D_2 \, a$:

$$m_1' (P \mid Q) = \begin{array}{ll} \textbf{if} & (\forall a : A \backslash A_1 \bullet P \, a = D_2 \, a = Q \, a) \\ \textbf{then} & m_1 \, (A_1 \triangleleft P \mid A_1 \triangleleft Q) \\ \textbf{else} & 0 \end{array}$$

Similarly for the construction of $m_2'$ out of $m_2$.

Note that the above definition of $R_1 \bowtie R_2$ has the same structure as the definition of the normal join, which we consider as a necessary condition in order to call our attempt successful (and the attempt by Choenni *et al.* fails in this respect). For both joins we find that only some pairs of the Cartesian product of $R_1 \times R_2$ will give rise to a row in the result, or more specifically, when "the combination" of a pair of rows exists, the pair contributes a row to the result, but when "the combination" does not exist, the rows are considered contradictory and the pair does not contribute to the result. Hence, for both joins, the size of the join $R_1 \bowtie R_2$ may be smaller than the size of the Cartesian product $R_1 \times R_2$.

## Discussion

Just presenting a set of formal well-formed definitions is in itself no guarantee that some practical problem has been solved. We do have the formal definitions, but we are not yet sure that they give the outcome that one wants to have in practice. For instance, does the initial relation *SHIP* (borrowed from Choenni [1]) make sense (*two* rows about *Maria*), and does *our* relation *SHIP$'$* makes sense? Is there something special about keys, like *Maria*? We do not know the final answers yet. We consider our work as an attempt improve the sketch by Choenni for encoding ignorance in the relational model.

## . REFERENCES

[1] Sunil Choenni, Henk Ernst Blok, and Maarten Fokkinga. Extending the Relational Model with Uncertainty and Ignorance. Technical Report TR-CTIT-04-29, Centre for Telematics and Information Technology, University of Twente, The Netherlands, July 2004.

[2] Sunil Choenni, Henk Ernst Blok, and Erik Leertouwer. Handling Uncertainty and Ignorance in Databases: A Rule to Combine Dependent Data. In *Proceedings of the 11th Intl. Conf. on Database Systems for Advanced Applications (DASFAA)*, April 2006.

[3] G. Shafer. *A Mathematical Theory of Evidence*. Princeton University Press, Princeton University, 1976.

## APPENDIX

### 18 Appendix: Interpretation of bpa

Referring to paragraph 8, we interpret the finite $D$ as an *exhaustive* set of *distinct* propositions about some topic: $D$ is the frame of discernment. We imagine that there are *infallible* agents that in some way or another may have obtained evidence for sets of propositions; more precisely, an agent $m$ may have for subsets $P$ of $D$ "evidence of amount $x$ supporting *just $P$*", meaning that agent $m$ is certain to degree $x$ that a proposition in $P$ is true. We denote this fact by:

$$m \, P = x$$

We assume that only the ratio between the various evidences matter, so that an amount of evidence is expressed as a number in $[0, 1]$. Moreover, in order that the agents can compare and combine their findings, we assume that each agent normalizes his evidences; thanks to the exhaustiveness of $D$ we can do it in such a way that the sum of all evidences is one: $\Sigma P \bullet m \, P = 1$. The sum is well-defined thanks to finiteness of $D$. Since set $D$ is exhaustive, it is impossible that there is evidence for the empty set of propositions: $m \, \varnothing = 0$. Since the propositions are *distinct*, different $p, p' \in D$ denote different propositions, so that there is no need for constraints that relate $m \{p\}$ to $m \{p'\}$. Such an agent $m$, then, is formally characterized by a bpa.

In practice, an agent will have evidences only for some subsets of $D$. In that case, we stipulate that the agent has evidence 0 for each missing $P \neq D$.

Note that an agent may have evidence $x$ supporting just $P$ and also another evidence supporting just a subset $P'$ of $P$;

in this way the agent may differentiate between the individual propositions of $D$. The following *indifference* principle is crucial (for the definition of combination in paragraph 10):

> "Evidence $x$ supporting *just $P$*" does not distinguish between the individual propositions in $P$: it allows for arbitrarily differentiated evidences for the individual propositions and subsets of $P$ by some other means.

In particular, $m(\{p_1, \ldots, p_n\}) = x$ does *not* imply or follow from $m(\{p_i\}) = \frac{x}{n}$ for $i = 1 \ldots n$; these two assertions will lead to distinct beliefs (defined below).

Further following the above interpretation, it is natural to say that for an agent $m$, the *belief* in $P$ is the sum of all evidences supporting parts of $P$. In addition, the plausibility in $P$ is the "un-belief" in the complement of $P$, and ignorance in $P$ is the difference between the plausibility and belief in $P$.

Note that the "weakest" set of propositions is $D$ itself, since evidence for the set $D$ gives no information at all; in particular, evidence 1 for $D$ (and consequently evidence 0 for every proper subset of $D$) signals complete ignorance.

## 19 Appendix: Intuition of Dempster's $\oplus$

Shafer "provides no conclusive *a priori* argument for Dempster's rule" but sees that "the rule does seem to reflect the pooling of evidence, provided..." [3, page 57]. He explains the rule by interpreting the definition for $m_1 \oplus m_2$ in a geometrical way as in Figure 2. Here we try to give "a motivated intuition" for Dempster's combination rule presented in paragraph 10.

Let $D$ be a frame of discernment. Consider two *independent* agents characterized by $m_1$ and $m_2$, respectively. In what way can the two agents further act as one, combining their evidences? For this, we make the following observations:

($i$) The infallibility and independence of the agents implies that each agent wants to combine *each* "piece of evidence" held by the other with *all* his own evidences, in such a way that it is done *proportionally to his own evidences*.

($ii$) The indifference principle implies that for both agents together the evidences supporting $P_1$ according to agent 1 and supporting $P_2$ according to agent 2 together support $P_1 \cap P_2$, provided this intersection is nonempty.
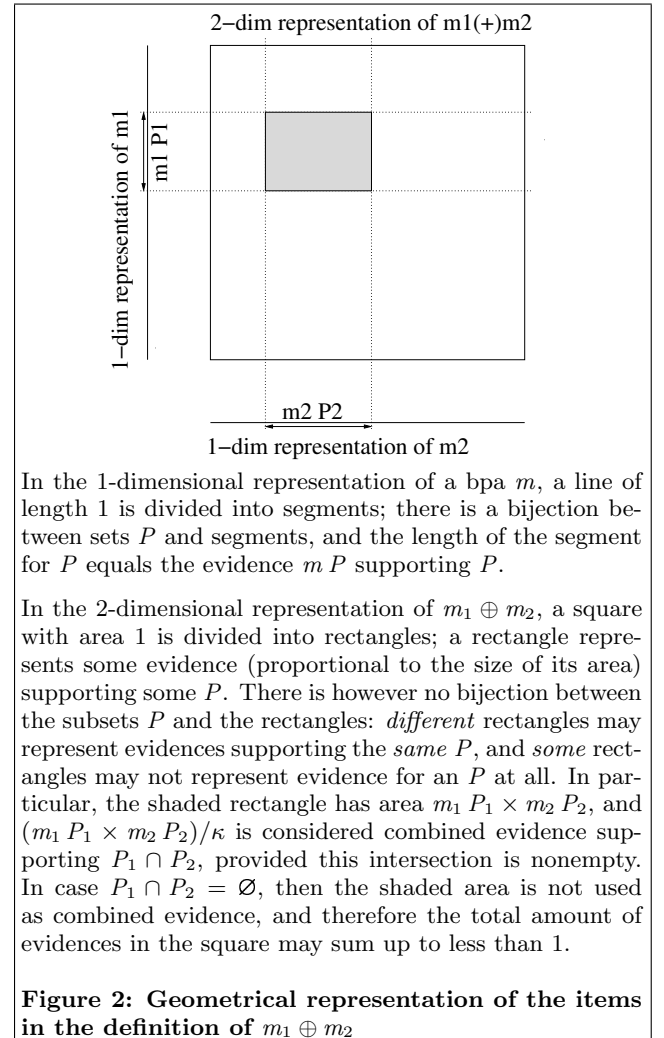
Elaborating this, we find the following:

- According to the notion of bpa, the empty set of propositions is not supported at all.

- According to ($i$) there exists a constant $c_1$ for agent 1 such that each "piece of evidence" $m_2\,P_2$ is 'combined' with each $P$ into evidence $c_1 \times m_1\,P \times m_2\,P_2$ that, according to ($ii$), supports $P \cap P_2$, provided this intersection is nonempty.

- Symmetrically, according to ($i$) there exists a constant $c_2$ for agent 2 such that each "piece of evidence" $m_1\,P_1$ is 'combined' with each $P$ into evidence $c_2 \times m_1\,P_1 \times$

$m_2\,P$ that, according to ($ii$), supports $P_1 \cap P$, provided this intersection is nonempty.

For the sake of symmetry we take $c_1 = c_2$, and define $\kappa = 1/c_1 = 1/c_2$, and we have established the defining equations for $m_1 \oplus m_2$.

Regarding the value of $\kappa$, we notice that in order that the combination is a bpa, all evidences must sum up to 1, and therefore $\kappa$ must be equal to the sum of all summands in the defining equations for $m_1 \oplus m_2$. It follows that when $\kappa = 0$, no subset $P$ of $D$ is supported by evidence of both $m_1$ and $m_2$. In that case the two agents cannot agree in accordance with the principles; the evidences that they hold are contradictory, and the combination of $m_1$ and $m_2$ simply does not exist. This concludes my intuition for the definition of $m_1 \oplus m_2$.

Figure 2 geometrically relates various items mentioned above to each other, and can be used to organize an actual computation of a bpa combination.



In the 1-dimensional representation of a bpa $m$, a line of length 1 is divided into segments; there is a bijection between sets $P$ and segments, and the length of the segment for $P$ equals the evidence $m\,P$ supporting $P$.

In the 2-dimensional representation of $m_1 \oplus m_2$, a square with area 1 is divided into rectangles; a rectangle represents some evidence (proportional to the size of its area) supporting some $P$. There is however no bijection between the subsets $P$ and the rectangles: *different* rectangles may represent evidences supporting the *same* $P$, and *some* rectangles may not represent evidence for an $P$ at all. In particular, the shaded rectangle has area $m_1\,P_1 \times m_2\,P_2$, and $(m_1\,P_1 \times m_2\,P_2)/\kappa$ is considered combined evidence supporting $P_1 \cap P_2$, provided this intersection is nonempty. In case $P_1 \cap P_2 = \varnothing$, then the shaded area is not used as combined evidence, and therefore the total amount of evidences in the square may sum up to less than 1.

**Figure 2: Geometrical representation of the items in the definition of $m_1 \oplus m_2$**

## 20 Appendix: Attack and defense

The following criticism on Dempster's combination $\oplus$ is well-known. Consider the following two agents in the context of $D = \{p_1, p_2, p_3\}$:

$$
\begin{aligned}
m_1 &= \{\{p_1\} \mapsto 0.9999, \ \{p_3\} \mapsto 0.0001\} \\
m_2 &= \{\{p_2\} \mapsto 0.9999, \ \{p_3\} \mapsto 0.0001\}
\end{aligned}
$$

So, each agent $i$ believes almost certainly in proposition $p_i$, and considers $p_3$ very unplausible. Yet, their combination gives full certainty to $p_3$, which might be considered counter-intuitive:

$$
m_1 \oplus m_2 \ = \ \{\{p_3\} \mapsto 1.0\}
$$

The defense is clear: in view of the infallibility of agents, proposition $p_3$ is the only one that can be true according to both agents together. Truth of an proposition with very low but positive plausibility is not inconsistent with an agents view of the world.

# Some formal definitions

We give some of the missing definitions; space limitations do not permit to give all.

## 21 Appendix: Schema

A pair $(A, D)$ of a set $A$ and a function $D$ that maps each $a \in A$ to a [finite] set, is called a [finite] *schema*.

Given a schema $(A, D)$, the notion of labeled products $\Pi_A D$ and $\Pi_A^{\mathbb{P}} D$ make sense, as explained in the next paragraph.

## 22 Appendix: Labeled products

Members of a product $D_1 \times \cdots \times D_n$ are called *tuples* and denoted $(x_1, \ldots, x_n)$. Unfortunately, for some manipulations the concepts of product and tuple, with the ellipses "...".-notation, do not work well (for example, the "union" and "join" are not easy to express). The formulas work out far more beautiful and manipulatable when we view a tuple $(x_1, \ldots, x_n)$ as a function $x = \{1 \mapsto x_1, \ldots, n \mapsto x_n\}$, so that $x\, i = x_i$. Correspondingly, $D$ is viewed a function $D = \{1 \mapsto D_1, \ldots, n \mapsto D_n\}$ with $D\, i = D_i$, and the role of $D_1 \times \ldots \times D_n$ is now taken over by 'the set of functions $x$ with $x\, i \in D\, i$', denoted by $\Pi_{1..n} D$. In short, we exploit the following isomorphism ($\approx$):

$$
(x_1, \ldots, x_n) \ \approx \ \{1 \mapsto x_1, \ \ldots, \ n \mapsto x_n\} \ = \ x
$$

$$
D_1 \times \cdots \times D_n \ \approx \ \begin{pmatrix} \text{the set of functions } x \\ \text{with } \forall\, i : 1..n \bullet x\, i \in D_i \end{pmatrix} \ = \ \Pi_{1..n} D
$$

Actually, we can now generalize a bit, and use arbitrary set $A$ instead of $1..n$ to label the components: the set $\Pi_A D$ is a *labeled product*, and a member of $\Pi_A D$ is an *A-labeled tuple over $D$*:

$$
\Pi_A D \ = \ \begin{aligned}&\text{the set of all functions } x \text{ with domain } A \\ &\text{satisfying } \ \forall\, a : A \bullet x\, a \in D\, a\end{aligned}
$$

*Example.* Take:

$$
\begin{aligned}
A &= \{ Name, && Age, && Sex &\} \\
D &= \{ Name \mapsto Text, && Age \mapsto Number, && Sex \mapsto \{\text{'F', 'M'}\}\} \\
r &= \{ Name \mapsto \text{'Alice'}, && Age \mapsto 13, && Sex \mapsto \text{'F'} &\} \\
r' &= \{ Name \mapsto \text{'Bill'}, && Age \mapsto 50, && Sex \mapsto \text{'M'} &\}
\end{aligned}
$$

Then, $r$ and $r'$ are $A$-labeled tuples over $D$, that is, $r, r' \in \Pi_A D$. Imposing an order on $A$, say $A = (Name, Age, Sex)$, and using the conventional tuple notation, these equations read:

$$
\begin{aligned}
A &= (Name, & Age, & \quad Sex & ) \\
D &= (Text, & Number, & \{\text{'F', 'M'}\}) \\
r &= (\text{'Alice'}, & 13, & \text{'F'} & ) \\
r' &= (\text{'Bill'}, & 50, & \text{'M'} & )
\end{aligned}
$$

Now $r, r' \in D_1 \times D_2 \times D_3$. The order on $A$, namely 'first *Name* then *Age* then *Sex*', is absent in the $A$-labeled tuples but essential in the conventional tuple notation. The conventional tuple notation *necessitates* an order on $A$ (thus forcing some overspecification), whereas the labeled products and tuples don't do so.
(*End of example.*)

Generalizing slightly, we also define labeled products and tuples that are *set-valued* (with $P_i \subseteq D_i$):

$$
(P_1, \ldots, P_n) \quad \approx \quad \{1 \mapsto P_1, \ldots, n \mapsto P_n\} \quad = \quad P
$$

$$
\mathbb{P} D_1 \times \cdots \times \mathbb{P} D_n \approx \begin{pmatrix} \text{the set of fcts } P \text{ with} \\ \forall\, i : 1..n \bullet P\, i \subseteq D_i \end{pmatrix} = \Pi_{1..n}^{\mathbb{P}} D
$$

So,

$$
\Pi_A^{\mathbb{P}} D \ = \ \begin{aligned}&\text{the set of all functions } P \text{ with domain } A \\ &\text{satisfying } \ \forall\, a : A \bullet P\, a \subseteq D\, a\end{aligned}
$$

## 23 Appendix: Domain restriction

Let $f$ be a function with domain $A$; then the *domain restriction* of $f$ to set $B$ is the function $\lambda\, a : A \cap B \bullet f\, a$, for which we introduce the abbreviation: $B \lhd f$.

Let $(A, D)$ be a finite schema, $m$ be a ct-bpa over $(A, D)$, and $B \subseteq A$. The *restriction of $m$ to $B$*, denoted $B \blacktriangleleft m$, is the ct-bpa over $(A, D)$ obtained from $m$ by changing in each entry the conditions and conclusions for $A \setminus B$ into $*$, and simultaneously also, if the change has effect on the conditions, the other conclusions of the entry:

$$
(B \blacktriangleleft m)(P \mid Q) = 0 \quad \text{if } \neg\, (\forall\, a : A \setminus B \bullet P\, a = * = Q\, a),
$$
$$
\text{else:}
$$
$$
\begin{aligned}
(B \blacktriangleleft m)(P \mid Q) = \Sigma\, P', Q' \\
| \quad (\forall\, a : B \bullet Q'\, a = Q\, a) \wedge \\
\textbf{if } Q' = Q \\
\textbf{then } (\forall\, a : B \bullet P'\, a = P\, a) \\
\textbf{else } P = \bar{*} \\
\bullet \quad m\, (P' \mid Q')
\end{aligned}
$$

For example, take $A = 1..3$ and $B = 1..2$ and consider, using the conventional tuple notation:

$$
\begin{aligned}
\{ &(a_0, b_0, * \mid a, b, *) \mapsto v, \\
&(a_0, b_0, c_0 \mid a, b, *) \mapsto w, \\
&(a_1, b_1, c_1 \mid a, b, c\ ) \mapsto x, \\
&(a_2, b_2, c_2 \mid a, b, c') \mapsto y \ \}
\end{aligned}
$$

This ct-bpa is mapped by $B\blacktriangleleft$ to the following ct-bpa:

$$
\begin{aligned}
\{ &(a_0, b_0, * \mid a, b, *) \mapsto v + w, \\
&(\, *\, ,\, *\, ,\, * \mid a, b, *) \mapsto x + y \ \}
\end{aligned}
$$