

PDF hosted at the Radboud Repository of the Radboud University Nijmegen

The following full text is a publisher's version.

For additional information about this publication click this link.

<http://hdl.handle.net/2066/113678>

Please be advised that this information was generated on 2017-12-06 and may be subject to change.

3826

ordinal
data
analysis:

biorder
representation
and
knowledge
spaces

mathieu
koppen

Ordinal Data Analysis :
Biorder Representation and Knowledge Spaces

Ordinal Data Analysis :
Biorder Representation and Knowledge Spaces

*Een wetenschappelijke proeve
op het gebied van de Sociale Wetenschappen*

Proefschrift

ter verkrijging van de graad van doctor
aan de Katholieke Universiteit te Nijmegen,
volgens besluit van het College van Decanen
in het openbaar te verdedigen op
maandag 21 augustus 1989 des namiddags
te 1.30 uur precies

door

Matheus Gerardus Marie Koppen

geboren op 16 augustus 1955 te Ospel

Promotores : Prof. dr. E. E. Ch. I. Roskam

Prof. dr. J.-C. Falmagne (University of California)

This thesis was written at New York University and the study was supported by DOD grant MDA903-87-K-0002 and AFOSR grant F49620-87-C-031 to Prof. Jean-Claude Falmagne at NYU. For the chapters 3 and 4, the research was also supported by the Netherlands Organisation for the Advancement of Pure Research, Z.W.O. grant 560-670-006 to Dr. Norman Verhelst.

The chapters 3, 4, 7, 8 and 9 have been written as self-contained papers. This is the cause of some amount of repetition in the introductory parts of various chapters. Chapter 3 has been published in the *Journal of Mathematical Psychology* and chapter 7 is to appear in the same journal. Chapters 4, 8 and 9 have been submitted for publication.

Met dank aan:

- **Norman Verhelst**, zonder wiens vasthoudendheid en persoonlijke inzet dit alles nooit van start zou zijn gegaan. Samenwerken met hem was altijd zeer plezierig en de sfeer werd niet slechter als er af en toe na de werkdag nog wat moest worden bijgepraat in de binnenstad.
- **Jean-Paul Doignon**, wiens werk de inspiratie vormde voor grote delen van deze dissertatie. Zijn invloed is onmiskenbaar, van het begin tot het eind. De trips van Utrecht en Arnhem naar Brussel waren voor mij zeer waardevol en na mijn overstap naar New York werd hij, dankzij de elektronische communicatie, niet minder bereikbaar.
- De afdeling OPD van het Cito te Arnhem voor het gastvrije onderdak.
- **Charlie Chubb** voor interessante discussies over diverse aspecten van de problemen betreffende biorde-representatie.
- **Maria Kambouri en Mike Villano** voor de prettige samenwerking in het “**knowledge assessment**” project onder leiding van Jean-Claude Falmagne.
- **Marleen**, voor onuitputtelijke kameraadschap.

CONTENTS

1. General Introduction	1
Part I : Biorde Representation	
2. From Guttman scale to biorde representation	13
3. On finding the bidimension of a relation	29
4. Finding minimal biorde extensions of a relation and maximal stable sets in the associated hypergraph	55
5. Discussion of Part I	103
Part II : Knowledge Spaces	
6. The basic theory of knowledge spaces	121
7. How to build a knowledge space by querying an expert	143
8. Extracting human expertise for constructing knowledge spaces : An algorithm	169
9. Surmise mappings and well graded knowledge spaces	193
10. Discussion of Part II	211
Summary	221
Samenvatting	225
Curriculum vitae	229

CHAPTER 1

GENERAL INTRODUCTION

Binary data are quite common in psychology and social sciences. They can come in various forms. In a mathematics test, a subject is presented with a mathematical problem and it is recorded whether she solves or fails it. In a signal detection task, a person has to indicate on each trial the presence or absence of some target stimulus. In a questionnaire on some subject matter, people are asked whether they agree or disagree with some proposition. Despite their diverse outward appearances, these examples seem to have something in common on a more abstract level. In each case, we may conceive the response to be “positive” (“correct”, “present”, “agree”) if and only if, on some underlying dimension, the person has “more” than is required by the item: a mathematical ability exceeding the (mathematical) difficulty of the problem; more “detecting” power than the “hiding” power of the stimulus; a more pronounced view on the subject matter than that expressed in the proposition.

There is a straightforward formalization of this idea. Given binary data of the above type, we want to assign numerical values to the subjects and the items, representing their respective positions on the underlying dimension. And we want to do this in such a way that if $f(a)$ is the value of subject a and $g(u)$ is the value of item u , then a “dominates” (gives a positive response to) u if and only if $f(a)$ is greater than $g(u)$. This is the classical scalogram analysis (Guttman, 1944). If a pair f, g of such scales exists, the data are said to be Guttman scalable. It must be clear that in such a situation this pair is not unique. In the above description we make only use of the relative ordering of the scale values, so any transformation of the scales that leaves this relative ordering intact yields an equally valid pair of scales. In other words, with these binary data we measure subjects and items only on an ordinal level: a Guttman scale is basically a joint ordering of subjects and items.

In this thesis, models will be considered that feature this Guttman scale as a fundamental building block. The fact is that only in very exceptional cases a set of binary data can be fully described by a single Guttman scale. Usually some extension or relaxation of the basic model is needed. A popular approach consists in dropping the deterministic character of the Guttman scale: the scale values $f(a)$ and

$g(u)$ of the subject and the item, respectively, only specify a probability for a positive response. In this way we obtain the typical psychometric models, where the functional dependence of the probability of a correct answer on the subject and item scale values generally takes a very specific form. The scale values are usually treated as if they are measured on a higher than ordinal level; only in the more interesting case of the one-parameter logistic model (Rasch, 1960) are there a formal justification for the assumption of measurement on interval level and ways of testing this assumption statistically on empirical data.

Actually, we will not be concerned with this approach here. Instead, we will study the relaxation of the model that is obtained by assuming that the situation is not governed by a single Guttman scale, but that multiple Guttman scales are involved. We will consider two elaborations of this idea and in both cases we will respect the ordinal character of binary data. That is, we will avoid making specific (parametric) assumptions, but rather investigate how far we can get with an analysis that is solely based on ordinal concepts. It must be realized that, until fairly recently, strong, numerical models were a kind of necessity, since these were the only models that could be analyzed in any detail. Such analyses are typically based on the methods of linear algebra and calculus; the corresponding algorithms consist of iterative numerical estimation procedures (or, if possible, closed formulae).

In a purely ordinal analysis, to the contrary, we have to resort to a more abstract algebra and to set-theoretic concepts. Accordingly, the algorithms are here of a combinatorial nature. Often, such algorithms are computationally very demanding. Only the most trivial cases can be solved by hand and modern, powerful computers may be needed even for moderate size problems. In both cases of multidimensional Guttman scales that we are going to consider in this thesis, we will indeed be dealing with problems whose algorithmic solution necessarily involves some sort of exhaustive search. The running time of such algorithms depends very strongly (typically, at least exponentially) on the size of the input problem, so it is critical to try and make the search as least exhaustive as possible by drawing inferences all the way. These may result in reducing markedly the number of alternatives that have to be considered at each step of the search procedure. This concern for deriving inference rules from the formal characteristics of the situation under study will be a recurring theme in this thesis.

As indicated, our mathematical tools derive mainly from the algebra of sets. Specific concepts will be defined where they are going to be used, but below we recall some fundamental notions that are at the heart of all that follows.

Some basic concepts and terminology

Given two sets X and Y , we can construct the *Cartesian product* of X and Y , denoted by $X \times Y$, which is the set of all ordered pairs (x, y) with $x \in X$ and $y \in Y$. When possible without creating confusion, an ordered pair (x, y) will more simply be written as xy . Since $X \times Y$ is just another set, we can look at subsets of this Cartesian product. Any $R \subseteq X \times Y$ is called a *relation between X and Y* (or: *from X to Y*). In case $Y = X$, that is, if $R \subseteq X \times X$, R is called a *relation on X* . We write, equivalently, xRy for $xy \in R$. (Strictly speaking, the above defines a *binary* relation, but since we will consider only binary relations, this adjective will be dropped.) The relevance of this relation concept for the analysis of binary data must be clear if we realize that any binary (say, 0-1) matrix defines a relation between the set indexing the rows and the set indexing the columns of the matrix: e.g., the relation consisting of the row-column pairs for which the corresponding entry in the matrix is a “1”.

Since relations are sets, the usual set relations (inclusion, equality, disjointness) and operations (union, intersection, set difference) are available. For instance, for any sets X and Y , the empty set \emptyset and the full product $X \times Y$ are two special relations between X and Y and any such relation includes the former and is included in the latter. If we have $R_1 \subseteq R_2$ for two relations between X and Y , we also say that R_2 *extends* (or: *is an extension of*) R_1 . The intersection of a number of relations between X and Y is the relation consisting of the pairs xy that are present in all of these relations, and so on.

The *complement (with respect to $X \times Y$)* of $R \subseteq X \times Y$ is the relation between X and Y consisting of the pairs xy that are *not* in R . If the implied full Cartesian product is clear from the context, the notation \bar{R} is often used; thus,

$$\bar{R} = (X \times Y) - R,$$

where “-” denotes set difference. (In the case of a 0-1 matrix, the relation induced by the “0” entries is the complement of the relation induced by the “1” entries.) The *converse* of $R \subseteq X \times Y$, denoted by R^{-1} , is the relation between Y and X containing the pairs yx such that $xy \in R$. Thus, $R^{-1} \subseteq Y \times X$ and

$$yR^{-1}x \quad \text{iff} \quad xRy$$

(*iff* is the usual abbreviation for “if and only if”). This notion is reminiscent of the notion of the inverse of a function, as is the adopted notation. In fact, the converse of a relation is a straightforward generalization of the inverse of a function. We have a similar generalization to relations of the concept of the composition of two functions. Again, the notation will be indicative of this connection. If R is a relation between X and Y and S is a relation between Y and Z , then the (*relative*) *product* of

R with S , denoted $R \circ S$ or simply RS , is the relation between X and Z defined by

$$xRsz \text{ iff there is } y \in Y \text{ such that } xRy \text{ \& } ySz.$$

As can be checked easily from this definition, the operation of taking the product of two relations is associative ($(RS)T = R(ST)$ whenever either side is defined) and thus we write unambiguously RST for the product of three relations, etcetera.

For a relation R on a set X there are a number of properties that will come up again and again in the following chapters. These properties have a very compact expression in terms of set relations between relations and the above notions of the converse of a relation and the product of a number of relations. We use the notation I for the *identity* relation; that is, $I = \{xx : x \in X\}$.

A relation R on X is called *reflexive* if $I \subseteq R$, that is, if for all $x \in X$

$$xRx;$$

it is *irreflexive* whenever its complement \bar{R} is reflexive. R is *symmetric* if $R = R^{-1}$, or, for all $x, y \in X$,

$$xRy \text{ iff } yRx;$$

it is *antisymmetric* if $R \cap R^{-1} = I$: for all $x, y \in X$

$$xRy \text{ \& } yRx \text{ implies } x = y,$$

and *asymmetric* if $R \cap R^{-1} = \emptyset$: for all $x, y \in X$

$$xRy \text{ implies not } yRx.$$

R is called *transitive* if $RR \subseteq R$; that is, if for all $x, y, z \in X$

$$xRy \text{ \& } yRz \text{ implies } xRz.$$

Finally, R is *complete* (or *connected*) if $R \cup R^{-1} = X \times X$, which means that

$$xRy \text{ or } yRx$$

for all $x, y \in X$.

Combinations of these properties lead to interesting classes of relations. Any relation that is reflexive, symmetric and transitive is called an *equivalence relation*. Such a relation partitions the underlying set X in a number of *equivalence classes*. A *quasi order* is any relation that is both reflexive and transitive. If a quasi order is also antisymmetric, it is a *partial order*; if it is complete, it is a *weak order*; if a quasi order is both antisymmetric and complete, it is a *linear order* or *total order*. We get *strict* versions of the partial and linear orders by considering only the irreflexive part of these relations, or, equivalently, by replacing the antisymmetry and reflexivity in the definitions by asymmetry. A strict linear order is also called a *simple order*.

Conceptually, a linear order corresponds to a simple ranking of all elements of X , without any ties. In a weak order, ties are allowed, but for any two untied elements one must still strictly precede the other. In a partial order, there are no ties, but there may be pairs of elements that are unordered. The generalization from partial orders to quasi orders is again obtained by allowing for ties. Clearly, any partial order is a quasi order and any linear order is a weak order; on the other hand, if we interpret a quasi order as a relation not on X , but on the equivalence classes of tied elements (i.e., we consider each collection of tied elements as one single element), then this relation becomes by definition antisymmetric, thus a partial order, and by this operation weak orders turn into linear orders. Because of this close connection, almost any result on partial or linear orders has a direct generalization to quasi and weak orders.

For any set X we can consider the collection of its subsets, called the *power set* of X and denoted by 2^X . Since this is again a set, we can of course consider its power set, the collection of all possible families of subsets of X , and so on, ad infinitum. We usually speak of a “collection” or “family” of sets instead of simply “set of sets”, but these are all synonyms. For any family F of subsets of X , that is, $F \subseteq 2^X$, the notations $\cap F$ and $\cup F$ denote the subsets of X that are the intersection and union, respectively, of all sets in F . For instance, if F is finite, $F = \{F_1, \dots, F_n\}$ say, then $\cap F = F_1 \cap \dots \cap F_n$. A family F of sets is said to be *closed under intersection* if for any subfamily $S \subseteq F$, $\cap S \in F$. (“Any intersection of members is again a member.”) The same remark applies, mutatis mutandis, to the term *closed under union*.

When we were discussing the various order relations on X above, we were of course, in effect, dealing with families of relations on X . As an illustration, let us apply the terminology and notation of the preceding paragraph to this case. Let $Q, P, W, L \subseteq 2^{X \times X}$ be the collections of quasi, partial, weak and linear orders on X , respectively. Then we have seen above that $L \subseteq P \subseteq Q$, $L \subseteq W \subseteq Q$ and, in fact, $L = P \cap W$. Also, it follows easily from definitions that Q and P are closed under intersection, while W and L are not (completeness is not preserved). The notation gives just a compact, but straightforward reformulation. The only difficulty that may arise when we work with collections of families of sets, etc., is that we have the same set relations and set operations on each level and we have to keep track at what level we are at each moment. For instance, note that with the binary “ \cap ” and “ \cup ” operators we remain on the same level, while their unary counterparts bring us one level down: $P \cap W = L \subseteq 2^{X \times X}$, but, e.g., $\cap Q = \cap P = I \in 2^{X \times X}$ (I is again the identity relation on X).

The compact notation is also convenient when we deal with the representation of a partial order as the intersection of a number of linear orders. Since P is closed

under intersection and $L \subseteq P$, it is in particular the case that the intersection of any number of linear orders is a partial order. It appears that we can also go the other way. Szpilrajn (1930) showed that (i) any partial order P on X can be extended to a linear order on X , and (ii) the intersection of all these linear extensions of P is P itself. In other words, if we define for $P \in \mathcal{P}$,

$$L(P) = \{L \in \mathcal{L} : P \subseteq L\},$$

then Szpilrajn's results are that $L(P) \neq \emptyset$ and

$$P = \bigcap L(P).$$

Thus, any partial order is obtained as the intersection of a number of linear orders. However, the above representation may be redundant in the sense that not all linear extensions of P are needed to obtain P as their intersection; a subcollection of $L(P)$ might suffice for this purpose. This led Dushnik and Miller (1941) to their classical definition of the *dimension* of a partial order as the minimum cardinality of a subfamily F of $L(P)$ such that $P = \bigcap F$. Clearly, the dimension of P equals 1 if and only if $P \in \mathcal{L}$. In this discussion, we may replace "partial order" by "quasi order" and "linear order" by "weak order" everywhere and then we can add the remark that we need only consider weak order extensions that have the same equivalence classes as the quasi order; other weak order extensions are bound to be redundant.

Not surprisingly, the various order relations we have considered here will play an important role in the ordinal analysis of binary data based on Guttman scales. After all, a Guttman scale corresponds to a joint ordering of subjects and items and since there are no reasons not to allow equivalence classes of items and subjects, it may be identified with a weak ordering on the union of these two sets. Especially the correspondence between a partial order and a collection of linear orders (a quasi order and a collection of weak orders) will arise in both Part I and Part II. It is to a summary description of these two parts that we now turn.

Part I : The Guttman scale and biorder representation

In Part I, we study two multidimensional extensions of the Guttman scale that were introduced in the psychological literature by Coombs and Kao in 1955 (see also Coombs, 1964). The binary response of a person to an item, say "pass" or "fail", is no longer determined by the relative position of the person and item scale values on one single dimension. Instead, a number of dimensions are invoked and the person and item have scale values on each of these dimensions. In the *conjunctive*

model the adopted response rule is such that the person passes the item if and only if on each dimension he is higher than the item; in the *disjunctive model* this is the case if and only if there is at least one dimension on which he is higher than the item.

Thus, in these models the latent structure consists of a number of Guttman scales, that is, a number of joint orderings of subjects and items, which are combined according to the above response rules. Analyzing a binary data matrix with such a model means finding a minimal number of joint subject and item orderings that, under the chosen response rule, will reproduce the data. In Coombs and Kao (1955) and in Coombs (1964) a procedure for doing this is suggested by way of a two-dimensional example. We argue that this "procedure" capitalizes on specific characteristics of the simple example, that it would in "almost every" data matrix meet with unresolved difficulties and that we need a more thorough mathematical theory in order to deal with the general problem. This mathematical underpinning was – finally – provided by Doignon, Ducamp and Falmagne in 1984. They define the *border* as the relation deriving from a Guttman scalable data matrix and recast the problem as that of finding, for an arbitrary relation (the observed data matrix), a representation as the intersection or union of a minimal number of borders.

Just the subproblem of deciding what this minimal number, the *border dimension* or *bidimension* of the relation, is, appears to be a very difficult problem in general. We present a solution for this problem in the form of a recursive procedure in which we minimize the amount of computation by applying, at each level of the recursion, a reduction mechanism to the subproblem at that level, before invoking the next recursive call. The reduction mechanism is based on the characterization by Doignon *et al.* of the bidimension as the *chromatic number* of some *hypergraph* associated with the relation. Next, we derive algorithms for producing actual border representations and we complete the algorithmic specification of the bidimension procedure. Again, the hypergraph approach plays an important role here. Finally, we discuss the close connection between the border representation problem and the problem of representing a partial order as the intersection of a number of linear orders. Doignon *et al.*'s hypergraph is a generalization of an idea of Trotter (1983) in the context of this latter problem. We also indicate some possible approaches for the case where we do not want a perfect, deterministic border representation *per se*, but are rather interested in a – low-dimensional – approximate or probabilistic solution. Problems arise mostly from the lack of uniqueness of solutions in this type of models.

Part II : The Guttman scale and knowledge spaces

Here we are also dealing with a collection of Guttman scales, but the situation is rather different. We consider a fixed set of items and a population of persons that we think of as being partitioned into a number of classes. For each class we have a Guttman scale, i.e., a joint ordering of the persons in that class and the items. Thus, the items are the same over all Guttman scales, and we have as many Guttman scales as there are classes in the population. The interpretation that we have in mind for this situation is that the persons are students, e.g. high school pupils, and the items are problems in some field of knowledge, e.g. a collection of problems in high school mathematics. In general there is some, but no complete freedom in the order in which the various problems can be mastered and the different classes of students correspond to different paths in acquiring the problems: the different Guttman scales represent the possible orderings.

We want to use this representation as the basis for knowledge assessment procedures. This kind of procedure is seen as an essential component of a computerized instruction system. Assuming that we know the different Guttman scales as far as the items are concerned, we want to use this knowledge about the possible orderings of the problems in order to determine most efficiently – i.e., by asking him a minimal number of the questions – the collection of problems an individual has mastered. This collection will be called his *knowledge state*. The family of Guttman scales places indeed restrictions on the sets of problems that are possible knowledge states, since a person has mastered *all* problems that are lower than her position on “her” Guttman scale and *none* of the problems that are higher. The collection of possible knowledge states, which is called the *knowledge structure* of the domain, is therefore given by the collection of lower sets of the various Guttman scales.

In Part II we start with an introduction to the knowledge assessment project that is based on this concept of knowledge structures (a more detailed overview can be found in Falmagne, Koppen, Villano, Johannesen and Doignon, 1989). We will discuss briefly two assessment procedures that have been developed in this framework by Falmagne and Doignon (1988a,b) and we will investigate two interesting special cases of knowledge structures. If the family of knowledge states is closed under union and intersection, it can alternatively be represented by a quasi order; if it is just closed under union, there is an alternative representation by *surmise mappings*, a generalization of quasi orders (Doignon and Falmagne, 1985). Closure under union appears to be a reasonable hypothesis in practice; a knowledge structure with this property is called a *knowledge space*. The assessment procedures start from a given knowledge structure and in Part II we focus on the problem of finding

out what this knowledge structure is in the first place. That is, we want to determine what the collection of Guttman scales is in the domain under consideration. We can think of two complementary methods; in both cases we will assume that the knowledge structure is a space. In the first instance we consult experts in the field who provide us with a knowledge space; in the presence of a sufficient amount of empirical data, this space can then be checked and possibly pruned down. A model for this second approach has been developed by Falmagne (1989); we will concentrate on the first problem, that of tapping human expertise for the purpose of constructing a knowledge space.

Since we cannot simply ask experts to give us the list of possible knowledge states – such a task would be impracticable – , we have to resort to an indirect method. We derive a second alternative representation for knowledge spaces that is well suited for this purpose. This time we will deal with quasi orders on the power set of the collection of problems. Next we present a procedure that is based on this new representation and that transforms the answers of an expert to a specific set of questions into the corresponding knowledge space. We also consider how special cases in the representation by surmise mappings relate to extra conditions on knowledge spaces. Of particular interest here is the case of a *well graded* knowledge space (Falmagne and Doignon, 1988b; Falmagne, 1989), where on all Guttman scales the items are totally ordered, instead of just weakly. We conclude with an overview of the various representations available and some remarks on on-going research and prospects in this project.

References

- Coombs, C. H. (1964). *A Theory of Data*. New York: Wiley.
- Coombs, C. H. & Kao, R. C. (1955). *Nonmetric factor analysis*. Engineering Research Bulletin No. 38. Ann Arbor: Univ. of Michigan Press.
- Doignon, J.-P., Ducamp, A. & Falmagne, J.-C. (1984). On realizable biorders and the biorder dimension of a relation. *Journal of Mathematical Psychology*, **28**, 73-109.
- Doignon, J.-P. & Falmagne, J.-C. (1985). Spaces for the assessment of knowledge. *International Journal of Man-Machine Studies*, **23**, 175-196.
- Dushnik, B. & Miller, E. W. (1941). Partially ordered sets. *American Journal of Mathematics*, **63**, 600-610.
- Falmagne, J.-C. (1989). A latent trait theory via a stochastic learning theory for a knowledge space. *Psychometrika*, in press.
- Falmagne, J.-C. & Doignon, J.-P. (1988a). A class of stochastic procedures for the assessment of knowledge. *British Journal of Statistical and Mathematical Psychology*, **41**, 1-23.
- Falmagne, J.-C. & Doignon, J.-P. (1988b). A Markovian procedure for assessing the state of a system. *Journal of Mathematical Psychology*, **32**, 232-258.

- Falmagne, J.-C., Koppen, M., Villano, M., Johannesen, L. & Doignon, J.-P. (1989). An introduction to knowledge spaces: How to build, test and search them. *Psychological Review*, to appear.
- Guttman, L. (1944). A basis for scaling qualitative data. *American Sociological Review*, **9**, 139-150.
- Rasch, G. (1960). *Probabilistic models for some intelligence and attainment tests*. Copenhagen: The Danish Institute of Educational Research.
- Szpilrajn, E. (1930). Sur l'extension de l'ordre partiel. *Fundamenta Mathematicae*, **16**, 386-389.
- Trotter, W. T., Jr. (1983). Graphs and partially ordered sets. In Beineke, L. W. & Wilson, R. J., *Selected Topics in Graph Theory 2*. London/New York: Academic Press.

Part I :

Biorder Representation

FROM GUTTMAN SCALE TO BIORORDER REPRESENTATION

1. The Guttman scale

For the case of binary data that consist of positive and negative responses of subjects to a number of "items" (which may be of various kinds) we can consider a model in which it is assumed that there is one underlying, latent dimension that mediates the responses. A positive response will be obtained if and only if the person has "more" of the quantity measured along this dimension than the item requires.

This idea is at the basis of Guttman's (1944) scalogram analysis. Let us pick the mental test situation as our typical concrete example. Then we may write aRu whenever subject a solves item u . That is, the collection of correct answers defines a binary relation R between the set A of subjects and the set D of items. The idea that a solves u if and only if a "dominates" u on the relevant dimension is then formalized as follows. We want to find scales for subjects and items, that is, mappings $f : A \rightarrow \mathbb{R}$ and $g : D \rightarrow \mathbb{R}$, where \mathbb{R} denotes the set of real numbers, such that for any $a \in A, u \in D$

$$aRu \quad \text{iff} \quad f(a) > g(u). \tag{1}$$

In our example, f would measure the subjects' ability and g the difficulty of the items. The model represented in (1) is a weak model in the sense that it is non-metric. It is clear that f and g in (1) are just ordinal scales: for any strictly increasing transformation ϕ of the reals, the pair of scales $(\phi \circ f, \phi \circ g)$ satisfies (1) whenever the pair (f, g) does. This means that a Guttman scale (f, g) is essentially a joint ordering of subjects and items. A subject solves an item if and only if the subject precedes the item in this ordering.

This shows in what sense the Guttman scale is a very strong model: the data must be accounted for by one single ordering. This puts severe restrictions on the data. Consider the data matrix where the rows are indexed by the subjects, the columns by the items, and where a "1" entry signifies that the subject solved the item, a "0" entry that he failed it. Now suppose that this matrix contains the 2x2 submatrix

$$\begin{array}{cc} & \begin{array}{c} u \ v \end{array} \\ \begin{array}{c} a \\ b \end{array} & \begin{array}{cc} 1 & 0 \\ 0 & 1 \end{array} \end{array} \tag{2}$$

and suppose that functions f and g exist, satisfying (1) for R the relation defined by the “1” entries of this matrix. Then we could derive:

$$f(a) > g(u) \geq f(b) > g(v) \geq f(a),$$

a clear absurdity. Obviously, (2) constitutes a forbidden submatrix for a Guttman scale. In fact, as Ducamp and Falmagne (1969) have shown, it is *the* forbidden submatrix: reordering the rows and columns of a binary matrix by their number of “1” entries leads to the triangular shape, typical for the Guttman scalogram, if and only if the pattern (2) does not appear.

Not surprisingly, empirical binary data typically do not display this perfect triangular shape. In such cases, a number of options are available. If the violations are minor in number and well localized (occurring in just a few response patterns having a low frequency) one might well postulate that the model is essentially valid. Even if there is an overall pattern of violations, model (1) can still be saved “in principle”. One may maintain that its only failure is that it is completely deterministic. Instead of deciding deterministically by (1) whether a subject solves an item, the scales f and g should only provide a probability for this event to occur. This is a sensible relaxation and it has become the virtually unanimous approach in test theory, leading to the heavily investigated field of Item Response Theory. (By the way, “determinism” is doomed, whenever it is judged to be a good idea to add deliberately noise to the situation, just for the sake of easy scoring: it is clear that the representation (1) does not stand a chance if the items are of the multiple choice type.)

We will, however, be concerned with a third option in the face of data that are not Guttman scalable. Incidentally, this option does not preclude – and may in fact just precede – the other two. We may well conclude that there is something more basically wrong with (1) as a model for our data; something that cannot be dealt with by just adding a random component. The central assumption of the model, that of *unidimensionality*, might be violated. If we accept that one ordering of subjects and items is not enough to explain the data, we can try to describe the data by two orderings, or by more. That is, we need a multidimensional extension of the Guttman scale. This is exactly what is provided by Coombs and Kao’s (1955) *conjunctive* and *disjunctive* models, which we describe in the next section. Their presentation was rather informal and in the following section we will discuss the rigorous mathematical treatment of this situation that was given only fairly recently by Doignon, Ducamp and Falmagne (1984). Their fundamental paper is at the basis of the investigations in the following chapters.

2. The conjunctive-disjunctive model

In a multidimensional extension of the Guttman scale, we no longer consider just one pair of subject-item scales, but rather scales $f_i : A \rightarrow \mathbb{R}$ and $g_i : D \rightarrow \mathbb{R}$ for the subjects and items, respectively, on each of n dimensions. The model is determined by the choice of a composition rule. This rule has to specify how the scale values $f_1(a), \dots, f_n(a), g_1(u), \dots, g_n(u)$ are combined to yield the observable score of subject a on item u . In the case of binary data, this score simply amounts to the decision “pass” or “fail”. Commonly used composition rules involve taking a weighted sum of the values on the underlying dimensions; this underlies, e.g., the classical linear models of analysis of variance, regression analysis and common factor analysis. This presumes, however, that measurements are at least on an interval level (otherwise, linear combinations do not make much sense). That assumption would be hard to justify in the case of binary data and a linear model would not be in the spirit of the non-metric Guttman scale.

We can conceive of an essentially different composition rule that leads to a straightforward extension of model (1). If we have n dimensions instead of just one, a reasonable assumption would be that a person solves an item if and only if she dominates it on every dimension. That is, the left hand side of (1) is true if and only if the right hand side is true for each dimension:

$$aRu \quad \text{iff} \quad (f_i(a) > g_i(u) \text{ for all } i = 1, \dots, n). \quad (3)$$

This is the *conjunctive* model as defined by Coombs and Kao (1955). On each dimension i we have a pair of scales (f_i, g_i) that is purely ordinal and that corresponds essentially to a joint ordering of subjects and items on that dimension. A subject solves an item whenever he precedes it in all of these n orderings.

This conjunctive rule readily suggests an alternative, in which an item is solved whenever the person dominates it on at least one dimension (precedes it in at least one of the n orderings). This is the *disjunctive* model

$$aRu \quad \text{iff} \quad (f_i(a) > g_i(u) \text{ for some } i = 1, \dots, n). \quad (4)$$

Although the disjunctive model has psychologically quite a different interpretation, it is, as Coombs and Kao note, formally isomorphic to the conjunctive model. Taking the negation of both sides of (4) shows that applying the disjunctive model is equivalent to applying the conjunctive model after flipping the binary data (considering the complement of R instead of R) and reversing the direction of the dimensions (orderings). Consequently, for the formal development only one model needs to be considered and this will be, rather arbitrarily, the conjunctive model.

Note that the conjunctive-disjunctive model is essentially distinct from a linear type of model in that it is non-compensatory. While in any linear model one can make up for a deficiency on one dimension by a surplus on another (and this in principle without limit), the conjunctive-disjunctive model involves, rather, thresholds on each dimension. And no matter how amply a threshold is surpassed on one dimension, this will, in the conjunctive model, not compensate for failing the threshold on another dimension. In many situations this seems, indeed, a more reasonable assumption than a (linear) compensating mechanism. Similarly, the disjunctive model formalizes the situation where failing the threshold badly on a number of dimensions does not hurt, as long as there is one dimension on which the threshold is passed, however narrowly. The use of logical (threshold) composition rules was already advocated by Johnson (1935) in the field of aptitude testing. He suggested a very general class of models, in which composition rules can consist of any logical combination of the dimensions. For instance, an item is solved whenever the person dominates it on the first dimension *or* on both the second *and* the third dimension. Coombs and Kao's conjunctive-disjunctive model is the only special case of this class that has actually been developed.

In Coombs and Kao (1955) – as in Coombs (1964), where chapter 12 covers essentially the same material – we find only a sketch of a procedure for constructing a conjunctive model representation of a binary data matrix; a sketch consisting in going through a small example. Before indulging in the mathematical analysis of Doignon, Ducamp and Falmagne (1984) in the next section, we illustrate here Coombs and Kao's approach on a tiny example indeed. The example is big enough, however, to show the kind of reasoning and the kind of difficulties involved and it will lead to an observation that formed the starting point for the investigations in Chapter 3.

Suppose we want to find a conjunctive model representation of the data in Table 1, concerning seven response patterns of subjects a to g on three items u, v, w . Of course, we want a representation in the minimum number of dimensions. Coombs and Kao's (1955) and Coombs' (1964) approach to achieving this can be summarized as follows: order the rows (response patterns) in increasing order of number of "0" entries (this has already been done in Table 1) and next consider each row in turn, from the top down, drawing inferences for orderings of the items if possible and introducing new dimensions if necessary. The obtained item orderings can then be extended to appropriate joint subject-item orderings.

Starting with row a , we see immediately that this does not tell us anything about an ordering of the items: the pattern simply does not discriminate between items. On any dimension we want to invoke, we will have to place subject a above all items. Since in pattern b only item w is failed, it is clear that we need a dimension

	<i>u</i>	<i>v</i>	<i>w</i>	orderings of items
<i>a</i>	1	1	1	—
<i>b</i>	1	1	0	$X_1 = [w (u, v)]$
<i>c</i>	0	1	1	$X_1 = [w (u, v)]$ $X_2 = [u (v, w)]$
<i>d</i>	1	0	0	$X_1 = [w v u]$ $X_2 = [u (v, w)]$
<i>e</i>	0	1	0	$X_1 = [w v u]$ $X_2 = [u (v, w)]$
<i>f</i>	0	0	1	$X_1 = [w v u]$ $X_2 = [u v w]$
<i>g</i>	0	0	0	$X_1 = [w v u]$ $X_2 = [u v w]$

Table 1. A hypothetical data matrix of 7 subjects on 3 items. The last column contains the gradual construction à la Coombs and Kao of the dimensions in a conjunctive model representation (items in parentheses are, as yet, unordered). See text.

(ordering) X_1 on which item w is above (precedes) the other two items. Subject b must come immediately below w on X_1 and can be put above all items on possible other dimensions. Similarly, row c implies the existence of a dimension on which item u is highest. This, clearly, cannot be X_1 , so we are forced to introduce a new dimension X_2 , with u highest on X_2 . Subject c can be put immediately below u on X_2 and above all items on X_1 and possible other dimensions. We can deal with row d without invoking another dimension. Since u is solved, d must be placed above u on X_2 . But then, u solved and v failed implies that v is above u on X_1 , with d inbetween. Next, we see that row e contains no further information: failing u and w , while solving v can already be accounted for by placing this subject immediately below w on X_1 and immediately below u on X_2 . The inferences for f are a mirror image of those for row d : with only w solved, f must be above w on X_2 and between v and w on X_1 , where v must be above w . Of course, row g is again uninformative. It can be given a place in any representation; everything is all right, for instance, as long as it is at the bottom of any one dimension. We thus arrive at the following two-dimensional solution:

$$\begin{aligned}
 X_1: & (a, c, f) w (b, e) v d u g \\
 X_2: & (a, b, d) u (c, e) v f w g,
 \end{aligned}
 \tag{5}$$

where the ordering between parentheses is arbitrary, and where g , rather arbitrarily, is placed at the bottom of *both* dimensions.

Manifestly, we cannot recover a complete strict ordering of the subjects on the separate dimensions. This is not surprising in view of the fact that such a strict ordering can only be obtained on the basis of distinct positions with respect to at least some item. Clearly, k items can only discriminate between $k+1$ classes of

	<i>u</i>	<i>v</i>	<i>w</i>
<i>a</i>	1	1	1
<i>b</i>	1	1	0
<i>c</i>	1	1	1
<i>d</i>	1	0	0
<i>e</i>	1	1	0
<i>f</i>	1	1	1
<i>g</i>	0	0	0

X_1

	<i>u</i>	<i>v</i>	<i>w</i>
<i>a</i>	1	1	1
<i>b</i>	1	1	1
<i>c</i>	0	1	1
<i>d</i>	1	1	1
<i>e</i>	0	1	1
<i>f</i>	0	0	1
<i>g</i>	0	0	0

X_2

Table 2. The two factor matrices of the matrix in Table 1 corresponding to the orderings X_1 and X_2 of Eq. (5).

subjects on each dimension. Here, we get on both dimensions complete strict orderings for the smaller set, the items, but that is by no means true in general. In Table 1, for instance, it can be seen that, in the absence of row *f*, we would not be able to decide on the ordering between items *v* and *w* on dimension X_2 . Elements between parentheses are considered as being equivalent (on that particular dimension). In this way, a dimension is identified with a *weak order* on the subjects and items, the generalization of a total (linear) order in which non-trivial equivalence classes are allowed (cf. the General Introduction).

Table 2 shows the two obtained orderings X_1 and X_2 , represented as ‘‘factor matrices’’. Note that both matrices display the triangular structure characterizing a Guttman scale and that the ‘‘observed’’ matrix in Table 1 is indeed the conjunction, that is, the element-by-element product, of the matrices in Table 2.

This looks all very neat. It seems like we have here a constructive scaling procedure for the conjunctive model that yields the minimum dimensionality as an automatically obtained by-product. As soon as one tries to turn the above sketch of a procedure into a general, explicit algorithm, however, one immediately runs into trouble. It appears that the example of Table 1 is a special, constructed example, like those in Coombs and Kao (1955) and Coombs (1964). In general, things are not so straightforward. Consider, for instance, the data in Table 3. This matrix is just a submatrix of Table 1 and it is not too difficult to find the two-dimensional solution that is the restriction of (5):

$$\begin{aligned} X_1: & f w e v d u \\ X_2: & d u e v f w . \end{aligned} \tag{6}$$

	<i>u</i>	<i>v</i>	<i>w</i>	orderings of items
<i>d</i>	1	0	0	$X_1 = [(v, w) u]$
<i>e</i>	0	1	0	$X_1 = [w v u]$ $X_2 = ["u \text{ above } v; w \text{ arbitrary}"]$
<i>f</i>	0	0	1	$X_1 = [w v u]$ $X_2 = [u v w]$

Table 3. Another hypothetical data matrix, submatrix of Table 1, with a possible "derivation" of a conjunctive model representation.

A possible way of arriving at this solution is sketched in Table 3 under the column "orderings of items". However, this solution results from applying heuristic rules, rather than some algorithm. This is obvious, once we notice the symmetry in Table 3. If (6) is a solution, then any permutation, simultaneously applied to the triples (*d, e, f*) and (*u, v, w*), must give an equally valid solution. Thus we find two more representations in two dimensions, essentially different from (6) and from each other:

$$\begin{aligned} X_1' &: e v d u f w \\ X_2' &: f w d u e v \end{aligned} \quad (7)$$

and

$$\begin{aligned} X_1'' &: e v f w d u \\ X_2'' &: d u f w e v. \end{aligned} \quad (8)$$

We would have found the representations (7) or (8) if we had applied the same "rules" as in Table 3, but processing the rows in a different order. We know that the dimensionality of 2 is minimal for Table 3 only because it is obvious that the matrix is not one-dimensional; it is not guaranteed by the method. In Table 3, for instance, we "concluded" from an inspection of the first row, *d*, that there must be a dimension on which item *u* is lowest. Solution (7), however, shows that this is not necessarily true. Here, this assumption did no harm, but in a more complex situation such an unwarranted conclusion might well be the decisive step in missing the minimum dimensionality.

The difficulties encountered in the analysis of Table 3 appear to be more typical than the apparent smoothness of the analysis of Table 1. In fact, in "almost every" binary data matrix (however one wants to define this notion) it will be the case that (i) the minimal dimensionality can only be found by an – in principle exhaustive – method of trial-and-error, and (ii) there will be various essentially different solutions in this minimum dimensionality. The minimum dimensionality problem is complex in a very precise, technical sense of the word. To establish this, we need the mathematical theory developed by Doignon, Ducamp and Falmagne (1984) that we

discuss in the next section. The above discussion and examples give an impression of the ideas behind the conjunctive-disjunctive model, but they also show how intuitive and informal the approach was until this paper appeared. One thing that emerged in the description above (although we did not stress it much until now), and that will not be directly visible in the next section, is the fact that there may be patterns (like row e in Table 1) that are non-trivial, but that, nevertheless, do not give us any additional information. This observation will be important in the sequel.

3. Representing a relation by biorders

A complete mathematical analysis of the conjunctive-disjunctive model was presented only in 1984, in a paper by Doignon, Ducamp and Falmagne. They considered the most general situation, where there are no restrictions on the cardinalities of the sets A and D . We will follow that approach here, since it shows clearly which characteristics and concepts are essential. In fact, the more general situation forces more direct arguments. Where needed we will give the specialization to the finite case that is ultimately of interest to us.

Doignon, Ducamp and Falmagne (1984) define a *biorder* between A and D to be a relation $B \subseteq A \times D$ such that

$$aBu \ \& \ b\bar{B}u \ \& \ bBv \ \text{implies} \ aBv, \quad (9)$$

where we introduce the notation \bar{B} for the complement of B : $\bar{B} = (A \times D) - B$. In terms of relative products of relations, this condition can be succinctly expressed as $B\bar{B}^{-1}B \subseteq B$. A slight rearrangement of (9) leads to

$$\text{not} (aBu \ \& \ bBv \ \& \ a\bar{B}v \ \& \ b\bar{B}u), \quad (10)$$

which makes it clear that a biorder is characterized by the forbidden submatrix (2). In other words, a binary data matrix is Guttman scalable if and only if the corresponding relation, defined by the positive answers, is a biorder. This was already noted by Ducamp and Falmagne (1969) in their axiomatization of the Guttman scale. By the way, relations defined by (9) have been introduced in a different context under the name of "Ferrers relations" (Riguet, 1951; see also the survey paper by Monjardet, 1978).

Since biorders correspond to Guttman scales, representing a binary data matrix according to the conjunctive model amounts to writing a relation R as the intersection of a number of biorders. Indeed, rewriting the right hand side of (3), using (1) with R replaced by B_i , yields

$$aRu \ \text{iff} \ (aB_iu \ \text{for all } i = 1, \dots, n),$$

or, more compactly,

$$R = \bigcap_{i=1}^n B_i.$$

Analogously, the disjunctive model corresponds to a representation of R as a union of biororders. Since the complement of a biorder is again a biorder (this is immediate from (10)), the set-theoretic equivalence

$$R = \bigcap_{i \in I} B_i \quad \text{iff} \quad \bar{R} = \bigcup_{i \in I} \bar{B}_i \quad (11)$$

is a restatement of the duality between the conjunctive and the disjunctive model. Here and in what follows, I is an index set of arbitrary cardinality, denoted by $|I|$. For the finite case, I may be identified with the set $\{1, \dots, |I|\}$. We will deal with the intersection representation that corresponds to the conjunctive model.

Since for any $au \in A \times D$ the relation $(A \times D) - \{au\}$ is a biorder (a violation of (10) requires at least two pairs in the complementary relation), we can write any relation R between A and D as the intersection of the biororders $(A \times D) - \{au\}$ with $au \notin R$. This justifies the definition of the *biorder dimension* or, shorter, the *bidimension* of R as the minimal number of biororders needed for such a representation. This number will be denoted by *Bidim* R :

$$\text{Bidim } R = \text{Min} \{ |I| : R = \bigcap_{i \in I} B_i, B_i \text{ biorder} \}. \quad (12)$$

Strictly speaking, we should call this the \cap -bidimension and define a \cup -bidimension in the same way, replacing intersection by union. By (11), however, any statement on the \cap -dimension translates into a dual statement on the \cup -dimension. Clearly, the \cap -bidimension (\cup -bidimension) of a relation R is precisely the minimum dimensionality in the conjunctive (disjunctive) model of the binary data represented by R . From the argument leading to definition (12) it follows that with finite A and D , *Bidim* R will also be a finite number.

We will now discuss a main result of Doignon, Ducamp and Falmagne (1984). Generalizing an approach of Trotter (1983) in the context of the order dimension of a partial order, they established an equivalence between the bidimension and the chromatic number of some hypergraph defined in terms of the relation. Let us first recall what is meant by these terms (see, e.g., Berge, 1973, for more). A *hypergraph* $H(V)$ is a set V of elements called *vertices*, together with a system of subsets of V called *edges*. Any edge has at least two elements. A hypergraph thus generalizes an ordinary, undirected graph, where all edges contain exactly two vertices. A subset of V is called *stable* in $H(V)$ if it does not contain an edge and a *coloring* of the hypergraph $H(V)$ is a partition of V into a number of stable sets, the *colors*. The *chromatic number* of $H(V)$, then, which is denoted by *Chrom* $H(V)$, is the minimal

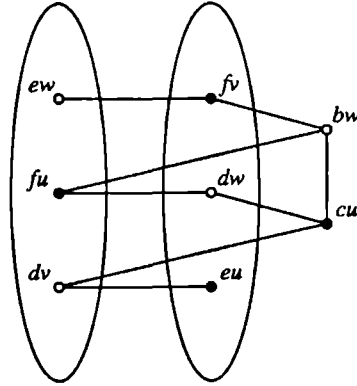


Figure 1. Example of hypergraph on eight vertices (the labeling is explained later in the text). There are eight 2-edges (the pairs of vertices connected by straight lines) and two 3-edges (the triples of vertices enclosed in ellipses). The chromatic number equals 2: the partitioning of the vertices in open and filled circles constitutes a minimal coloring.

number of colors needed in a coloring of $H(V)$. Figure 1 illustrates these concepts.

We will give here our own derivation of this equivalence, which differs somewhat – in wording, not in essence – from the presentation in Doignon, Ducamp and Falmagne (1984). We begin with a couple of rather trivial rewritings of (12). First, we use (11) and the fact that the complement of a biorder is a biorder to turn (12) into

$$\text{Bidim } R = \text{Min} \{ |I| : \bar{R} = \bigcup_{i \in I} B_i, B_i \text{ biorder} \}. \quad (13)$$

Now call $C \subseteq \bar{R}$ a *feasible* set if and only if C is the subset of a biorder contained in \bar{R} . Then (13) is equivalent to

$$\text{Bidim } R = \text{Min} \{ |I| : \bar{R} = \bigcup_{i \in I} C_i, C_i \text{ feasible} \}. \quad (14)$$

Indeed, any biorder in (13) is contained in \bar{R} , thus a feasible set; on the other hand, if we have the representation in (14), we let B_i be a biorder such that $C_i \subseteq B_i \subseteq \bar{R}$ and we obtain

$$\bar{R} = \bigcup_{i \in I} C_i \subseteq \bigcup_{i \in I} B_i \subseteq \bar{R}.$$

The results do not seem very impressive so far, but feasible sets have one property biorders do not have: any subset of a feasible set is again feasible. This means that we might just as well assume that the feasible sets in (14) are disjoint:

$$\text{Bidim } R = \text{Min} \{ |I| : \bar{R} = \sum_{i \in I} C_i, C_i \text{ feasible} \}, \quad (15)$$

where the sigma sign indicates taking the union over mutually disjoint sets. Clearly, the representation in (15) is a special case of (14). From (14) we can get to (15), however, by fixing some simple order $<$ on I and defining, for $i \in I$:

$$C'_i = C_i - (\cup_{j < i} C_j).$$

Any C'_i is a subset of C_i , thus feasible; all C'_i are mutually disjoint and there are no more sets C'_i than there were sets C_i .

Now the right hand side of the equality in (15) is the definition of the chromatic number of a hypergraph, if we replace "feasible" by "stable". That is, we have obtained

$$Bidim R = Chrom H(\bar{R}), \tag{16}$$

where $H(\bar{R})$ is the hypergraph with vertex set \bar{R} and where the stable sets are the subsets of \bar{R} that can be extended to a biorder contained in \bar{R} . The maximal stable sets are, thus, the maximal biorders in \bar{R} and the edges of $H(\bar{R})$ are the (minimal) subsets of \bar{R} that do not have such an extension.

This equivalence is not very helpful, however, unless we can find a more intrinsic characterization of the collection of stable (feasible) sets, or, equivalently, the collection of edges. To this end, we define the following relation Γ_R on \bar{R} . For $au, bv \in \bar{R}$,

$$au \Gamma_R bv \quad \text{iff} \quad bRu.$$

The situation $au \Gamma_R bv$ corresponds to the following submatrix:

$$\begin{array}{c} a \\ b \end{array} \begin{array}{|c|} \hline u \quad v \\ \hline 0 \\ \hline 1 \quad 0 \\ \hline \end{array}$$

where the open entry is arbitrary. An n -cycle in Γ_R is a sequence $a_1u_1, a_2u_2, \dots, a_nu_n$ in \bar{R} , such that

$$a_1u_1 \Gamma_R a_2u_2 \Gamma_R \dots \Gamma_R a_nu_n \Gamma_R a_1u_1.$$

A 4-cycle in Γ_R , for instance, corresponds to the following submatrix:

$$\begin{array}{c} a \\ b \\ c \\ d \end{array} \begin{array}{|c|} \hline u \quad v \quad w \quad x \\ \hline 0 \qquad \qquad \qquad 1 \\ \hline 1 \quad 0 \\ \hline \qquad 1 \quad 0 \\ \hline \qquad \qquad 1 \quad 0 \\ \hline \end{array}$$

Again, open entries are arbitrary. The special case $n=2$ shows that cycles in Γ_R are

related to the biorder concept:

$$\begin{array}{c} \\ a \\ b \end{array} \begin{array}{c} u \quad v \\ \hline 0 \quad 1 \\ 1 \quad 0 \end{array}$$

Comparing this with the forbidden submatrix (2), we immediately conclude:

$$(\Gamma 1) \quad R \text{ is a biorder} \quad \text{iff} \quad \Gamma_R \text{ has no 2-cycle.}$$

It is, however, not difficult to see that if Γ_R has an n -cycle, then it must have an $(n-1)$ -cycle or a 2-cycle. (In the above 4x4 submatrix, the pair aw , for instance, is either in R or in \bar{R} .) Thus, by induction,

$$(\Gamma 2) \quad R \text{ is a biorder} \quad \text{iff} \quad \Gamma_R \text{ is acyclic.}$$

It can be established that if M is a maximal subset of \bar{R} such that Γ_R has no cycle in M , then M is a biorder. This is equivalent to the following crucial generalization of ($\Gamma 2$):

$$(\Gamma 3) \quad C \subseteq \bar{R} \text{ is feasible} \quad \text{iff} \quad \Gamma_R \text{ is acyclic on } C.$$

(We get ($\Gamma 2$) from ($\Gamma 3$) by taking $C = \bar{R}$.) This is the characterization of the stable sets of $H(\bar{R})$: a subset of \bar{R} is stable if and only if it does not contain a cycle in Γ_R . In other words, the hypergraph $H(\bar{R})$ in (16) has as vertex set \bar{R} and as collection of edges the cycles in Γ_R .

Note that, while ($\Gamma 1$) shows that for biorderhood it is sufficient to check for 2-cycles in Γ_R , this is no longer the case for feasible sets. If we assume that in the above 4x4 submatrix all open entries are zero, then the subset of vertices $\{au, bv, cw, dx\}$ consists of a 4-cycle, but does not contain any 3- or 2-cycle. From this example it is clear how to construct, for arbitrary n , a relation R and a subset of n vertices of $H(\bar{R})$ consisting of an n -cycle in Γ_R while not containing any k -cycle for $k < n$.

With the above characterization of $H(\bar{R})$ in mind, it may be checked that the hypergraph of Fig. 1 is exactly the hypergraph associated with the relation defined in Table 1. The pairs gu, gv, gw are left out of the hypergraph of Fig. 1, since they are not contained in any edge and thus irrelevant to the chromatic number. Note that the coloring in Fig. 1 corresponds exactly to the decomposition of the matrix into the two factor matrices of Table 2. This reflects the special character of the Table 1 matrix. Here, there is essentially – that is, disregarding the trivial zero pattern g – just one minimal coloring and in that coloring the color classes are themselves already maximal biorders in \bar{R} , so there is just one pair of biorders covering \bar{R} . Consequently, the complementary biorders constitute the (essentially) unique two

dimensional solution for R that is given in (5). In general, there may be several distinct minimal colorings, and color classes may have various extensions to maximal stable sets. This means there are multiple ways of covering \bar{R} by a minimal number of biorders in \bar{R} and, thus, by complementation, multiple ways of writing R as the intersection of a minimal number of biorders.

4. The following chapters

The reinterpretation of the bidimension of a relation as the chromatic number of some explicitly defined hypergraph is an important constructive result. This does not mean, however, that the problem of finding the bidimension of a relation – thus the minimum dimensionality for the conjunctive model – is easy. To the contrary; we have a very precise and rather disappointing result in terms of computational complexity theory. As was already noted by Cogis (1982), finding the bidimension is polynomially equivalent to finding the usual order dimension of a partial order. By a rather recent result of Yannakakis (1982), the latter problem is NP -complete for a dimension exceeding 2. (See, e.g., Garey and Johnson, 1979, or Papadimitriou and Steiglitz, 1982, for a detailed exposition of the above concepts; we discuss some of the consequences below.)

The NP -completeness of the bidimension problem is theoretically bad news: it means that, in all probability, there will be no algorithm computing the bidimension in an order of time that is a polynomial function of input size. The practical consequences of this fact are, however, not always as clear-cut. For data of moderate size, the differences may not be dramatic. Moreover, polynomial time algorithms may, in practice, be not very efficient whenever the degree of the polynome exceeds, say, 2 or 3. A good example is the Ellipsoid algorithm developed by Shor, Judin and Nemirovskii for non-linear optimization, but which Khachiyan proved to provide a polynomial algorithm for the Linear Programming problem (see, e.g., Papadimitriou and Steiglitz, 1982, for more details and references). This was a major theoretical break-through in that it showed that the Linear Programming problem is solvable in polynomial time. It is, however, hardly applied, since the old Simplex algorithm, though not polynomial and with a terrible worst case performance, is, for all practical purposes, considerably faster.

The NP -completeness does indicate that the practical efficiency of an algorithm for the bidimension will depend strongly on the size of the input and that the solution will involve some kind of exhaustive search. In devising such an algorithm it is thus essential to maximally reduce the input size and to maximally prune the search tree that is traversed in the algorithm. Chapter 3 describes a procedure for doing exactly

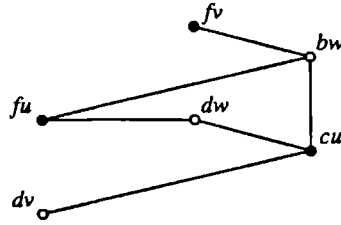


Figure 2. The subhypergraph (in this case, a simple graph), induced by removing vertices ew and eu from the hypergraph in Fig. 1.

this. Starting point was an observation in the informal discussion of the conjunctive model. In the example of Table 1, we saw that there was a non-trivial pattern (row e) that did not provide any new information, but was, in a sense, already implied by other patterns. In terms of the data matrix, the implication is that the dimensionality of the matrix is the same as that of the submatrix with row e removed. Or, for the relation: the bidimension of the relation is the same as the bidimension of the restriction of this relation to $(A - \{e\}) \times D$. With the help of the equivalence (16) we get, finally, a version in terms of hypergraphs. The hypergraph corresponding to Table 1 – that is, the hypergraph in Fig. 1 – has the same chromatic number as the subhypergraph that is induced by removing all vertices that have the subject e as a member and all edges that contain such vertices. (Figure 2 shows the resulting subhypergraph in this case.) The idea of Chapter 3 is to describe the most general conditions under which hypergraphs can be reduced to subhypergraphs without changing the chromatic number. Removal of patterns like row e of Table 1 will appear to be just a very special case of the kind of reduction that is possible. These reductions can be used throughout in a recursive procedure for computing the chromatic number of the original hypergraph, thus, by (16), the bidimension of the original relation.

Chapter 4 deals with obtaining actual representations in the minimum dimensionality. Algorithms are developed for generating minimal border extensions of a relation R . (Any border B_i in the representation $R = \bigcap B_i$ clearly contains R and the minimal extensions are those that are closest to R , i.e., for which the difference $B_i - R$ is minimal.) This is directly related to the question of generating maximal stable sets of the hypergraph $H(\bar{R})$. This is then generalized to an algorithm for maximal stable sets that contain some specified stable set. This is interesting, since it can be combined with the procedure for determining the bidimension, as described in Chapter 3, to produce not just the bidimension but, at the same time, some representations of R in Bidim R borders. The point is that this procedure is constructive; that is, once the bidimension is known, we are in fact

given some minimal coloring of the hypergraph. This allows us to use this last algorithm to produce maximal stable set extensions of the various color classes in this minimal coloring. As we have seen, maximal stable sets of $H(\bar{R})$ are biorders contained in \bar{R} and if we have an extension for each color class, then these biorders clearly cover \bar{R} . As a consequence, R can be written as the intersection of the complement biorders. Last, but of course not least, a version of the algorithms in Chapter 4 is used to compute the bidimension in the first place. In Chapter 3, a recursion formula is presented for the computation of the bidimension, but it is not specified how to actually construct the implied search tree. In Chapter 4 is described how this can be achieved, thus completing the description of the bidimension algorithm.

References

- Berge, C. (1973). *Graphs and Hypergraphs*. Amsterdam: North-Holland.
- Cogis, O. (1982). On the Ferrers dimension of a digraph. *Discrete Mathematics*, **38**, 47-52.
- Coombs, C. H. (1964). *A Theory of Data*. New York: Wiley.
- Coombs, C. H., & Kao, R. C. (1955). *Nonmetric factor analysis*. Engineering Research Bulletin No. 38. Ann Arbor: Univ. of Michigan Press.
- Doignon, J.-P., Ducamp, A., & Falmagne, J.-C. (1984). On realizable biorders and the biorder dimension of a relation. *Journal of Mathematical Psychology*, **28**, 73-109.
- Ducamp, A., & Falmagne, J.-C. (1969). Composite measurement. *Journal of Mathematical Psychology*, **6**, 359-390.
- Garey, M. R., & Johnson, D. S. (1979). *Computers and Intractability: A Guide to the Theory of NP-completeness*. San Francisco: Freeman.
- Guttman, L. (1944). A basis for scaling qualitative data. *American Sociological Review*, **9**, 139-150.
- Johnson, H. M. (1935). Some neglected principles in aptitude testing. *American Journal of Psychology*, **47**, 159-165.
- Monjardet, B. (1978). Axiomatiques et propriétés des quasi-ordres. *Mathématiques et Sciences Humaines*, **63**, 51-82.
- Papadimitriou, C. H., & Steiglitz, K. (1982). *Combinatorial Optimization : Algorithms and Complexity*. Englewood Cliffs NJ: Prentice-Hall.
- Riguet, J. (1951). Les relations de Ferrers. *Comptes rendus des Séances de l'Académie des Sciences (Paris)*, **232**, 1729-1730.
- Trotter, W. T., Jr. (1983). Graphs and partially ordered sets. In Beineke, L. W., & Wilson, R. J., *Selected Topics in Graph Theory 2*. London/New York: Academic Press.
- Yannakakis, M. (1982). The complexity of the partial order dimension problem. *SIAM Journal on Algebraic and Discrete Methods*, **3**, 351-358.

CHAPTER 3

ON FINDING THE BIDIMENSION OF A RELATION

(Journal of Mathematical Psychology, 31, 1987)

On Finding the Bidimension of a Relation

M. G. M. KOPPEN

University of Utrecht, Utrecht, The Netherlands

A method is presented for evaluating the bidimension of a finite binary relation, i.e., the number of borders (Guttman relations) needed to yield the relation as their intersection. In case the relation is induced by a binary data matrix, the bidimension equals the minimal number of dimensions needed for a representation of the data matrix according to the conjunctive model of C. H. Coombs and R. C. Kao (*Nonmetric factor analysis*, Engineering Research Bulletin No 38, Univ. of Michigan Press, Ann Arbor, 1955). Central to the evaluation of the bidimension is its characterization, provided by J.-P. Doignon, A. Ducamp, and J.-C. Falmagne (*Journal of Mathematical Psychology*, **28**, 73-109, 1984), as the chromatic number of some associated hypergraph. A procedure is described to reduce hypergraphs of this kind to subhypergraphs with the same chromatic number. This reduction can be used throughout in applying a recurrence relation that expresses the chromatic number of a hypergraph in terms of the chromatic numbers of some of its subhypergraphs. © 1987 Academic Press, Inc.

1. INTRODUCTION

In many settings of psychological research and testing the situation is so complex that it is highly unrealistic to expect that the behaviour of the subjects can be explained satisfactorily by one single dimension along which they and the experimental stimuli ("items") vary. In such cases, a multidimensional model is called for. An all-important aspect of such a model is its composition rule, by which positions on the separate, unobserved dimensions combine to give a "score" that is directly related to the observed behaviour. The most widely used composition rule consists in mapping positions on the separate dimensions into the real number line and expressing the observed score as a linear combination (weighted sum) of these scores on the separate dimensions (e.g., factor analysis, analysis of variance, regression analysis). A model with this kind of composition rule can be considered as a compensatory one: a deficiency on one dimension can be compensated for (and without limit) by a surplus on another dimension. In many situations such a rule is, again, not very realistic, especially in the case of binary data that are instances of the general type: subject solves/fails item. In such cases an essentially different composition rule is conceivable in which "a person solves an item" (these are all used as generic terms) if and only if on each of the separate dimensions the position of the

This research was supported by the Netherlands Organisation for the Advancement of Pure Research (Grant No 560-670-006). The author is grateful to J.-P. Doignon for his useful comment on a previous draft and to a referee for supplying the data set of Section 4.

0022-2496/87 \$3.00

Copyright © 1987 by Academic Press, Inc.
All rights of reproduction in any form reserved.

person dominates the position of the item. In this multidimensional generalization of Guttman's (1944) idea there is no compensating mechanism involved; rather, an item corresponds to a threshold on each dimension that has to be surpassed by a person in order to solve the item. The observed score, then, is not seen as an arithmetical sum of the scores on the separate dimensions, but rather as the logical product of these scores. This explains the name conjunctive model (Coombs and Kao, 1955; Coombs, 1964) for models with this composition rule.

The conjunctive composition rule immediately suggests another one, in which a person solves an item if and only if the position of the item on at least one dimension. With this rule the observed score corresponds to the logical sum of the scores on the separate dimensions, so this model is called disjunctive. Thanks to the logical equivalence

$$(A_1 \text{ or } \cdots \text{ or } A_n) \text{ iff not } ((\text{not } A_1) \text{ and } \cdots \text{ and } (\text{not } A_n)),$$

the disjunctive model needs no separate consideration. By flipping the binary data and reversing the direction of each dimension it is directly translated into the conjunctive form.

An informal discussion of how to construct for some given data matrix a representation according to the conjunctive model can be found in Coombs and Kao (1955) and Coombs (1964). In order to develop an explicit general algorithm, however, one has to face two important problems. The first one is that of finding the minimum dimensionality needed for such a representation (note that for the present we are considering deterministic models only) and the second has to do with the uniqueness of an obtained representation in the minimum dimensionality. This paper will be concerned with the problem having logical priority, finding the minimum dimensionality.

That this is not a trivial problem at all is shown by Doignon, Ducamp, and Falmagne (1984). In their fundamental paper they cast the problem in terms of representing the relation between the set of persons and the set of items defined by the binary data matrix as the intersection of a minimal number of biorders (Guttman relations), this number being its so-called bidimension. Some of their results were, for the finite case, independently obtained by Cogis (1980, 1982). Regarding computational complexity (see, e.g., Garey and Johnson, 1979), Cogis (1982) showed that the problem of finding the bidimension of a relation is polynomially equivalent to that of determining the dimension of a partial order (i.e., the minimal number of linear orders to yield the partial order as their intersection; see Dushnik and Miller, 1941), and Yannakakis (1982) showed the latter problem to be NP-complete for a dimension greater than 2 (the two-dimensional case is polynomially decidable). Doignon *et al.* give a characterization of this bidimension as the chromatic number of a certain hypergraph, generalizing a similar approach by Trotter (1983) for the usual order dimension. It is this equivalence that we will use for the computation of the bidimension. So first we summarize in the next section the relevant part of the theory developed in Doignon

et al. (1984). Next, in Section 3, we describe some principles enabling us to reduce a hypergraph of the considered type to a subhypergraph having the same chromatic number. This reduction will appear to include as a special case the collapsing of the data matrix into a submatrix, its so-called core, as is obtained by Chubb (1986) in a somewhat more general context (i.e., in principle not restricted to biorder representations). Furthermore, it generalizes the restriction in the case of a partial order to the so-called "non-forced" pairs (Trotter, 1983). In Section 4 this reduction process is illustrated, using an empirical data set. In Section 5 we give a recursive formula for the chromatic number of a hypergraph in terms of the chromatic number of some of its subhypergraphs and in Section 6 we show how the reduction and recursion we developed can be combined to compute the bidimension of another empirical data set. In the last section we discuss the prospects for integrating the findings of the preceding sections in a really explicit and reasonably efficient algorithm. In this context our second problem, that of uniqueness, also comes into view: Is it possible to combine our approach of determining the bidimension with a computation or a characterization in some sense of all possible representations in that dimensionality?

2. BASIC THEORY

First we will fix some general notation and definitions. For two sets A and D , $R \subseteq A \times D$ is called a relation between A and D ; if $R \subseteq A \times A$ it is a relation on A . We will write ad for the ordered pair (a, d) and $ad \in R$ or aRd , equivalently. $\bar{R} = (A \times D) - R$ denotes the complement of R . The cardinality of a set A is denoted $|A|$. Since we aim to applying their theory to real data, we will, in contrast to Doignon *et al.*, throughout assume A and D to be finite. So $A = \{a_1, \dots, a_N\}$ and $D = \{d_1, \dots, d_K\}$ for some natural numbers N and K . As a consequence, a relation R between A and D , as well as its complement \bar{R} , will always be finite.

For interpretative purposes the elements of A can be thought of as persons, those of D as items, and R can be regarded as a dominance relation, $a_i R d_j$, meaning: person a_i solves correctly item d_j . Because of finiteness we can represent a relation R in an $N \times K$ $(0, 1)$ -matrix, called the data matrix, having a 1 in cell (i, j) if $a_i R d_j$ and a 0 if $a_i \bar{R} d_j$. We will denote this matrix as $[R]$; $[R]_{.i}$ is its i th row and is sometimes called the pattern of a_i , $[R]_{.j}$ is its j th column or the pattern of d_j , and $[R]_{ij}$ is the value in cell (i, j) .

There are partial orders on the rows and columns of $[R]$ (corresponding directly to the quasi-orders R_A and R_D of Doignon *et al.*) defined by

$$[R]_{.i} \leq [R]_{.j} \quad \text{iff} \quad (\text{for } k = 1, \dots, K, [R]_{.ik} \leq [R]_{.jk})$$

and

$$[R]_{.i} \leq [R]_{.j} \quad \text{iff} \quad (\text{for } n = 1, \dots, N, [R]_{ni} \leq [R]_{nj}).$$

Now we will summarize the findings of Doignon *et al.* (1984) as far as they are relevant to our more practical goal. Where useful we will specialize their statements

to the finite case and give a translation in terms of the data matrix. In the sequel we may assume that A and D are disjoint; if they are not, we pass to disjoint copies A' and D' of A and D , respectively, and for any relation $R \subseteq A \times D$ we consider its isomorphic image $R' \subseteq A' \times D'$, results there can be directly translated back to the original $R \subseteq A \times D$ (see Doignon *et al*).

The central concept, already defined in Ducamp and Falmagne (1969), is that of a *border*

$R \subseteq A \times D$ is called a *border* between A and D iff for all $a, b \in A, d, e \in D$ we have if $(aRd$ and $bRe)$ then $(aRe$ or $bRd)$

An equivalent, more symmetrical formulation of this condition reads

not $(aRd$ and bRe and $a\bar{R}e$ and $b\bar{R}d)$,

from which it is immediately clear that \bar{R} is a border iff R is one. For the matrix $[R]$ this definition means that R is border iff $[R]$ has no 2×2 submatrix (permutations of rows and columns allowed) of the form

$$\begin{array}{cc} 0 & 1 \\ 1 & 0 \end{array} \quad (21)$$

which implies there are permutations of the rows and the columns of $[R]$ that bring the matrix in triangular form (the matrix is then said to have *triangular structure*).

Of course not every relation is a border, but Doignon *et al* show that any relation R between A and D is, in a trivial way, the intersection of $|\bar{R}|$ borders, which in our case is a finite cardinal. This leads to the concept of *bidimension*, which for the finite case reads

The *bidimension* of a relation $R \subseteq A \times D$, denoted $Bidim R$, is the smallest number q for which there is a collection of q borders $B_i \subseteq A \times D, i = 1, \dots, q$, with $R = \bigcap_{i=1}^q B_i$.

The relevance of the concept of bidimension for a representation of R according to the conjunctive model lies in the following proposition (where \mathbb{R} denotes the set of real numbers)

$Bidim R$ is the smallest number q for which there are two mappings $f = (f_1, \dots, f_q) : A \rightarrow \mathbb{R}^q$ and $g = (g_1, \dots, g_q) : D \rightarrow \mathbb{R}^q$ such that for all $a \in A, d \in D$ aRd iff $f_i(a) \leq g_i(d)$ for $i = 1, \dots, q$.

So we have the practical problem of computing the bidimension of a relation. To that end we will use the equivalence, derived by Doignon *et al*, between the bidimension and the chromatic number of a certain hypergraph (for hypergraphs, see Berge, 1973)

A *hypergraph* $H = \langle V, E \rangle$ is a set V of elements called *vertices* together with a collection E of subsets of V , called *edges*. A subset of V is called

stable in H iff it includes no edge of H and the *chromatic number* of H , denoted $Chrom H$, is the smallest number q for which there is a q -colouring of H , that is, a partition of V into q stable sets, the *colour classes*; in other words, in a colouring no edge is "monochromatic."

We will not consider hypergraphs having singleton edges; in this case the trivial partition into one-element classes is a colouring, so the chromatic number is well-defined and is finite whenever the number of vertices is.

The definition of the hypergraph in question is based on the following generalization of the violating case for a biorder:

An n -alternating cycle of \bar{R} is a sequence $(a_i d_i)_{i=0}^{n-1}$ of elements in \bar{R} such that $a_{i+1} R d_i$ for all i taken modulo n .

This definition is ours; Doignon *et al.* define the corresponding notion as a sequence of elements of A and D , alternately, that "induces" the sequences in \bar{R} and R of the definition above. (The sequence $(a_{i+1} d_i)$ in R is in turn an n -alternating cycle of R , in reverse order.) Since it is the sequence in \bar{R} that is important in the sequel, we think our definition is more direct.

So by definition R is a biorder iff \bar{R} has no 2-alternating cycle, but it follows more generally:

R is a biorder iff \bar{R} has no alternating cycle.

In terms of the matrix $[R]$ we see that the existence of, say, a 4-alternating cycle of \bar{R} implies the existence of a 4×4 submatrix of $[R]$ (permutations of rows and columns allowed) of the form

$$\begin{matrix} 0 & x & x & 1 \\ 1 & 0 & x & x \\ x & 1 & 0 & x \\ x & x & 1 & 0 \end{matrix} \tag{2.2}$$

where the x 's are arbitrary.

Now the *hypergraph associated with R* , which we will, in a slight departure from Doignon *et al.*, denote $H(\bar{R})$, is defined as follows:

$H(\bar{R})$ is the hypergraph whose vertices are the elements of \bar{R} and whose edges are the alternating cycles of \bar{R} , interpreted as sets.

Then we have the promised equivalence

$$\text{Bidim } R = \text{Chrom } H(\bar{R}),$$

on which our method of finding the bidimension of a relation will be based.

3 REDUCTION OF THE HYPERGRAPH

The significance of the equivalence between the bidimension of a relation and the chromatic number of the associated hypergraph for the practical purpose of computing the bidimension depends on the extent to which this approach gives way to more efficient algorithms. By the computational equivalence, established by Cogis (1982), of the bidimension problem for an arbitrary relation and the order dimension problem for a partial order, combined with a result of Yannakakis (1982) for the complexity of the latter problem, we know that for $k > 2$ deciding whether k is an upper bound for the bidimension of a given relation is an NP-complete problem. So we cannot expect to find a theoretically efficient, i.e., polynomial time-bounded algorithm. The efficiency will, in practice, strongly depend on the size of the input. Hence, when we translate the problem in terms of finding the chromatic number of a hypergraph, it will be of importance to take care that this hypergraph be as small as possible. To that end we will search for conditions under which we can reduce the hypergraph associated with a relation to one having fewer vertices and edges but the same, yet unknown, chromatic number.

In order to formulate our principle of reducing the hypergraph, we first need some definitions.

3.1 DEFINITION Let $\langle V, E \rangle$ be a hypergraph. For $V^* \subseteq V$ we define $E(V^*) = \{U \in E \mid U \subseteq V^*\}$. So $\langle V^*, E(V^*) \rangle$ is the subhypergraph obtained from $\langle V, E \rangle$ by restricting the set of vertices to V^* and the set of edges to those included in V^* .

Obviously, if $V_1 \subseteq V_2 \subseteq V$, then $\text{Chrom} \langle V_1, E(V_1) \rangle \leq \text{Chrom} \langle V_2, E(V_2) \rangle$.

3.2 DEFINITION For a hypergraph $\langle V, E \rangle$ and subsets $V_1 \subseteq V_2 \subseteq V$, let φ_1 (φ_2) be a colouring of $\langle V_1, E(V_1) \rangle$ ($\langle V_2, E(V_2) \rangle$). Then φ_2 is said to be an *extension* of φ_1 iff for all $C \in \varphi_1$ there is $C' \in \varphi_2$ with $C \subseteq C'$.

So extensions to V_2 do not break up the classes already formed in V_1 . Note that we may obtain non-trivial extensions while $V_2 = V_1$ by merging different colour classes into one new.

3.3 DEFINITION For a hypergraph $H = \langle V, E \rangle$ and some $V^* \subseteq V$, let φ be a colouring of $\langle V^*, E(V^*) \rangle$. Then $\text{Chrom}_\varphi H$ is defined to be the chromatic number of H under the restriction to colourings that are extensions of φ .

As can be checked easily from Definition 3.2, the extension relation on colourings of subhypergraphs of H is transitive (in fact, it is a partial order): an extension of an extension of some colouring φ is again an extension of φ . In this way we see that for two colourings φ_1 and φ_2 of subhypergraphs of H , φ_2 being an extension of φ_1 implies

$$\text{Chrom} H \leq \text{Chrom}_{\varphi_1} H \leq \text{Chrom}_{\varphi_2} H,$$

simply because, going from left to right in the above inequalities, the minimum is taken over a decreasing collection of colourings of H . Clearly both inequalities may be replaced by equalities iff φ_2 (and thus φ_1) can be extended to a minimal colouring of H .

The reduction of the hypergraph we are going to consider here is based on the following property of a vertex

3.4 DEFINITION Let $H = \langle V, E \rangle$ be a hypergraph and $v, w \in V$. Then v is said to be *dominated in H by w* iff for any edge U of H that contains v the set $(U - \{v\}) \cup \{w\}$ is non-stable in H .

More generally, for $v \in V$ and $T \subseteq V$, stable in H , we will say that v is *dominated in H by T* iff for any edge U of H that contains v the set $(U - \{v\}) \cup T$ is non-stable in H .

This definition of being dominated is equivalent to saying that whenever adding a dominated vertex to a stable set in H would turn it into a non-stable set, adding the dominating vertex (subset) instead would have the same effect.

Now we are ready to formulate, for hypergraphs in general, the following

3.5 REDUCTION PRINCIPLE Let $H = \langle V, E \rangle$ be a hypergraph and let v be an element of V that is dominated in H by another vertex $w \in V$. Then any colouring of $H^* = \langle V - \{v\}, E(V - \{v\}) \rangle$ can be extended to a colouring of H by giving v the colour of w .

In particular, $\text{Chrom } H = \text{Chrom } H^* = \text{Chrom}_{\{v, w\}} H$

The proof of this reduction principle really is immediate from Definition 3.4, but we will deduce it from the following more general version 3.5' that makes some assumption on the colouring of the subhypergraph. The proof of 3.5 will again appear to be an easy consequence of Definition 3.4.

3.5' REDUCTION PRINCIPLE (generalized, conditional version) Let $H = \langle V, E \rangle$ be a hypergraph and let v be an element of V that is dominated in H by some stable $T \subseteq V - \{v\}$. Then any colouring φ^* of $H^* = \langle V - \{v\}, E(V - \{v\}) \rangle$ that is an extension of $\{T\}$, i.e., in which T is monochromatic, can be extended to a colouring φ of H by adding v to the colour class of T .

In particular, $\text{Chrom}_{\varphi^*} H = \text{Chrom}_{\varphi^*} H^* = \text{Chrom}_{\varphi^*} H$

Proof Any monochromatic edge U of H induced by adding v to the colour class of T in φ^* must contain v and so the set $(U - \{v\}) \cup T$, which is a subset of $V - \{v\}$, is monochromatic in φ^* . However, since v is dominated by T , this set is non-stable and this contradicts the assumption that φ^* is a colouring of H^* . ■

Proof of 3.5 Put $T = \{w\}$ and apply 3.5', noting that any colouring of H^* is an extension of $\{\{w\}\}$. ■

Note that the reduction principle (in both versions) is constructive in the sense

that it gives a way to construct a minimal colouring of the greater hypergraph from a minimal colouring of the smaller one

By reduction principle 3.5', a vertex that is dominated by a (disjoint) stable subset of vertices needs no separate consideration in the colouring process of a hypergraph, *provided that a dominating subset is monochromatic*. In that case any colour that is good for the dominating subset will do for the dominated vertex as well. As for the chromatic number of the hypergraph, we can split the collection of colourings into those in which the dominating subset is monochromatic and those in which it is not. In finding the chromatic number over the former subcollection we may discard the dominated vertex from the hypergraph. The version 3.5 is just the special case where there is a singleton dominating subset, being monochromatic in any colouring. In this case the dominated vertex can be discarded unconditionally. In the rest of this paper we will in fact use only this special version, in Section 5 a procedure is sketched that allows this restriction, but we may think of alternatives for that procedure that use the more general, conditional version 3.5'.

All this is very fine, but from a practical point of view two intrinsically related questions naturally arise. Will there, in the type of hypergraph defined in the preceding section, be any dominated vertices and how are we going to find them? In order to establish that vertex v is dominated by vertex u , we have, by Definition 3.4, to enumerate all edges containing v , replace v by u in each such edge, and determine whether the resulting set contains an edge. This seems like a lot of work, even when we use the observation that we need not really check *all* edges, but only those that are minimal (i.e., that do not strictly include another edge), which observation by the way introduces the problem of checking whether an edge is minimal. In our type of hypergraph a vertex will generally be part of many (minimal) edges of varying sizes. Moreover, this kind of hypergraph is only implicitly given in the data matrix: its edges are not directly visible, but must be detected by completing alternating cycles.

To get an idea of the intricacy of the hypergraphs considered here, suppose we have tracked down in our data matrix a minimal 4-edge. It must be contained in a submatrix as in (2.2) in which all x 's are zeros (for otherwise the 4-edge is not minimal, as can be checked easily). This means, however, that this submatrix contains not just one minimal 4-edge, but six of them. In general, a minimal n -edge implies an $n \times n$ submatrix with just one 1-entry in each row and each column. In such a matrix there are in fact $(n-1)!$ minimal n -edges, more generally, it contains, for $k = 2, \dots, n$, $n!/(n-k)!/k$ minimal k -edges. These counts can be justified by noting that each such k -edge is determined by the sequence of k 1-entries, used as "stepping-stones" in the alternating cycle. Clearly there are $n!/(n-k)!$ sequences of k out of n 1-entries. We have to divide by a factor k because a k -edge is invariant under cyclic permutations of the k 1-entries involved. The reader may verify that the matrix of (2.2) (with x 's equal to zero) contains six 4-edges, eight 3-edges, and six 2-edges.

In view of this, prospects for applying the reduction principle 3.5 (or 3.5') seem rather poor. Therefore the following proposition is crucial, which states that for this

special type of hypergraph, dominated vertices can be found by inspecting 2-edges only We first need one more definition

3.6 DEFINITION We call $ad, a^*d^* \in \bar{R}$ enemies in \bar{R} iff $ad^* \in R$ and $a^*d \in R$ The (simple, unoriented) graph of the so defined symmetric enemy relation on \bar{R} is denoted as $G(\bar{R})$

We see that two elements of \bar{R} are enemies iff they constitute a 2-edge of $H(\bar{R})$, which means that they can never be in one colour class of $H(\bar{R})$ In the matrix two zeros are enemies iff they are the zeros of a 2×2 submatrix as in (2.1) The graph $G(\bar{R})$ is the partial (hyper)graph of $H(\bar{R})$, obtained by discarding all edges with more than two elements

For $V \subseteq \bar{R}$ denoting by $H(V)$ the subhypergraph $\langle V, E(V) \rangle$ of $H(\bar{R})$, we can now state

3.7 PROPOSITION Let $R \subseteq A \times D, V \subseteq \bar{R}, ad, a'd' \in V$ Then ad is dominated in $H(V)$ by $a'd'$ if it is dominated in $G(\bar{R})$ by $a'd'$

We will again deduce Proposition 3.7 from a more general version 3.7', which is in terms of being dominated by a subset of vertices

3.7' PROPOSITION Let $R \subseteq A \times D, V \subseteq \bar{R}, ad \in V$, and $T \subseteq V$, stable in $H(V)$ Then ad is dominated in $H(V)$ by T if it is dominated in $G(\bar{R})$ by T

Proof Let U be an edge of $H(V)$ containing $ad, U = \{ad, a_1d_1, \dots, a_nd_n\}$, say, where the ordering of the elements corresponds to the underlying alternating cycle In order to show that $(U - \{ad\}) \cup T$ contains an edge of $H(V)$ if T dominates ad in $G(\bar{R})$, we consider the pair a_1d_1 If it is in R , the subset $\{a_1d_1, \dots, a_nd_n\} = U - \{ad\}$ is an edge of $H(V)$ and we are finished If a_1d_1 is in \bar{R} , then ad and a_1d_1 are enemies in \bar{R} Since T dominates ad in $G(\bar{R})$, $T \cup \{a_1d_1\}$ contains an edge of $G(\bar{R})$, which means that a_1d_1 has some enemy $a^*d^* \in T$ This, however, implies that $\{a^*d^*, a_1d_1, \dots, a_nd_n\} \subseteq (U - \{ad\}) \cup T$ is an edge of $H(V)$ ■

Proof of 3.7 Put $T = \{a'd'\}$ and apply 3.7' ■

By Propositions 3.7 and 3.7' we may, when applying 3.5, resp 3.5', to (sub)hypergraphs of our special type, replace the phrase "dominated in H " by "dominated in $G(\bar{R})$ " The fact that a vertex ad is dominated in $G(\bar{R})$ by $a'd'$ simply means that any enemy of ad in \bar{R} is an enemy of $a'd'$ as well, likewise, ad being dominated in $G(\bar{R})$ by a subset T means that any enemy of ad in \bar{R} has an enemy in T So questions concerning being dominated in $G(\bar{R})$ can easily be settled from inspection of the tableau of the binary enemy relation on \bar{R} In this way, finding dominated vertices is in some sense reduced to inspecting 2-edges only Notice, however, the appearance of $G(\bar{R})$ in 3.7 and 3.7' instead of the perhaps more expected $G(V)$ (with the obvious meaning for this notation) While it is clear that being dominated in $G(V)$ is necessary for being dominated in $H(V)$ ($G(V)$ being a partial hypergraph of $H(V)$), it is not sufficient The reason for this can be seen in the

proof of 3.7': there we call on a pair $a_i d_n$ that, in the non-trivial case, is in \bar{R} , but in no way needs to be in V . On the other hand, being dominated in $G(\bar{R})$, while sufficient, is not necessary for being dominated in $H(V)$: there is no "only if" in 3.7 or 3.7'. This means that by applying 3.7 or 3.7' we may not detect all dominated vertices in $H(V)$. The danger of missing some dominated vertices certainly depends on the discrepancy between the sufficient $G(\bar{R})$ and the necessary $G(V)$, that is, the discrepancy between V and \bar{R} . (In case $V = \bar{R}$ we trivially do have "only if's" in 3.7 and 3.7'.) In this respect it is important to realize that whenever $V \subseteq A' \times D'$ for some $A' \subseteq A$, $D' \subseteq D$, then $V \subseteq \bar{R}'$, where R' is the restriction of R to $A' \times D'$, and we may use \bar{R}' instead of \bar{R} in 3.7 and 3.7'.

There is an important special application of Proposition 3.7 in which dominated vertices can be detected from inspection of the data matrix itself, that is, without even explicitly considering the enemy relation on the vertices (the zeros in the matrix). For this special kind of dominated vertices there is, moreover, some intuitive justification that indeed they have no bearing on the dimensionality of a representation of a relation R according to the conjunctive model. The subset of vertices in question can best be characterized in terms of the data matrix $[R]$, in which we need to find a colouring for all zeros.

Consider two ordered patterns, of persons a_i and a_j say; assume $[R]_i \leq [R]_j$. Because of this ordering, in any column in which $[R]_i$ has a zero, $[R]_j$ has one, too. Let k be such a column and let $a_i d_k$ and $a^* d^*$ be enemies. Then we know that $a^* d_k \in R$ and $a_i d^* \in R$; by the ordering the latter implies $a_i d^* \in R$ and together with the assumption $a_i d_k \in \bar{R}$ it follows that $a_i d_k$ and $a^* d^*$ are enemies. So, by Proposition 3.7, the vertex $a_i d_k$ is dominated in $H(\bar{R})$ by $a_i d_k$.

There is a completely analogous version for the case of two ordered column patterns. From this we see that in a pattern we have left to be coloured only those zeros for which there is no higher-ordered pattern with a zero in the same position. A direct consequence is that in determining the bidimension, a pattern in which all occurring zeros can be "explained" in this way, that is, a pattern that is the conjunction of a number of higher-ordered patterns, can be discarded altogether. Intuitively one might feel that such a pattern does not offer any new information regarding the multidimensional representation. Suppose, for instance, we have a representation for two persons, one of which failed on item 1 only, the other on item 2 only. Then, precisely because we work in the conjunctive model, there must already be in that representation a place for a person failing on just the items 1 and 2. According to the conjunctive model the (possible) occurrence of the last pattern is, actually, implied by the occurrence of the first two. Again there is an analogous intuitive argument with the roles of persons and items reversed. Patterns, rows or columns, that for this reason can be removed from the data matrix will be called *implied patterns*. Separate zeros that, in the sense defined above, can be explained by a zero a in a higher-ordered pattern will be called *implied zeros*. Note that repetitions in the data matrix of one and the same pattern constitute a special, trivial case of implied patterns. So, such repetitions can be discarded without changing the bidimension, a fact that is intuitively self-evident.

Although it is quite easy to present data matrices in which reduction of the hypergraph according to 3.5 and 3.7 has no effect at all, the special case of implied zeros shows that under rather mild conditions (occurrence of ordered rows or columns in the data matrix) the associated hypergraph is bound to contain dominated vertices. On the other hand, because the conjunctive model predicts in some sense, their occurrence, as we have seen above, the amount of implied zeros in a data matrix may be considered as some sort of measure of confirming evidence for the conjunctive model as the operative one in producing the data.

There exist, as was pointed out by Doignon (personal communication), close connections between the notions of implied patterns and implied zeros and the work of Chubb (1986) and Trotter (1983), respectively.

In fact, the special application of Proposition 3.7 consisting of removing implied rows and columns gives exactly the restriction of R as described by Chubb. Our non-implied rows (columns) constitute his minimal row (column) \cap -generating set and the corresponding restriction of R is its so-called \cap -core. We see that by the present reduction we can, in addition, "remove" separate zeros and by fully using Proposition 3.7 this may well be more zeros than just the implied ones. In particular this means that we may, possibly, remove more rows or columns, thereby obtaining a still smaller "core" of R .

In the special case where $D = A$ and the relation R is a (reflexive) partial order on A , its bidimension equals its order dimension (Doignon *et al.*, 1984) and then the non-implied zeros correspond exactly to the "non-forced pairs," the subset of incomparable pairs of R on which Trotter (1983) defines his hypergraph. It can be shown that in this case all dominated zeros are implied zeros and in this way Proposition 3.7 (together with 3.5) is equivalent to a result of Maurer, Rabinovitch, and Trotter (1980), implying that for the dimension of a partial order consideration can be confined to the set of non-forced pairs. In the description of the implied zeros above we have seen a way to find this set.

4 AN ILLUSTRATION OF THE REDUCTION PROCESS

Here we will demonstrate the potential significance of the reduction of a hypergraph according to Proposition 3.7. We will use a data set from Chubb (1986) that originates, in turn, from data of Stouffer *et al.* (1950) on six polytome items administered to a set of American World War II GI's and intended to assess their "readiness to enter into battle." By appropriately dichotomizing the responses Chubb obtained the data matrix of Table 4.1(a), where capitals A to F denote the six items and the rows represent the observed response patterns. Of course, the response patterns occurred with varying frequencies, but for a deterministic analysis of the data these frequencies are irrelevant and it suffices to consider the reduced data matrix without duplicate rows or columns. (When searching for approximate representations of the data in a lower dimensionality it is natural to take the frequencies into account.)

TABLE 4 1

(a) Data Matrix from Chubb (* Indicates a Row-Implied Zero Non-Implied Rows Are Numbered),
 (b) Submatrix of Non-Implied Patterns (* for an Implied Zero)

(a)	A	B	C	D	E	F	(b)	A	B	C	D	E	F
	1	1	1	1	1	1	1	1	1	0	1	1	1
(1)	1	1	0	1	1	1	2	1	0	1	1	1	1
(2)	1	0	1	1	1	1	3	0	1	1	1	1	1
(3)	0	1	1	1	1	1	4	*	1	1	1	1	0
(4)	*	1	1	1	1	0	5	*	1	1	0	1	1
(5)	*	1	1	0	1	1	6	1	1	*	0	0	1
	*	*	1	1	1	1	7	*	*	*	1	0	*
	1	*	*	1	1	1							
	*	*	1	*	1	1							
	*	*	*	1	1	1							
	*	1	*	*	1	1							
(6)	1	1	*	0	0	1							
	*	*	1	*	1	*							
	*	*	*	*	1	1							
	*	1	*	*	*	1							
	*	*	*	*	*	1							
(7)	*	*	*	1	0	*							
	*	*	*	*	*	*							

In Table 4 1(a) we have already indicated which zeros are row-implied a "*" denotes a zero for which there is a non-implied zero, denoted as "0," in the same column in some higher-ordered pattern (i.e., higher in the partial order on the rows of the matrix) By the special application of Proposition 3 7 discussed in the preceding section we may confine our attention to the subset of non-implied rows These are numbered in Table 4 1(a) and the corresponding submatrix is given separately in Table 4 1(b) There are no extra column-implied zeros only one row (6) has more than one "living" zero, but the corresponding columns (*D* and *E*) are unordered So here ends the special application of Proposition 3 7 consisting of removing implied zeros from the hypergraph, and the submatrix of Table 4 1(b) (without the distinction between implied and non-implied zeros) is the "core," the restriction that is obtained by Chubb

In order to fully use Proposition 3 7 we now consider the enemy relation between the vertices of the reduced hypergraph (the non-implied zeros) and those of the full hypergraph (all zeros in Table 4 1(b)) This enemy relation is represented in the 8×15 matrix of Table 4 2(a), where 1-entries indicate that the vertices in the corresponding row and column constitute a 2-edge We see that the vertex $3A$ is dominated by $4F$, $6E$ is dominated by $1C$, and $7E$ by $4F$ (for instance) In Table 4 2(a) the rows of dominated vertices are starred and by Proposition 3 7 the problem reduces to the situation in Table 4 2(b) There we see that this reduction has turned the vertex $6D$ into a dominated one and its removal leads us to the

TABLE 4.2

(a) Enemy Relation for the Hypergraph of Table 4.1(b) (Starred Rows Indicate Dominated Vertices);
 (b) Resulting Submatrix and Enemy Relation after Removing Dominated Vertices,
 (c) Final Reduction

(a)	1C	2B	3A	4F	5D	6D	6E	7E	4A	5A	6C	7A	7B	7C	7F
1C	.	1	1	1	1	.	.	.	1	1
2B	1	.	1	1	1	1	1	.	1	1	1
*3A	1	1	.	.	.	1	1	.	.	.	1
4F	1	1	.	.	1	1	1	.	.	.	1
5D	1	1	.	1	.	.	.	1	1	1	1
6D	.	1	1	1	1	.	.	1	1	.	1
*6E	.	1	1	1	1	1
*7E	1

(b)	B	C	D	F	1C	2B	4F	5D	6D	6C
1	1	0	1	1	1C	1	1	1	.	.
2	0	1	1	1	2B	1	.	1	1	1
4	1	1	1	0	4F	1	1	.	1	1
5	1	1	0	1	5D	1	1	1	.	.
6	1	*	0	1	*6D	.	1	1	.	.

(c)	B	C	D	F	1C	2B	4F	5D
1	1	0	1	1	1C	1	1	1
2	0	1	1	1	2B	1	.	1
4	1	1	1	0	4F	1	1	1
5	1	1	0	1	5D	1	1	1

hypergraph consisting of four vertices and the 4×4 submatrix of Table 4.2(c). In fact, here reduced and full hypergraphs are the same and after inspecting the enemy relation on its vertices we see that now Proposition 3.7 has spent itself: there are no more dominated vertices. But in this case, there is no more problem, either! For we have shown that the chromatic number of the full hypergraph corresponding to Table 4.1(a) is equal to the chromatic number of the subhypergraph consisting of the four vertices that are present in Table 4.2(c). But from Table 4.2(c) it is clear that no two of these vertices can have the same colour: they are all mutual enemies. So we will need exactly four colours and we may conclude that the bidimension of the original data matrix equals 4.

This means that the relation underlying the data matrix can be represented as the intersection of four biorders. But these biorders are by no means unique: many collections of four biorders have the given relation as their intersection. In other words, there are many different collections of four dimensions (Guttman scales) that are, under the conjunctive composition rule, compatible with the observed data matrix. By inspecting the content of the four items present in Table 4.2(c),

Chubb manages to attach a verbal label to each dimension on which one of these items scores highest and to construct four dimensions that reflect reasonably well the interpretation suggested by these labels. There is, however, definitely some arbitrariness involved here, which seems inevitable when embedding six items in 4-dimensional space on the basis of binary data. Generally, representations in the minimum dimensionality will be far from unique and we clearly see the need for "best approximate" representations in some lower dimensionality with a higher degree of uniqueness. (Of course, of greatest help would be a strong psychological theory about the data that allows focusing on specific aspects in the collection of representations.)

Returning to our more basic problem of finding the dimensionality of the observed data, we see that by applying Proposition 3.7 we were able to reduce from a hypergraph with 51 vertices based on a 18×6 matrix to one with four vertices based on a 4×4 matrix. This shows the potential power of the reduction principle described in the preceding section. Here we were very lucky indeed, in general, however, there will be a non-trivial problem left when no more reduction is possible. In the next section we discuss a possible approach in that case. This approach will then be illustrated in Section 6, using another data set. There, again, we will make use of Proposition 3.7 whenever we can. Its effects will not be as dramatic as was the case here, but it still will turn out to be very useful.

5. A RECURSION FORMULA FOR THE CHROMATIC NUMBER

With all possible reductions carried out, there will come a moment when the real work has got to start. The problem's being NP-complete suggests that at such a point some sort of exhaustive search will be inevitable. We will describe here some such search in which we use the notion of a *maximal stable set* in a hypergraph, that is, a stable set that is not contained in any other stable set.

For maximal stable sets in an arbitrary hypergraph we can establish the following properties:

5.1 PROPOSITION (i) *For any vertex v of a hypergraph H there is a minimal colouring of H in which the colour class of v is maximal.*

(ii) *For any maximal stable set M in a hypergraph $\langle V, E \rangle$ we have $\text{Chrom}_{\{M\}} \langle V, E \rangle = 1 + \text{Chrom} \langle V - M, E(V - M) \rangle$.*

Proof. (i) Let C be the colour class of v in a minimal colouring of H . If there is a vertex w , not in C , such that $C \cup \{w\}$ still is stable in H , then the transfer of w to C clearly gives a colouring of H without introducing new colours. This process can be repeated until there are no more candidates, that is, until the class of v is maximal.

(ii) Because M is stable, any $(q-1)$ -colouring of $\langle V - M, E(V - M) \rangle$ can be extended to a q -colouring of $\langle V, E \rangle$ by adding M as a colour class. This proves $\text{Chrom}_{\{M\}} \langle V, E \rangle \leq 1 + \text{Chrom} \langle V - M, E(V - M) \rangle$. For the reverse inequality,

consider a q -colouring of $\langle V, E \rangle$ in which M is monochromatic Then M , being maximal, must itself be a colour class and we are in fact given a $(q - 1)$ -colouring of $\langle V - M, E(V - M) \rangle$ ■

From these properties we easily obtain

5.2 COROLLARY For any vertex v of a hypergraph $\langle V, E \rangle$ we have

$$\text{Chrom} \langle V, E \rangle = 1 + \min_{M \in \text{MAX}(v)} \{ \text{Chrom} \langle V - M, E(V - M) \rangle \},$$

where $\text{MAX}(v)$ denotes the collection of maximal stable sets in $\langle V, E \rangle$ containing the vertex v

Proof By Proposition 5.1(i), $\text{Chrom} \langle V, E \rangle = \min_{M \in \text{MAX}(v)} \{ \text{Chrom}_{\{M\}} \langle V, E \rangle \}$ and applying Proposition 5.1(ii) completes the proof ■

We see that Corollary 5.2 gives a recursive formula for the chromatic number of a hypergraph in terms of the chromatic numbers of a collection of strictly smaller hypergraphs So, when applied to the finite, special type of hypergraph that turns up in the bidimension problem, this recursion will give the chromatic number and thereby the bidimension in a finite number of steps The big question, of course, is whether computation of the bidimension according to this recursion will be feasible in practice In this context it is worth noting some additional properties of maximal stable sets in a hypergraph associated with a relation

5.3 PROPOSITION For $R \subseteq A \times D$ the following two properties hold

- (i) Any maximal stable set in $H(\bar{R})$ is a biorder between A and D
- (ii) If M is a maximal stable set in $H(\bar{R})$, then for some $a_0 \in A$, $\bar{R} \cap (\{a_0\} \times D) \subseteq M$ and for some $d_0 \in D$, $\bar{R} \cap (A \times \{d_0\}) \subseteq M$

Proof (i) (Also in Doignon *et al.*, 1984, p 95) Let M be a maximal stable set in $H(\bar{R})$ and suppose M has a violation of the biorder property $ad, be \in M$ and $ae, bd \in \bar{M}$ We are going to derive a contradiction by showing that M must contain an edge of $H(\bar{R})$ If both ae and bd are in R , then M contains the 2-edge $\{ad, be\}$ If one, say ae , is in R and the other, bd , is in \bar{R} , then by maximality of M there is an edge included in $M \cup \{bd\}$, which of course contains the vertex bd Let $bd, a_1d_1, \dots, a_nd_n$ be the corresponding alternating cycle, then the sequence $be, ad, a_1d_1, \dots, a_nd_n$ is an alternating cycle in M If, eventually, both ae and bd are in \bar{R} , then, by the same argument, there are alternating cycles $bd, a_1d_1, \dots, a_nd_n$ and $ae, a'_1d'_1, \dots, a'_md'_m$ in $M \cup \{bd\}$ and $M \cup \{ae\}$, respectively But then the sequence $ad, a_1d_1, \dots, a_nd_n, be, a'_1d'_1, \dots, a'_md'_m$ is an alternating cycle in M

(ii) Consider two elements a_i, a_j in A If there are $d_1, d_2 \in D$ such that $a_id_1 \in M$ and $a_jd_2 \in \bar{M}$, while $a_jd_1 \in \bar{M}$ and $a_id_2 \in M$, then M certainly is no biorder So if M is a maximal stable set in $H(\bar{R})$ it is, by part (i), a biorder between A and D , consequently for any pair a_i, a_j we have either $a_id \in M$ implies $a_jd \in M$ for all

$d \in D$, or $a, d \in M$ implies $a, d \in M$ for all $d \in D$. In this way M induces a weak ordering on A (transitivity is easily checked) and since A is finite we can find an element a_0 in A that is maximal in this ordering, i.e., for which for any $a' \in A$ and any $d \in D$, $a'd \in M$ implies $a_0d \in M$. We will show $\bar{R} \cap (\{a_0\} \times D) \subseteq M$ for this a_0 by showing that the set $M \cup (\bar{R} \cap (\{a_0\} \times D))$ is a biorder, which is an equivalent assertion since M is maximal. If $M \cup (\bar{R} \cap (\{a_0\} \times D))$ is not a biorder, then, since M itself is, any violation of the biorder property must involve the element a_0 . In particular there must then be $a' \in A$ and $d' \in D$ such that $a'd'$ is in $M \cup (\bar{R} \cap (\{a_0\} \times D))$ and a_0d' is not. The former, however, implies $a'd' \in M$ (a_0 and a' clearly being distinct) and thus, by our choice of a_0 , $a_0d' \in M \subseteq M \cup (\bar{R} \cap (\{a_0\} \times D))$, a contradiction. The proof of the other half of (ii), with the roles of A and D reversed, is completely analogous. ■

In terms of the data matrix, Proposition 5.3(i) states that a maximal stable set of zeros has triangular structure in the matrix. In this perspective 5.3(ii) really is obvious. It asserts that in a row (column) that contains, compared to other rows (columns), a maximal number of elements of the set in question, *all* zeros belong to the set. If not, they could be added without disturbing the triangular structure (i.e., biorderhood) of the set, which consequently would not have been maximal.

Corollary 5.2 poses the problem of finding the relevant collection of maximal stable sets. We want to use the corollary for subhypergraphs $H(V)$, obtained after maximally reducing a hypergraph $H(\bar{R})$. If \bar{R}' is the smallest restriction of \bar{R} that contains V , we may apply Corollary 5.2 to the hypergraph $H(\bar{R}')$, in which case Proposition 5.3(i) gives a translation in terms of maximal biorders contained in \bar{R}' and containing some specified vertex. When applying the corollary to the reduced hypergraph $H(V)$ itself, this equivalence may still be useful if we notice that any maximal stable set in $H(V)$ is the restriction to V of some maximal stable set in $H(\bar{R}')$.

Obviously, the amount of work implied by the recursion formula in 5.2 will depend on the reduction obtained before invoking it, but we must realize that at that moment the reduction process of Section 3 is not simply set aside. Rather, at each step of the recursion, by deleting a maximal stable set from the hypergraph under consideration, we get not only a strictly smaller hypergraph, but also one that is again susceptible of reduction. The removal of a maximal stable set of $H(V)$ leads to the removal of at least one row and one column from the underlying matrix. Any maximal stable set of $H(V)$ is the restriction to V of some maximal stable set of $H(\bar{R}')$, and by Proposition 5.3(ii) the latter "includes" at least one column and one row of $[\bar{R}']$. So, in general, before applying the next step in the recursion, we can further shrink the hypergraph and possibly also the submatrix that it is based on.

Further remarks on the question of turning the results of this section and Section 3 into a reasonably practical algorithm will be made in the last section.

6. EXAMPLE OF FINDING THE BIDIMENSION ON EMPIRICAL DATA

Here we will show how the findings of Sections 3 and 5 can be combined to compute the bidimension of an empirical data matrix. The data set we use is borrowed from Marcovici (1981, p. 122). The context is a signal-detection experiment with several conditions of induced colour-blindness. The obtained data matrix is reproduced in Table 6.1(a). Columns refer to figures to be detected and rows of the matrix correspond to subject \times condition combinations (there were 3 subjects and 7 conditions). Since our only objective here is to compute the dimensionality of this data matrix, we will in the sequel pay no attention to its construction and interpretation; it will simply be considered a given 21×10 (0, 1)-matrix.

We will compute the bidimension of the data matrix as the chromatic number of its associated hypergraph and we will start therefore by reducing this hypergraph as much as possible according to Proposition 3.7. In Table 6.1(a) we have already

TABLE 6.1

(a) Data Matrix from Marcovici (1981) (* Indicates a Row-Implied Zero, Non-Implied Row Patterns are Numbered), (b) Matrix of Non-Implied Row Patterns (Extra Column-Implied Zeros Found, Here Marked by an Underscore), (c) Further Reduction by Removal of Implied Column Patterns (All Implied Zeros Are Detected and Displayed as *)

(a)	A	B	C	D	E	F	G	H	I	J	(b)	A	B	C	D	E	F	G	H	I	J
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	2	1	1	1	0	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	3	0	1	1	*	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	4	1	1	1	1	<u>*</u>	*	1	1	1	0
(1)	1	1	1	1	1	0	1	1	1	1	5	1	1	0	*	<u>*</u>	*	1	1	1	1
(2)	1	1	1	0	1	1	1	1	1	1	6	1	1	1	*	*	*	1	1	0	1
(3)	0	1	1	*	1	1	1	1	1	1	7	0	1	0	1	*	*	0	1	1	*
	*	1	1	*	1	1	1	1	1	1	8	1	1	0	1	*	*	<u>*</u>	0	0	*
	*	1	1	*	1	1	1	1	1	1	9	1	1	1	*	*	<u>*</u>	0	*	*	
(4)	1	1	1	1	0	*	1	1	1	0	10	*	0	*	*	*	*	1	1	*	
(5)	1	1	0	*	0	*	1	1	1	1	11	*	*	*	*	*	*	0	1	*	
(6)	1	1	1	*	0	*	1	1	0	1											
(7)	0	1	0	1	*	*	0	1	1	*											
(8)	1	1	0	1	*	*	0	0	0	*	(c)	A	B	C	D	F	G	H	I	J	
(9)	1	1	1	*	*	*	0	0	*	*	1	1	1	1	0	1	1	1	1	1	
	1	1	*	*	*	*	*	*	*	*	2	1	1	1	0	1	1	1	1	1	
(10)	* 0	*	*	*	*	*	*	1	1	*	3	0	1	1	*	1	1	1	1	1	
	*	*	*	*	*	*	*	1	*	*	4	1	1	1	1	*	1	1	1	0	
(11)	* *	*	*	*	*	*	*	0	1	*	5	1	1	0	*	*	1	1	1	1	
	* 1	*	*	*	*	*	*	*	*	*	6	1	1	1	*	*	1	1	0	1	
	* *	*	*	*	*	*	*	*	*	*	7	0	1	0	1	*	0	1	1	*	
											8	1	1	0	1	*	*	0	0	*	
											9	1	1	1	*	*	*	0	*	*	
											10	*	0	*	*	*	*	1	1	*	
											11	*	*	*	*	*	*	0	1	*	

indicated which zeros in the matrix are row-implied and we see that we can restrict our attention to the 11×10 submatrix in Table 6.1(b). Finding here some new column-implied zeros we obtain the 11×9 matrix of Table 6.1(c).

Next, to further exploit Proposition 3.7, we consider the enemy relation in this submatrix for the 15 vertices of the reduced hypergraph. This relation is represented in Table 6.2(a), where, in order to save space, enemies that are no longer contained in the hypergraph are given in listed form. Inspecting the partial order on the rows we find 5 dominated vertices and the data matrix reduces to that of Table 6.2(b). The tableau of the corresponding enemy relation reveals another 3 dominated vertices and we can further reduce to the situation of Table 6.2(c). Inspection shows that now there are no more dominated vertices and so the reduction process stops here.

Proposition 3.7 has enabled us to reduce the original hypergraph containing 82 vertices and based on a 16×10 matrix (not counting perfect one- or zero-patterns) to one having 7 vertices and based on a 7×6 matrix. But in this case we are still left with a problem. From inspection of Table 6.2(c) it is not at all clear which value the chromatic number of the resulting hypergraph has. So we have to invoke Corollary 5.2, that is, we have to choose a vertex of the reduced hypergraph, generate all maximal stable sets in the reduced hypergraph that contain the chosen vertex, and compute successively the chromatic number of the hypergraphs obtained by removing each such maximal stable set from the reduced hypergraph. The chromatic number we are searching for will be the minimum over this set of numbers, increased by one.

We choose, for instance, the vertex $3A$. From Table 6.2(c) we see that there are only two vertices, $2D$ and $7A$, which are not enemies of $3A$. Since $2D$ and $7A$, on their turn, are mutual enemies, it is clear that there are just two maximal stable sets in the reduced hypergraph that contain the vertex $3A$, $\{3A, 2D\}$ and $\{3A, 7A\}$. So, first we will remove the vertices $3A$ and $2D$ from the reduced hypergraph. The resultant submatrix and its enemy relation are displayed in Table 6.3(a).

The important point to be noticed is that, by removing a maximal stable set from the reduced hypergraph, at least one row and one column can be removed from the data matrix and we get a hypergraph in which, again, some vertices may be dominated ones. So we see in Table 6.3(a) that we may remove vertex $8H$, next in Table 6.3(b) $7A$ may be discarded, and as a result we are left with the three vertices in the 3×3 matrix of Table 6.3(c). These vertices, all being mutual enemies, clearly need three different colours. An alternative way of putting this is that any vertex in Table 6.3(c) constitutes a maximal stable set by itself and so, by repeatedly applying Corollary 5.2, we can remove the vertices one after the other, each time increasing the chromatic number by one, until there are no more vertices left. Anyway, the conclusion is the same: if the vertices $3A$ and $2D$ are to have the same colour, we need four colours for the hypergraph represented in Table 6.2(c).

Now we can do the same, taking $\{3A, 7A\}$ as the maximal stable set to start with. The results of this choice are given in Table 6.4.

Apparently we must conclude: if the vertices $3A$ and $7A$ are to have the same

TABLE 62

(a) Enemy Relation for the Hypergraph of Table 61(c) (Removed Enemies Are Listed, Dominated Rows Are Starred), (b), (c) Further Reductions until No More Dominated Vertices Appear

(a)	1F	2D	3A	4J	5C	6I	7A	7C	7G	8C	8H	8I	9H	10B	11H
* 1F		1	1												3D
2D	1			1		1	1	1	1	1	1	1			4F, 7F, 7J, 8F, 8G, 8J
3A	1			1	1	1				1	1	1	1		4F, 5F, 6F, 8F, 8G, 8J, 9F, 9G, 9I, 9J
4J		1	1		1	1									3D, 5D, 6D
5C			1	1		1							1		9G, 9I, 9J
6I			1	1	1		1	1	1					1	7J, 10A, 10C, 10G, 10J, 11A, 11B, 11C, 11G, 11J
7A		1				1					1	1	1		5D, 6D, 9D, 9I
7C		1				1							1		3D, 6D, 9D, 9I
* 7G		1				1									3D, 5D, 6D
8C		1	1												3D, 6D, 9D
8H		1	1				1							1	3D, 5D, 6D, 10A, 10D
8I		1	1				1							1	3D, 5D, 6D, 10A, 10D, 11A, 11B, 11D
* 9H			1		1		1	1						1	10A, 10C
*10B						1					1	1	1		9I
*11H						1									—

(b)	A	C	D	H	I	J	2D	3A	4J	5C	6I	7A	7C	8C	8H	8I	
2	1	1	0	1	1	1	2D		1			1	1	1	1	1	7J, 8J
3	0	1	*	1	1	1	3A		1	1	1			1	1	1	8J
4	1	1	1	1	1	0	4J	1	1		1	1					3D, 5D, 6D
5	1	0	*	1	1	1	5C		1	1		1					—
6	1	1	*	1	0	1	6I		1	1	1		1	1			7J
7	0	0	1	1	1	*	7A	1							1	1	5D, 6D
8	1	0	1	0	0	*	*7C	1									3D, 6D
							*8C	1	1								3D, 6D
							8H	1	1								3D, 5D, 6D
							*8I	1	1								3D, 5D

(c)	A	C	D	H	I	J	2D	3A	4J	5C	6I	7A	8H	
2	1	1	0	1	1	1	2D			1		1	1	7C, 7J, 8C, 8I, 8J
3	0	1	*	1	1	1	3A			1	1	1	1	8C, 8I, 8J
4	1	1	1	1	1	0	4J	1	1		1	1		3D, 5D, 6D
5	1	0	*	1	1	1	5C		1	1		1		—
6	1	1	*	1	0	1	6I		1	1	1		1	7C, 7J
7	0	*	1	1	1	*	7A	1				1	1	5D, 6D, 8I
8	1	*	1	0	*	*	8H	1	1				1	3D, 5D, 6D

TABLE 63

(a) Data Matrix and Corresponding Enemy Relation after Removing Vertices 3A and 2D from Hypergraph of Table 62(c), (b), (c) Reductions by Removing Dominated Vertices

(a)	A	C	H	I	J	4J	5C	6I	7A	8H
4	1	1	1	1	0	4J	1	1		—
5	1	0	1	1	1	5C	1	1		
6	1	1	1	0	1	6I	1	1	1	7C, 7J
7	0	*	1	1	*	7A		1		1 8I
8	1	*	0	*	*	*8H			1	—

(b)	A	C	I	J	4J	5C	6I	7A
4	1	1	1	0	4J	1	1	—
5	1	0	1	1	5C	1	1	—
6	1	1	0	1	6I	1	1	1 7C, 7J
7	0	*	1	*	*7A		1	--

(c)	C	I	J	4J	5C	6I
4	1	1	0	4J	1	1 -
5	0	1	1	5C	1	1
6	1	0	1	6I	1	1 -

TABLE 64

(a) Data Matrix and Corresponding Enemy Relation after Removing Vertices 3A and 7A from Hypergraph of Table 62(c), (b), (c) Reductions by Removing Dominated Vertices

(a)	C	D	H	I	J	2D	4J	5C	6I	8H
2	1	0	1	1	1	2D	1			1 8C, 8I, 8J
4	1	1	1	1	0	4J	1	1	1	5D, 6D
5	0	*	1	1	1	5C	1	1		—
6	1	*	1	0	1	6I	1	1		--
8	*	1	0	*	*	*8H	1			5D, 6D

(b)	C	D	I	J	2D	4J	5C	6I
2	1	0	1	1	*2D	1		
4	1	1	1	0	4J	1	1	1 5D, 6D
5	0	*	1	1	5C	1	1	—
6	1	*	0	1	6I	1	1	

(c)	C	I	J	4J	5C	6I
4	1	1	0	4J	1	1 —
5	0	1	1	5C	1	1 —
6	1	0	1	6I	1	1 —

colour, we need four colours for the hypergraph of Table 6.2(c). By Corollary 5.2, the chromatic number of this hypergraph is the minimum of the values obtained from these two computations and, by Proposition 3.7 and the reduction principle 3.5, this number equals the chromatic number of the hypergraph in Table 6.1(a) we started with. By the fundamental theorem of Doignon *et al.*, then, we may conclude that the bidimension of the original data matrix equals 4.

7 DISCUSSION

We have seen that determining the dimensionality needed for the representation of a data matrix according to the conjunctive model is, in general, a hard, indeed an NP-hard problem, whereas from a casual inspection of Coombs' (1964) example it first appeared to be an automatically obtained by-product of a constructive scaling procedure. On the basis of the equivalence, given by Doignon *et al.*, between this bidimension and the chromatic number of a certain hypergraph, we have gathered in Sections 3 and 5 some results which seem to make computation of the bidimension feasible in practice, at least for data sets of moderate size (As we have seen, the computations of Sections 4 and 6 were rather easily carried out with paper and pencil.) So the first thing to do is to combine these findings in an explicit, reasonably efficient algorithm. As for the reduction process, this does not seem to offer many problems: it is already rather explicitly described and illustrated in Section 3 and in the examples of Sections 4 and 6.

The recursion formula of Corollary 5.2 will need more consideration. It poses the problem of computing the collection of all maximal stable sets containing a certain vertex. In the example of the preceding section this problem was easily solved "by inspection," but the general case will be NP-hard (Without appealing to any reduction mechanism, the recursion still solves the NP-hard problem of finding the bidimension.) So it will be important to devise a practical algorithm for this sub-problem and Proposition 5.3(1) may turn out to be useful in this context. At least as important, however, will be trying to avoid needless exhaustive execution of this algorithm. Returning, for instance, to our example in the preceding section, we can see that consideration of the second computation, based on the maximal stable set $\{3A, 7A\}$, was in fact pointless. For, as a result of the first computation, we had established 4 as an upper bound for the chromatic number. In this computation, however, the vertices 3A, 4J, 5C, and 6I were put in different maximal stable sets and in Table 6.2(c) we can easily check that these four vertices are all mutual enemies (they are said to form a 4-clique in $G(\bar{R})$). So in any colouring they must have different colours and this establishes 4 as a lower bound and thereby solves the problem.

In general, then, we may use the sequence of maximal stable sets formed in a computation for detecting maximal such cliques from the last set backwards. In this way the search tree induced by the recursion formula of Corollary 5.2 may be pruned considerably. Suppose, for instance, we have a "solution" in five maximal

stable sets and suppose we have traced a 3-clique of enemies in the sets at the last three levels. Then we know that, given our present choice of the maximal stable set at level 2, we can, at the next levels, do no better than we did. Hence we need not consider any alternative choices at the levels 5, 4, or 3, instead we may backtrack to the next alternative at the second level. We can, moreover, try to exploit the freedom we have in choosing a vertex in Corollary 5.2 in order to keep the branching of this tree to a minimum in the first place. Heuristically it seems reasonable to choose a vertex with a maximum number of enemies in the subhypergraph in question, expecting such a vertex to have a minimum number of "surrounding" maximal stable sets. We must remark, after all this, that application of Corollary 5.2 is just one possible way of tackling the problem. Alternative and possibly more efficient procedures may exist, waiting to be developed.

Just computing the dimensionality of a representation will not be very useful, however. What we really want are the very representations. This presents the second problem alluded to in the introduction, that of uniqueness of solutions.

It may be noticed that the manner of determining the bidimension of a relation R as sketched in the previous sections is constructive in the sense that it is easy to derive from the computations at least one possible representation of R in the minimum dimensionality, i.e., we can construct at least one minimal collection of biorders having R as their intersection. For the bidimension is found by completing sequences of alternately reducing a hypergraph to a subhypergraph of non-dominated vertices and removing from this subhypergraph a maximal stable set, starting with the original hypergraph $H(R)$ and ending when there are no more vertices left. At the moment the bidimension is known, we have executed at least one such computation that annihilates $H(\bar{R})$ in a minimal number of steps. In such a computation any removed dominated vertex has left a dominating vertex in the subhypergraph as a representative, so ultimately any vertex has a representative in one of the maximal stable sets of this sequence. Hence by adding each removed dominated vertex to a set containing a representative for it, we obtain a minimal colouring of the original hypergraph. Now, from a q -colouring of $H(\bar{R})$ we can derive at least one representation of R as the intersection of q biorders: consider for each colour class the collection of biorders contained in \bar{R} and containing that class. This collection is non-empty, since any colour class can be expanded to a maximal stable set and these are biorders by Proposition 5.3(1). Clearly any combination of such biorder extensions of the different colour classes covers \bar{R} , hence the intersection of the complementary biorders is R . In this way one can find some of the generally many distinct representations in the minimum dimensionality. If one should really want to have them all, the only way seems to be an—in principle—exhaustive trial of all combinations of Bidim R biorders containing R (or their complements contained in \bar{R}).

For interpretative purposes we may think of a reasonable reduction of the uniqueness problem in that we do not want all representations of R , but only those in which the biorders are minimal. Suppose $R = \bigcap B_i$ for biorders B_i , ($i = 1, \dots, q$). Now for each i we can choose a biorder $B'_i \subseteq B_i$ that still includes R and that is minimal

in this respect. Then still $R = \cap B_i'$ and we may prefer the latter representation because "each B_i is more like R than B_i' is" ($B_i' - R$ is a subset of $B_i - R$). In terms of matrices, a biorder corresponds to a $(0, 1)$ -matrix having triangular structure that represents the hypothetical data matrix on one of the latent dimensions. Thus, writing R as the intersection of a number of biorders is equivalent to writing $[R]$ as the direct logical product of the same number of matrices having triangular structure. Now the restriction to representations in minimal biorders amounts to the fact that if we have a 0 in the observed data matrix, then we assume a 0 in the corresponding position in each hypothetical factor matrix, unless the triangular structure of that matrix forces a 1. To give an extreme example, if an observed 0 can be "explained" by 0's in all dimensions we are not going to explain it by a 0 in one dimension and 1's in all other dimensions. If both representations are possible, one can argue that the former is more likely. (Another possible restriction, one that would be less severe, is to representations that are "obedient" as defined by Chubb (personal correspondence). Among other things this means restriction to biorders that are compatible with the partial orders on the rows and columns of the data matrix.)

If we restrict the class of solutions for R to the set of all representations in minimal biorders, then, in the procedure sketched above for obtaining some representations from a minimal colouring of $H(\bar{R})$, we need only consider expansions of the colour classes to maximal stable sets. By taking complements we obtain from such a collection a representation of R in minimal biorders.

Even with the above restrictions in mind we may expect that in some cases finding all solutions is not practically feasible. More important, however, are the theoretical problems posed by the occurrence of multiple solutions. Can we single out any one from these as "the right one"? If we can appeal to some psychological theory underlying the data this could be possible. For cases in which this does not apply we can try to find a formal rationale to grade the various solutions, thus obtaining a "best one." Another approach could be trying to capture in some explicit characterization the essence common to large classes of solutions. The last two questions will in particular be important when we do not really want a solution of the deterministic model, but instead our ultimate goal is to find a "best fitting" solution in some lower dimensionality of a probabilistic version of the conjunctive model.

REFERENCES

- BERGE, C. (1973) *Graphs and hypergraphs*. Amsterdam: North-Holland.
- CHUBB, C. (1986) Collapsing binary data for algebraic multidimensional representation. *Journal of Mathematical Psychology*, **30**, 161-187.
- COGIS, O. (1980) *La dimension Ferrers des graphes orientés*. These, Université Pierre et Marie Curie, Paris.
- COGIS, O. (1982) On the Ferrers dimension of a digraph. *Discrete Mathematics*, **38**, 47-52.
- COOMBS, C. H. (1964) *A theory of data*. New York: Wiley.

- COOMBS C H, & KAO R C (1955) *Nonmetric factor analysis* Engineering Research Bulletin No 38 Ann Arbor Univ of Michigan Press
- DOIGNON, J-P, DUCAMP, A, & FALMAGNE J-C (1984) On realizable borders and the border dimension of a relation *Journal of Mathematical Psychology* **28** 73-109
- DUCAMP A, & FALMAGNE, J-C (1969) Composite measurement *Journal of Mathematical Psychology* **6** 359-390
- DUSHNIK B, & MILLER, E W (1941) Partially ordered sets *American Journal of Mathematics* **63** 600-610
- GARY M R, & JOHNSON D S (1979) *Computers and intractability. A guide to the theory of NP-completeness* San Francisco Freeman
- GUTTMAN I (1944) A basis for scaling qualitative data *American Sociological Review* **9** 139-150
- MARCOVICI, S (1981) *The multidimensional analysis of partially ordered data. Theory, method and applications* Ph D thesis, New York University New York
- MAURER S B, RABINOVITCH I & TROTTER W T JR (1980) Large minimal realizers of a partial order II *Discrete Mathematics* **31** 297-313
- STOLFFER, S A, GUTTMAN L, SUCHMAN F A, LAZARFIELD P F, STARR S A, & CLAUSEN J A (1950) *Measurement and prediction, Studies in social psychology in World War II* Princeton NJ Princeton Univ Press
- TROTTER W T JR (1983) Graphs and partially ordered sets. In L W Beineke & R J Wilson *Selected topics in graph theory 2* London New York Academic Press
- YANNAKAKIS M (1982) The complexity of the partial order dimension problem *SIAM Journal on Algebraic and Discrete Methods* **3** 351-358

RECEIVED February 20, 1986

CHAPTER 4

FINDING MINIMAL BIORORDER EXTENSIONS OF A RELATION AND MAXIMAL STABLE SETS IN THE ASSOCIATED HYPERGRAPH

(Submitted)

Finding minimal biorder extensions of a relation and maximal stable sets in the associated hypergraph

Mathieu Koppen

New York University

In this paper a number of closely related algorithms are developed that are relevant to the *biorder representation* problem (Doignon, Ducamp & Falmagne, 1984, *JMP*, 28, 73-109). One of the algorithms described here completes, on an algorithmic level, the procedure for computing the *biorder dimension* of a relation as presented in Koppen (1987, *JMP*, 31, 155-178). (In case the relation is a partial order, this biorder dimension coincides with the usual order dimension.)

1. INTRODUCTION

For any binary relation R between two sets A and D we can consider the *biorder representation problem* (Doignon, Ducamp and Falmagne, 1984): write R as the intersection of a minimal number of *biorders* between A and D . (There is a dual problem with intersection replaced by union.) A relation B is a *biorder* between A and D iff for all $a, b \in A$, $d, e \in D$ we have:

$$\text{not } (aBd \ \& \ bBe \ \& \ a\bar{B}e \ \& \ b\bar{B}d),$$

where we introduce the notation \bar{B} for the complement of B . In general, $\bar{R} = A \times D - R$ denotes the complement of a relation R relative to $A \times D$. Any four elements a, b, d, e such that aRd , bRe , $a\bar{R}e$ and $b\bar{R}d$ are said to be a *violation of the biorder property* for the relation R .

Doignon *et al.* show that any relation has a biorder representation and they introduce the notion of the *bidimension* of a relation R as the minimal number of biorders needed for such a representation. If the sets A and D are finite, which we

This work was supported by AFOSR grant F49620-87-C-0131 to New York University and, in an earlier stage, by grant 560-670-006 of the Netherlands Organisation for the Advancement of Pure Research. Address comments and requests for reprints to M. Koppen, Dept. of Psychology NYU, 6 Washington Place, New York, NY 10003.

will assume throughout, any $R \subseteq A \times D$ will have a finite bidimension. Clearly, any biorder B involved in a biorder representation of R contains the relation R ; it is called a *biorder extension* of R . If no strict subset of B is a biorder extension of R , B is called a *minimal biorder extension* of R . These are the biorder extensions that are closest to R (there are no biorders “inbetween”). Any biorder B_j in a representation $R = \bigcap_i B_i$ may be replaced by a minimal biorder B_j' such that $R \subseteq B_j' \subseteq B_j$, and clearly we will still have $R = \bigcap_i B_i'$. This shows that for any representation of R in a minimal number of biorders, we have such a representation where the biorders are minimal extensions of R . This, in itself, makes it interesting to have the collection of minimal biorder extensions of a relation R available, and the next section is devoted to the algorithmic problem of producing this collection for an arbitrary relation R . Another reason why this collection is interesting is given in the following paragraphs.

Doignon *et al.* give a characterization of the bidimension as the chromatic number of some hypergraph $H(\bar{R})$ that is associated with R and Koppen (1987) describes a procedure for determining the bidimension of R by computing this chromatic number $\text{Chrom } H(\bar{R})$. The procedure is constructive in the sense that, once $\text{Chrom } H(\bar{R})$ is known, we have completed at least one minimal coloring of $H(\bar{R})$. The color classes are, by definition, *stable sets* of the hypergraph $H(\bar{R})$ (i.e., they do not contain an edge of $H(\bar{R})$); consequently they may be extended to *maximal stable sets* of $H(\bar{R})$. The essential point now is that such a maximal stable set corresponds directly to a minimal biorder extension of R . Thus, we can obtain some representations of R in minimal biorder extensions simply as a by-product of computing the bidimension of R . Any collection of maximal extensions of the different color classes of a minimal coloring of $H(\bar{R})$ leads directly to such a representation.

This is the motivation for Sections 3 and 4. In Section 3 we first state the correspondence between minimal biorder extensions of R and maximal stable sets of $H(\bar{R})$ and next we describe how the results of Section 2 can be transformed to yield an algorithm for generating the collection of all such maximal stable sets. However, for the above stated purpose we do not need all maximal stable sets, but rather the subcollection of these that are extensions of the constructed colors. Accordingly, we derive in Section 4 how the algorithm of Section 3 can be modified to produce, for a given stable subset of \bar{R} , the collection of maximal stable sets containing this stable set. This algorithm can be used to compute the collections of maximal stable extensions for all colors in a minimal coloring of $H(\bar{R})$, and thus, by the correspondence of Section 3, a number of minimal biorder representations of R .

It appears that some variation of the algorithm of Section 4 is needed to compute $\text{Chrom } H(\bar{R})$ in the first place. In the procedure described in Koppen (1987), the

following recurrence relation is used:

$$\text{Chrom } H(V) = 1 + \min \{ \text{Chrom } H(V-M) : M \text{ is a maximal stable set of } H(V) \text{ containing the vertex } ad \}, \quad (1.1)$$

where V is a subset of \bar{R} , $H(V)$ is the subhypergraph of $H(\bar{R})$ obtained by restricting the set of vertices to V and the set of edges to those included in V , and ad is an element of V , arbitrarily chosen. Obviously, any singleton $\{ad\} \subseteq \bar{R}$ is stable in $H(\bar{R})$ and thus (1.1) poses the problem dealt with in Section 4: produce all maximal stable extensions of the stable set $\{ad\}$. There is however one difference: we want in (1.1) maximal stable sets of a subhypergraph $H(V)$ instead of the full $H(\bar{R})$. In some practical sense this makes the task easier, since it appears that there can be no more maximal stable sets in the subhypergraph than there are in the full hypergraph. Theoretically, the change from $H(\bar{R})$ to $H(V)$ in (1.1) introduces some difficulties, since for maximal stable sets in $H(V)$ the direct correspondence with biorders is lost. These issues are addressed in Sections 5 and 6, where we discuss how the results of the previous sections can be adapted to the case of subhypergraphs of $H(\bar{R})$ and where an alternative algorithm is derived, directly in terms of the subhypergraph $H(V)$. This leads to algorithmic solutions for applying the recurrence relation (1.1), thereby completing the specification of the procedure for computing the bidimension of a relation given in Koppen (1987). As shown in Doignon *et al.* (1984), the bidimension generalizes the notion of the (order) dimension of a partial order, so for the particular case where $D = A$ and the relation R is a partial order on A , the procedure for computing the bidimension of R computes in fact the usual order dimension of R .

Let us introduce here some more concepts and notation, used in subsequent sections. We consider two finite sets A and D , fixed throughout, and relations between A and D , that is, subsets of the Cartesian product $A \times D$. As we have seen above, we write ad for the ordered pair $(a, d) \in A \times D$ and aRd or $ad \in R$, equivalently. A relation $R \subseteq A \times D$ induces subsets of D and A , respectively, in the following way:

$$aR = \{d \in D : aRd\}, \quad a \in A$$

and

$$Rd = \{a \in A : aRd\}, \quad d \in D.$$

For arbitrary $\alpha \subseteq A$, $\delta \subseteq D$ and $R \subseteq A \times D$ we use $R[\alpha, \delta]$ to denote the restriction of R to $\alpha \times \delta$, that is, $R[\alpha, \delta] = R \cap (\alpha \times \delta)$. From this definition we can easily deduce some useful identities, such as

$$R[\alpha_1 * \alpha_2, \delta] = R[\alpha_1, \delta] * R[\alpha_2, \delta],$$

$$R[\alpha, \delta_1 * \delta_2] = R[\alpha, \delta_1] * R[\alpha, \delta_2],$$

where $*$ denotes any of the set operations intersection, union or difference, and

$$R[\alpha_1, \delta_1] \cap R[\alpha_2, \delta_2] = R[\alpha_1 \cap \alpha_2, \delta_1 \cap \delta_2].$$

Any pair of subsets $\alpha \subseteq A$, $\delta \subseteq D$ partitions a relation $R \subseteq A \times D$ into four restrictions:

$$R = R[\alpha, \delta] + R[\alpha, D - \delta] + R[A - \alpha, \delta] + R[A - \alpha, D - \delta].$$

(We will use the plus sign to denote taking the union over mutually disjoint sets.) Clearly, $R_1 \subseteq R_2$ is equivalent to $R_1[X, Y] \subseteq R_2[X, Y]$ for $X \in \{\alpha, A - \alpha\}$, $Y \in \{\delta, D - \delta\}$.

Considering restrictions of relations between A and D it is important to stress that the notion of complement is always meant with respect to the full Cartesian product $A \times D$ and never with respect to some restriction. So $\overline{R[\alpha, \delta]}$ denotes the complement of the restriction of R to $\alpha \times \delta$ while $\overline{R}[\alpha, \delta]$ is the restriction to $\alpha \times \delta$ of the complement of R , and these are not the same; instead we have:

$$\overline{R[\alpha, \delta]} = \overline{R}[\alpha, \delta] + \overline{\alpha \times \delta}.$$

Regarding biorderhood of relations there are some properties that will be used repeatedly in the sequel and that are collected in the following lemma. Here $H(\overline{R})$ denotes the hypergraph that Doignon *et al.* (1984) associated with a relation R . Its vertices are the elements of \overline{R} and its edges the sets consisting of sequences (a, d_i) of elements of \overline{R} such that $a_{i+1}d_i \in R$ (cyclically).

1.1. LEMMA. *Let $R \subseteq A \times D$, $\alpha \subseteq A$, $\delta \subseteq D$.*

- (i) \overline{R} is a biorder if R is.
- (ii) $R[\alpha, \delta]$ is a biorder if R is.
- (iii) $R[\alpha, \delta] + \overline{\alpha \times \delta}$ is a biorder if R is.
- (iv) Any biorder contained in \overline{R} is stable in $H(\overline{R})$.
- (v) Any maximal stable set of $H(\overline{R})$ is a biorder.

That the complement of a biorder is a biorder (i) is immediate from the definition. A restriction of a biorder is a biorder (ii) since any violation of the biorder property for the restriction would constitute a violation for the unrestricted relation. Property (iii) follows from (i) and (ii) since $R[\alpha, \delta] + \overline{\alpha \times \delta}$ is the complement of $\overline{R}[\alpha, \delta]$. Properties (iv) and (v), finally, are not obvious; proofs can be found in Doignon *et al.*, p. 92 and p. 95, respectively.

2. MINIMAL BIORDER EXTENSIONS

Let us first recall the precise definition of a minimal biorder extension of a relation:

2.1. DEFINITION. Let $R, B \subseteq A \times D$. We call B a *minimal biorder extension* of R iff:

- (i) B is a biorder,
- (ii) $R \subseteq B$,
- (iii) if B_0 is a biorder and $R \subseteq B_0 \subseteq B$, then $B_0 = B$.

The collection of all minimal biorder extensions of a relation R will be denoted by $\mathbf{B}(R)$.

In this section we consider the problem of generating, for arbitrary $R \subseteq A \times D$, the collection $\mathbf{B}(R)$. To that end the following definition and lemmas are useful.

2.2. DEFINITION. Let $R \subseteq A \times D$. We call $d_0 \in D$ *minimal in R* iff for all $d \in D$, $Rd \subseteq Rd_0$ implies $Rd = Rd_0$.

2.3. LEMMA. If B is a biorder and $d_0 \in D$ is minimal in B , then $Bd_0 \subseteq Bd$ for any $d \in D$.

Proof. Let $d \in D$. If neither $Bd_0 \subseteq Bd$, nor $Bd \subseteq Bd_0$ we would have a violation of the biorder property. By the minimality of d_0 the case $Bd \subseteq Bd_0$ is equivalent to $Bd = Bd_0$, which means that we always have $Bd_0 \subseteq Bd$. ■

2.4. LEMMA. Let $B, R \subseteq A \times D$, $B \in \mathbf{B}(R)$ and let $Rd_0 = \emptyset$ for some $d_0 \in D$. Then $Bd_0 = \emptyset$.

Proof. $B[A, D - \{d_0\}]$ is a biorder (Lemma 1.1(ii)) and $R = R[A, D - \{d_0\}] \subseteq B[A, D - \{d_0\}] \subseteq B$, so the minimality of B implies $B = B[A, D - \{d_0\}]$, which means that $Bd_0 = \emptyset$. ■

2.5. LEMMA. Let $B, R \subseteq A \times D$, $B \in \mathbf{B}(R)$, and let $d_0 \in D$ be minimal in B . Then d_0 is minimal in R and $Bd_0 = Rd_0$.

Proof. Let $Rd_1 \subseteq Rd_0$ for some $d_1 \in D$. We need to show that $Rd_1 = Rd_0 = Bd_0$. Define $B' = B[A, D - \{d_1\}] + R[A, \{d_1\}]$. Then, clearly, $R \subseteq B' \subseteq B$ and, moreover, B' is a biorder. For, because B' equals the biorder B except for d_1 , the element d_1 must be involved in any violation. But we have $B'd_1 = Rd_1 \subseteq Rd_0 \subseteq Bd_0 \subseteq Bd$ for any $d \in D$ (the last inclusion by Lemma 2.3) and since $Bd = B'd$ for any $d \neq d_1$, we

obtain $B'd_1 \subseteq B'd$ for any $d \in D$, which excludes this possibility. B being minimal, we conclude $B' = B$ and thus $Rd_0 \subseteq Bd_0 \subseteq Bd_1 = B'd_1 = Rd_1 \subseteq Rd_0$, from which the desired equality $Rd_1 = Rd_0 = Bd_0$ immediately follows. ■

The next two propositions will be at the base of our algorithm for generating $\mathbf{B}(R)$.

2.6. PROPOSITION. *Let $B, R \subseteq A \times D$. Then $B \in \mathbf{B}(R)$ iff there is $d_0 \in D$, minimal in R , such that $B[Rd_0, D] = Rd_0 \times D$ and $B[\bar{R}d_0, D] = B'$, with $B' \in \mathbf{B}(R')$, where $R' = R[\bar{R}d_0, D] = R[\bar{R}d_0, D - D_0]$ and $D_0 = \{d \in D : Rd = Rd_0\}$.*

Proof. (If.) Suppose we have d_0 and B' with the above properties. B' is a biorder, so, by Lemma 1.1(iii), $B'[\bar{R}d_0, D] + (A \times D - \bar{R}d_0 \times D) = B' + Rd_0 \times D = B$ is a biorder. Furthermore, $R = R[Rd_0, D] + R[\bar{R}d_0, D - D_0] \subseteq Rd_0 \times D + B' = B$. To show the minimality of B , suppose $R \subseteq B_0 \subseteq B$ for some biorder B_0 . This implies $R' \subseteq B_0[\bar{R}d_0, D] \subseteq B'$ and because $B' \in \mathbf{B}(R')$ we obtain $B_0[\bar{R}d_0, D] = B[\bar{R}d_0, D] = B'$. On the other hand, using $B'd_0 = \emptyset$, which follows by Lemma 2.4 from the obvious fact $R'd_0 = \emptyset$, and using Lemma 2.3, we get for d_1 , minimal in B_0 (because of finiteness such a d_1 can always be found): $Rd_1 \subseteq B_0d_1 \subseteq B_0d_0 = B_0[Rd_0, D]d_0 + B'd_0 = B_0[Rd_0, D]d_0 \subseteq Rd_0$. Since d_0 is minimal in R this implies $Rd_1 = B_0d_1 = B_0d_0 = Rd_0$ and because d_1 is minimal in B_0 we obtain, by Lemma 2.3, $Rd_0 = B_0d_1 \subseteq B_0d$ for any $d \in D$. In other words, $B_0[Rd_0, D] = Rd_0 \times D = B[Rd_0, D]$ and we conclude that $B_0 = B$.

(Only if.) Suppose $B \in \mathbf{B}(R)$. Choose $d_0 \in D$, minimal in B . Then, by Lemmas 2.5 and 2.3, $Rd_0 = Bd_0 \subseteq Bd$ for any $d \in D$, which means that $B[Rd_0, D] = Rd_0 \times D$. So we are left to show that $B' \in \mathbf{B}(R')$ for $B' = B[\bar{R}d_0, D]$ and $R' = R[\bar{R}d_0, D - D_0]$. By Lemma 1.1(ii), B' is a biorder and since $R[\bar{R}d_0, D - D_0] = R[\bar{R}d_0, D] \subseteq B[\bar{R}d_0, D]$, the only point of concern is the minimality of B' . Suppose there is a biorder B'_0 such that $R' \subseteq B'_0 \subseteq B'$. Define $B_0 = Rd_0 \times D + B'_0$; by Lemma 1.1(iii), B_0 is a biorder. Since $R = R[Rd_0, D] + R' \subseteq Rd_0 \times D + B'_0 = B_0 \subseteq Rd_0 \times D + B' = B$, it follows from the minimality of B that $B_0 = B$ and thus, in particular, that $B'_0 = B_0[\bar{R}d_0, D] = B[\bar{R}d_0, D] = B'$. ■

2.7. PROPOSITION. *Let $B, R \subseteq A \times D$, $B \in \mathbf{B}(R)$, such that for $d_0, d'_0 \in D$ we have $B[Rd_0, D] = Rd_0 \times D$ and $B[Rd'_0, D] = Rd'_0 \times D$. Then both d_0 and d'_0 are minimal in R and $Rd_0 = Rd'_0$.*

Proof. Choose d_1 minimal in B (possible since D is finite). According to Lemma 2.5, d_1 is minimal in R and $Bd_1 = Rd_1$. From $B[Rd_0, D] = Rd_0 \times D$ and

$B[Rd'_0, D] = Rd'_0 \times D$, respectively, it follows that $Rd_0 \subseteq Bd_1$ and $Rd'_0 \subseteq Bd_1$. Using the minimality of d_1 in R we obtain $Rd_0 = Bd_1 = Rd_1 = Rd'_0$ and we see that both d_0 and d'_0 are minimal in R . ■

2.8. ALGORITHM. The Propositions 2.6 and 2.7 can be used to generate $\mathbf{B}(R)$ for a given relation $R \subseteq A \times D$. More precisely, Proposition 2.6 shows that the following recursive procedure MBE generates exactly $\mathbf{B}(R)$ and Proposition 2.7 shows that each element is generated just once:

```

procedure MBE ( $\alpha \subseteq A, \delta \subseteq D, B \subseteq A \times D$ ):
if  $\delta = \emptyset$  then output ( $B$ )
else
  enumerate all pairs of sets ( $A_i, D_i$ ) where
     $A_i = R[\alpha, \delta]d$  for some  $d \in \delta$ , minimal in  $R[\alpha, \delta]$ ,
    and where  $D_i = \{d \in \delta : R[\alpha, \delta]d = A_i\}$ ;
  for each pair ( $A_i, D_i$ )
  do
    MBE ( $\alpha - A_i, \delta - D_i, B + A_i \times \delta$ )
  od
fi.

```

Using this procedure, we have the following algorithm for generating $\mathbf{B}(R)$ for an arbitrary relation $R \subseteq A \times D$:

```

INITIALIZE ( $A, D, R$ );
MBE ( $A, D, \emptyset$ ).

```

2.9. We see that the only problem in MBE is to compute the collection of pairs (A_j, D_j) . In order to decide whether $d \in \delta$ is minimal in $R[\alpha, \delta]$ we need to compute $R[\alpha, \delta]d$; this means that we will have to compute the sets A_j corresponding to each $d \in \delta$ in order to select the ones corresponding to minimal elements of δ . So suppose that in MBE (α, δ, B) we have computed all pairs (A_j, D_j) where $A_j = R[\alpha, \delta]d$ for some $d \in \delta$ and where $D_j = \{d \in \delta : R[\alpha, \delta]d = A_j\}$ and have marked the pairs that correspond to elements of δ that are not minimal in $R[\alpha, \delta]$. If (A_i, D_i) is an unmarked pair, we will have the call MBE ($\alpha - A_i, \delta - D_i, B + A_i \times \delta$) and we will have to compute all pairs (A'_h, D'_h) for the relation $R[\alpha - A_i, \delta - D_i]$. These, however, can easily be obtained from the collection computed in the calling procedure: each A'_h equals $A_j - A_i$ for some A_j of the calling procedure and each D'_h equals the union over those D_l of the calling procedure for which $A_l - A_i$ equals A'_h . In this way we can get all collections of pairs needed in subsequent calls of MBE, once we have established the collection of pairs for the original relation R .

2.10. Any biorder between A and D is a subset of $A \times D$ and as such has an $O(|A| + |D|)$ representation (enumeration of its elements). Since a biorder is a special kind of relation, however, it allows a more economical $O(|A| + |D|)$ representation, namely as a sequence

$$A_1, D_1, \dots, A_k, D_k,$$

where $\{A_i\}$ is a partition of A (possibly $A_1 = \emptyset$) and $\{D_i\}$ is a partition of D (possibly $D_k = \emptyset$). An element ad of $A \times D$ belongs to the biorder iff the class of a precedes the class of d , that is, iff, in the above representation, $a \in A_i$ and $d \in D_j$ imply $i \leq j$. So, (A_1, D_1) being the first pair in the representation means that the biorder consists of $A_1 \times D$ plus a biorder restricted to $(A - A_1) \times (D - D_1)$. In other words, this representation mirrors exactly the way in which the biorders are constructed by Algorithm 2.8. Such a biorder can thus be identified with the sequence of pairs (A_i, D_i) for which $\text{MBE}(\alpha - A_i, \delta - D_i, B + A_i \times \delta)$ was subsequently called to produce it and by the procedure described in 2.9 we can, starting from the collection of pairs (A_i, D_i) corresponding to the original relation R , construct all possible such sequences.

2.11. EXAMPLE. Suppose we have sets $A = \{a, b, c, d, e, f\}$ and $D = \{u, v, w, x, y, z\}$ and the following relation R between these sets:

	u	v	w	x	y	z
a	R	—	—	R	—	—
b	—	R	R	R	R	R
c	R	R	—	—	R	R
d	R	—	R	—	—	R
e	—	—	R	R	R	R
f	—	R	—	R	—	R

From this matrix we can establish the following pairs of sets for this relation R (for simplicity we denote sets without braces or separators, that is, we write abc for the set $\{a, b, c\}$):

acd u
 bcf v
 bde w
 $abef$ x
 bce y
 $bcdef$ z

where only the last pair belongs to an element of D that is not minimal in R (Rz

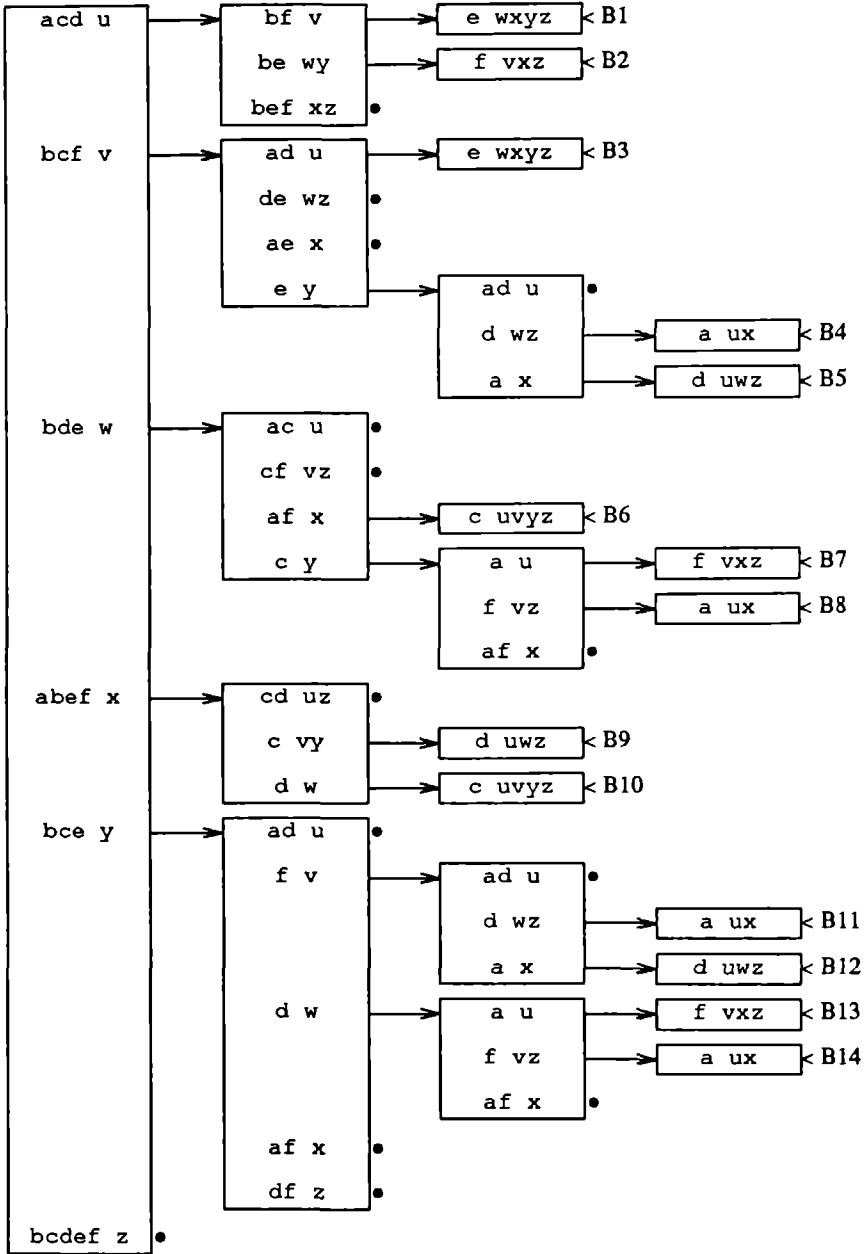


Figure 2.1. Generating $B(R)$ for Example 2.11. For explanation, see text.

	u	v	w	x	y	z
a	1	!	!	1	!	!
b	.	1	1	1	1	1
c	1	1	!	!	1	1
d	1	!	1	!	!	1
e	.	.	1	1	1	1
f	.	1	!	1	!	1

B1

	u	v	w	x	y	z
a	1	!	!	1	!	!
b	.	1	1	1	1	1
c	1	1	!	!	1	1
d	1	!	1	!	!	1
e	.	!	1	1	1	1
f	.	1	.	1	.	1

B2

	u	v	w	x	y	z
a	1	.	!	1	!	!
b	!	1	1	1	1	1
c	1	1	!	!	1	1
d	1	.	1	!	!	1
e	.	.	1	1	1	1
f	!	1	!	1	!	1

B3

	u	v	w	x	y	z
a	1	.	.	1	.	.
b	!	1	1	1	1	1
c	1	1	!	!	1	1
d	1	.	1	!	.	1
e	!	.	1	1	1	1
f	!	1	!	1	!	1

B4

	u	v	w	x	y	z
a	1	.	!	1	.	!
b	!	1	1	1	1	1
c	1	1	!	!	1	1
d	1	.	1	.	.	1
e	!	.	1	1	1	1
f	!	1	!	1	!	1

B5

	u	v	w	x	y	z
a	1	!	.	1	!	!
b	!	1	1	1	1	1
c	1	1	.	.	1	1
d	1	!	1	!	!	1
e	!	!	1	1	1	1
f	!	1	.	1	!	1

B6

	u	v	w	x	y	z
a	1	!	.	1	.	!
b	!	1	1	1	1	1
c	1	1	.	!	1	1
d	1	!	1	!	!	1
e	!	!	1	1	1	1
f	.	1	.	1	.	1

B7

	u	v	w	x	y	z
a	1	.	.	1	.	.
b	!	1	1	1	1	1
c	1	1	.	!	1	1
d	1	!	1	!	!	1
e	!	!	1	1	1	1
f	!	1	.	1	.	1

B8

	u	v	w	x	y	z
a	1	!	!	1	!	!
b	!	1	1	1	1	1
c	1	1	!	.	1	1
d	1	.	1	.	.	1
e	!	!	1	1	1	1
f	!	1	!	1	!	1

B9

	u	v	w	x	y	z
a	1	!	!	1	!	!
b	!	1	1	1	1	1
c	1	1	.	.	1	1
d	1	!	1	.	.	1
e	!	!	1	1	1	1
f	!	1	!	1	!	1

B10

	u	v	w	x	y	z
a	1	.	.	1	.	.
b	!	1	1	1	1	1
c	1	1	!	!	1	1
d	1	.	1	!	.	1
e	!	!	1	1	1	1
f	!	1	!	1	.	1

B11

	u	v	w	x	y	z
a	1	.	!	1	.	!
b	!	1	1	1	1	1
c	1	1	!	!	1	1
d	1	.	1	.	.	1
e	!	!	1	1	1	1
f	!	1	!	1	.	1

B12

	u	v	w	x	y	z
a	1	!	.	1	.	!
b	!	1	1	1	1	1
c	1	1	!	!	1	1
d	1	!	1	!	.	1
e	!	!	1	1	1	1
f	.	1	.	1	.	1

B13

	u	v	w	x	y	z
a	1	.	.	1	.	.
b	!	1	1	1	1	1
c	1	1	!	!	1	1
d	1	!	1	!	.	1
e	!	!	1	1	1	1
f	!	1	.	1	.	1

B14

Figure 2.2. Matrix representation of the minimal border extensions constructed in Fig. 2.1. Elements of R are denoted by "1"; elements of \bar{R} , added to obtain the border, by "!".

includes strictly Rv , for instance). Generating $B(R)$ by Algorithm 2.8 according to the specifications given in 2.9 and 2.10 is displayed in the tableau of Figure 2.1. Here, collections of pairs computed in one call of MBE are grouped together in boxes, the leftmost box containing the pairs belonging to the original relation R . Pairs of sets that correspond to non-minimal elements are marked by a bullet; unmarked pairs lead to a recursive call of MBE which is represented by an arrow from that pair pointing to a box to the right. The content of this box can directly be derived from that of the parent box: according to 2.9 we must copy all other pairs of the parent box (i.e., apart from the pair that originated the call), while discarding from their left members the elements of the left member of the originating pair and next joining pairs with identical left members. Finally, the different left members are compared in order to detect non-minimal elements to be marked. This process of choosing a pair in the current box and constructing from it a box to the right stops when there are no more other pairs to be copied (the next box would be empty), that is, when we have a box containing just one pair. At this point we have finished another biorder, represented by the concatenation of the pairs along the path from the first column to this endpoint, and we go backtracking to find the first occasion of completing a next biorder. We see in Fig. 2.1 that in this case $B(R)$ has 14 elements which we have labeled $B 1$ to $B 14$ and the matrix representation of which is given separately in Figure 2.2.

2.12. A STACK IMPLEMENTATION. Following the discussion in 2.9 and 2.10, we can think of a non-recursive version of Algorithm 2.8 that operates on some sort of stack. A typical element of this stack is the collection of pairs of sets (A_i, D_i) belonging to some restriction $R[\alpha, \delta]$ that turns up in the computation. That is, the elements of the stack are boxes as in Fig. 2.1. On this stack the following operations are defined:

LOAD (R):

- computes the collection of pairs (Rd_i, d_i) where d_i runs through D ;
- pushes this collection as the top box on the empty stack.

POP: pops the top box off the stack.

PUSH: pushes a copy of the top box on the stack.

ACTIVATE (A_0, D_0):

- chooses an unmarked pair (A_0, D_0) in the top box as the "active" pair;
- marks the active pair (just as MARK marks some pairs).

TRIM (A_0, D_0):

- removes in top box the active pair (A_0, D_0) ;
- removes from all left members of top box tripartitions the elements of A_0 , the

left member of the active pair.

MERGE :

merges in the top box pairs with identical left members into one pair by taking the union over their right members.

MARK :

marks in the top box pairs with non-minimal left members (non-minimal compared to other left members in the box).

In addition the following tests can be made:

EMPTY : returns TRUE iff the stack is empty.

EXHAUSTED :

returns TRUE iff the stack is not empty and all pairs in the top box are marked.

SINGLETON : returns TRUE iff the top box contains just one pair.

Finally we assume a procedure

OUTPUT-RESULT :

yields the sequence of active pairs, starting at the bottom of the stack and ending with the (!) pair of the top box.

2.13. A NON-RECURSIVE ALGORITHM. With the specifications of 2.12, the following algorithm is an iterative elaboration of Algorithm 2.8:

```

INITIALIZE (A, D, R);
LOAD (R);
MERGE;
MARK;
while not EMPTY
do
  while not SINGLETON
  do
    ACTIVATE (A0, D0);
    PUSH;
    TRIM (A0, D0);
    MERGE;
    MARK
  od;
  OUTPUT-RESULT;
  while POP; EXHAUSTED do od
od.

```

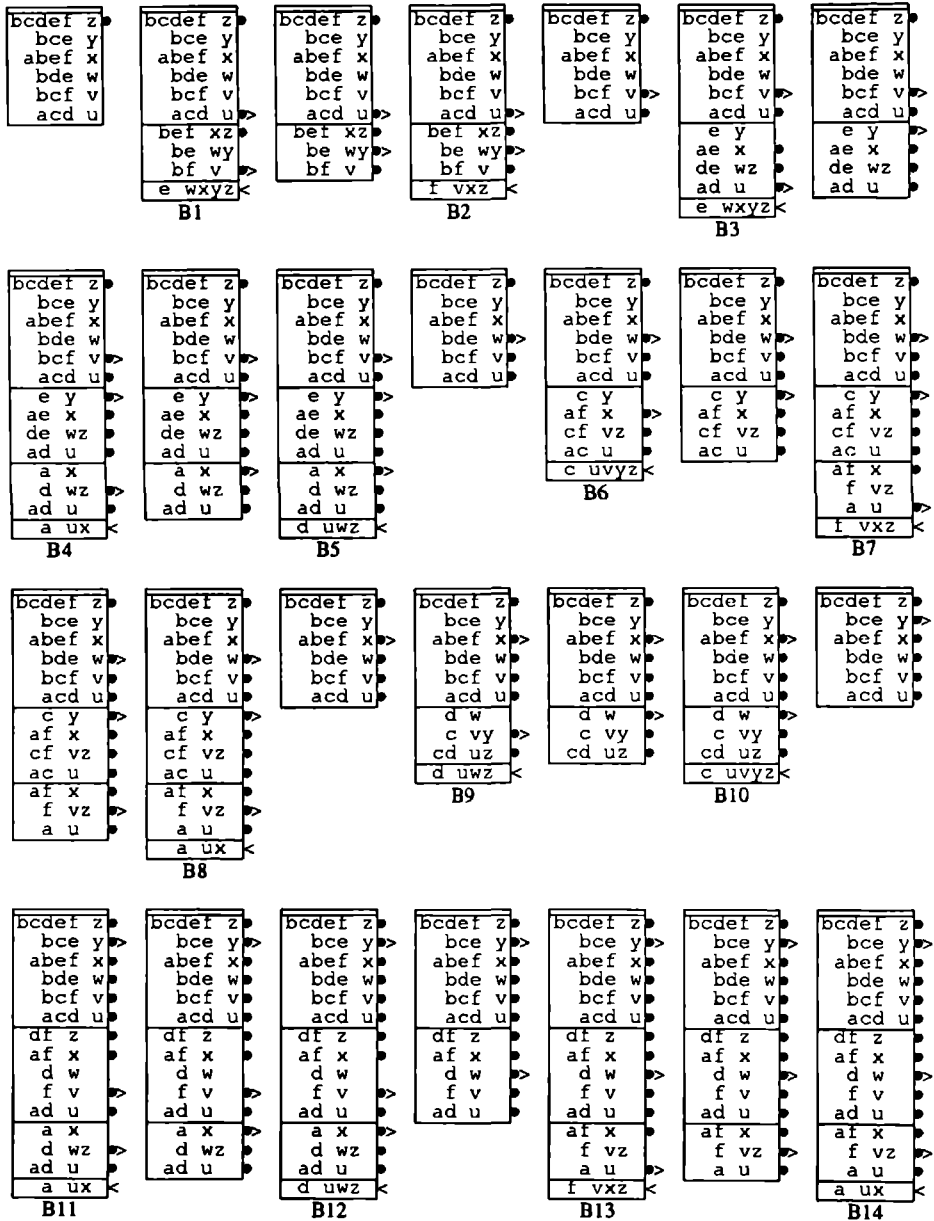


Figure 2.3. Algorithm 2.13 working on Example 2.11. For explanation, see text.

2.14. Figure 2.3 shows the history of the stack when this algorithm is applied to the relation of Example 2.11. We start in the upper left corner with the state of the stack just before entrance of the outer loop and move to the right and down until we end up with the empty stack. The stack grows downwards; “•” is the marking symbol and “>” denotes an active pair. ACTIVATE chooses the lowest unmarked pair in the top box as the active pair. The connection with Fig. 2.1 is obvious.

3. MAXIMAL STABLE SETS

In this section we consider the problem of finding all maximal stable sets of the hypergraph $H(\bar{R})$, for $R \subseteq A \times D$. We recall the definition of a maximal stable set:

3.1. DEFINITION. Let $M, R \subseteq A \times D$. We call M a *maximal stable set* of the hypergraph $H(\bar{R})$ iff

- (i) $M \subseteq \bar{R}$,
- (ii) M is stable in $H(\bar{R})$ and
- (iii) if M_0 is stable in $H(\bar{R})$ and $M \subseteq M_0 \subseteq \bar{R}$, then $M_0 = M$.

The collection of all maximal stable sets of the hypergraph $H(\bar{R})$ will be denoted by $S(\bar{R})$.

The next lemma shows that there exists a close connection between minimal border extensions of R and maximal stable sets of $H(\bar{R})$.

3.2. LEMMA. Let $M, R \subseteq A \times D$. Then $M \in S(\bar{R})$ iff $\bar{M} \in B(R)$.

Proof. By Lemma 1.1(v), any element of $S(\bar{R})$ is a border and by Lemma 1.1(iv) we see that $S(\bar{R})$ coincides with the collection of maximal borders contained in \bar{R} . Clearly, the complement of any border contained in \bar{R} is a border extension of R and vice versa (Lemma 1.1(i)); maximal borders in \bar{R} correspond in this way to minimal border extensions of R . ■

3.3. Lemma 3.2 entails that finding $S(\bar{R})$ really is the same problem as finding $B(R)$. In fact, replacing the statement “output(B)” in the procedure MBE of Algorithm 2.8 by the statement “output(\bar{B})” turns this algorithm into one for generating $S(\bar{R})$. However, since we will, in subsequent sections, be interested not so much in generating all of $S(\bar{R})$ but rather in generating certain subsets thereof, it will be convenient to have versions of the results of the preceding section in terms of subsets of \bar{R} instead of the complementary sets and versions of Algorithms 2.8 and

2.13 that use the $O(|A| \cdot |D|)$ representation of maximal stable sets as real subsets of $A \times D$, rather than the possible $O(|A| + |D|)$ representation of them as biorders. Therefore we give here the complementary versions of Propositions 2.6 and 2.7 and Algorithm 2.8.

3.4. PROPOSITION. *Let $R \subseteq A \times D$, $M \subseteq \bar{R}$. Then $M \in \mathbf{S}(\bar{R})$ iff for some $d_0 \in D$, minimal in R , and with $D_0 = \{d \in D : \bar{R}d = \bar{R}d_0\}$, we have $M = M[\bar{R}d_0, D] = \bar{R}d_0 \times D_0 + M'$, with $M' \in \mathbf{S}(\bar{R}[\bar{R}d_0, D - D_0])$.*

Proof. Directly from Proposition 2.6 by Lemma 3.2. ■

3.5. PROPOSITION. *Let $M \in \mathbf{S}(\bar{R})$ for $R \subseteq A \times D$ such that for $d_0, d'_0 \in D$ we have $M = M[\bar{R}d_0, D] = M[\bar{R}d'_0, D]$. Then both d_0 and d'_0 are minimal in R and $\bar{R}d_0 = \bar{R}d'_0$.*

Proof. Directly from Proposition 2.7 by Lemma 3.2. ■

3.6. ALGORITHM. $\mathbf{S}(\bar{R})$ is produced by the following recursive procedure MSS:

```

procedure MSS ( $\alpha \subseteq A$ ,  $\delta \subseteq D$ ,  $M \subseteq A \times D$ ):
  if  $\delta = \emptyset$  then output ( $M$ )
  else
    enumerate all pairs of sets ( $A_i, D_i$ ) where
       $A_i = \bar{R}[\alpha, \delta]d$  for some  $d \in \delta$ , minimal in  $R[\alpha, \delta]$ ,
      and where  $D_i = \{d \in \delta : \bar{R}[\alpha, \delta]d = A_i\}$ ;
    for each pair ( $A_i, D_i$ )
      do
        MSS ( $A_i, \delta - D_i, M + A_i \times D_i$ )
    od
  fi.

```

Using this procedure the algorithm

```

INITIALIZE ( $A, D, R$ );
MSS ( $A, D, \emptyset$ )

```

generates $\mathbf{S}(\bar{R})$ for an arbitrary relation $R \subseteq A \times D$. Correctness follows either from Propositions 3.4 and 3.5 or, by Lemma 3.2, from the correctness of Algorithm 2.8, if we note that the sets, generated by each procedure in turn, are complementary.

3.7. As with MBE, the only problem in MSS consists in computing the collection of pairs (A_j, D_j). As before, we will have to compute the sets A_j corresponding to each $d \in \delta$ in order to select the ones corresponding to minimal elements of δ . (An element of $D_j \subseteq \delta$ is minimal iff the corresponding A_j is not (strictly) included in

any other set A_i .) Suppose in $MSS(\alpha, \delta, M)$ we have computed all pairs (A_j, D_j) with $A_j = \bar{R}[\alpha, \delta]d$ for some $d \in \delta$ and $D_j = \{d \in \delta : \bar{R}[\alpha, \delta]d = A_j\}$, and have marked the pairs that correspond to elements of δ that are not minimal in $\bar{R}[\alpha, \delta]$. If (A_i, D_i) is an unmarked pair we will have the call $MSS(A_i, \delta - D_i, M + A_j \times D_j)$ and we will have to compute all pairs (A_h', D_h') in $\bar{R}[A_i, \delta - D_i]$. These, however, are again easily obtained from the collection computed in the calling procedure: this time each A_h' equals $A_j \cap A_i$ for some A_j of the calling procedure and again each D_h' equals the union over those D_l of the calling procedure for which $A_l \cap A_i$ equals A_h' . So, as before, we can get all collections of pairs needed in subsequent calls of MSS once we have established the collection of pairs for the original relation R .

3.8. Clearly each maximal set constructed in Algorithm 3.6 is completely determined by the sequence of calls of MSS that produced it, that is, by the sequence of pairs (A_i, D_i) that were chosen in each subsequent call after the invocation $MSS(A, D, \emptyset)$. Now what are the consequences of choosing a particular pair for the maximal sets based upon this choice? Suppose that, in constructing a maximal set M_0 , we have had the call $MSS(\alpha, \delta, M)$ and that there we have the candidate pair (A_i, D_i) , corresponding to some element of δ that is minimal in $\bar{R}[\alpha, \delta]$. First we must note that being in $MSS(\alpha, \delta, M)$ we can add only elements of $\bar{R}[\alpha, \delta]$ to M . This means that outside of $\bar{R}[\alpha, \delta]$ the issue of membership of M_0 has already been decided: elements from there included in M will be included in M_0 and elements from there not included in M will not be present in M_0 . It is just for the set $\bar{R}[\alpha, \delta]$ that decisions about admission to M_0 have to be made. With respect to the pair (A_i, D_i) to be considered in $MSS(\alpha, \delta, M)$, we have following partition of $\bar{R}[\alpha, \delta]$:

$$\bar{R}[\alpha, \delta] = \bar{R}[A_i, D_i] + \bar{R}[A_i, \delta - D_i] + \bar{R}[\alpha - A_i, D_i] + \bar{R}[\alpha - A_i, \delta - D_i],$$

where $\bar{R}[\alpha - A_i, D_i]$ is empty by definition. Since to (A_i, D_i) corresponds the call $MSS(A_i, \delta - D_i, M + A_i \times D_i)$, we see that this choice implies the following for the set M_0 in construction:

- inclusion in M_0 of $\bar{R}[A_i, D_i] = A_i \times D_i$,
- exclusion from M_0 of $\bar{R}[\alpha - A_i, \delta - D_i] = \bar{R}[\alpha - A_i, \delta]$,
- further search for M_0 in $\bar{R}[A_i, \delta - D_i]$.

Thus, each pair (A_i, D_i) computed in $MSS(\alpha, \delta, M)$ corresponds to a tripartition of the elements of \bar{R} not yet decided upon into a set for which, when MSS is called with this pair as an argument, the decision is positive (we will call this the P -class of the tripartition), one for which the decision is negative (the N -class) and one for which the decision is postponed (the U -class, U for undecided). Denoting the P -class (etc.) of the tripartition corresponding to the pair (A_i, D_i) by $P(A_i, D_i)$ (etc.),

we have

$$P(A_i, D_i) = \bar{R}[A_i, D_i] = A_i \times D_i,$$

$$N(A_i, D_i) = \bar{R}[\alpha - A_i, \delta - D_i] = \bar{R}[\alpha - A_i, \delta],$$

$$U(A_i, D_i) = \bar{R}[A_i, \delta - D_i].$$

3.9. As we saw in 3.7 that the collection of pairs needed in a call of MSS can be computed from the collection in the calling procedure, the same must hold for the corresponding tripartitions. In fact, suppose that (A_h', D_h') is a pair occurring in MSS $(A_i, \delta - D_i, M + A_i \times D_i)$ when this procedure is called by MSS (α, δ, M) and let l be an index such that $A_h' = A_l \cap A_i$, or, when appropriate, let l run through all such indices. From 3.7 we know that this index set will be non-empty. Now the following straightforward computations yield the tripartition belonging to (A_h', D_h') in terms of the tripartitions in the calling procedure:

$$\begin{aligned} P(A_h', D_h') &= \bar{R}[A_h', D_h'] = \bar{R}[A_h', \cup_l D_l] = \cup_l \bar{R}[A_h', D_l] \\ &= \cup_l \bar{R}[A_l \cap A_i, D_l] = \cup_l (\bar{R}[A_l, D_l] \cap \bar{R}[A_i, D_l]) \\ &= \cup_l (\bar{R}[A_l, D_l] \cap \bar{R}[A_i, \delta - D_i]) \quad (\text{since } D_l \subseteq \delta - D_i) \\ &= (\cup_l \bar{R}[A_l, D_l]) \cap \bar{R}[A_i, \delta - D_i] \\ &= (\cup_l P(A_l, D_l)) \cap U(A_i, D_i); \\ N(A_h', D_h') &= \bar{R}[A_i - A_h', \delta - D_i] = \bar{R}[A_i - (A_l \cap A_i), \delta - D_i] \\ &= \bar{R}[A_i \cap (\alpha - A_l), \delta - D_i] = \bar{R}[\alpha - A_l, \delta - D_i] \cap \bar{R}[A_i, \delta - D_i] \\ &= \bar{R}[\alpha - A_l, \delta] \cap \bar{R}[A_i, \delta - D_i] \quad (\text{since } \delta - D_i \subseteq \delta) \\ &= N(A_l, D_l) \cap U(A_i, D_i); \\ U(A_h', D_h') &= \bar{R}[A_h', \delta - D_i - D_h'] = \bar{R}[A_h', \delta - D_i - (\cup_l D_l)] \\ &= \cap_l \bar{R}[A_h', \delta - D_i - D_l] = \cap_l \bar{R}[A_l \cap A_i, \delta - D_i - D_l] \\ &= \cap_l \bar{R}[A_l \cap A_i, (\delta - D_i) \cap (\delta - D_l)] \\ &= \cap_l (\bar{R}[A_l, \delta - D_i] \cap \bar{R}[A_i, \delta - D_l]) \\ &= (\cap_l \bar{R}[A_l, \delta - D_l]) \cap \bar{R}[A_i, \delta - D_i] \\ &= (\cap_l U(A_l, D_l)) \cap U(A_i, D_i). \end{aligned}$$

3.10. LEMMA. *Let $R \subseteq A \times D$, $\alpha \subseteq A$, $\delta \subseteq D$, $d_0 \in \delta$, $A_0 = \bar{R}[\alpha, \delta]d_0$ and $D_0 = \{d \in \delta : \bar{R}[\alpha, \delta]d = A_0\}$. Then d_0 is minimal in $R[\alpha, \delta]$ iff $N(A_0, D_0)$ is minimal compared to other $N(A_j, D_j)$, for pairs (A_j, D_j) computed in $MSS(\alpha, \delta, M)$.*

Proof. It suffices to show $Rd_j \subseteq Rd_0$ iff $N(A_j, D_j) \subseteq N(A_0, D_0)$. We have following equivalences: $N(A_j, D_j) \subseteq N(A_0, D_0)$ iff $\bar{R}[\alpha - A_j, \delta] \subseteq \bar{R}[\alpha - A_0, \delta]$ iff $\bar{R}[\alpha, \delta] - \bar{R}[\alpha - A_j, \delta] \supseteq \bar{R}[\alpha, \delta] - \bar{R}[\alpha - A_0, \delta]$ iff $\bar{R}[A_j, \delta] \supseteq \bar{R}[A_0, \delta]$ iff $\bar{R}[\bar{R}d_j, \delta] \supseteq \bar{R}[\bar{R}d_0, \delta]$. Now first assume $Rd_j \subseteq Rd_0$; this means $\bar{R}d_j \supseteq \bar{R}d_0$ and consequently $\bar{R}[\bar{R}d_j, \delta] \supseteq \bar{R}[\bar{R}d_0, \delta]$. For the reverse implication note that by definition we have $\bar{R}d_0 \times D_0 \subseteq \bar{R}[\bar{R}d_0, \delta]$. So $\bar{R}[\bar{R}d_0, \delta] \subseteq \bar{R}[\bar{R}d_j, \delta]$ implies $\bar{R}d_0 \times D_0 \subseteq \bar{R}[\bar{R}d_j, \delta]$ and this is only possible if $\bar{R}d_0 \subseteq \bar{R}d_j$, that is, if $Rd_j \subseteq Rd_0$. ■

3.11. A STACK IMPLEMENTATION. By 3.8 and 3.10 we have a translation of Algorithm 3.6 in terms of tripartitions; it can be viewed as generating all sequences of tripartitions where each next tripartition has as its domain the U -class of its predecessor and has a minimal N -class compared to alternative tripartitions at this point. Each such sequence starts with a tripartition corresponding to the original relation and ends with one having an empty U -class. In 3.9 we saw how to compute the collection of tripartitions at any level from the collection at the previous level and this suggests again a stack implementation of Algorithm 3.6, very similar to the one described in 2.12 and 2.13 for Algorithm 2.8. Here the stack elements will be boxes consisting of tripartitions (P_i, N_i, U_i) instead of pairs of sets (A_i, D_i) . From 2.12 we take over the procedures POP, PUSH and EMPTY and also ACTIVATE and EXHAUSTED, on the understanding that the elements of the top box are now tripartitions instead of pairs of sets. We need following new versions of the procedures LOAD, TRIM, MERGE, MARK and OUTPUT-RESULT which have an apparent connection to the old versions:

LOAD (R) :

computes the collection of tripartitions (P_i, N_i, U_i) , where, for some $d_i \in D$,
 $P_i = \bar{R}d_i \times \{d_i\}$, $N_i = \bar{R}[Rd_i, D]$ and $U_i = \bar{R}[\bar{R}d_i, D - \{d_i\}]$;
 pushes this collection as the top box on the empty stack.

TRIM (P_0, N_0, U_0) :

removes from all classes of the tripartitions in the top box elements that are not in U_0 .

MERGE :

removes tripartitions with empty P -classes;
 merges tripartitions with identical N -classes into one by taking the union over their P -classes and the intersection over their U -classes.

MARK :

marks tripartitions with non-minimal N -classes (non-minimal compared to other N -classes in the box).

OUTPUT-RESULT :

yields the set of elements present in the P -classes of the active tripartitions at the different levels.

The procedures TRIM and MERGE do the computations of 3.9 (for the removal of tripartitions with empty P -classes in MERGE, see below) and MARK is based on Lemma 3.10. Finally we need a new termination condition, testing whether another maximal stable set has been completed:

DONE (P_0, N_0, U_0):

returns TRUE iff U_0 , the U -class of the active tripartition, is empty.

3.12. A NON-RECURSIVE ALGORITHM. Using the procedures described in 3.11, we have following iterative version of Algorithm 3.6:

```

INITIALIZE ( $A, D, R$ );
LOAD ( $R$ );
MERGE;
MARK;
while not EMPTY
do
  while
    ACTIVATE ( $P_0, N_0, U_0$ );
    not DONE ( $P_0, N_0, U_0$ )
  do
    PUSH;
    TRIM ( $P_0, N_0, U_0$ );
    MERGE;
    MARK
  od;
  OUTPUT-RESULT;
  while EXHAUSTED do POP od
od.

```

The only aspect of the algorithm still needing justification seems to be why we are allowed to discard tripartitions with empty P -classes in the procedure MERGE. (Note that, after TRIM, this includes the active tripartition.) Though it is clear that such a tripartition cannot introduce any elements in the maximal sets in construction, it could in Algorithm 3.12 in principle be necessary to ACTIVATE such a tripartition in order to fulfill the termination condition DONE. This is, however, not the case, since (i) the collection of P -classes of the LOADED tripartitions clearly is a

partition of \bar{R} , and (ii) by the definitions of PUSH, TRIM and MERGE it follows from (i) that at any time the collection of P -classes in the top box constitutes a partition of the U -class of the ACTIVATED tripartition that generated the top box. As a consequence we can always ACTIVATE a tripartition with a non-empty P -class, thereby strictly reducing the cardinality of the U -class tested in DONE.

3.13. EXAMPLE. Let us demonstrate Algorithm 3.12 by computing $S(\bar{R})$ for R the relation of Example 2.11. The elements of the maximal stable sets of $H(\bar{R})$ are ordered pairs pq , with $p \in \{a, b, c, d, e, f\}$ and $q \in \{u, v, w, x, y, z\}$ for which $p\bar{R}q$. For ease of notation we will number these pairs hexadecimally from 1 to F and in the sequel the pairs are denoted by their ordinal number. To show the chosen numbering scheme we recall the matrix of R , where now “~” entries indicate that R holds and a hexadecimal entry gives the ordinal number of a pair in \bar{R} :

	u	v	w	x	y	z
a	~	1	2	~	3	4
b	5	~	~	~	~	~
c	~	~	6	7	~	~
d	~	8	~	9	A	~
e	B	C	~	~	~	~
f	D	~	E	~	F	~

Figure 3.1 shows the search tree constructed by Algorithm 3.12 and Figure 3.2 shows the history of the stack during the computation. In these figures, a tripartition is represented as a sequence, indexed by 1 to F, of symbols where a plus sign, a minus sign and a dot indicate that the corresponding element of \bar{R} is respectively in the P -, N - or U -class of the tripartition. A blank indicates that the corresponding element is not in the partition since it is in the P - or N -class of a “parent” tripartition. Fig. 3.1 may be compared to Fig. 2.1 and consulting Fig. 2.2 one can easily check that maximal stable set $M1$ is the complement of minimal border extension $B1$, etcetera. Fig. 3.1 makes it clear that stable sets and border extensions are constructed simultaneously: they correspond to the cumulative P - and N -classes, respectively. In other words, by modifying OUTPUT-RESULT of 3.11 so that the cumulative N -class is put out instead of the P -class, we get an algorithm for $B(R)$ in which borders are represented as subsets of $A \times D$ instead of sequences over $A \cup D$, as was the case in Algorithm 2.13. In Figure 3.2 we start in the upper left corner with the state of the stack before entrance of the outer loop and we end with the empty stack. Again, the stack grows downwards, “•” is used as the marking symbol and “>” denotes an active tripartition. As in 2.13, ACTIVATE chooses the lowest unmarked tripartition in the top box as the active element.

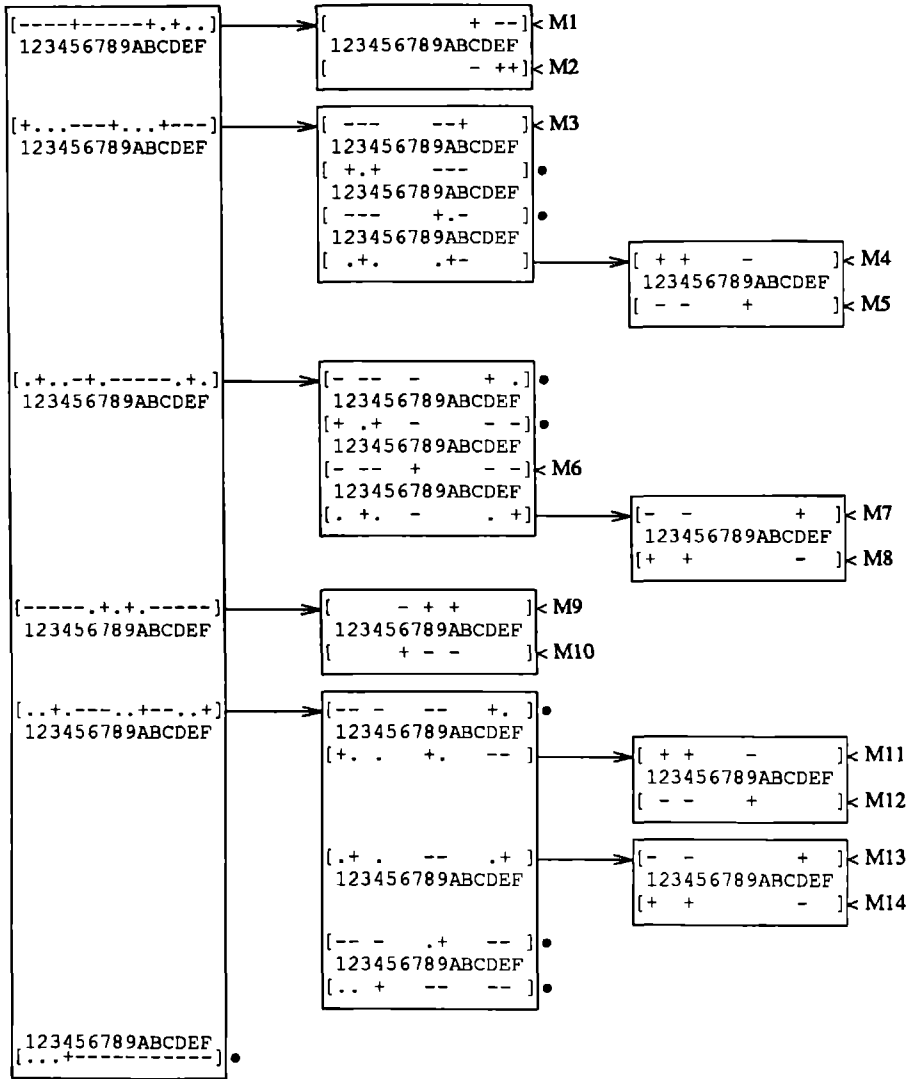


Figure 3.1. Generating $S(\bar{R})$ for Example 2.11. For explanation, see text.

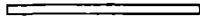
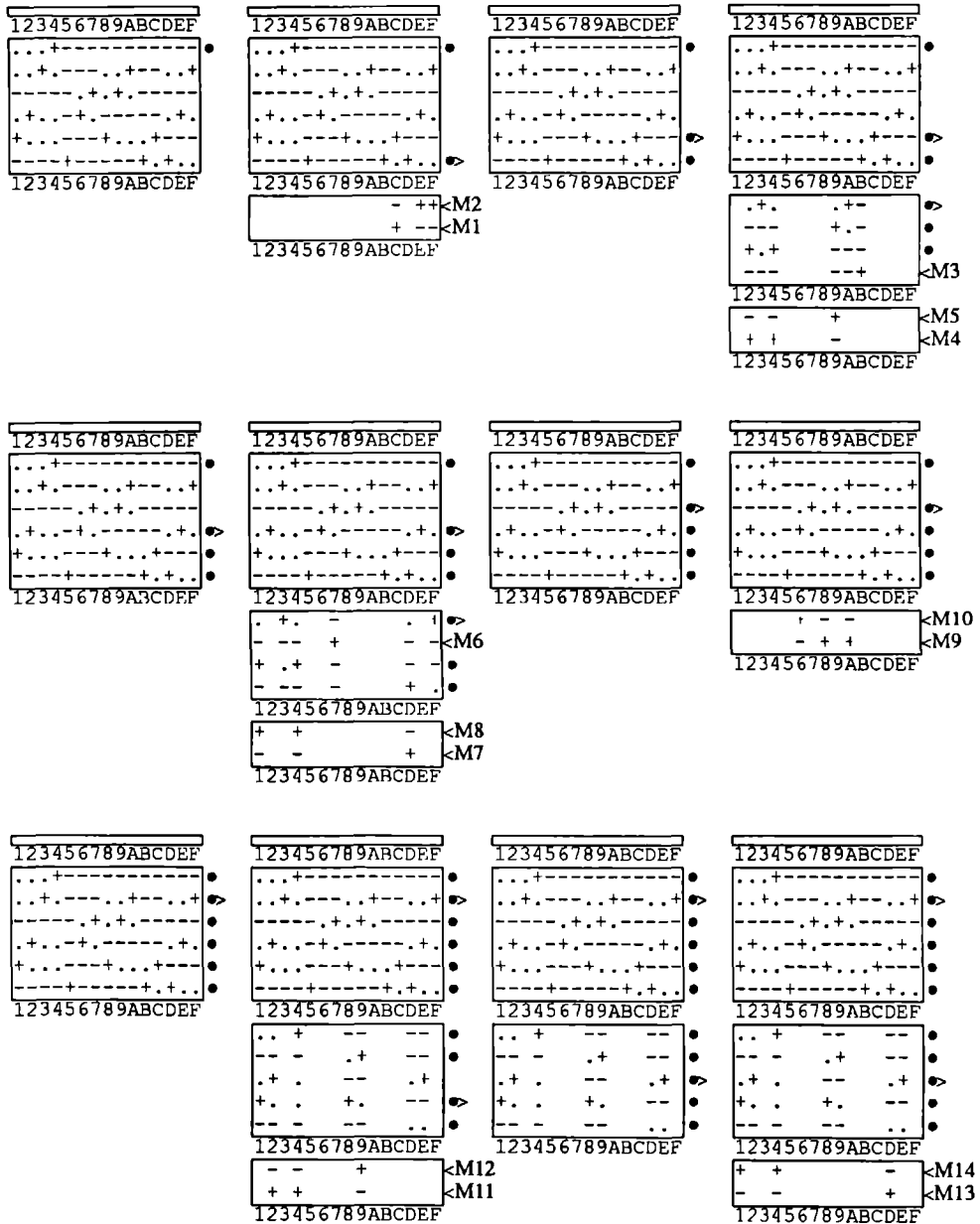


Figure 3.2. Algorithm 3.12 working on Example 2.11. For explanation, see text.

4. MAXIMAL STABLE SETS CONTAINING SOME SPECIFIED SUBSET

Here we show how the algorithm for $S(\bar{R})$ of the preceding section can be adapted to produce only those maximal stable sets that contain some specified subset of vertices. Such a version can be used to obtain representations of R in a (minimal) number of minimal biorder extensions from a (minimal) coloring of $H(\bar{R})$. Indeed, any combination of maximal stable extensions of the different color classes will cover \bar{R} and thus, by Lemma 3.2, the collection of complements of these sets will constitute a representation of R by minimal biorder extensions.

We begin with the special case where we want the subcollection of maximal stable sets containing a particular single vertex. The general case will follow easily from this special case, which is also of interest in itself. As indicated in the Introduction, it appears in the recurrence formula (1.1) that has to be applied in the computation of the bidimension; the next section will be concerned with this aspect. Here, we begin with a proposition telling us that imposing the restriction to produce only maximal stable sets containing some specified vertex does not really change the problem. We use the notation $S(ad; \bar{R})$ for the subcollection of $S(\bar{R})$ consisting of the sets that contain the vertex ad .

4.1. PROPOSITION. *Let $R \subseteq A \times D$ and $ad \in \bar{R}$. Define a second relation $R' \subseteq A \times D$ by $R' = R + \bar{R}[Rd, aR]$. That is, R' is obtained from R by adding all elements of \bar{R} that together with ad constitute a 2-edge of $II(\bar{R})$. (Equivalently, \bar{R}' is obtained from \bar{R} by discarding just these elements.) Then we have: $S(ad; \bar{R}) = S(\bar{R}')$.*

Proof. Suppose $M \in S(ad; \bar{R})$. Since M is stable in $II(\bar{R})$ and $ad \in M$ we must have $M \cap R' = \emptyset$, that is, $M \subseteq \bar{R}'$. Furthermore, M is a biorder (being maximal in $II(\bar{R})$), so it certainly is stable in $H(\bar{R}')$. Finally, if $M \subseteq M_0 \subseteq \bar{R}' \subseteq \bar{R}$ for some $M_0 \in S(\bar{R}')$, then M_0 is a biorder, thus stable in $H(\bar{R})$ and because M is maximal in $H(\bar{R})$ we have $M_0 = M$; we conclude that M is maximal in $H(\bar{R}')$. This proves $S(ad; \bar{R}) \subseteq S(\bar{R}')$. For the reverse inclusion, suppose $M \in S(\bar{R}')$. We have $M \subseteq \bar{R}' \subseteq \bar{R}$ and M is a biorder (being maximal in $II(\bar{R}')$), thus stable in $II(\bar{R})$. To prove that M contains ad it suffices to show that $M \cup \{ad\}$ is stable in $II(\bar{R}')$ (obviously, $ad \in \bar{R}'$). Suppose to the contrary that an edge of $II(\bar{R}')$ is included in $M \cup \{ad\}$: there exists a sequence $ad, a_1d_1, \dots, a_nd_n$ with $a_1d_1, \dots, a_nd_n \in M \subseteq \bar{R}'$ and $ad_n, a_1d, a_2d_1, \dots, a_nd_{n-1} \in R'$. The definition of R' makes it clear that $a_1d, ad_n \notin R' - R$. So $a_1d, ad_n \in R$ and, again by definition of R' , it follows that $a_1d_n \in R'$. But in this case the subsequence a_1d_1, \dots, a_nd_n forms an edge of $H(\bar{R}')$ included in M , which is impossible for $M \in S(\bar{R}')$. So $M \cup \{ad\}$ is stable in $II(\bar{R}')$

and by the maximality of M it must be that $ad \in M$. To show that M is maximal in $H(\bar{R})$ we note that for any M_0 that is stable in $H(\bar{R})$ and for which $M \subseteq M_0 \subseteq \bar{R}$ it follows that $ad \in M_0$. Being stable in $H(\bar{R})$, then, M_0 cannot contain any element of $R' - R$, so we have in fact $M \subseteq M_0 \subseteq \bar{R}' \subseteq \bar{R}$ and $M \in S(\bar{R}')$ implies $M_0 = M$. We conclude that M is maximal in $H(\bar{R})$ and have thus established $S(\bar{R}') \subseteq S(ad; \bar{R})$. ■

4.2. ALGORITHM. By the previous proposition we see that generating $S(ad; \bar{R})$ for some $ad \in \bar{R}$ is the same as generating the collection of maximal stable sets for some other relation. So the only adjustment we have to make, compared to Algorithm 3.6, is assigning this new value to the variable R before calling the procedure MSS. The algorithm for generating $S(ad; \bar{R})$ thus becomes:

```
INITIALIZE (A, D, R, ad);
R := R +  $\bar{R}$  [Rd, aR];
MSS (A, D,  $\emptyset$ )
```

where MSS is as defined in 3.6.

4.3. Of course we can treat Algorithm 3.12 the same way in order to turn it into an algorithm that computes $S(ad; \bar{R})$ for some $ad \in \bar{R}$. There we need only replace the fragment

```
INITIALIZE (A, D, R);
LOAD (R)
```

by

```
INITIALIZE (A, D, R, ad);
R := R +  $\bar{R}$  [Rd, aR];
LOAD (R).
```

Here, however, there exists an alternative solution that does not need a new relation R' to be computed explicitly. We can, instead, compute the collection of tripartitions belonging to this R' directly from the collection belonging to R . (That is, we can appropriately modify the top box of the stack after it has been LOADED according to the original relation R .) To see how this can be done, let $\{(P_i, N_i, U_i)\}$ be the collection of tripartitions belonging to a relation $R \subseteq A \times D$, where for each i there are $A_i \subseteq A$ and $D_i \subseteq D$ such that

$$\begin{aligned} P_i &= \bar{R}[A_i, D_i] = A_i \times D_i, \\ N_i &= \bar{R}[A - A_i, D - D_i] = \bar{R}[A - A_i, D], \\ U_i &= \bar{R}[A_i, D - D_i]. \end{aligned}$$

Now, for some $ad \in \bar{R}$, define $R' = R + \bar{R}[Rd, aR]$. Suppose $ad \in A_0 \times D_0$, where (A_0, D_0) is a pair of sets belonging to the relation R and let (P_0, N_0, U_0) be the corresponding tripartition (that is, $ad \in P_0$). Consider a pair of sets (A_j, D_j) belonging to the relation R . Then there are two exclusive possibilities: (i) $D_j \cap aR = \emptyset$ or (ii) $D_j \subseteq aR$.

In case (i), $(A_j \times D_j) \cap \bar{R}' = (A_j \times D_j) - (Rd \times aR) = A_j \times D_j$, and we see that (A_j, D_j) is also a pair of sets belonging to R' and consequently the tripartition (P_j, N_j, U_j) is a tripartition belonging to the new relation. (Perhaps the pair (A_j, D_j) is, in R' , part of a "bigger" pair (A_h, D_h) where $A_h = A_j$ and $D_h \supseteq D_j$, but this is irrelevant for the argument; it only means that the tripartition (P_j, N_j, U_j) can, in the new relation R' , be merged with another tripartition.)

In case (ii), we see that

$$(A_j \times D_j) \cap \bar{R}' = (A_j \times D_j) - (Rd \times aR) = (A_j - Rd) \times D_j,$$

which means that we have a new pair of sets (A_j', D_j') , belonging to R' , defined by

$$A_j' = A_j - Rd = A_j \cap A_0 \quad \text{and} \quad D_j' = D_j.$$

Consequently there is a new tripartition (P_j', N_j', U_j') corresponding to this pair and

$$\begin{aligned} P_j' &= A_j' \times D_j' = (A_j \cap A_0) \times D_j = (A_j - (A - A_0)) \times D_j \\ &= (A_j \times D_j) - ((A - A_0) \times D) \\ &= P_j - N_0; \\ N_j' &= \bar{R}[A - A_j', D] = \bar{R}[A - (A_j \cap A_0), D] = \bar{R}[(A - A_j) \cup (A - A_0), D] \\ &= \bar{R}[A - A_j, D] \cup \bar{R}[A - A_0, D] \\ &= N_j \cup N_0; \\ U_j' &= \bar{R}[A_j', D - D_j'] = \bar{R}[A_j \cap A_0, D - D_j] = \bar{R}[A_j - (A - A_0), D - D_j] \\ &= \bar{R}[A_j, D - D_j] - \bar{R}[A - A_0, D - D_j] = \bar{R}[A_j, D - D_j] - \bar{R}[A - A_0, D] \\ &= U_j - N_0. \end{aligned}$$

Finally it appears that, for a particular pair of sets (A_j, D_j) , we can decide between the cases (i) and (ii) on the basis of the tripartition (P_j, N_j, U_j) : (ii) is the case iff aRd_j , for $d_j \in D_j$, which by definition is equivalent to $a \in A - A_j$. We conclude that $D_j \subseteq aR$ iff $ad \in \bar{R}[A - A_j, D]$, that is, iff $ad \in N_j$.

4.4. A NON-RECURSIVE ALGORITHM. Based on the discussion in 4.3 we can make following adjustment to Algorithm 3.12 to turn it into an algorithm for producing $S(ad; \bar{R})$ for some $ad \in \bar{R}$. Define the procedure

ADD (ad):

determines the tripartition (P_0, N_0, U_0) for which $ad \in P_0$;
 changes each tripartition (P_j, N_j, U_j) for which $ad \in N_j$ into the tripartition $(P_j - N_0, N_j \cup N_0, U_j - N_0)$

and insert this procedure into Algorithm 3.12 after the LOAD instruction:

```
INITIALIZE (A, D, R, ad);
LOAD (R);
ADD (ad);
MERGE;
MARK;
while not EMPTY
do

od.
```

Except for ADD, the procedures are as defined in 3.11; the code “...” is identical to the body of the outer loop in Algorithm 3.12.

4.5. EXAMPLE. Let us illustrate the working of the algorithm in 4.4 by generating $S(fy; \bar{R})$ for the relation R of Example 2.11. Figure 4.1(a) shows the state of the stack after the LOAD instruction. This is the same as the starting position in Fig. 3.2(a). Now the instruction ADD (fy) is executed. (We use the numbering of the vertices as given in 3.13, so fy corresponds to the hexadecimal number F.) This results in the modified top box of Fig. 4.1(b) and from here on the stack is operated upon in the same way as in 3.13. We see that successively the maximal stable sets $M_2, M_7, M_8, M_{11}, M_{12}, M_{13}$, and M_{14} are produced and in Fig. 3.1 or Fig. 2.2 it can be checked that these are indeed precisely the maximal stable sets that contain the vertex fy .

4.6. Since Proposition 4.1 brings us, in some sense, back to the old situation, it can be applied iteratively, that is, to describe the subcollection of $S(\bar{R})$ consisting of the sets containing the element a_1d_1 , as well as a_2d_2 , etc. We will denote this subcollection by $S(a_1d_1, a_2d_2, \dots; \bar{R})$. Alternatively, for $C \subseteq \bar{R}$, $S(C; \bar{R})$ will denote the subcollection of maximal stable sets of $H(\bar{R})$ that contain the set C . Then we can formulate the following generalization of Proposition 4.1:

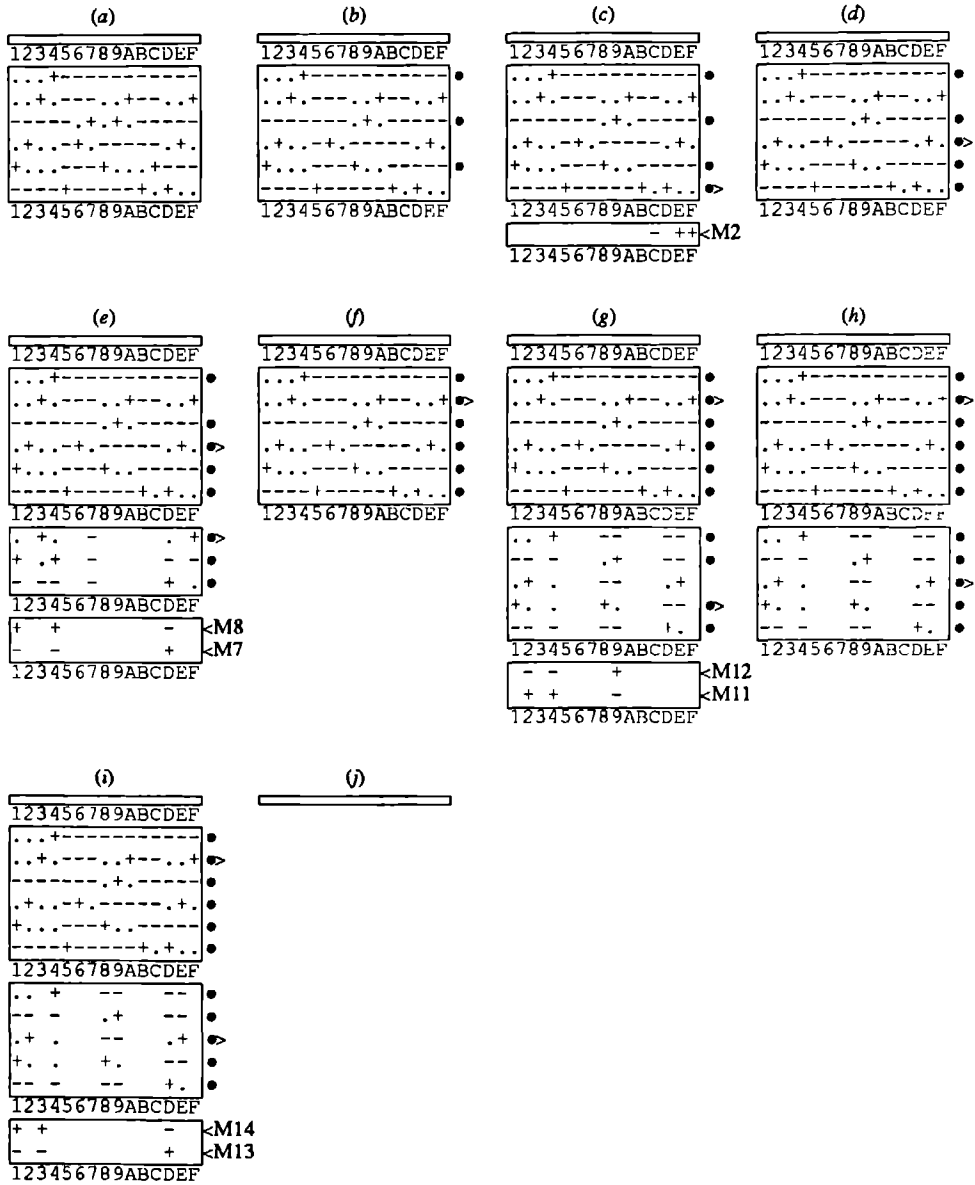


Figure 4.1. Generating $S(fy; \bar{R})$ with Algorithm 4.4. For explanation, see text.

4.7. PROPOSITION. Let $R \subseteq A \times D$ and let $C = \{a_1 d_1, \dots, a_n d_n\} \subseteq \bar{R}$ be stable in $H(\bar{R})$. Define relations $R_i \subseteq A \times D$ by $R_0 = R$ and for $i = 1, \dots, n$,

$$R_i = R_{i-1} + \bar{R}_{i-1}[R_{i-1}d_i, a_i R_{i-1}].$$

Then we have $S(C; \bar{R}) = S(\bar{R}_n)$.

Proof. Suppose k is the least index for which $a_k d_k \in R_{k-1}$. By Proposition 4.1 we have for $i = 1, \dots, k-1$:

$$S(a_i d_i, \dots, a_n d_n; \bar{R}_{i-1}) = S(a_{i+1} d_{i+1}, \dots, a_n d_n; \bar{R}_i).$$

Then $S(C; \bar{R}) = S(a_1 d_1, \dots, a_n d_n; \bar{R}_0) = S(a_k d_k, \dots, a_n d_n; \bar{R}_{k-1}) \subseteq S(a_k d_k; \bar{R}_{k-1}) = \emptyset$, because by definition all elements of $S(a_k d_k; \bar{R}_{k-1})$ are subsets of \bar{R}_{k-1} , so they cannot contain $a_k d_k \in R_{k-1}$. We see that if $a_i d_i \in R_{i-1}$ for some $i = 1, \dots, n$, then there is no maximal stable set of $H(\bar{R})$ that contains $C \subseteq \bar{R}$, which means that C is not stable in $H(\bar{R})$. Since C is stable, we must have $a_i d_i \in \bar{R}_{i-1}$ for all $i = 1, \dots, n$ and thus, by Proposition 4.1, $S(a_i d_i, \dots, a_n d_n; \bar{R}_{i-1}) = S(a_{i+1} d_{i+1}, \dots, a_n d_n; \bar{R}_i)$ for all $i = 1, \dots, n$. ■

4.8. REMARK. The proof of Proposition 4.7 describes a procedure for deciding whether a (finite) subset $C = \{a_1 d_1, \dots, a_n d_n\}$ is stable. If, in constructing the relations R_1, \dots, R_n successively, $C \cap R_i \neq \emptyset$ for some i , then C is not stable and if $C \cap R_n = \emptyset$, then C is stable in $H(\bar{R})$.

4.9 REMARK. For the relation R_n in Proposition 4.7 we will in general *not* have $R_n = R + \cup_i \bar{R}[Rd_i, a_i R]$. This ‘‘straightforward’’ generalization of Proposition 4.1 will not work. The left side does include the right side, but in general the inclusion will be strict. That is, we do *not* get \bar{R}_n by discarding from \bar{R} just all vertices that constitute a 2-edge with any of the elements in C .

4.10. ALGORITHM. From Proposition 4.7 we directly obtain a generalization of Algorithm 4.2 that produces the subset of $S(\bar{R})$ in which each element includes some stable set C :

```
INITIALIZE (A, D, R, C = {a_1 d_1, ..., a_n d_n});
for i := 1 to n do R := R + \bar{R}[Rd_i, a_i R] od;
MSS (A, D, \emptyset)
```

where MSS is as defined in 3.6.

4.11. A NON-RECURSIVE ALGORITHM. Again we may, comparing Algorithm 4.10 to Algorithm 3.6, deduce that Algorithm 3.12 may be turned into an algorithm for $S(C; \bar{R})$ by replacing the fragment

```
INITIALIZE (A, D, R);
LOAD (R)
```

in 3.12 by

```
INITIALIZE (A, D, R, C = {a1d1, . . . , andn});
for i := 1 to n do R := R +  $\bar{R}$ [Rdi, a, R] od;
LOAD (R).
```

Alternatively we may take the approach described in 4.3 and then it is clear that we turn Algorithm 4.4 into an algorithm for $S(C; \bar{R})$ by repeating the ADD call for each element of C . To be explicit, the following is an algorithm for $S(C; \bar{R})$, generalizing Algorithm 4.4:

```
INITIALIZE (A, D, R, C = {a1d1, . . . , andn});
LOAD (R);
for i := 1 to n
do
    ADD (a, di);
    MERGE
od;
MARK;
while not EMPTY
do
    . . .
od.
```

ADD is defined in 4.4, the other procedures in 3.11; the code “. . .” is identical to the body of the outer loop in Algorithm 3.12.

Note that the repeated ADDing works correctly only if C is stable. If C is not stable, then for some $k \leq n$, ADD ($a_k d_k$) will fail to find the tripartition whose P -class contains $a_k d_k$, simply because there is no longer any such tripartition: by a previous ADD operation the element $a_k d_k$ has been moved from the P -class to the N -class of “its” tripartition. In this way, the procedure for deciding whether $C \subseteq \bar{R}$ is stable that was described in Remark 4.8 may be translated in operations on tripartitions.

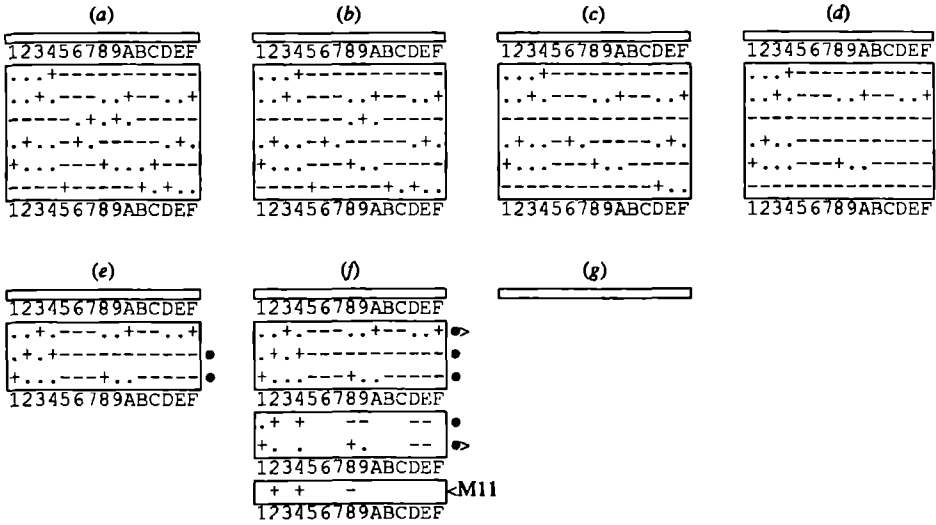


Figure 4.2. Generating $S(fy, aw, dv; \bar{R})$ with Algorithm 4.11. For explanation, see text.

4.12. EXAMPLE. In Figure 4.2, Algorithm 4.11 is illustrated, computing $S(C; \bar{R})$ where R is the relation of Example 2.11 and $C = \{fy, aw, dv\}$. We use again the numbering of 3.13 for the vertices, so the elements of C correspond respectively to the hexadecimal numbers F, 2 and 8 in Fig. 4.2. In Fig. 4.2(a) the situation after LOAD is given and in (b) the vertex fy (number F) is ADDED (compare with the situation in Fig. 4.1(a) and (b)). Now in Fig. 4.2(c) and (d) the vertices aw (number 2) and dv (number 8) are ADDED respectively. For reasons of clarity we have chosen in Fig. 4.2 not to MERGE after each ADD, as in Algorithm 4.11. Instead, after all the ADDing has been done the top box is MERGED and MARKED; the result is given in 4.2(e). In this state the main loop is entered and in 4.2(f) and (g) we see that in this case $S(C; \bar{R})$ consists of one maximal stable set only. In Fig. 3.1 or Fig. 2.2 it can be checked that $M11$ is indeed the only maximal stable set of R that contains the vertex fy , as well as aw , as well as dv .

4.13. REMARK. If we had for the subset of $S(\bar{R})$ consisting of the elements *not* containing some specified vertex ad of $II(\bar{R})$ a similar result as in Proposition 4.1 – that is, if, denoting this subset by $S(\neg ad; \bar{R})$, we would have $S(\neg ad; \bar{R}) = S(\bar{R}^*)$, for some relation R^* –, then we would have a very simple recursive procedure for generating $S(\bar{R})$ (and thus $S(C; \bar{R})$):

$$\begin{aligned}
 S(\bar{R}) &= S(ad; \bar{R}) + S(\neg ad; \bar{R}) \quad (ad \in \bar{R} \text{ arbitrary}) \\
 &= S(R') + S(R'')
 \end{aligned}$$

for appropriate $R', R'' \subseteq A \times D$.

But here the situation is not so simple: for instance, it is generally not the case that $S(\neg ad; \bar{R}) = S(\bar{R} - \{ad\})$. We do have inclusion from left to right, but, conversely, an element V of $S(\bar{R} - \{ad\})$ is not maximal in \bar{R} unless the set $V + \{ad\}$ contains a 2-edge of $H(\bar{R})$ (involving ad) and this is not necessarily so. This makes iterative application rather awkward: $S(\neg a_1d_1, \neg a_2d_2, \dots; \bar{R})$ is the subset of $S(\bar{R} - \{a_1d_1, a_2d_2, \dots\})$ consisting of the elements that contain a vertex that is in a 2-edge with a_1d_1 , one that is in a 2-edge with a_2d_2 , etc. (these vertices may or may not be different from each other and there may or there may not exist such sets in $S(\bar{R})$). Such a description does not seem very promising for algorithmic purposes.

5. MAXIMAL STABLE SETS OF A SUBHYPERGRAPH

In Koppen (1987) a procedure is described for computing $\text{Chrom } H(\bar{R})$ (and thus finding the bidimension of a relation R). It consists of alternately reducing a subhypergraph of $H(\bar{R})$ to a smaller subhypergraph with the same chromatic number and next applying the recurrence relation (1.1) for the chromatic number of a hypergraph to the reduced subhypergraph. In Koppen (1987) the reduction procedure is proved correct and described in detail, but there is no explicit procedure given how to evaluate the right hand side of (1.1) and this is the problem we will consider here. Since we are now interested in maximal stable sets of a subhypergraph $H(V)$ rather than $H(\bar{R})$, we first generalize some previously introduced notions.

5.1. DEFINITION. For $R \subseteq A \times D$ and $V \subseteq \bar{R}$, let $S(V)$ denote the collection of maximal stable sets of the subhypergraph $H(V)$ and let $S(\bar{R} | V)$ be the collection of stable sets of $H(V)$ that are restrictions to V of maximal stable sets of $H(\bar{R})$. So $S(\bar{R} | V) = \{M \cap V : M \in S(\bar{R})\}$.

For $ad \in V$, let $S(ad; V)$ and $S(ad; \bar{R} | V)$ denote the subcollections of $S(V)$ and $S(\bar{R} | V)$, respectively, consisting of the sets containing ad .

Since any set in $S(V)$ is stable in $H(\bar{R})$ it can be extended to a maximal stable set of $H(\bar{R})$ and by the maximality in $H(V)$ this is obtained by adding elements of $\bar{R} - V$ only. So $S(V) \subseteq S(\bar{R} | V)$. On the other hand, each set in $S(\bar{R} | V)$ is stable in $H(V)$ and thus we see that $S(V)$ is the collection of maximal elements of $S(\bar{R} | V)$. In the same way $S(ad; V)$ is the collection of maximal elements of $S(ad; \bar{R} | V)$.

With the help of these definitions, we can write the recurrence relation (1.1) for $H(V)$, given in the Introduction, as

$$\text{Chrom } H(V) = 1 + \min \{ \text{Chrom } H(V-M) : M \in S(ad; V) \}, \quad (5.1)$$

where ad is an arbitrary element of V .

Inspired by Proposition 4.1 we may think that the problems of generating $S(ad; V)$ and $S(ad; \bar{R} | V)$ can be reduced to the problems of generating $S(V)$ and $S(\bar{R} | V)$, respectively. This is indeed the case:

5.2. PROPOSITION. *Let $R \subseteq A \times D$, $V \subseteq \bar{R}$ and $ad \in V$. Define a second relation $R' \subseteq A \times D$ by $R' = R + \bar{R}[Rd, aR]$ and define $V' = V \cap \bar{R}'$. Then we have:*

(i) $S(ad; \bar{R} | V) = S(\bar{R}' | V')$ and (ii) $S(ad; V) = S(V')$.

Proof. (i) Using definitions and Proposition 4.1 we can write $S(ad; \bar{R} | V) = \{M \cap V : M \in S(ad; \bar{R})\} = \{M \cap V : M \in S(\bar{R}')\} = \{(M \cap \bar{R}') \cap V : M \in S(\bar{R}')\} = \{M \cap V' : M \in S(\bar{R}')\} = S(\bar{R}' | V')$. (ii) Follows directly from (i), since $S(ad; V)$ and $S(V')$ are the collections of maximal elements of $S(ad; \bar{R} | V)$ and $S(\bar{R}' | V')$, respectively. ■

5.3. Let us first show that by the results of the previous sections we already have some solution for the problem of evaluating (5.1). If $V = \bar{R}$ there would be no problem, since in Section 4 we have seen how to produce $S(ad; \bar{R})$. This gives one way of evaluating (5.1): replacing V by \bar{R} in (5.1), we get

$$\text{Chrom } H(\bar{R}) = 1 + \min \{ \text{Chrom } H(\bar{R}-M) : M \in S(ad; \bar{R}) \} \quad (5.2)$$

and since in the procedure of Koppen (1987) $H(V)$ is such that

$$\text{Chrom } H(\bar{R}) = \text{Chrom } H(V),$$

we can evaluate (5.2) instead of (5.1). In fact, we can replace \bar{R} by \bar{R}^* in (5.2), where R^* is the smallest restriction of R such that $V \subseteq \bar{R}^*$. Since the subhypergraph $H(\bar{R}^*)$ of $H(\bar{R})$ is the full hypergraph associated with the relation R^* we can find $S(ad; \bar{R}^*)$ by the methods of Section 4.

5.4. Although using \bar{R}^* instead of V does solve the problem, it is not quite satisfactory a solution, since we do not fully use the reduction to $H(V)$. Without loss of generality we will in the sequel assume that $R^* = R$; this simplifies notation. Now it is clear that $S(ad; \bar{R})$ contains more (not fewer) elements than $S(ad; V)$: any maximal stable set of $H(V)$ can be extended to a maximal stable set of $H(\bar{R})$ and in this way different elements of $S(ad; V)$ give rise to different elements of $S(ad; \bar{R})$ (only elements of $\bar{R}-V$ are added to obtain the extensions). So, solving

(5.2) instead of (5.1) obliges us to generate more maximal stable sets and since each maximal stable set M of $H(\bar{R})$ induces a recursive "call" of (5.2) for the (reduced subhypergraph of the) hypergraph $H(\bar{R}-M)$, we introduce in doing so an extra number of branches at each level of the recursion.

This problem can be remedied by computing $S(ad; V)$ from $S(ad; \bar{R})$. As we have seen above, $S(ad; V)$ is a subset of $S(ad; \bar{R} | V)$. On the other hand, the restriction to V of a maximal stable set of $H(\bar{R})$ need not be maximal in $H(V)$ and different maximal stable sets of $H(\bar{R})$ may, when restricted to V , yield the same stable set of $H(V)$. So an algorithm for $S(ad; \bar{R})$ can be changed into an algorithm for $S(ad; V)$ in the following way: maintain a list of restrictions to V of generated elements of $S(ad; \bar{R})$ and when a new element of $S(ad; \bar{R})$ is produced, add its restriction to V to the list, provided that this is not a subset of an element already in the list (this includes the case of duplicates), and remove from the list any elements that are subsets of the newly inserted element. Given the resulting list of $S(ad; V)$ we can apply (5.1) instead of (5.2), thereby reducing the number of recursive calls.

5.5. In the above procedure we still need to compute all of $S(ad; \bar{R})$ while we are only interested in subsets of V . It would be nice if the generating process, from the start, was restricted to subsets of V . Since $S(ad; V)$ consists precisely of the maximal elements of $S(ad; \bar{R} | V)$, it is clear that

$$\begin{aligned} \min \{Chrom H(V-M) : M \in S(ad; V)\} = \\ \min \{Chrom H(V-M') : M' \in S(ad; \bar{R} | V)\} \end{aligned} \quad (5.3)$$

and, consequently, generating $S(ad; \bar{R} | V)$ is an alternative solution to our problem.

The point is that Algorithm 4.4 produces in fact $S(ad; \bar{R} | V)$ rather than $S(ad; \bar{R})$, provided it is supplied with the appropriate "input". The only provision we have to take is that in the LOADING of the stack (with LOAD as defined in 3.11) we do not push the tripartitions belonging to the relation R on the stack, but instead the tripartitions of V that we get by intersecting each class of each tripartition with V . That is, the collection $\{(P_i, N_i, U_i)\}$ is changed to the collection $\{(P'_i, N'_i, U'_i)\}$, where

$$P'_i = P_i \cap V, \quad N'_i = N_i \cap V, \quad U'_i = U_i \cap V.$$

With this modification the algorithm can run as before; to establish that it correctly produces $S(ad; \bar{R} | V)$ we need to show that the MERGE and MARK operations, when controlled by the new tripartitions, do not suppress any element of $S(ad; \bar{R} | V)$. This is guaranteed by the following lemma.

5.6. LEMMA. *Let $R \subseteq A \times D$, $V \subseteq \bar{R}$, $d_0, d_1 \in D$. Then $V[Rd_0, D] \supseteq V[Rd_1, D]$ iff, for any $M \subseteq \bar{R}$, $M \in \mathcal{S}(\bar{R}[\bar{R}d_0, D])$ implies $M \cap V \subseteq \bar{R}[\bar{R}d_1, D]$.*

Proof. (Only if.) Suppose $V[Rd_0, D] \supseteq V[Rd_1, D]$ or, equivalently, $V[\bar{R}d_0, D] \subseteq V[\bar{R}d_1, D]$. If $M \subseteq \bar{R}[\bar{R}d_0, D]$, then $M \cap V \subseteq \bar{R}[\bar{R}d_0, D] \cap V = V[\bar{R}d_0, D] \subseteq V[\bar{R}d_1, D] \subseteq \bar{R}[\bar{R}d_1, D]$. (If.) Suppose $ad \in V[\bar{R}d_0, D] - V[\bar{R}d_1, D]$. This means $ad \in V$, $a\bar{R}d_0$ and aRd_1 . Put M any element of $\mathcal{S}(\bar{R}[\bar{R}d_0, D])$ that contains ad . Then $ad \in M \cap V$ while aRd_1 , so $M \cap V \not\subseteq \bar{R}[\bar{R}d_1, D]$. ■

Interpreting this lemma in the context of Algorithm 3.12, that is, letting R stand for some restriction $R[\alpha, \delta]$ implicitly computed there, we note that $V[Rd_0, D] = \bar{R}[Rd_0, D] \cap V = N_0 \cap V$, where N_0 is the N -class of the tripartition of \bar{R} corresponding to d_0 ; likewise, $V[Rd_1, D]$ is the restriction to V of the N -class of the tripartition corresponding to d_1 . On the other hand, $M \in \mathcal{S}(\bar{R}[\bar{R}d_0, D])$ means that M is an element of $\mathcal{S}(\bar{R})$ that, according to Algorithm 3.12, is generated when the tripartition corresponding to d_0 is ACTIVATED. So in terms of Algorithm 3.12, the only-if-part of Lemma 5.6 conveys that if $N_0 \supseteq N_1$, then for any element of $\mathcal{S}(\bar{R})$ generated by choosing the tripartition corresponding to d_0 , there is an element of $\mathcal{S}(\bar{R})$ generated by choosing the tripartition corresponding to d_1 with the same restriction to V . As we have seen in the previous section, the step from Algorithm 3.12 to Algorithm 4.4 is trivial. We see that, with all classes restricted to V , tripartitions with non-minimal N -classes do not generate “new” elements of $\mathcal{S}(ad; \bar{R} | V)$ (i.e., MARK still works well) and in particular, tripartitions with identical N -classes generate the same subcollection of $\mathcal{S}(ad; \bar{R} | V)$ and we do not lose any element of this subcollection by choosing these tripartitions one after the other (i.e., MERGE is still okay). The if-part of Lemma 5.6 implies that we risk missing some elements of $\mathcal{S}(ad; \bar{R} | V)$ if not all tripartitions with minimal N -classes are ACTIVATED.

Note that, while Lemma 5.6 assures that the modification of Algorithm 4.4 as described in 5.5 produces each element of $\mathcal{S}(ad; \bar{R} | V)$, it does not guarantee that each element is generated just once. In fact, duplicates may occur: in the “restricted version” of Algorithm 4.4 it is possible that one and the same element is produced by two tripartitions with incomparable N -classes. In the original version for $\mathcal{S}(ad; \bar{R})$ this is impossible (cf. Proposition 3.5).

5.7. As we have seen in 5.5 and 5.6, a minor modification of Algorithm 4.4 allows us to work with $\mathcal{S}(ad; \bar{R} | V)$ rather than $\mathcal{S}(ad; \bar{R})$. At this point, we have the same alternatives as sketched in 5.3 and 5.4. The first possibility is to work with $\mathcal{S}(ad; \bar{R} | V)$ as generated by the “restricted version” of Algorithm 4.4 (that is, with

possible duplicates); this means, we compute the right hand side of (5.3). The second possibility is to select the maximal elements of $S(ad; \bar{R} | V)$, giving us $S(ad; V)$ by which the left hand side of (5.3) can be evaluated. Which of these alternatives will be more efficient may be very implementation and data dependent. The latter alternative saves us a number of recursive “calls” at the cost of comparing the elements of $S(ad; \bar{R} | V)$ in order to select the maximal elements. If V is “almost like” \bar{R} , we will do a lot of comparing with little or no reduction, while if V is “considerably smaller” than \bar{R} , the reduction obtained may be substantial.

5.8. The two alternatives sketched above evoke a third possibility, an approach that in fact seems to be the obvious one, when one is confronted with the problem of evaluating (5.1): find an adaptation of Algorithm 4.4 that directly produces $S(ad; V)$ rather than $S(ad; \bar{R})$. This is indeed possible, though not by a “minor modification”. The problem is that in generating maximal stable sets of a subhypergraph $H(V)$ the connection with generating minimal biorder extensions is lost: a maximal stable set of $H(V)$ is, generally, not a biorder. It is exactly this connection (Lemma 3.2) on which the recursive Algorithm 3.6 is based that in turn is at the origin of Algorithms 3.12 and 4.4. In the next section we will develop an algorithm for directly producing $S(V)$ for a subset V of \bar{R} . By Proposition 5.2, this is equivalent to an algorithm for $S(ad; V)$, $ad \in V$. The algorithm for $S(V)$ will not use biorder properties, but is solely based on the definition of the hypergraph $H(\bar{R})$ and its subhypergraphs.

6. AN ALTERNATIVE APPROACH : GENERAL TRIPARTITIONS

In this section we take an alternative approach to the problem of generating $S(V)$, with $V \subseteq \bar{R}$ for some relation R . We start with something that is apparently unrelated to this issue, namely the notion of general tripartitions on a (finite) set. We single out one of the three classes, we consider the problem of finding in certain collections of such tripartitions the elements for which this class is maximal, and we formulate an algorithm that solves this problem. Then we come back to our hypergraph problem, in two steps. First, we consider a relation Q on a set E and show that the problem of finding the maximal subsets S of E such that the restriction of Q to S is acyclic can be seen as an instance of the above tripartition problem. Next, our original problem of finding the maximal stable sets in a subhypergraph $H(V)$ is interpreted as an instance of this maximal acyclic restriction problem. In

sum, the algorithm for the tripartition problem can be applied to find $S(V)$ for $V \subseteq \bar{R}$. We finally compare the ensuing algorithm with the ones obtained in the preceding sections. One distinction is that in the algorithm of this section the notion of biorder does not appear.

6.1. GENERAL TRIPARTITIONS. Consider an arbitrary finite set E and let T_1, T_2 be two tripartitions of E . We denote the classes of a tripartition by capitals P , N and U , respectively (why not?). So,

$$T_i = (P(T_i), N(T_i), U(T_i)),$$

where $P(T_i) + N(T_i) + U(T_i) = E$; often we will write more briefly $T_i = (P_i, N_i, U_i)$. For two tripartitions T_1, T_2 of E we define $T_1 \circ T_2$, the *composition* of T_1 with T_2 , as

$$T_1 \circ T_2 = (P_1 + (P_2 \cap U_1), N_1 + (N_2 \cap U_1), U_2 \cap U_1).$$

It is easily checked that $T_1 \circ T_2$ is again a tripartition of E and that the so defined composition operation \circ is associative. Let τ be a collection of tripartitions of E that is closed under this composition operation. Then we define for a tripartition $T_0 \in \tau$ the *offspring* of T_0 in τ , denoted by $T_0 \circ \tau$, as

$$T_0 \circ \tau = \{T_0 \circ T : T \in \tau\}.$$

A P -class of a tripartition in τ is called *maximal* in τ iff it is not strictly included in the P -class of another element of τ . In the following lemma some properties are collected that anticipate the problem of finding maximal P -classes in τ .

6.2. LEMMA. Let τ be a collection of tripartitions of a set E , closed under the composition operation \circ .

- (i) For any $T'_0 \in T_0 \circ \tau$ we have $P'_0 \supseteq P_0, N'_0 \supseteq N_0, U'_0 \subseteq U_0$.
- (ii) For any $T_0 \in \tau$, $T_0 \in T_0 \circ \tau$ and if $U_0 = \emptyset$, then $T_0 \circ \tau = \{T_0\}$.
- (iii) For any $T'_0 \in T_0 \circ \tau$, $T'_0 \circ T_0 = T'_0$.
- (iv) If $N_1 \subseteq N_2$ for $T_1, T_2 \in \tau$, then for any $T'_2 \in T_2 \circ \tau$ there is $T'_1 \in T_1 \circ \tau$ such that $P'_2 \subseteq P'_1$.
- (v) If $N_1 = N_2$ for $T_1, T_2 \in \tau$, then $T_1 \circ T_2 = T_2 \circ T_1$ and for any $T \in T_1 \circ \tau$ (or $T \in T_2 \circ \tau$) there is $T^* \in T_1 \circ T_2 \circ \tau$ such that $P \subseteq P^*$.
- (vi) If $P_2 \cap U_1 = \emptyset$ for $T_1, T_2 \in \tau$, then $N(T_1 \circ T_2 \circ T_0) \supseteq N(T_1 \circ T_0)$ for any $T_0 \in \tau$.
- (vii) If $N_2 \cap P_1 \neq \emptyset$ for $T_1, T_2 \in \tau$, then there is no $T'_2 \in T_2 \circ \tau$ such that $P'_2 \supseteq P_1$.
- (viii) If $U_1 = \emptyset$ and $N_1 \not\subseteq N_2$ for $T_1, T_2 \in \tau$, then there is no $T'_2 \in T_2 \circ \tau$ such that $P'_2 \supseteq P_1$.

- (ix) If τ is such that $\cup \{P(T) : T \in \tau\} = E$, then for any $T_0 \in \tau$ and any $s \in U_0$ there is $T'_0 \in T_0 \circ \tau$ such that $s \in P'_0$.
- (x) If τ is such that $\cup \{P(T) : T \in \tau\} = E$, then P_0 is maximal in τ only if $U_0 = \emptyset$.

Proof. (i) to (iii): Immediate from the definitions.

(iv): If $N_1 \subseteq N_2$, then $N(T_1 \circ T_2) = N_1 + (N_2 \cap U_1) = (N_1 \cup N_2) \cap (N_1 \cup U_1) = N_2 - P_1$. Suppose $T'_2 = T_2 \circ T_0$ for some $T_0 \in \tau$; then, for $T'_1 = T_1 \circ T_2 \circ T_0 \in T_1 \circ \tau$, we obtain: $N(T'_1) = N(T_1 \circ T_2) + (N_0 \cap U(T_1 \circ T_2)) = (N_2 - P_1) + (N_0 \cap U_1 \cap U_2) \subseteq N_2 + (N_0 \cap U_2) = N(T'_2)$ and $U(T'_1) = U_1 \cap U_2 \cap U_0 \subseteq U_2 \cap U_0 = U(T'_2)$. Since both T'_1 and T'_2 are tripartitions of E , it follows that $P(T'_1) \supseteq P(T'_2)$.

(v): $T_1 \circ T_2 = T_2 \circ T_1$ because $U(T_1 \circ T_2) = U_1 \cap U_2 = U(T_2 \circ T_1)$ and, by the proof of (iv), $N(T_1 \circ T_2) = N_2 - P_1 = N_1 - P_1 = N_1 = N_2 = N_2 - P_2 = N_1 - P_2 = N(T_2 \circ T_1)$. Let $T = T_1 \circ T_0$ for some $T_0 \in \tau$; then for $T^* = T_1 \circ T_2 \circ T_0$ we have: $N(T^*) = N(T_1 \circ T_2) + (N_0 \cap U(T_1 \circ T_2)) = N_1 + (N_0 \cap U_1 \cap U_2) \subseteq N_1 + (N_0 \cap U_1) = N(T)$ and $U(T^*) = U_1 \cap U_2 \cap U_0 = U(T)$, so $P(T^*) \supseteq P(T)$.

(vi): Clearly, $U(T_1 \circ T_2 \circ T_0) \subseteq U(T_1 \circ T_0)$ and $P(T_1 \circ T_2) = P_1$, if $P_2 \cap U_1 = \emptyset$. But then $P(T_1 \circ T_2 \circ T_0) = P(T_1 \circ T_2) + P_0 \cap U(T_1 \circ T_2) = P_1 + (P_0 \cap U_1 \cap U_2) \subseteq P_1 + (P_0 \cap U_1) = P(T_1 \circ T_0)$.

(vii): For any $T'_2 \in T_2 \circ \tau$, $N'_2 \cap P_1 \supseteq N_2 \cap P_1 \neq \emptyset$, so $P'_2 \not\subseteq P_1$.

(viii): $U_1 = \emptyset$ and $N_1 \not\subseteq N_2$ imply $N_2 \cap P_1 \neq \emptyset$, so apply (vii).

(ix): Take $T'_0 = T_0 \circ T_1$, where T_1 is such that $s \in P_1$.

(x): Immediate from (ix). ■

6.3. FINDING MAXIMAL P -CLASSES. Now we consider the problem of finding the maximal P -classes in τ , where τ is a collection of tripartitions of E that is closed under composition. We denote by MAX_τ the collection of P -classes that are maximal in τ ; for $C \subseteq \tau$ we define

$$MAX_\tau(C) = MAX_\tau \cap \{P(T') : T' \in T \circ \tau, T \in C\},$$

which means that $MAX_\tau(C)$ is the subcollection of MAX_τ that is in the offspring of some tripartition in C .

We will describe a procedure for finding MAX_τ , starting from what may be called a *base* for τ , that is, a subset $\tau_0 \subseteq \tau$ such that τ is the closure of τ_0 under composition. In other words, a tripartition T is in τ iff there is a number k and there are tripartitions $T_1, \dots, T_k \in \tau_0$ such that

$$T = T_1 \circ \dots \circ T_k.$$

The minimum number k for which $T \in \tau$ can be "factorized" in this way into elements of τ_0 will be called the *degree* of T (over τ_0). Since $\cup \{T \circ \tau : T \in \tau_0\} = \tau$

(any element of τ is in the offspring of some element of τ_0), it is clear that

$$MAX_{\tau}(\tau_0) = MAX_{\tau}.$$

We will, moreover, assume that $\cup \{P(T) : T \in \tau\} = E$, which is equivalent to $\cup \{P(T) : T \in \tau_0\} = E$. This is no real restriction in the context of finding maximal P -classes in τ : if $\cup \{P(T) : T \in \tau\} = E' \subseteq E$, we consider the collection τ' of tripartitions of E' that is obtained by taking the restriction to E' of each class of the tripartitions in τ . Since the P -classes remain the same, the collections of maximal P -classes in τ and τ' clearly are identical.

In the next proposition we state some operations that may be performed on a subset C of τ without changing $MAX_{\tau}(C)$; these operations will be used in an algorithm for finding MAX_{τ} .

6.4. PROPOSITION. *With τ and τ_0 as in 6.3, let C, C' be subsets of τ . Then, in all of the cases (i) to (iii), we have $MAX_{\tau}(C) = MAX_{\tau}(C')$:*

- (i) C' is obtained by forming in C classes of tripartitions with identical N -classes and replacing each class by the composition of its elements.
- (ii) C' is the subset of C consisting of the tripartitions with minimal N -classes (N -classes not strictly including the N -class of another element in C).
- (iii) $C' = C_0 + \{T \circ t : T \in C_1, t \in \tau_0, U(T) \cap P(t) \neq \emptyset\}$, where $C_0, C_1 \subseteq C$ with $C_0 = \{T \in C : U(T) = \emptyset\}$ and $C_1 = C - C_0$.

Proof. (i) Lemma 6.2.(v) tells us that the composition is independent of the order of its constituents (so C' is well defined) and that by this operation we are not throwing away any maximal P -classes in the offspring of elements of C . (ii) Lemma 6.2.(iv) shows that all elements in $MAX_{\tau}(C)$ are also in $MAX_{\tau}(C')$. (iii) By Lemma 6.2.(ii), $C_0 = \{T \circ t : T \in C_0, t \in \tau_0\}$, so $C' = C^* - \{T \circ t : T \in C_1, t \in \tau_0, U(T) \cap P(t) = \emptyset\}$, where $C^* = \{T \circ t : T \in C, t \in \tau_0\}$. For any $T \in \tau$ we clearly have $T \circ \tau = \cup \{T \circ t \circ \tau : t \in \tau_0\}$. Consequently, $\cup \{T \circ \tau : T \in C\} = \cup \{T \circ t \circ \tau : T \in C, t \in \tau_0\}$, which implies $MAX_{\tau}(C) = MAX_{\tau}(C^*)$. By Lemma 6.2.(ix) there is for any $T_0 \in C_1$ some $t_0 \in \tau_0$ for which $U(T_0) \cap P(t_0) \neq \emptyset$. This $T_0 \circ t_0$ is "saved" in C' and Lemmas 6.2.(vi) and (iii) together imply that any maximal P -class found in the offspring of $T_0 \circ t$ where $U(T_0) \cap P(t) = \emptyset$, is also found in the offspring of $T_0 \circ t_0 \in C'$. It follows that $MAX_{\tau}(C^*) = MAX_{\tau}(C')$. ■

6.5. ALGORITHM. In view of cases (i) and (ii) of Proposition 6.4 we define a procedure

COMPRESS ($C \subseteq \tau$):

- replaces in C classes of tripartitions with identical N -classes by the composition of their elements;
- discards tripartitions with non-minimal N -classes.

Note that by applying the transformation of 6.4(i) we get a collection where all tripartitions have different N -classes; by combining this with the transformation of 6.4(ii), as is done in COMPRESS, the end result is a collection where all N -classes are mutually incomparable. Using the above procedure COMPRESS, we have following algorithm for computing MAX_τ :

```

M := ∅; C1 := τ0;
i := 1;
while Ci ≠ ∅
do
  COMPRESS (Ci);
  Ci+1 := ∅;
  for T ∈ Ci
  do
    if U(T) = ∅ then M += P(T)
    else
      for t ∈ τ0
      do
        if U(T) ∩ P(t) ≠ ∅ then Ci+1 += T ∘ t
      od
    fi
  od;
  i += 1
od.

```

This algorithm terminates because $|E|$ is an upper bound for the cardinality of a U -class in C_1 and, by construction, any U -class in C_{i+1} is strictly smaller than some U -class in C_i . Consequently, for some $n \leq |E| + 1$ all U -classes in C_n will be empty, which implies that for $i = n+1$ the termination condition of the main loop will be fulfilled.

Before entering the main loop we have

$$M + MAX_\tau(C_i) = \emptyset + MAX_\tau(C_1) = MAX_\tau(\tau_0) = MAX_\tau$$

Within the main loop $M + MAX_\tau(C_i)$ is not changed by COMPRESSING C_i (cases (i) and (ii) of Proposition 6.4). We have seen above that after COMPRESSING all

N -classes in C_i are incomparable; by Lemma 6.2(viii), no P -class of a tripartition in C_i with empty U -class will be included in any P -class produced in the offspring of the tripartitions in C_{i+1} . In other words: at the moment a P -class is added to M , it is known to be an element of MAX_τ . Using 6.4(iii), we see that $M + MAX_\tau(C_i)$ before entrance of the second level loop equals $M + MAX_\tau(C_{i+1})$ at termination of this loop. (In particular, the two constituent collections are indeed still disjoint.) In sum, $M + MAX_\tau(C_i) = MAX_\tau$ is a main loop invariant that will hold at termination; combining with the termination condition we obtain $M = MAX_\tau$. We conclude that on termination of the above algorithm M contains the collection of maximal P -classes in τ .

6.6. REMARK. We can mention two obvious ways of modifying Algorithm 6.5 to improve its performance. Looking at the innermost loop (controlled by "for $t \in \tau_0$ "), we can easily see that if $N(t_1) = N(t_2)$ for $t_1, t_2 \in \tau_0$, then for any $T \in \tau$, $N(T \circ t_1) = N(T \circ t_2)$. This means that in the next passage through the main loop $T \circ t_1$ and $T \circ t_2$ will be merged into $(T \circ t_1) \circ (T \circ t_2) = (T \circ t_2) \circ (T \circ t_1) = T \circ (t_1 \circ t_2) = T \circ (t_2 \circ t_1)$. So we might just as well do the merging here, that is, work with the composition $t_1 \circ t_2$ instead of t_1 and t_2 separately. Applying to the collection τ_0 the operation 6.4(i) of merging tripartitions with identical N -classes, we get a collection τ_1 that may be used instead of τ_0 in Algorithm 6.5 (it is easy to see that also the initialization $C_1 := \tau_0$ may be replaced by $C_1 := \tau_1$). Replacing τ_0 by τ_1 speeds up the algorithm because $|\tau_1| \leq |\tau_0|$ and, in addition, any U -class in τ_1 is a subset of corresponding U -classes in τ_0 . The second obvious refinement is to do the COMPRESSing of C_{i+1} not at the beginning of the next passage through the main loop, but rather while creating C_{i+1} in the innermost loop. Instead of unconditionally adding to C_{i+1} each $T \circ t$ for which $U(T) \cap P(t) \neq \emptyset$, we can compare $N(T \circ t)$ with the N -classes of the tripartitions already collected in C_{i+1} . If there is T' in C_{i+1} such that $N(T') \subseteq N(T \circ t)$, we leave C_{i+1} unchanged in case $N(T') \neq N(T \circ t)$ and we replace T' by $T' \circ T \circ t$ in C_{i+1} if $N(T') = N(T \circ t)$; otherwise we add $T \circ t$ to C_{i+1} . Clearly, any thus constructed C_{i+1} is COMPRESSED, and initializing $C_1 = \tau_1$ we start with a COMPRESSED collection. Doing the COMPRESSing while creating the collection minimizes the number of comparisons that have to be made.

6.7. ACYCLIC RESTRICTIONS OF A RELATION. Why would we want to find maximal P -classes in a collection of tripartitions of E ? One possible application is in the following situation. We consider an arbitrary relation Q on E . A Q -cycle in E is a sequence of elements $x_i \in E$, $i = 1, \dots, n$, such that

$$x_1 Q x_2 Q \cdots Q x_n Q x_1.$$

Q is called *acyclic* on E iff there is no Q -cycle in E . An arbitrary $Q \subseteq E \times E$ need not be acyclic on E , of course, but we can pose the problem of finding subsets, and in particular *maximal* subsets $S \subseteq E$, for which $Q \cap (S \times S)$, the restriction of Q to S , is acyclic. We are going to characterize such subsets in terms of tripartitions of E . To each $x \in E$ we attribute a tripartition $T(x)$ of E , defined in the following way:

$$T(x) = (\{x\}, xQ, E - \{x\} - xQ),$$

where, in line with our conventions, xQ denotes the set of $y \in E$ for which xQy . We define τ_0 to be the collection of all such tripartitions $T(x)$:

$$\tau_0 = \{T(x) : x \in E\}$$

and we let τ denote the closure of τ_0 under composition. Now we can establish following connection between these tripartitions and acyclic restrictions of a relation.

6.8. PROPOSITION. *Let E be a finite set and let Q be a relation on E . Then, for $S \subseteq E$, $Q \cap (S \times S)$ is acyclic iff S is the P -class of some tripartition in τ , where τ is as defined in 6.7. In particular, the collection of maximal subsets of E on which Q is acyclic is MAX_{τ} , for τ defined as in 6.7.*

Proof. (If.) Consider $T_0 = (P_0, N_0, U_0) \in \tau$. To show that P_0 does not contain a Q -cycle, we use induction on the degree of T_0 . If this degree equals 1, the result is immediate, since then P_0 is singleton. If T_0 has degree $k > 1$ we can write $T_0 = T(x) \circ T'_0$ for some $x \in E$ and some $T'_0 \in \tau$ with degree $k-1$. By definition $U(x) \cap xQ = \emptyset$, so the element x cannot be involved in any Q -cycle in $P_0 = \{x\} + (P'_0 \cap U(x))$; on the other hand, by the induction hypothesis there is no Q -cycle in P'_0 , let alone in $P'_0 \cap U(x)$. The conclusion must be that there is no Q -cycle in P_0 .

(Only if.) Let S be such that $Q \cap (S \times S)$ is acyclic. To show that S is the P -class of some tripartition in τ , we proceed by induction on the cardinality of S . If $|S| = 1$, $S = \{x\}$ say, then $S = P(x)$ for the P -class of $T(x) \in \tau_0 \subseteq \tau$. If $|S| > 1$, let $x \in S$ be such that $S \cap xQ = \emptyset$. (Since S is finite, there has to be such an element if $Q \cap (S \times S)$ is acyclic.) Defining $S_0 = S - \{x\}$ we see that $S_0 \subseteq E - \{x\} - xQ$ and since $|S_0| = |S| - 1$, we may, by the induction hypothesis, assume $S_0 = P_0$ for the P -class of some $T_0 \in \tau$. But then $S = \{x\} + S_0 = \{x\} + (S_0 \cap (E - \{x\} - xQ)) = P(x) + (P_0 \cap U(x))$ is the P -class of $T(x) \circ T_0 \in \tau$. ■

6.9. APPLICATION TO HYPERGRAPHS $H(V)$. Proposition 6.8 gives us a translation of the problem of finding maximal acyclic restrictions of a relation in terms of finding maximal P -classes in some collection of tripartitions and via this result we can return to our original context of a relation R between two finite sets A and D . To show the connection of our excursion in the preceding paragraphs to the subhypergraphs $H(V)$ of $H(\bar{R})$, we define a relation Γ on \bar{R} in the following way:

$$ad \Gamma be \quad \text{iff} \quad bRd \quad (ad, be \in \bar{R}).$$

If we compare this definition with the definition of the hypergraph $H(\bar{R})$, given in the Introduction, we see that the edges of $H(\bar{R})$ coincide with the Γ -cycles in \bar{R} . The following equivalence is an immediate consequence:

$$S \subseteq \bar{R} \text{ is stable in } H(\bar{R}) \text{ iff } \Gamma \cap (S \times S) \text{ is acyclic.}$$

In fact, since, by the definition of the subhypergraph $H(V)$ of $H(\bar{R})$, a subset S of V is stable in $H(V)$ iff S is stable in $H(\bar{R})$, we have more generally:

$$S \subseteq V \subseteq \bar{R} \text{ is stable in } H(V) \text{ iff } \Gamma \cap (S \times S) \text{ is acyclic.}$$

In this way, stable sets in $H(V)$ are seen as a special instance of acyclic restrictions of a relation and from Proposition 6.8 we may obtain a translation in terms of tripartitions.

We define for $ad \in \bar{R}$ a tripartition $T(ad; \bar{R})$ of \bar{R} by

$$T(ad; \bar{R}) = (\{ad\}, ad\Gamma, \bar{R} - \{ad\} - ad\Gamma)$$

and collections of tripartitions

$$\tau_0(\bar{R}) = \{T(ad; \bar{R}) : ad \in \bar{R}\}$$

and $\tau(\bar{R})$, the closure of $\tau_0(\bar{R})$ under the composition operation \circ . For $V \subseteq \bar{R}$ and $ad \in V$, we consider the tripartition $T(ad, V)$ of V that results from restricting to V each class of $T(ad; \bar{R})$:

$$T(ad, V) = (\{ad\}, V \cap ad\Gamma, V - \{ad\} - ad\Gamma)$$

and we use the expected notations

$$\tau_0(V) = \{T(ad; V) : ad \in V\}$$

and $\tau(V)$ for the closure of $\tau_0(V)$ under composition. With these definitions the equivalence between stable sets in $H(V)$ and acyclic restrictions of Γ allows the following application of Proposition 6.8:

6.10. COROLLARY. *Let $R \subseteq A \times D$, $V \subseteq \bar{R}$. Then $S \subseteq V$ is stable in $H(V)$ iff S is the P -class of some tripartition in $\tau(V)$, where $\tau(V)$ is as defined in 6.9. In particular, $S(V) = \text{MAX}_{\tau(V)}$. ■*

6.11. A REINTERPRETATION OF ALGORITHM 3.12. Corollary 6.10 shows that Algorithm 6.5 is an algorithm for producing $S(V)$ when we define $\tau_0 = \tau_0(V)$ (or $\tau_0 = \tau_1(V)$, see Remark 6.6). In the special case $V = \bar{R}$, both Algorithm 6.5 and Algorithm 3.12 can be used to generate $S(\bar{R})$ and we may compare their ways of doing this. Reinterpreting Algorithm 3.12, we see that it starts with a collection of tripartitions

$$\{(\bar{R}d_i \times \{d_i\}, \bar{R}[Rd_i, D], \bar{R}[\bar{R}d_i, D - \{d_i\}])\}$$

that is “somewhere between” $\tau_0(\bar{R})$ and $\tau_1(\bar{R})$. From the definition of Γ it is clear that $a_1d_i\Gamma = a_2d_i\Gamma$ whenever $a_1d_i, a_2d_i \in \bar{R}$. In other words, $N(T(a_1d_i; \bar{R})) = N(T(a_2d_i; \bar{R}))$ and consequently tripartitions in $\tau_0(\bar{R})$ belonging to one and the same $d_i \in D$ may be replaced by their composition. Thus we obtain the collection of tripartitions that is LOADED in Algorithm 3.12: for any $d_i \in D$ we have

$$\begin{aligned} \cup \{P(T(ad_i; \bar{R})) : ad_i \in \bar{R}\} &= \cup \{\{ad_i\} : ad_i \in \bar{R}\} = \bar{R}d_i \times \{d_i\}; \\ N(T(ad_i; \bar{R})) &= ad_i\Gamma = \{be \in \bar{R} : bRd_i\} = \bar{R}[Rd_i, D]; \\ \cap \{U(T(ad_i; \bar{R})) : ad_i \in \bar{R}\} &= \bar{R} - \cup \{P(T(ad_i; \bar{R}))\} - N(T(ad_i; \bar{R})) \\ &= \bar{R} - (\bar{R}d_i \times \{d_i\}) - \bar{R}[Rd_i, D] \\ &= \bar{R}[\bar{R}d_i, D - \{d_i\}]. \end{aligned}$$

The MERGE operation that follows the LOADING in Algorithm 3.12 turns this collection of tripartitions into $\tau_1(\bar{R})$ and it is this collection on which subsequent computations are based. New tripartitions are computed from old ones according to the rules given in 3.9. The connection between these rules and the composition operation defined in 6.1 is obvious. Considering from now on the “cumulative” tripartitions in Algorithm 3.12 rather than the tripartitions of the current U -class, we see that these new tripartitions are obtained by taking compositions of already computed tripartitions. That is to say, just as Algorithm 6.5, Algorithm 3.12 computes members of $\tau(\bar{R})$, starting with $\tau_0(\bar{R})$ (or $\tau_1(\bar{R})$) and ending with tripartitions having a maximal P -class. In fact, the collection of unmarked tripartitions computed at the i -th level by Algorithm 3.12 (in the i -th “column” of Fig. 3.1) will equal the collection of tripartitions in the COMPRESSED C_i of Algorithm 6.5. The two algorithms differ in the way and the order in which these collections are generated. In Algorithm 3.12, a sort of depth-first search is used,

which is possible since comparing the N -classes can be done *locally*, that is, within the boxes. This may be seen as a consequence of Lemma 3.5, that is based on the biorderhood of maximal stable sets of $H(\bar{R})$. The same fact can be represented in terms of Algorithm 6.5, as a corollary of the following proposition:

6.12. PROPOSITION. *If Algorithm 6.5 is applied with $\tau_0 = \tau_0(\bar{R})$ (or $\tau_0 = \tau_1(\bar{R})$), then, for any i and any distinct tripartitions T_1, T_2 in a COMPRESSED C_i , we have $N_2 \cap P_1 \neq \emptyset$.*

Proof. As we have seen in 6.6, we may assume that we work with $\tau_0 = \tau_1(\bar{R})$. For $i = 1$ we know that the COMPRESSED C_i equals $\tau_1(\bar{R})$. So then

$$T_j = (\bar{R}d_j \times D_j, \bar{R}[Rd_j, D], \bar{R}[\bar{R}d_j, D - D_j])$$

where $D_j = \{d \in D : \bar{R}d = \bar{R}d_j\}$ for $j = 1, 2$ and some $d_1, d_2 \in D$. Consequently, $N_2 \cap P_1 = \bar{R}[Rd_2, D] \cap (\bar{R}d_1 \times D_1) = (\bar{R}d_1 \cap Rd_2) \times D_1$. We see that $N_2 \cap P_1 = \emptyset$ implies $\bar{R}d_1 \cap Rd_2 = \emptyset$, by which $N_2 - N_1 = \bar{R}[Rd_2, D] - \bar{R}[Rd_1, D] = \bar{R}[Rd_2 \cap \bar{R}d_1, D] = \emptyset$. This contradicts the fact that T_1 and T_2 are distinct elements of a COMPRESSED C_i . Now consider $T_1 \circ t_1$ and $T_2 \circ t_2$, two distinct tripartitions in C_{i+1} , where T_1 and T_2 are from the COMPRESSED C_i and $t_1, t_2 \in \tau_1(\bar{R})$. If T_1 and T_2 are distinct, we may by induction hypothesis assume that $N_2 \cap P_1 \neq \emptyset$. But then clearly $N(T_2 \circ t_2) \cap P(T_1 \circ t_1) \supseteq N_2 \cap P_1 \neq \emptyset$. So suppose $T_1 = T_2 = T$. Since $T \in C_i$ is the composition of a number of tripartitions in $\tau_1(\bar{R})$, we know that $U(T)$ is the intersection of a number of U -classes in $\tau_1(\bar{R})$. So $U(T) = \bar{R}[\alpha, \delta]$ for some $\alpha \subseteq A$ and $\delta \subseteq D$. Let t_j be the tripartition in $\tau_1(\bar{R})$ corresponding to $d_j \in D$, $j=1, 2$. If $T \circ t_j$ is to be a candidate for C_{i+1} we know that $P(t_j) \cap U(T) \neq \emptyset$. But then $D_j \subseteq \delta$, since the only alternative, $D_j \cap \delta = \emptyset$, would lead to $P(t_j) \cap U(T) = (\bar{R}d_j \times D_j) \cap \bar{R}[\alpha, \delta] = \emptyset$. Thus, $N(T \circ t_2) \cap P(T \circ t_1) = (N(T) + (N(t_2) \cap U(T))) \cap (P(T) + (P(t_1) \cap U(T))) = N(t_2) \cap P(t_1) \cap U(T) = \bar{R}[Rd_2, D] \cap ((\alpha \cap \bar{R}d_1) \times (\delta \cap D_1)) = (\alpha \cap \bar{R}d_1 \cap Rd_2) \times D_1$. So $N(T \circ t_2) \cap P(T \circ t_1) = \emptyset$ implies $\alpha \cap \bar{R}d_1 \cap Rd_2 = \emptyset$ and in this case $N(T \circ t_2) - N(T \circ t_1) = (N(T) + (N(t_2) \cap U(T))) - (N(T) + (N(t_1) \cap U(T))) = (N(t_2) - N(t_1)) \cap U(T) = (\bar{R}[Rd_2, D] - \bar{R}[Rd_1, D]) \cap \bar{R}[\alpha, \delta] = \bar{R}[\alpha \cap \bar{R}d_1 \cap Rd_2, \delta] = \emptyset$. This means that $T \circ t_1$ and $T \circ t_2$ will not survive as two distinct members of the COMPRESSED C_{i+1} . ■

6.13. By Lemma 6.2(vii) we may conclude from Proposition 6.12 that when $S(\bar{R})$ is generated by Algorithm 6.5, for any COMPRESSED C_i all N -classes of tripartitions in the offspring of different elements of C_i will be mutually incomparable and the collections of maximal P -classes generated in the offspring of different elements

will be disjoint. For $V \subseteq \bar{R}$, that is, with $\tau_0 = \tau_1(V)$, Proposition 6.12 is no longer valid. So, for generating $S(V)$ we have to resort to the breadth-first search of Algorithm 6.5, comparing globally all N -classes at each level.

6.14. Returning to the problem at the base of this section, that of evaluating (5.2), we see that we have established a third alternative beside the two mentioned in 5.7. Proposition 5.2(ii) shows that Algorithm 6.5, applied with $\tau_0 = \tau_0(V')$ for the proper V' , is an algorithm for directly generating $S(ad; V)$. This is not to say that it will be the best alternative in practice. We can summarize the two alternatives of 5.7 as:

- (i) generate $S(ad; \bar{R} | V)$ and work with this collection, or
- (ii) generate $S(ad; V)$ by generating $S(ad; \bar{R} | V)$ and applying one COMPRESS operation to the corresponding collection of tripartitions.

In view of the discussion in 6.11 to 6.13, we can describe the third alternative, computing $S(V')$ by Algorithm 6.5, as:

- (iii) generate $S(ad; V)$ by generating $S(ad; \bar{R} | V)$ while applying a COMPRESS operation to the collections of tripartitions produced at each level in the generating procedure.

In this perspective, alternative (iii) is one further step in the same direction after the step from (i) to (ii). It would seem that (iii) will not be preferable to (i) unless (ii) is preferable to (i) and the relative merits of (ii) and (iii) will again depend on the peculiarities of the situation, in particular the discrepancy between V and \bar{R} .

REFERENCES

- Doignon, J.-P., Ducamp, A. & Falmagne, J.-C. (1984). On realizable biorders and the biorder dimension of a relation. *Journal of Mathematical Psychology*, **28**, 73-109.
- Koppen, M. G. M. (1987). On finding the bidimension of a relation. *Journal of Mathematical Psychology*, **31**, 155-178. (Chapter 3.)

CHAPTER 5

DISCUSSION OF PART I

Here we reconsider the results of the preceding chapters from a more general perspective. The chapter is divided in two sections. First we have another look at the translation of the biorder representation problem in terms of a hypergraph and we relate its significance to the close link between biorders and linear orders. This first section is more theoretical and involves primarily a reinterpretation of the results of Chapter 3. In the second section we shift our attention to more practical issues and we investigate the possibilities for applying the obtained results. Two major problems are addressed, the uniqueness of solutions and the deterministic character of the developed methods. In this section the emphasis is more on the results of Chapter 4.

1. The connection between bidimension and order dimension

Central to the developments in the preceding chapters was Doignon, Ducamp and Falmagne's (1984) characterization of the biorder dimension of a relation as the chromatic number of a particular hypergraph. Actually, this way of putting it is a bit misleading, since it is not just the rewriting of the bidimension as a chromatic number that does the job. It appears that, for the results of the preceding chapters, the specific structure of the hypergraph that is involved is essential. As we saw in Chapter 2, Eqs. (12) to (16), the hypergraph concept is so general that it entered the discussion almost sneakily, on a very general level. Let us first show that, indeed, a chromatic number characterization is available for other notions of the "dimension" of a relation. Next, we will give an example of such an alternative definition, in which the reformulation in terms of a hypergraph does not appear to be of much help for the purpose of computing this dimension. The rest of this section is devoted to a discussion of properties that are specific to the biorder dimension and the associated hypergraph; they will be intimately related to the theory of partial orders.

1.1. The general representation problem for relations.

In the preceding chapters we were concerned with the collection \mathbf{B} of biorders. More generally, we consider here an arbitrary specified collection \mathbf{T} of relations between A and D . (No finiteness conditions are used.) For such a collection \mathbf{T} , we can investigate the problem of representing any relation $R \subseteq A \times D$ as the union of a number of members of \mathbf{T} :

$$R = \cup \tau, \quad \tau \subseteq \mathbf{T}. \quad (1)$$

There is, of course, a dual representation problem with "union" replaced by "intersection":

$$R = \cap \tau, \quad \tau \subseteq \mathbf{T}. \quad (2)$$

In this context, it is of special interest to find *minimal* families τ for which (1) or (2) holds. If the representation (1) is possible, the minimum cardinality of such a family τ is called the *union \mathbf{T} -dimension* of R . Similarly, (2) leads to the definition of the *intersection \mathbf{T} -dimension*.

We confine here our attention to the union representation problem (1) and the union \mathbf{T} -dimension. This is no real restriction, since any intersection representation problem can be formulated as a union representation problem. Denoting the collection of complements (with respect to $A \times D$) of the relations in a collection τ by $\bar{\tau}$, i.e.,

$$T \in \tau \quad \text{iff} \quad (A \times D) - T \in \bar{\tau},$$

we immediately obtain:

$$(R = \cup \tau, \tau \subseteq \mathbf{T}) \quad \text{iff} \quad (\bar{R} = \cup \bar{\tau}, \bar{\tau} \subseteq \bar{\mathbf{T}}).$$

Consequently, any intersection representation problem for R with respect to the collection \mathbf{T} is a union representation problem for \bar{R} with respect to $\bar{\mathbf{T}}$ and the intersection \mathbf{T} -dimension of R equals the union $\bar{\mathbf{T}}$ -dimension of \bar{R} . (In the case of biorders, $\bar{\mathbf{B}} = \mathbf{B}$ – the complement of a biorder is again a biorder – and the above equivalence reduces to the duality between the conjunctive and disjunctive models.)

Accordingly, let R be a relation for which the representation (1) exists. (This is the case for *any* R if and only if all singleton relations $\{au\}$, $a \in A$, $u \in D$, are in \mathbf{T} ; e.g., the union biorder dimension is defined for any relation because any singleton relation is a biorder.) Then the (union) \mathbf{T} -dimension of R is defined as:

$$\mathbf{T}\text{-dim } R = \text{Min} \{ |\tau| : R = \cup \tau, \tau \subseteq \mathbf{T} \}. \quad (3)$$

Now we can mimick the operations that brought us in Chapter 2 from Eq. (13) to (16). Calling a subset C of R *feasible* (for \mathbf{T} with respect to R) if and only if there is $T \in \mathbf{T}$ such that $C \subseteq T \subseteq R$, we obtain from (3):

$$\mathbf{T}\text{-dim } R = \text{Min } \{ |\gamma| : R = \cup \gamma, \gamma \text{ consists of feasible sets} \}. \quad (4)$$

Since subsets of feasible sets are again feasible, we can, as in Chapter 2, replace any collection of feasible sets by a collection of mutually disjoint feasible sets without changing the union. Consequently:

$$\mathbf{T}\text{-dim } R = \text{Min } \{ |\gamma| : R = \sum \gamma, \gamma \text{ consists of feasible sets} \}. \quad (5)$$

The sigma sign means union over disjoint sets (γ is a partition of R) and the right hand side of (5) is the definition of the chromatic number of a hypergraph. We conclude that whenever the union \mathbf{T} -dimension of a relation R is defined, it equals the chromatic number of the hypergraph with vertex set R and where maximal stable sets consist of the maximal members of \mathbf{T} contained in R . The edges of this hypergraph are, thus, the (minimal) subsets of R that cannot be extended to a member of \mathbf{T} contained in R . The characterization of these edges clearly depends on the chosen collection \mathbf{T} and this characterization determines the prospects of success for the reduction method described in Chapter 3.

1.2. Another example: the matching dimension.

That these prospects may differ from case to case can be seen from another example that has been described in the literature. Doignon and Falmagne (1984) call a relation M between A and D a *matching relation* if, for all $a, b \in A, u, v \in D$,

$$aMu \ \& \ bMu \ \& \ bMv \ \text{implies} \ aMv,$$

which means that a matching relation is any relation M for which there exist functions $f : A \rightarrow \mathbb{R}$ and $g : D \rightarrow \mathbb{R}$ such that

$$aMu \ \text{iff} \ f(a) = g(u).$$

Doignon and Falmagne consider the problem of representing a relation R as the union of a minimal number of matching relations. Since any singleton relation $\{au\}$ is a matching relation, this minimal number, the *matching dimension*, is defined for any relation R and the approach of the previous subsection is available.

Doignon and Falmagne give indeed a reformulation of this matching dimension as the chromatic number of a hypergraph and they characterize the edges of this hypergraph. That is, everything is set for applying the reduction methods of Chapter 3 to this hypergraph. Then, however, we are in for a disappointment. Not only are the minimal edges in practice hard to find (this shows the importance of a result like Proposition 3.7 of Chapter 3), but, typically, no reduction of the hypergraph is obtained at all. This shows how critically the structure of the associated hypergraphs can vary: while the biorder hypergraph generally contains a number of irrelevant vertices (irrelevant for the chromatic number that is: the *dominated* vertices), this is

not the case for the matching hypergraph. The occurrence of dominated vertices in the biorder hypergraph can, at least in part, be related to the distinction between *forced* and *non-forced* pairs in the theory of partial orders (see, e.g., Trotter, 1983; Maurer, Rabinovitch and Trotter, 1980). Before investigating this relationship in the subsections 1.5 and 1.6, we first describe the close connection between the biorder dimension of a relation and the (linear) order dimension of a partial order.

1.3. The bidimension as an order dimension.

If we consider the special case where $D=A$, it is easy to check that any linear order on A is a biorder on A . An immediate consequence is that for any quasi order Q on A the bidimension of Q cannot exceed the order dimension of Q as defined by Dushnik and Miller (1941) (see the General Introduction). On the other hand, as Doignon, Ducamp and Falmagne (1984) show, any *minimal* biorder extension of a quasi order is in fact a weak order. From this observation they derive easily that the two numbers must coincide: for any quasi order Q on A we have

$$\text{Bidim } Q = \text{Dim } Q, \quad (6)$$

where $\text{Dim } Q$ denotes the classic order dimension of Q . Thus, the biorder dimension, defined for any relation, may be considered as a generalization of the order dimension as defined for quasi orders.

It appears, however, that we can also go the other way: for any relation R between (possibly different) sets A and D , $\text{Bidim } R$ can be interpreted as the order dimension of a distinguished quasi order associated with R . Using again the notations

$$aR = \{u \in D : aRu\} \quad \text{and} \quad Ru = \{a \in A : aRu\},$$

we define, for any $R \subseteq A \times D$, relations R_A on A , R_D on D and R_X from D to A as follows: for any $a, b \in A$, $u, v \in D$,

$$aR_A b \quad \text{iff} \quad aR \supseteq bR, \quad (7a)$$

$$uR_D v \quad \text{iff} \quad Ru \subseteq Rv, \quad (7b)$$

$$uR_X a \quad \text{iff} \quad Ru \times aR \subseteq R. \quad (7c)$$

Definitions (7a) and (7b) show that R_A and R_D are quasi orders on A and D , respectively. In terms of the corresponding $(0,1)$ -matrix, $aR_A b$ reflects the fact that the row indexed by a dominates the row indexed by b ; $uR_D v$ means that the column indexed by u is dominated by that indexed by v ; finally, $uR_X a$ holds whenever the submatrix consisting of the rows dominating u and the columns that are dominated by a contains only "1" entries.

These four relations R , R_A , R_D and R_X can be regarded as a partitioning of a relation Q_R on $A \cup D$:

$$Q_R = R_A + R + R_X + R_D. \quad (8)$$

This relation Q_R , the original construction of which is due to Bouchet (1971), defines a quasi order on $A \cup D$; it is in fact the maximal quasi order on $A \cup D$ whose restriction to $A \times D$ coincides with R (Doignon *et al.*, 1984). With respect to this quasi order we have the following fundamental result:

$$\text{Bidim } R = \text{Dim } Q_R. \quad (9)$$

This was already established by Cogis (1980, 1982) in the finite case; Doignon *et al.* (1984) give a proof for the unrestricted case. Together, the results (6) and (9) show why computing the bidimension of a relation and computing the order dimension of a quasi order are two equivalent problems.

1.4. An alternative proof.

It may be interesting to note that the reduction theory for the hypergraph $H(R)$ that was developed in Chapter 3 can be used to obtain a simple proof for the equality in (9). From the definition of Q_R it follows that for $a \in A$:

$$Q_R a = \bigcap_{u \in aR} Q_R u \quad (10)$$

and for $u \in D$:

$$u Q_R = \bigcap_{a \in Ru} a Q_R. \quad (11)$$

To establish (10), for instance, we note that if $p \in A$, then, by (7a),

$$p Q_R a \text{ iff } pR \supseteq aR \text{ iff } (aRu \text{ implies } pQ_R u);$$

if instead $p \in D$, then, by (7c) and (7b),

$$\begin{aligned} p Q_R a \text{ iff } Rp \times aR \subseteq R \text{ iff } (aRu \text{ implies } Rp \subseteq Ru) \\ \text{iff } (aRu \text{ implies } pQ_R u). \end{aligned}$$

There is a similar proof for (11). In Chapter 3, we defined an ‘‘implied’’ column as one that could be obtained as the conjunction of a number of other columns of the matrix representation of a relation, and we observed there that implied columns could be removed without changing the bidimension. Now (10) means that, in the matrix for Q_R , any column indexed by an element of A is an implied column; similarly, according to (11), any row indexed by an element of D is an implied row. Removing these columns and these rows reduces the matrix of Q_R to that of R ,

which implies that

$$\text{Bidim } Q_R = \text{Bidim } R,$$

and combining this with (6) we obtain (9). (Although in Chapter 3 we restricted ourselves to the finite case, it can easily be checked that the above proof works in the general case, i.e. that the removal of any (infinite) collection of implied patterns does not change the bidimension.)

1.5. Trotter's original hypergraph.

The hypergraph defined by Doignon *et al.* (1984) was a generalization of the hypergraph that Trotter (1983) associates with a partial order. This would imply that for this special case, when $D=A$ and the relation R is a partial order, the two definitions should coincide. This is indeed the case, except for one detail: while Doignon *et al.*'s hypergraph for a partial order P would have the full set \bar{P} as vertex set, Trotter defined his hypergraph on the subset of "non-forced" pairs in \bar{P} . A pair $xy \in \bar{P}$ is called *non-forced* if and only if the relation $P \cup \{xy\}$ is a partial order. A non-forced pair must certainly be an incomparable pair for P , that is, $xy \in \bar{P} \cap \bar{P}^{-1}$, since if $yx \in P$, then $P \cup \{xy\}$ is not antisymmetric ($xy \in \bar{P}$ implies $x \neq y$). But not all incomparable pairs are non-forced: it might well be that $P \cup \{xy\}$ is no longer transitive and that to restore this property more pairs would have to be added. Clearly, if there are no incomparable pairs there can be no non-forced pairs and Trotter's hypergraph has no vertices. To avoid trivialities, we assume in the sequel that P is not a linear order, so there are incomparable pairs. Then there must be non-forced pairs among these incomparable pairs and Trotter's hypergraph is well-defined.

Apparently, Trotter's hypergraph is a subhypergraph of that of Doignon *et al.*, the subhypergraph induced by restricting the vertex set to the collection of non-forced pairs in \bar{P} . However, by the result (6), both must have the same chromatic number. This can be seen to follow directly from the reduction theorems of Chapter 3: it appears that in applying these to Doignon *et al.*'s hypergraph we would immediately reduce it to Trotter's subhypergraph on the non-forced pairs (and the reduction might go beyond this point). Let us state the relevant result more precisely:

1.6. Proposition.

Let P be a partial order on A and $xy \in \bar{P}$. Then $P \cup \{xy\}$ is a partial order if and only if xy is not a dominated vertex of $H(\bar{P})$.

Recall that, according to Proposition 3.7 and Definition 3.6 of Chapter 3, the vertex xy of $H(\bar{P})$ is dominated by $x'y'$ if and only if for any 2-edge of $H(\bar{P})$ containing xy

we get another 2-edge by replacing xy by $x'y'$. (The pairs xy and uv constitute a 2-edge whenever $x\bar{P}y$, $u\bar{P}v$, xPy and uPy .)

Proof of the Proposition. We first show that $P \cup \{xy\}$ is not a partial order if xy is a dominated vertex. This is immediate if $yx \in P$ (antisymmetry would be violated), so we may suppose that both $xy, yx \in \bar{P}$. A partial order is a reflexive relation, so $xx, yy \in P$. In other words, $\{xy, yx\}$ is a 2-edge of the hypergraph $H(\bar{P})$. If xy is dominated in $H(\bar{P})$ by $x'y'$, then $\{x'y', yx\}$ is also a 2-edge of $H(\bar{P})$, which yields $x'Px$ and yPy' . For any partial order P' containing both P and xy we can then write $x'P'xP'yP'y'$, so P' contains also $x'y' \notin P \cup \{xy\}$.

Conversely, if $P \cup \{xy\}$ is not a partial order, either the antisymmetry or the transitivity must be violated. The former is the case if and only if yPx . But then uPy and xPv would imply uPv , which shows that we cannot have any 2-edge $\{xy, uv\}$. Thus xy is dominated (for instance by a vertex of the incomparable pair that we suppose to be there). $P \cup \{xy\}$ violates transitivity if and only if there is $x'y' \notin P \cup \{xy\}$, such that $x'Px$ and yPy' . But then uPy implies uPy' and xPv implies $x'Pv$, which shows that for any 2-edge $\{xy, uv\}$ we have a 2-edge $\{x'y', uv\}$: xy is dominated by $x'y'$. ■

(It can readily be checked that, at the end of the proof, the vertex xy is not only dominated, but even "implied" by $x'y'$ in the sense of the definition of an implied zero in Chapter 3, Section 3. Indeed, if $x' \neq x$, then $x'\bar{P}y$ and $xP \subseteq x'P$; if $y' \neq y$, then $x\bar{P}y'$ and $Py \subseteq Py'$. This illustrates the remark in the last paragraph of Section 3 of Chapter 3 to the effect that here any dominated zero is an implied zero. However, it must be noted that this refers only to the non-trivial situation where xy is contained in any 2-edge at all; this proviso should have been added in Chapter 3.)

The above proposition shows that, in the case of a partial order, the methods of Chapter 3 reduce Doignon *et al.*'s hypergraph to Trotter's hypergraph. For the special kind of quasi order Q_R on $A \cup D$ as defined by (8) and (7), all non-dominated vertices of $H(\bar{Q}_R)$ are within $A \times D$, as we saw in 1.4. Thus, according to Proposition 1.6, Trotter's hypergraph for Q_R is a hypergraph on \bar{R} right from the start, and as such it must be a subhypergraph of Doignon *et al.*'s $H(\bar{R})$. This is another illustration of the basic equality (9), which establishes the bidimension as an instance of the classical order dimension in the sense of Dushnik and Miller (1941).

2. Prospects for applications

2.1. The uniqueness problem.

The results of the preceding chapters are constructive in nature. In Chapter 4 we developed algorithms for the construction, for any relation $R \subseteq A \times D$, of all minimal biorders containing R . The combination of results in the Chapters 3 and 4 leads to an effective procedure for computing the bidimension of R , together with a number of biorder representations of R in this dimensionality. So, in principle, we have found a satisfactory solution for the problem we started with in Chapter 2, finding a conjunctive or disjunctive model representation of a binary data matrix in the minimum dimensionality. However, one major drawback in applying these models is the lack of uniqueness of an obtained solution.

This problem is already apparent in the method where we compute the bidimension via the equivalence with the hypergraph. Suppose we have a coloring of $H(\bar{R})$ in a minimal number of color classes $C_i \subseteq \bar{R}$, $i = 1, \dots, q$. Denote by \mathbf{B}_i the collection of biorders B_i such that $C_i \subseteq B_i \subseteq \bar{R}$. Since each C_i is stable, these collections will be non-empty and \mathbf{B}_i will in fact contain more than one element unless the color class C_i is itself a maximal biorder in \bar{R} . It is clear that any arbitrary combination of elements B_i from the different \mathbf{B}_i leads to a distinct biorder representation for R :

$$\bar{R} = \bigcup_{i=1}^q C_i \subseteq \bigcup_{i=1}^q B_i \subseteq \bar{R},$$

and thus

$$R = \bigcap_{i=1}^q \bar{B}_i.$$

Consequently, this coloring alone produces

$$\prod_{i=1}^q |\mathbf{B}_i|$$

distinct solutions in the minimum dimensionality. In general there will be a number of minimal colorings and each of these may add new combinations of biorder extensions of the colors. The examples in Coombs and Kao (1955), Coombs (1964) and Table 1 of Chapter 2 have an (essentially) unique solution because there is (essentially) just one minimal coloring of the corresponding hypergraph, in which the color classes themselves are maximal biorders in \bar{R} . These examples were carefully constructed; they do not represent the typical situation.

The above discussion shows that when we compute the bidimension by the methods of Chapters 3 and 4, we obtain a number of minimal representations almost

as a by-product. If we are interested in the practical problem of finding all minimum dimension solutions for R , only a straightforward approach seems available: compute the bidimension k of R , compute the collection of minimal biorder extensions of R and test for every k -combination of this collection whether their intersection equals R . For fixed k , the number of such combinations is given by a k th degree polynomial in the number of minimal biorder extensions (provided that k is less than half this latter number). Consequently, to the extent that this approach is feasible at all (that is, to the extent that the collection of minimal biorder extensions is manageable), it will only be so for small values of k .

There are, however, deeper reasons why we should not try to find all solutions when the bidimension turns out to be relatively high. We are referring here to another obvious drawback of the procedures as developed so far: they are completely deterministic. One "error" in the data, one "0" that "really" should have been a "1", or vice versa, can easily change the bidimension of a data matrix. Since the binary data in question typically are noisy, such a deterministic analysis will generally reveal a larger number of dimensions than we want to retain. This introduces the problem of finding low dimensional representations that "best" fit the observed high dimensional data matrix. In the following two subsections we explore some possible first steps in this direction. Unfortunately, the above mentioned uniqueness problem will appear to be our constant companion on these excursions.

2.2. Approximate biorder representations.

When computing the bidimension of a data matrix results in a number that we deem unacceptably high, it is natural to inquire how bad the situation really is. Is the high dimensionality a stable characteristic of the data, overall, or is it perhaps caused by a few deviating response patterns with low frequencies of occurrence? More precisely, how many of the responses would have to be adjusted (a "0" converted to a "1" or vice versa) in order to obtain an acceptable bidimension? It is important to realize that now the frequency of observed patterns plays a role, while in the deterministic computations and constructions of the preceding chapters only the occurrence or non-occurrence of a pattern was significant. Another way of making the same point: we now have to work – at least conceptually – with the full data matrix or relation, with all replications of identical row or column patterns present.

The above considerations lead to the formulation of the following approximate biorder representation problem, which was already suggested by Charles Chubb (personal communication):

Given a relation R and a positive integer k (the intended low dimensionality), find the minimal number m_k for which there exist biorders B_1, \dots, B_k such that, with $L = \bigcap_{i=1}^k B_i$,

$$|L - R| + |R - L| = m_k. \quad (12)$$

Here L is the low dimensional approximation to R and the distance ("stress") m_k between L and R is, according to (12), measured by the size of their symmetric difference (the number of entries on which they disagree).

In this form, the problem is obviously extremely difficult. It is true that for the above L we need to consider only borders that are closest to R , i.e., which have a minimal symmetric difference with R . Indeed, replacing one of the borders constituting L by one whose symmetric difference with R is a subset of the first can never increase the size of the symmetric difference of the intersection with R . But there will generally be many borders that are closest to R in terms of the symmetric difference. At the two extremes this collection includes the minimal border extensions of R (for which the second term in (12) will vanish) as well as the complements of those of \bar{R} (with first term in (12) equal to zero), but there will be many more "inbetween".

It might be interesting, then, to try and simplify the problem. The first step in this direction will actually consist in complicating the issue. Let us think, very informally, of a "model" through which the "latent", "true" relation L gives rise to the observed relation R . One way to account for the discrepancies between L and R is to imagine that the positions of persons on the postulated dimensions are not fixed, but can have some disturbances round their average point. Such disturbances can cause a person to be momentarily higher on a dimension than an item, even if his average position is below the item on that dimension. In other words, a latent "0" score on a dimension can be turned into a "1" score. A similar process can take place in the opposite direction. It is convenient to abuse some terminology from the unrelated field of signal detection theory and call a *miss* a "1" entry in the latent relation L that appears as a "0" in the observed R ; conversely, a "1" in R corresponding to a "0" in L is designated as a *false positive*. Now, according to the conjunctive model, a "1" in R is produced by a conjunction of "1"s on the separate dimensions and a "0" by a disjunction of "0"s. Therefore, a miss may be produced by a significant negative disturbance on *any one* dimension, while a false positive does not obtain unless there is a significant disturbance in the positive direction on *all* dimensions on which the person is, on average, below the item, *and* on *all* dimensions on which the average person position is above the item there is no significant disturbance in the negative direction.

All of this suggests that the status of "0" and "1" entries in R is not symmetrical. Consequently, in computing a distance between L and R , false positives and misses should be weighted equally only in the one-dimensional case; with more dimensions, false positives should have a greater weight, the ratio of the weights being an increasing function of the number of dimensions. If we normalize

by giving unit weight to a miss, we obtain the following generalization of the above approximation problem:

Given a relation R , an integer $k > 1$ and a weight $\omega_k > 1$, find the minimal number m_k for which there exist biorders B_1, \dots, B_k such that, with

$$L = \bigcap_{i=1}^k B_i,$$

$$|L - R| + \omega_k |R - L| = m_k. \quad (13)$$

Obviously we have made the problem more difficult; we have in fact created an infinity of approximation problems, one for each choice of ω_k . However, at this point a simplification is possible. In the above discussion we derived a miss from a disjunction of disturbances in a negative direction on k dimensions, while a false positive requires a conjunction of disturbances in specific directions on the k dimensions. Assuming independence of the disturbances on the various dimensions (at least to some approximation) it is then tempting to conclude that, for any $k > 1$, a false positive is much less likely to occur than a miss. In terms of the weights this implies that ω_k is really large: $1 \ll \omega_k$. This means effectively that in the above generalized version of the approximation problem the minimum m_k will only be obtained for approximations L for which $|R - L| = 0$. To the extent that this is true, the general approximation problem is equivalent to the following simpler version:

Given a relation R and an integer $k > 1$, find the minimal number m_k for which there exist biorders $B_1, \dots, B_k \supseteq R$ such that, with $L = \bigcap_{i=1}^k B_i$,

$$|L - R| = m_k. \quad (14)$$

It is clear that in this version we need to consider only *minimal* biorder extensions; the problem is reduced to approximating R by k minimal biorder extensions. All such extensions are computed by an algorithm presented in Chapter 4 and we face again the problem of trying all k -combinations to find a best approximation. An upper bound on the minimum number m_k in (14) can be obtained from any minimal coloring of $H(\bar{R})$, by considering all combinations of minimal biorder extensions of k out of the $\text{Bidim}(R)$ different colors.

It must be clear that looking for best approximations does only aggravate the uniqueness problem. Not only can there be multiple solutions L for the minimal m_k in (14), but for each such L we face again the problem discussed in the preceding subsection: there will in general be a number of decompositions of L as the intersection of k biorders.

2.3. Towards a probabilistic model

Note that the above formulated approximation problems are still essentially non-probabilistic (although we used some quasi-probabilistic argument to obtain the last, simple version). In order to evaluate how well the data can be described by a low-dimensional approximation we need a truly probabilistic model. So let us finally sketch a possible – but, as we will see, not quite satisfactory – approach in this direction.

In a probabilistic model, the observed relation $R \subseteq A \times D$ is considered the realization of some random variable. This random variable is a matrix of $|A| \times |D|$ jointly distributed scalar random variables, any relation between A and D being determined by this number of entries. Random variables corresponding to different persons (elements of A) may naturally assumed to be mutually independent; for random variables representing the results for one person on different items, the principle of local stochastic independence is invoked. Thus, all the variables are independent and the joint distribution follows from a specification of the distributions of the separate variables.

The random variable corresponding to the pair au can take two values: we may observe aRu (person a solves item u) or $a\bar{R}u$ (a fails u). The probabilities for these two events depend on the “true”, situation. In the model that follows, which was suggested by Jean-Claude Falmagne (personal communication), we assume that there is a latent $T \subseteq A \times D$ representing the true relation between a and u : aTu if person a masters item u , $a\bar{T}u$ if a does not master u . The random component in R is now introduced by allowing for a correct answer when the item is not really mastered and an incorrect response when the item is mastered. More specifically, we postulate for every item u a careless error parameter α_u and a lucky guess parameter β_u , which results in the following distribution:

$$\left. \begin{array}{l} \mathbb{P}(aRu) = 1 - \alpha_u \\ \mathbb{P}(a\bar{R}u) = \alpha_u \end{array} \right\} \text{if } aTu$$

$$\left. \begin{array}{l} \mathbb{P}(aRu) = \beta_u \\ \mathbb{P}(a\bar{R}u) = 1 - \beta_u \end{array} \right\} \text{if } a\bar{T}u \quad (15)$$

Notice that this error model is different from the “model” that was suggested, very informally, in the subsection on approximate solutions. In particular, here the error parameters do not depend on the relative positions of items and persons on the separate dimensions. Indeed, the errors are not related at all to the underlying dimensions. Clearly, some refinement would be in order here, but we will continue with the above approach, accepting it as a first, maybe crude version of a probabilistic model.

In sum, then, we have three sets of parameters: $\{\alpha_u\}_{u \in D}$ and $\{\beta_u\}_{u \in D}$, all with values in the interval $[0, 1]$, and $T \subseteq A \times D$, with values in the collection \mathbf{R}_k of relations between A and D having bidimension not exceeding k . Here k is the maximum dimensionality we are willing to accept for the true, latent structure. Given the specification (15), we can readily write down the likelihood of the observed relation R as a function of these parameters. Using local stochastic independence it follows:

$$\begin{aligned} L(R) &= \prod_{au \in T \cap R} L(aRu) \prod_{au \in T \cap \bar{R}} L(a\bar{R}u) \prod_{au \in \bar{T} \cap R} L(aRu) \prod_{au \in \bar{T} \cap \bar{R}} L(a\bar{R}u) \\ &= \prod_{u \in D} (1 - \alpha_u)^{h_u} \cdot \alpha_u^{m_u} \cdot \beta_u^{f_u} \cdot (1 - \beta_u)^{c_u}, \end{aligned}$$

where (with Tu as usual denoting $\{a \in A : aTu\}$ and similarly for Ru)

$$\begin{aligned} h_u &= |Tu \cap Ru|, & m_u &= |Tu \cap \bar{R}u|, \\ f_u &= |\bar{T}u \cap Ru|, & c_u &= |\bar{T}u \cap \bar{R}u| \end{aligned} \quad (16)$$

are the numbers of ‘‘hits’’, ‘‘misses’’, ‘‘false positives’’ and ‘‘correct negatives’’, respectively (in R with respect to T). (To avoid heavy notation, we do not explicitly show the dependence of these numbers on T and R , which is obvious from the definition.) It is trivial to solve the likelihood equations for the parameters α_u and β_u in terms of the parameter T :

$$\log L(R) = \sum_{u \in D} h_u \log(1 - \alpha_u) + m_u \log \alpha_u + f_u \log \beta_u + c_u \log(1 - \beta_u), \quad (17)$$

and thus,

$$\begin{aligned} \frac{\partial \log L(R)}{\partial \alpha_u} &= \frac{-h_u}{1 - \alpha_u} + \frac{m_u}{\alpha_u} = 0 \quad \text{iff} \quad \alpha_u = \frac{m_u}{h_u + m_u} \\ \frac{\partial \log L(R)}{\partial \beta_u} &= \frac{f_u}{\beta_u} + \frac{-c_u}{1 - \beta_u} = 0 \quad \text{iff} \quad \beta_u = \frac{f_u}{f_u + c_u} \end{aligned} \quad (18)$$

We see that the maximum likelihood estimates of α_u and β_u are exactly as we would expect. The careless error parameter α_u is estimated by the proportion of misses on u by persons that have mastered item u ($h_u + m_u = |Tu|$), the lucky guess parameter β_u by the number of false positives on u by persons that have not mastered item u ($f_u + c_u = |\bar{T}u|$). On substituting (18) in (17) we obtain an expression for the log likelihood to be maximized solely in terms of T :

$$\log L(R) = \sum_{u \in D} h_u \log \frac{h_u}{h_u + m_u} + m_u \log \frac{m_u}{h_u + m_u} + f_u \log \frac{f_u}{f_u + c_u} + c_u \log \frac{c_u}{f_u + c_u}$$

$$\begin{aligned}
&= \sum_{u \in D} (h_u \log h_u + m_u \log m_u - (h_u + m_u) \log (h_u + m_u) + \\
&\quad f_u \log f_u + c_u \log c_u - (f_u + c_u) \log (f_u + c_u)). \quad (19)
\end{aligned}$$

This may all look very fine, but now we must finally face the fact that the likelihood is a function of the discrete parameter T with values in \mathbf{R}_k . Even for small k this collection of relations with bidimension bounded by k may be huge, so the practical problem of efficiently searching this collection for an element that maximizes (19) is a formidable combinatorial task. The situation is not quite as bad as this description suggests, since it can be shown that we need to consider only biorders whose symmetric difference with R is minimal: the maximum for (19) can always be found within the subcollection of \mathbf{R}_k consisting of the intersections of k such biorders.

Restricting the search to this subcollection of \mathbf{R}_k makes sense for another reason: it will prevent us from obtaining some ridiculous solutions. Note the symmetry in (19) between h_u and m_u , as well as that between f_u and c_u . Nothing in the likelihood expression indicates that hits and correct negatives are "good", and that misses and false positives are "bad". To put it in other words, the likelihood expression is totally blind for whether it describes a fit for R or for \bar{R} , since changing from an observed relation to its complement amounts to switching the roles of h_u and m_u , as well as those of f_u and c_u , simultaneously. This undesirable symmetry in the likelihood expression is due to the simple way in which (15) incorporates the random components in our model.

We encountered the subcollection of biorders having minimal symmetrical difference with R already in the case of the approximate representation problem. There we noted that this collection will still be huge and we went on to alleviate the problem by imposing some simplification that brought the results of Chapter 4 to bear on the situation. Here we can do something similar. Under some rather general circumstances (open ended items) it is reasonable to assume that the probability of a correct guess is negligibly small. Then we may set $\beta_u = 0$ for all $u \in D$. This implies that in (19) $f_u = 0$ and $h_u = |R u|$, independent of T , and consequently (19) reduces to

$$\log L(R) = \sum_{u \in D} |R u| \log |R u| + m_u \log m_u - (|R u| + m_u) \log (|R u| + m_u), \quad (20)$$

where the domain of T is restricted to the elements of \mathbf{R}_k that contain R ($f_u = |\bar{T}u \cap R u| = 0$). Thus, T now ranges over the intersections of k or less biorder extensions of R . As with (19), it is the case that the maximum in (20) can be found by considering only the biorder extensions that have a minimal symmetric

difference with R , and these are of course the minimal biorder extensions of R , which can be generated by the methods of Chapter 4.

Thus, with this restriction in the model, the practical problem consists again of inspecting k -combinations of minimal biorder extensions. This might be feasible in some practical cases, and we might obtain maximum likelihood estimates for approximations in k dimensions. Barring for the moment a number of important statistical issues, like how "stable" (in some sense yet to be defined) such an estimate from a discrete set will be, we must recognize that it is only an estimate for T , not for a particular set of biorders whose intersection is T . Thus, we are still confronted with the uniqueness problem: there may be (and in general there will be) multiple collections of biorders whose intersection is T and the probabilistic model considered here cannot give any clue as to which of these collections to choose. This is another consequence of the fact that the error component of this model, as given in (15), is defined solely in terms of the intersection T : this implies that different collections of biorders with identical intersection are indistinguishable in the model. To obtain a really interesting probabilistic model that estimates the best collection of biorders representing some data matrix, we would have to find a way of defining a sensible and testable error model that relates (the likelihood of) observed events to the relative positions of items and persons on the separate dimensions.

References

- Bouchet, A. (1971). *Etude combinatoire des ordonnés finis. Applications*. Thèse, Université Scientifique et Médicale, Grenoble.
- Cogis, O. (1980). *La dimension Ferrers des graphes orientés*. Thèse, Université Pierre et Marie Curie, Paris.
- Cogis, O. (1982). On the Ferrers dimension of a digraph. *Discrete Mathematics*, **38**, 47-52.
- Coombs, C. H. (1964). *A Theory of Data*. New York: Wiley.
- Coombs C. H. & Kao, R. C. (1955). *Nonmetric factor analysis*. Engineering Research Bulletin No. 38. Ann Arbor: Univ. of Michigan Press.
- Doignon J.-P. & Falmagne J.-C. (1984). Matching relations and the dimensional structure of social choices. *Mathematical Social Sciences*, **7**, 211-229.
- Doignon J.-P., Ducamp, A. & Falmagne, J.-C. (1984). On realizable biorders and the biorder dimension of a relation. *Journal of Mathematical Psychology*, **28**, 73-109.
- Dushnik, B. & Miller, E. W. (1941). Partially ordered sets. *American Journal of Mathematics*, **63**, 600-610.
- Maurer S. B., Rabinovitch, I. & Trotter W. T., Jr. (1980). Large minimal realizers of a partial order, II. *Discrete Mathematics*, **31**, 297-313.
- Trotter, W. T., Jr. (1983). Graphs and partially ordered sets. In Beineke, L. W. & Wilson R. J. (Eds.), *Selected Topics in Graph Theory 2*. London/New York: Academic Press.

Part II :

Knowledge Spaces

THE BASIC THEORY OF KNOWLEDGE SPACES

1. From Guttman scales to knowledge structure

At the basis of the investigations in the following chapters is the problem of developing well-defined, efficient procedures for assessing the knowledge of an individual with respect to a specified domain of information. The need for such procedures is apparent in view of the development of more advanced computerized instruction systems. We assume that the domain of knowledge in question can be represented as a fixed, finite set X of *problems* (or: *items*) for which a binary response is possible: correct or incorrect. The idea is to make the assessment efficient by using the amount of structure that is present in this collection of items. Usually, at least, the items will not all be independent. In very special circumstances it may even be the case that the order in which they can be acquired is completely fixed: we can label the elements of X as x_1, x_2, \dots, x_n in such a way that any student learning the material in X has to master x_1 before she can attack x_2 , x_2 before x_3 , and so on. In such a case, the population of students we consider and the collection X of items are clearly Guttman scalable: they can all be given a position on one dimension such that a student is either below x_1 , or above x_n , or between x_{i-1} and x_i for some $i = 2, \dots, n$.

However, the situation where there is just one possible order of mastering the items in X is very exceptional. Even in a field like elementary arithmetic, which may be regarded as highly structured, we can easily come up with examples where this order is not fully determined. Consider, for instance, a problem in X dealing with subtraction (of 3-digit numbers, with borrowing, say) and one involving multiplication of two single digit numbers. Then we can imagine that there are students solving one and failing the other, either way. This would imply that we need at least two Guttman scales, one in which the subtraction problem precedes ("is lower than") the multiplication problem and one in which it is the other way around. These two orderings of the items divide the population of students into two classes, depending on which of the two items is mastered first. In one group we may observe that subtraction is solved, while multiplication is failed, in the other group the reverse pattern may appear. If a student solves both items or if he fails both, we

cannot tell from which group he was drawn.

This simple example prompts two remarks. First, although we keep talking about Guttman “scales” it may be clear that we are here not interested at all in obtaining actual “scale values” for the items (which values would be determined only up to a monotonic transformation anyway). As before, we are only concerned with the possible orderings of the items in a population of students; as discussed in the General Introduction, this ordinal aspect is the essence of the notion of a Guttman scale.

As a second remark, we want to emphasize the wording “we can imagine” in the above example. The possibility of two different orders is there because *a priori*, as outsiders, we can find no compelling *logical* reasons why one of the two orders would be dictated. However, whether or not both orders appear is an *empirical* matter. It may well be, for instance, that *in practice* the multiplication problem is only mastered after the subtraction problem. In such a case, an experienced teacher would infer from a student’s failure on the subtraction that this student would also fail the multiplication and he would not ask her that problem. It must be stressed that our concern is with the practical situation. Our goal is not some kind of “cognitive analysis” of the domain X , but rather a factual description of the possible orderings of the items in the field, conditional on some chosen target population of students. This description will constitute a key parameter of our assessment procedures. We can think of two sources for arriving at such a description: information from experts in the field, here experienced teachers and tutors, and, when available, extensive empirical data from a sample of student from our population. The dependence of the possible orderings on the chosen population is something to keep in mind: switching from one population to another may certainly change the relative frequency of the classes induced by the different orderings and if this involves a change from a positive to a zero proportion or vice versa, the collection of possible orderings itself has changed.

Returning to our example of two orderings, we see that it readily generalizes to the case of an arbitrary number of Guttman scales: we assume that a specified domain of knowledge is characterized by the collection of possible orderings of its items. More formally, we endow X with a family G of weak orders on X . We consider the general case of weak orders, because we want to allow for equivalent items in the different orderings. The situation where all weak orders are linear orders (no equivalent items) will appear to be an interesting special case; we will come back to this issue in the next section. This family G places restrictions on which subsets of items a person in the population under consideration can have mastered. After all, such a person is supposed to have a position on one of the Guttman scales created by the weak orders on the items, which means he knows all

the items that are lower on this scale and he fails all the problems that are higher. (We restrict ourselves here to the deterministic case, where knowing implies giving a correct answer and not knowing giving an incorrect answer; the possibility of making careless errors and lucky guesses will come in later on.) If $W \in \mathbf{G}$ is the order that is valid for this person, then we know that a correct answer to item x would imply a correct answer to any item y for which $y W x$ ("y precedes x in the order W") holds. With a finite number of items there must be a last item the student has acquired and thus, with the usual notation

$$Wx = \{y \in X : y W x\},$$

it must be the case that the student knows nothing or there is some $x \in X$ such that the student knows exactly the subset Wx of items.

In the framework of knowledge assessment procedures, this subset of items that an individual is capable of solving is obviously of central importance. It is called the *knowledge state* of the individual. The collection of all possible knowledge states is called a *knowledge structure* on X (this collection depends in principle on the intended population of students). Thus, a knowledge structure on X is a family of subsets of X and in the preceding paragraph we saw how such a family is derived from a collection of weak orders on X representing the possible orderings of the items. If \mathbf{G} is this collection of weak orders, then the corresponding knowledge structure \mathbf{K} is given by

$$\mathbf{K} = \{Wx : W \in \mathbf{G}, x \in X\} \cup \{\emptyset\}. \quad (1)$$

Note that thus derived knowledge structures have one particular property: they always contain the null state \emptyset and the perfect state X . This is, however, the only thing special to \mathbf{K} in (1): for any family \mathbf{F} of subsets of X , the knowledge structure $\mathbf{F} \cup \{\emptyset, X\}$ can be obtained by (1) from a collection of weak orders. Including the null and perfect states in any knowledge structure seems reasonable: we always want to cover the case of students knowing nothing and students knowing everything in the domain of choice.

Obviously, not every combination of $W \in \mathbf{G}$ and $x \in X$ yields a distinct knowledge state in (1). That is to say, for $K \in \mathbf{K}$ we may have $K = W_1 x = W_2 y$ with $W_1 \neq W_2$ and/or $x \neq y$. Consequently, from knowing a student's knowledge state we cannot infer on which Guttman scale she is. (Indeed, a student knowing nothing or one knowing all of the items may be on any Guttman scale.) The only thing we know (and that is of real interest to us) is which of the items she has already mastered and which not yet. More generally, given a knowledge structure, the collection of weak orders in (1) is not uniquely determined. The structure \mathbf{K} is certainly produced by the family of all weak orders W such that $\{Wx : x \in X\} \subseteq \mathbf{K}$,

but various subfamilies of this collection may suffice to generate \mathbf{K} . Of special interest in this respect is the subfamily of *minimal* elements in this collection, since they correspond to the “learning paths” in \mathbf{K} . For weak orders, $W_1 \subseteq W_2$ means that W_1 is a *refinement* of W_2 ; it must resolve equivalence classes of W_2 into a number of ordered smaller classes; thus, minimal elements in a collection of weak orders are the ones with the smallest equivalence classes.

To see how these are detectable from the knowledge structure, we consider “chains” in the power set of X . A family \mathbf{C} of subsets of X is called a *chain* if it is completely ordered for the inclusion: for any elements A, B of \mathbf{C} , $A \subseteq B$ or $B \subseteq A$ holds. The chain \mathbf{C} contained in a family of subsets \mathbf{K} is *maximal* in \mathbf{K} if no other element of \mathbf{K} can be added to \mathbf{C} without destroying the chain property. Such a chain may be identified with a *learning path* in \mathbf{K} , a sequence of steps in which the various items are acquired, bringing a person over time from the null state \emptyset into the perfect state X . Any chain \mathbf{C} in 2^X induces a weak order $W_{\mathbf{C}}$ on X by the definition

$$xW_{\mathbf{C}}y \text{ iff } (y \in A \text{ implies } x \in A \text{ for all } A \in \mathbf{C}) \quad (2)$$

and it is readily seen that in this correspondence inclusions are reversed: adding a subset to a chain \mathbf{C} amounts to an extra restriction on the pairs (x, y) in $W_{\mathbf{C}}$, so the corresponding weak order will be smaller. Thus, maximal chains generate by (2) minimal (most refined) weak orders. The collection \mathbf{G} of weak orders on X that correspond to the learning paths in \mathbf{K} ,

$$\mathbf{G} = \{W_{\mathbf{C}} : \mathbf{C} \text{ is a maximal chain in } \mathbf{K}\}, \quad (3)$$

where $W_{\mathbf{C}}$ is defined in (2), does indeed generate \mathbf{K} by (1).

To sum up, starting from the conceptualization of a domain of knowledge by a finite set consisting of all the items or problems in this domain, we have two ways of describing the structure of such a field: by collecting the possible orders in which the various items may be acquired or by collecting the possible knowledge states in the so-called knowledge structure of the field. These two representations are closely related: any collection of weak orders induces a unique knowledge structure by (1) and any knowledge structure (including \emptyset and X) can be represented by the collection of weak orders, given by (3), which correspond to its maximal chains.

1.1 Example.

Let us illustrate the introduced notions on a miniature example, where the domain of knowledge X consists of just 5 items: $X = \{a, b, c, d, e\}$. The simple example of this subsection will be used throughout this chapter for illustrative purposes. Suppose that there are 10 possible orders in which these items can be acquired. More specifically, let the collection \mathbf{G} consist of the following linear orders:

$$\begin{aligned}
 L_1: a < b < d < c < e; & L_6: c < a < b < d < e; \\
 L_2: a < b < c < d < e; & L_7: c < a < b < e < d; \\
 L_3: a < b < c < e < d; & L_8: c < b < a < d < e; \\
 L_4: a < c < b < d < e; & L_9: c < b < a < e < d; \\
 L_5: a < c < b < e < d; & L_{10}: c < b < d < a < e.
 \end{aligned}
 \tag{4}$$

(Note that we have here the special case where all orders are linear orders: no equivalent items in any ranking.) As we have described, such a collection induces knowledge states Lx with L in G and x in X . For instance, L_1 generates the knowledge states;

$$L_1 a = \{a\}, L_1 b = \{a, b\}, L_1 d = \{a, b, d\}, L_1 c = \{a, b, c, d\}, L_1 e = X.$$

Also, $\{a, b, c\} = L_8 a = L_5 b = L_2 c$ is a state. It can easily be checked that the complete knowledge structure K deriving from G according to (1) has 12 states and is given by

$$\begin{aligned}
 K = \{ \emptyset, \{a\}, \{c\}, \{a, b\}, \{a, c\}, \{b, c\}, \\
 \{a, b, c\}, \{a, b, d\}, \{b, c, d\}, \{a, b, c, d\}, \{a, b, c, e\}, X \}.
 \end{aligned}
 \tag{5}$$

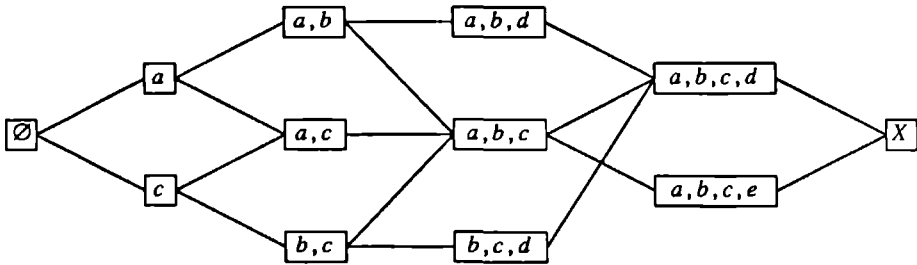


Figure 1. Graph of the knowledge structure of Eq. (5).

Representing this knowledge structure graphically, as is done in Figure 1, clearly shows the maximal chains in K . In this graph, two states are joined by lines if and only if the state to the left is a subset of the state to the right. (In technical terms, Fig. 1 presents the *Hasse diagram* of the knowledge structure K , partially ordered by the inclusion relation.) Any maximal chain in K corresponds to a path from the null to the perfect state. For instance, we can detect in Fig. 1 the maximal chain

$$\emptyset \subseteq \{c\} \subseteq \{c, b\} \subseteq \{c, b, a\} \subseteq \{c, b, a, d\} \subseteq \{c, b, a, d, e\}.$$

Along this path, the items are mastered in the order $c < b < a < d < e$, which

coincides with L_8 in (5). Indeed, there are 10 different paths from \emptyset to X in Fig. 1 and the correspondence with the elements of G as defined by (3) is easily established. It appears that the collection G in (4) was constructed with the knowledge structure K of (5) in mind: (4) collects exactly all learning paths in this structure. Various other collections of weak orders would have produced the same knowledge structure. For example, the collection G' where L_1 is replaced by

$$W_1: \{a, b, d\} < \{d, e\}$$

(elements between braces denote equivalent items in the weak order), but this means replacing a finer ordering by a coarser one (clearly, $L_1 \subseteq W_1$). Or, the collection $G'' = \{L_1, L_5, L_{10}\}$, but this entails the assumption that only three of the 10 possible learning paths in Fig. 1 are actually in use. In practical situations, only the knowledge structure can be observed directly (at least in principle, disregarding for the moment the question of noisy observations), so this kind of assumptions are in general to be avoided. In this sense, the collection in (4) (generally: the collection of weak orders defined by (3)) gives the most refined and complete representation of the knowledge structure K in (5).

2. Discriminating and well graded knowledge structures

Since we are dealing with collections of weak orders, there may be orders in which two distinct items are equivalent, that is, along that ordering one is mastered at the same time as the other. For instance, if we have the two orders

$$W_1: a < \{b, c\} \quad \text{and} \quad W_2: b < c < a$$

on the three item set $\{a, b, c\}$, then b and c are equivalent in W_1 , but not in W_2 . The corresponding knowledge structure in this case is

$$\{\emptyset, \{a\}, \{b\}, \{b, c\}, \{a, b, c\}\}. \quad (6)$$

Now consider the special case where there are two items x and y in X that are equivalent in all orders in a collection G . It is easy to check that then, and only then, a state of the knowledge structure K defined from G by (1) contains x if and only if it contains y . In other words, the items x and y are *indistinguishable* in K . Any student who can solve one can solve the other. For instance, if we replace W_2 in the above example by the extension

$$W_2': \{b, c\} < a,$$

then b and c are equivalent in both orders and indeed they are indistinguishable in

the corresponding knowledge structure, which is now given by

$$\{ \emptyset, \{a\}, \{b, c\}, \{a, b, c\} \}. \quad (7)$$

If we take the knowledge structure seriously as a complete description of the possible knowledge states, a natural interpretation is that indistinguishable items in a structure test in fact one and the same abstract notion. For example, we can think of many problems that are all equivalent versions of the notion “adding of two 2-digit numbers with repeated carrying” in elementary arithmetic. Accordingly, we define technically a *notion* in \mathbf{K} as a maximal collection of indistinguishable items in \mathbf{K} . Thus, in (7) the knowledge structure has the two notions $\{a\}$ and $\{b, c\}$.

We call a knowledge structure *discriminating* when there are no indistinguishable elements: any notion consists of one single item. The structure in (6) is discriminating, while that in (7) is not. Notice, however, that it is possible to consider any equivalence class of items, i.e., any notion, as a single element and interpret a knowledge structure as defined on these notions, instead of the separate items. In this way, a knowledge structure becomes discriminating by definition. The structure in (7) defined on the elements $\{a\}$ and $\{b, c\}$ is discriminating: it corresponds to the full power set on these elements. Since being discriminating is just a technical requirement that can always be met by applying this standard reduction operation, we may, whenever convenient, simply assume that this operation has been carried out.

An interesting special class of knowledge structures is obtained by restricting attention to collections of linear orders, instead of weak orders in general. A knowledge structure \mathbf{K} that can be derived, according to (1), from a collection \mathbf{G} of linear orders, is called *well graded*. This notion of well-gradedness can be given a number of equivalent reformulations. A linear order on X corresponds, by definition (2), to a chain that is maximal in 2^X , that is, an increasing sequence of subsets of X , starting with \emptyset , terminating with X , and such that any set except the first contains exactly one more item than its predecessor. If such an $|X| + 1$ element chain is contained in \mathbf{K} it is called a *gradation* in \mathbf{K} and \mathbf{K} is well graded if it is the union of its gradations, that is, if any state in \mathbf{K} is contained in at least one gradation. Still equivalently, \mathbf{K} is well graded if for any $K \in \mathbf{K} - \{ \emptyset, X \}$ there are $x \in K$ such that $K - \{x\} \in \mathbf{K}$ and $y \notin K$ such that $K + \{y\} \in \mathbf{K}$.

The intuitive idea is that in a well graded knowledge structure we have only to deal with learning paths in which the items are acquired one by one. The transition from one state to the next is always obtained by adding one item; we never have to make a jump of two or more. This seems to make sense pedagogically, at least when we replace “item” with “notion”. Obviously, a knowledge structure that is not discriminating cannot be well graded, but in such a case the reduced structure,

defined on the notions, may or may not be well graded, and it is here where well-gradedness seems to be an interesting and reasonable additional condition on knowledge structures.

For example, neither the knowledge structure in (6), nor the one in (7) is well graded. In (6), the state $\{a\}$ is not contained in a gradation (from $\{a\}$ we have to make a 2-item jump to the next state, $\{a, b, c\}$). In the knowledge structure of (7) there is no gradation at all, since this structure is not discriminating. However, the reduced structure on the two notions $\{a\}$ and $\{b, c\}$ is well graded. Also, the knowledge structure \mathbf{K} in (5) is well graded, since it is derived from the linear orders in (4). Each of these produces one of the ten gradations in \mathbf{K} .

3. Surmise relations

Its knowledge structure characterizes a field of information X in a very concrete way and this representation will actually be used in the assessment procedures. However, since for only a moderate number of items (say, 30) such a structure can easily contain hundreds, if not thousands of states, the description of the cognitive organization of X by the enumeration of all subsets that are possible knowledge states is in general not very enlightening. In search for a more succinct representation we may turn to the collection of different orderings of the items that is implied by the knowledge structure.

Let \mathbf{K} be a knowledge structure on X and let \mathbf{G} be the collection of weak orders corresponding to \mathbf{K} according to (3). Now suppose that for some $x, y \in X$ it is the case that x precedes y in all elements of \mathbf{G} . Then it clearly must hold that any state in \mathbf{K} containing y must also contain x . On the other hand, if the assumption is not true, that is, if there is some order in \mathbf{G} in which x does not precede y , then y precedes x (a weak order is complete) and thus this weak order induces a state of \mathbf{K} containing y while not containing x . Since the fact that x precedes y in all orders of \mathbf{G} is formally represented as $xy \in \bigcap \mathbf{G}$, we have established the following equivalence:

$$xy \in \bigcap \mathbf{G} \quad \text{iff} \quad (y \in K \text{ implies } x \in K \text{ for all } K \in \mathbf{K}), \quad (8)$$

when \mathbf{K} and \mathbf{G} are related by Eqs. (1) and (3).

As we have seen in the General Introduction, whenever \mathbf{G} is a collection of weak orders, $\bigcap \mathbf{G}$ is a quasi order. Thus, any knowledge structure \mathbf{K} defines a collection of weak orders \mathbf{G} and through this a quasi order $Q = \bigcap \mathbf{G}$, which satisfies (8). This quasi order associated with \mathbf{K} is called the *surmise relation* of \mathbf{K} , since, by (8), xQy may be interpreted as "if a student can solve y , it may be surmised that

this same student can also solve x ". In the sequel we will abbreviate this as " x may be surmised from y ". Note that $\cap G$ is antisymmetric whenever no two distinct items are equivalent in all orders in G , that is, whenever \mathbf{K} is discriminating. Consequently, in the same way that we may assume a knowledge structure to be discriminating (by considering the notions, see the previous section), we may also assume that the surmise relation is in fact a partial order.

The surmise relation gives such an easily interpretable and concise description of properties of \mathbf{K} that it is tempting to investigate if it determines the knowledge structure uniquely, that is, if we can recover \mathbf{K} from the knowledge of its surmise relation Q . There is a classical mathematical result telling us that this is possible only in special cases:

3.1. Theorem. (Birkhoff, 1937.)

For any set X , the formula

$$xQy \text{ iff } (y \in K \text{ implies } x \in K \text{ for all } K \in \mathbf{K}) \quad (9)$$

defines a 1-1 correspondence between the set of all quasi orders Q on X and the set of all families \mathbf{K} of subsets of X that are closed under union and intersection.

Recall from the General Introduction that a family \mathbf{K} of subsets of X is closed under union if $(\cup \mathbf{K}) \in \mathbf{K}$, that is, any union of states is again a state. Similarly for closure under intersection. If we restrict ourselves to a finite set of items X (note that Birkhoff's Theorem is also valid in the infinite case), it is sufficient to check for pairs of states: \mathbf{K} is closed under union (intersection) if and only if for any $K_1, K_2 \in \mathbf{K}$ we have also $K_1 \cup K_2 \in \mathbf{K}$ ($K_1 \cap K_2 \in \mathbf{K}$).

Note that (9) defines a surmise relation Q for any knowledge structure \mathbf{K} . If we then use (9) to obtain a knowledge structure \mathbf{K}^* from Q (by allowing any subset K as a state for which $y \in K$ and $x Q y$ imply $x \in K$), we know by Theorem 3.1 that \mathbf{K}^* is closed under union and intersection. If the structure \mathbf{K} we started with does not have this property, it follows that $\mathbf{K}^* \neq \mathbf{K}$. It turns out that in this situation always $\mathbf{K}^* \supseteq \mathbf{K}$; more specifically: \mathbf{K}^* is the closure of \mathbf{K} under union and intersection, that is, it contains all the states in \mathbf{K} , plus all unions and intersections of states in \mathbf{K} . See Monjardet (1970) (also Doignon and Falmagne, 1985) for a reformulation of Birkhoff's result in terms of such *closure operators*.

Note finally that for a partial order P the knowledge structure derived by (9) from P and the one derived by (1) from the collection of linear order extensions of P coincide: these are two methods for obtaining the largest structure having P as surmise relation. In particular, deriving from linear orders, such a structure is well graded by the definition of the previous section. Thus, by Theorem 3.1, any

knowledge structure closed under union and intersection is well graded, after forming the equivalence classes of indistinguishable items.

3.2. Example.

In the situation of Example 1.1, with G and K given by (4) and (5), respectively, it can be seen that, for instance, b precedes d in all orders L_1 to L_{10} . Consequently, $bd \in \cap G$ and indeed any state of K in (5) containing d also contains b . Since all members of G are linear orders, the relation $P = \cap G$ is a partial order and it can easily be checked that P , except for the loops aa to ee that are trivially there, consists of the pairs ae , bd , be and ce . That is, P is the partial order whose Hasse diagram is given in Figure 2. In such a diagram, a pair xy is in the partial order whenever there is a path from x to y by ascending lines.

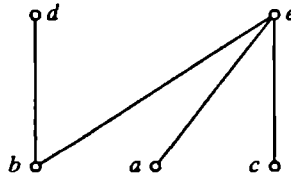


Figure 2. Hasse diagram of the partial order $P = \cap G$, where G is given in (4).

Now we may ask, does this partial order P describe the knowledge structure K completely? In other words, if we use P to generate a knowledge structure by the formula (9), would we get our original K back? The answer to these questions must be negative, since any knowledge structure constructed by (9) is necessarily closed under union and intersection, while our K of Eq. (5) is not. Indeed, $\{a, b\} \cap \{b, c\} = \{b\} \notin K$ and also $\{a, b, d\} \cap \{b, c, d\} = \{b, d\} \notin K$. According to (9), a subset A of X is a state of the knowledge structure induced by P whenever $y \in A$ and xPy together imply $x \in A$. In terms of Fig. 2, a state is any set containing with an item all items that can be reached from this item by descending lines. (Such sets are called the “lower sets” of the partial order P on X .) Clearly, both $A = \{b\}$ and $A = \{b, d\}$ satisfy this criterion and it can easily be checked that P generates the structure $K \cup \{\{b\}, \{b, d\}\}$, which is indeed closed under union and intersection.

4. Surmise mappings

Representing knowledge structures by quasi orders seems an attractive option because of its simplicity and economy, but it has one important drawback: any

problem $x \in X$ has one unique set of prerequisites. Here, “a set of prerequisites” for x means a minimal set of problems that must have been mastered before problem x can be tackled. In terms of the surmise relation Q , this unique set of prerequisites for x is given by $Qx - \{x\}$, the set of items $y \neq x$ that can be surmised from x . Although this condition may describe an interesting special case, it is too strict in general: it is not uncommon that there are different ways of solving a problem, which rely on different sets of notions. A knowledge structure in which an item has multiple sets of prerequisites is not fully described by a quasi order and, according to Theorem 3.1, it will not be closed under union and intersection.

4.1. Example.

In 3.2 we saw that the example knowledge structure K of Eq. 5 (Fig. 1) is not represented by its surmise relation P of Fig. 2. The partial order P corresponds to the following restrictions:

- (i) from d we can surmise b ;
- (ii) from e we can surmise a AND b AND c .

According to the partial order we cannot surmise anything from a , b and c . In the knowledge structure K , however, there is some restriction on b , simply because the singleton $\{b\}$ is not a state. Inspecting the states in K containing b , we see that while there is not a single other item contained in all these states, it is true that in each of them at least one of the items a and c is present. In other words:

- (iii) from b we can surmise a OR c .

Such an inference in disjunctive form cannot be captured by a surmise relation; it can only deal with conjunctions (cf. (ii) above). Similarly, (i) above does not represent all restrictions on d in K since it allows the state $\{b, d\}$ that is absent from K . An inspection of the states containing d shows indeed an additional restriction, again involving a disjunction. In sum we have:

from d we can surmise b AND (a OR c),

which we can write in the equivalent form

- (iv) from d we can surmise (b AND a) OR (b AND c).

It is clear that once disjunctions enter in the description, we are dealing with multiple sets of prerequisites. Thus, (iii) means that b has the two sets of prerequisites $\{a\}$ and $\{c\}$, and according to (iv) for d we have the two sets $\{a, b\}$ and $\{b, c\}$.

The preceding Example shows that to fully describe the knowledge structure of Example 1.1 we need something like an “AND/OR graph” (a tool in use in some parts of artificial intelligence, see e.g. Nilsson, 1971), while, in these terms, a surmise relation only represents an “AND graph”. More precisely, we want to define a generalization of the notion of a surmise relation where for each item there

is not just one collection of items that can be surmised from it, but a number of alternative collections. That is, with any $x \in X$ we want to associate a non-empty collection $\sigma(x)$ of subsets of X representing the idea that from the mastery of x we can surmise the mastery of *all* items in *at least one* of the elements of $\sigma(x)$. (The “all” refers to the AND-component, the “at least one” to the OR-component in this representation.) Such a mapping σ is called a *surmise mapping* on X (*space-like surmise mapping* in Doignon and Falmagne, 1985) and the elements of $\sigma(x)$, which are subsets of X , are called the *clauses* for x . From its interpretation we can deduce several properties we would like such a surmise mapping to have. For one thing, since from x we can always trivially surmise x itself, it must be the case that any clause for x contains x . In formula:

$$C \in \sigma(x) \text{ implies } x \in C. \quad (10)$$

Thus, any clause for x consists of x plus a possible set of prerequisites for x ; or, any clause for x collects the items in some minimal learning path leading to the mastery of x . Such a path for x defines a (minimal) subpath for any of the items on this path, and this implies that a clause for x must include at least one clause for each of its elements:

$$y \in C \ \& \ C \in \sigma(x) \text{ implies } C' \subseteq C \text{ for some } C' \in \sigma(y). \quad (11)$$

Finally, since a set of prerequisites is a minimal set of problems that must have been mastered before x can be tackled, it is clear that it does not make sense to have two clauses for x where one is included in the other. Or, in terms of surmising: if $A_1 \subseteq A_2 \subseteq X$, then to surmise A_1 OR A_2 is equivalent to surmising A_1 . So we have

$$C, C' \in \sigma(x) \ \& \ C \subseteq C' \text{ implies } C = C'. \quad (12)$$

To summarize, a surmise mapping on X is any mapping from X into the power set (minus \emptyset) of the power set of X that satisfies (10), (11) and (12).

At this moment we might want to check that with this definition a surmise mapping is indeed a generalization of a surmise relation. Such a surmise relation Q corresponds to a single clause Qx for any x , so we must show that $\sigma(x) = \{Qx\}$ defines a surmise mapping. This follows easily: (12) is trivial and (10) and (11) amount to the reflexive and transitive properties, respectively, of a surmise relation (i.e., quasi order).

As with a surmise relation (Formula (8) or (9)), any knowledge structure \mathbf{K} on X induces a surmise mapping. This is obtained by letting the clauses for x consist of the minimal states in \mathbf{K} containing the item x . (Note that with finite X there must always be such states.) Conditions (10) and (12) are trivial, and, for (11), a minimal state for x containing y includes, by definition, a *minimal* state containing y .

On the other hand, any surmise mapping σ on X induces in a natural way a knowledge structure on X . A subset K of X is a state in this structure whenever it contains some clause of $\sigma(x)$ for any $x \in K$. (A state contains with each element at least one set of prerequisites for that element.) This amounts to saying that the knowledge structure induced by σ consists precisely of the unions of clauses of σ .

The interesting question is again: for which class of knowledge structures are the above correspondences one-to-one? Or, equivalently, which knowledge structures are characterized by their surmise mapping? On the one hand this class must contain the structures closed under union and intersection, by Theorem 3.1 and the fact that the surmise mapping generalizes the surmise relation; on the other hand, we have seen in the preceding paragraph that any knowledge structure induced by a surmise mapping is still closed under union. In this light, the following theorem, proved by Doignon and Falmagne (1985), should not come as a complete surprise.

4.2. Theorem. (Doignon and Falmagne, 1985.)

For any finite set X , the formula

$$\sigma(x) = \hat{K}_x, \quad (13)$$

where \hat{K}_x denotes the subcollection of minimal states in K containing x , defines a 1-1 correspondence between the set of all surmise mappings σ on X and the set of all families K of subsets of X that are closed under union.

For the finite case, this theorem provides a generalization of Birkhoff's classical result (Theorem 3.1). Apparently, the relaxation from surmise relation to surmise mapping, that is, allowing items to have multiple sets of prerequisites, is obtained by dropping the requirement of closure under intersection for a knowledge structure. Such a knowledge structure, which is still closed under union, is christened a *knowledge space* by Doignon and Falmagne (1985).

4.3. Example.

From the "AND/OR graph" description in 3.1 we can easily deduce what the surmise mapping σ for the knowledge structure K of Example 1.1 should be (cf. (ii), (iii) and (iv) in 3.1):

$$\begin{aligned} \sigma(a) &= \{ \{a\} \} & \sigma(d) &= \{ \{a, b, d\}, \{b, c, d\} \} \\ \sigma(b) &= \{ \{a, b\}, \{b, c\} \} & \sigma(e) &= \{ \{a, b, c, e\} \} \\ \sigma(c) &= \{ \{c\} \} \end{aligned}$$

In Fig. 1 it is easily checked that, for $x = a, \dots, e$, $\sigma(x)$ does indeed consist of the minimal states of K containing x . If, conversely, we start with the above surmise mapping σ and construct the induced knowledge structure, that is, the collection of

all unions of clauses in σ , we get exactly our original K back. (Note that the null state is obtained as the union of *zero* clauses in σ .) This is as predicted by Theorem 4.2, since the knowledge structure K in Fig. 1 is closed under union – it is a knowledge *space* – and thus completely determined by its surmise mapping.

For knowledge spaces it is no longer true that they are well graded whenever they are discriminating, as was the case with structures closed under union and intersection. For instance, the knowledge structure

$$\{ \emptyset, \{a, b\}, \{a, c\}, \{b, c\}, \{a, b, c\} \}$$

on $\{a, b, c\}$ is closed under union and discriminating, but not well graded. It is true that the criterion for well-gradedness simplifies in the case of a space: a knowledge space is well graded if and only if any non-null state K has an element x such that $K - \{x\}$ is a state.

To summarize this section, we noted that modeling knowledge structures by surmise relations (quasi orders) had the undesirable consequence that each item was forced to have a unique set of prerequisites. This restriction was overcome by generalizing the notion of a surmise relation to that of a surmise mapping; the corresponding class of knowledge structures consists of the families closed under union, and a knowledge structure in this class is called a knowledge space.

5. A sketch of two assessment procedures

In the following chapters we will be dealing with the question of how we can construct a knowledge space for a particular domain in practice. Let us illustrate here why we are interested in this question in the first place. We want to show how the representation of a field of information by a knowledge structure (or, more specifically, space) leads to the design of efficient knowledge assessment procedures. The results of this section do not play any role in the sequel; they are presented here only to provide the motivation for what follows.

Given the collection of possible knowledge states and a student picked from the target population, the general problem for a knowledge assessment procedure is to determine the knowledge state of this particular student by asking him a minimal number of the problems of the domain. A general scheme for doing this is presented in Figure 3, which has been adapted from Falmagne, Koppen, Villano, Johannesen and Doignon (1989). At the start of each trial of the procedure, the information obtained from the student's responses to previously posed problems is summarized

by a *plausibility function*, assigning plausibility values to the various states. These values are used by a *questioning rule* to determine a subset of most informative problems to ask on this trial (“most informative” according to some criterion) and one of these problems is chosen at random. The student’s response to this problem is based on his knowledge state through a *response rule*. This response is then processed by an *updating rule* to recompute the plausibility function which will be the start of the next trial. In the absence of any information the procedure starts with all states equally plausible and it terminates when, according to some criterion, the plausibility function provides enough evidence for singling out the student’s state. We will first describe a straightforward application of this scheme in a deterministic framework; next follow two variations on this theme that have been worked out by Falmagne and Doignon (1988a,b) for the practical case where the information may be noisy.

5.1. The deterministic case.

In the simplest case we assume that the response to a problem is completely determined by the student’s knowledge state: a correct answer is obtained if and only if the problem is in this state. Then there is a straightforward procedure of uncovering this state. Note that any problem x divides any collection F of states in two: the subcollection F_x of states in F containing x and the subcollection $F_{\bar{x}}$ of states not containing x . In some cases, this partitioning may be trivial (one of the two classes being empty), but as long as F consists of more than one state there is always an item x that precludes this situation, i.e., such that both F_x and $F_{\bar{x}}$ are strictly smaller than F .

These observations lead to the following procedure. At the start of trial n , we have a collection $M^{(n)}$ of states that are still plausible at this moment (the “marked” states). This corresponds to a binary plausibility function with, for instance, a “1” value for a plausible and a “0” value for a non-plausible state. At the start, all states are plausible, that is, $M^{(1)} = K$, the knowledge structure for the domain under investigation. While $M^{(n)}$ contains more than one state, we choose an item x that strictly partitions $M^{(n)}$ and record the student’s answer to this problem. Under the assumption of error-free conditions, a correct response of the student entails that his state must contain x , and, accordingly, we set $M^{(n+1)} = M_x^{(n)}$; if the response is incorrect, we can draw the opposite conclusion and we set $M^{(n+1)} = M_{\bar{x}}^{(n)}$. This is the updating rule, which in terms of the binary plausibility function amounts to converting a “1” into a “0” for any state incompatible with the observed response. The procedure terminates when there is only one marked state left (only one “1” entry in the plausibility function).

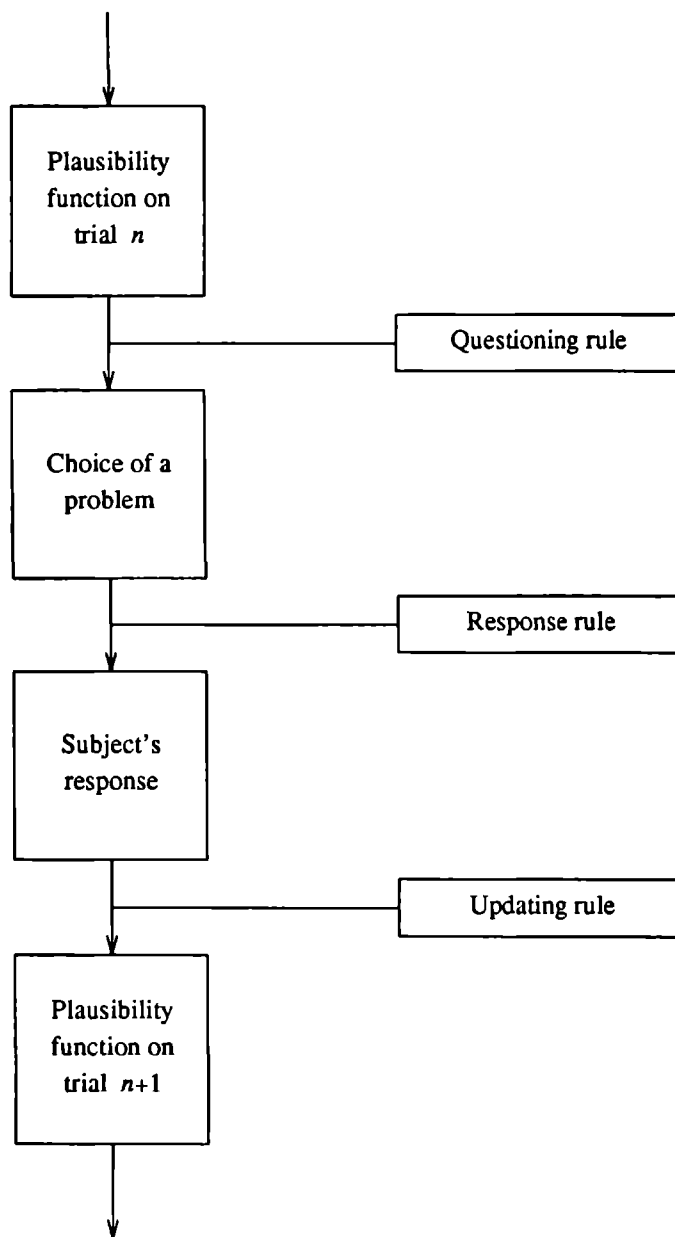


Figure 3. A general scheme for knowledge assessment procedures (adapted from Falmagne, Koppen, Villano, Johannesen and Doignon, 1989).

The maximum number of trials is n (when we have asked all items, we have complete information). This maximum will, of course, always be needed in the trivial case where every subset of items is a state ($|K| = 2^n$), but even with certain small knowledge structures it may, for some states, be the best we can do. (Consider locating a student knowing all n problems in a structure consisting, except for the null and the perfect states, of all $(n-1)$ -subsets.) However, if the domain is structured (the number of states is only a small fraction of the number of subsets), we will, averaged over all states, be able to do considerably better than asking all items. Imagine the special case where at trial k there happens to be an item x that divides the collection $\mathbf{M}^{(k)}$ evenly. Then, no matter what the response of the student to x will be, we will reduce the number of marked states by one half by choosing x . If we could find such an item at each trial, we would be able to determine the student's state in $\log_2 \mathbf{K} \leq n$ steps. This is the optimal over-all bound, and we can try to approximate it as closely as possible by adopting a "split-half" rule. This questioning rule designates as a most informative item at trial n any item x that partitions $\mathbf{M}^{(n)}$ most evenly, that is, one that minimizes (over all items) the function $||\mathbf{M}_x^{(n)}| - |\mathbf{M}_{\bar{x}}^{(n)}||$.

5.2. A discrete Markovian procedure.

It is clear that the above procedure cannot deal with random fluctuations in the student's performance. For instance, if the student would make a careless error on a problem that is in fact in his knowledge state, then, in the above procedure the true state, being inconsistent with the observed response, would be discarded from the list of plausible states, forever. The same would happen when the student, by making a lucky guess, would give a correct answer to a problem that he does not really master. Falmagne and Doignon (1988a,b) developed two elaborations of the deterministic procedure that can handle this kind of noise that will inevitably be present in practical assessment situations. We give only a sketch of the basic ideas; for details the reader is referred to the original papers.

The first of these procedures we are going to describe here, Falmagne and Doignon (1988b), is in fact only a minor variation on the deterministic procedure. (We consider here the special case that is of practical importance, which in the original paper appears in a broader context.) This procedure assumes that the knowledge structure we are dealing with is well graded. First the deterministic procedure is completed, at some point leaving us with just one marked state, say $\mathbf{M}^{(N)} = \{K\}$. Now, because of the presence of noise, we cannot be sure that K is the true state of the student (he may have made a careless error under way, or a lucky guess), but we may assume that the remaining state has much in common with this true state.

So, once there is only one state left, the procedure continues, but under modified questioning and updating rules. We want to compare a single remaining state locally, with states that are closest to it; in a well graded structure the neighboring states of K are of the form $K \pm \{x\}$ (where we assume that the “+” alternative is used for $x \notin K$, the “-” for $x \in K$). The modified questioning rule no longer applies to the collection of marked states, but rather to the single marked state plus its neighboring states: an item is a most informative one if it splits this collection most evenly. Notice, however, that “most evenly” has a very restricted meaning here: the smaller class of the partition will contain one state at most. Indeed, if $x \in K$, all neighboring states will contain x , with the possible exception of one, $K - \{x\}$, while, for $x \notin K$, $K + \{x\}$ is the only possibility for a neighboring state containing x . Thus, the questioning rule changes, in effect, to: choose randomly between the items x such that either $K \neq K - \{x\} \in \mathbf{K}$ or $K \neq K + \{x\} \in \mathbf{K}$. With one marked state, trials amount to a test of the current marked state K against an alternative $K \pm \{x\}$ and the updating rule is modified to follow the result of this test directly. We keep K as marked state if the response to the chosen problem x is consistent with K , otherwise we replace K by the new marked state $K \pm \{x\}$; in both cases we are ready to start a new trial with a single marked state.

5.3. A continuous Markovian procedure.

While in the preceding subsection we dealt with the noise in the situation by supplementing the deterministic procedure with a part in which we can recover from wrongly discarding the true state, Falmagne and Doignon (1988a) describe an alternative solution in which this kind of recovery is not needed since no state is ever really discarded. In this version of an assessment procedure the plausibility function is no longer binary, indicating marked and unmarked states, but now it takes the form of a likelihood function. At the start of trial n , the plausibility of each state K is given by a likelihood $L_K^{(n)} > 0$, such that $\sum_{K \in \mathbf{K}} L_K^{(n)} = 1$. Thus, at every trial we are given a probability distribution over the states that is nowhere zero.

As questioning rule, there is again a kind of half-split rule available. This time it is not a collection of marked sets that is split in two by an item, but rather the total mass of the likelihood function. The likelihood that the student under investigation will solve problem x is represented by $\sum_{x \in K} L_K^{(n)}$, the total of the mass on states containing x and the likelihood he will fail x by $\sum_{x \notin K} L_K^{(n)}$, the mass on states not containing x . A problem may be judged to be most informative when the likelihoods for a correct and an incorrect response are closest (and thus both closest to one half). Thus, the questioning rule selects x as a most informative item if it minimizes (over all items) $|\sum_{x \in K} L_K^{(n)} - \sum_{x \notin K} L_K^{(n)}|$, or, equivalently, $|\sum_{x \in K} L_K^{(n)} - 0.5|$. (Falmagne and Doignon, 1988a, considered also an alternative

questioning rule, in which an item is most informative at trial n if it minimizes the (expected) entropy of the updated likelihood function at trial $n+1$.)

For the updating rule to make sense, it must increase the likelihood of states compatible and decrease the likelihood of states incompatible with the latest observed response. Different ways of doing this are conceivable, but an obvious possibility is a Bayesian rule (in Falmagne and Doignon, 1988a, this is called the *multiplicative* updating rule). Here the posterior (updated) likelihood of a state is proportional to the probability of the observed response given this state times the prior likelihood of the state. In formula, if x is the problem at trial n ,

$$L_K^{(n+1)} \propto \mathbb{P}(\text{response on } x | K) \cdot L_K^{(n)}.$$

The conditional probabilities of the responses are parameters of the assessment procedure. More precisely, with each item x are associated a “careless error” parameter $\beta_x = \mathbb{P}(\text{incorrect on } x | x \in K)$ and a “lucky guess” parameter $\gamma_x = \mathbb{P}(\text{correct on } x | x \notin K)$. Accordingly, there are probabilities $1 - \beta_x$ for a correct response if $x \in K$, and $1 - \gamma_x$ for an incorrect response if $x \notin K$. The $\mathbb{P}(\text{response on } x | K)$ in the above formula is one of these values β_x , $1 - \beta_x$, γ_x or $1 - \gamma_x$, depending on whether the observed response was correct or incorrect and whether $x \in K$ or $x \notin K$. Of course, this updating rule makes only sense if all $\beta_x < 0.5$ and all $\gamma_x < 0.5$.

6. The following chapters

The assessment procedures described in the preceding section are all based on a fixed, predetermined knowledge structure of the domain under investigation. This structure is supposed to give an adequate specification of the possible knowledge states in the domain, so that any student from the target population can be assessed, if not perfectly, then at least to a good approximation. It is natural, then, to wonder how, in practice, we can arrive at such a knowledge structure representation of particular domains. Two sources of information may be available: experts in the field (in our case, experienced teachers and tutors) and, sometimes, extensive empirical data.

We can envision a two stage process for constructing the knowledge structure. A first sketch is obtained from systematically consulting a number of experts. We would like this first sketch to be conservative in the sense that it certainly contains all relevant states, even if this means that a number of superfluous states are still there. In the presence of experimental data, such a version can then be used as a model making predictions regarding observed response patterns in the relevant population.

However, such data are typically noisy and therefore the obtained knowledge structure has to be embedded in a probabilistic model. Then a statistical test of the model structure becomes feasible and possible superfluous states can be removed by formulating appropriate restricted versions of the model and testing these through standard likelihood ratio methods. Falmagne (1989) has developed a probabilistic learning model that deals precisely with these issues of the second stage; in this model it is assumed that the knowledge structure we work with is well graded. This leaves us with the problem of the first stage, of how to construct a knowledge structure from expert opinions, and this will largely be the subject of the following chapters.

In a straightforward approach we would supply an expert teacher with some domain (a collection of problems) and we would ask her to list all possible knowledge states. This approach is not feasible in practice. First, even for a very moderate number of items, say 30, this list may contain many states (a few thousand, say). Second, the concept of a knowledge state, in all its concreteness, is certainly not familiar to experts in the field. Most probably they do not have access to an internal representation of the required list of states; the knowledge structure is only implicitly there. So, we want to apply indirect methods to make this implicit structure explicit and for this purpose we have to rely on alternative representations for knowledge structures.

In Section 3 we presented such an alternative: a class of knowledge structures appeared to allow a characterization by quasi orders (surmise relations) and this representation would be very suitable for questioning experts. Unfortunately, we had to conclude that surmise relations constitute too strict a model and in Section 4 we described the generalization to surmise mappings as a representation for a broader, acceptable class of knowledge structures, the knowledge spaces. Surmise mappings, however, are not easy to use as a base for questioning experts about knowledge spaces. In Chapter 7 we develop a second alternative representation for knowledge spaces, one that is fit for this purpose. We will again be dealing with quasi orders, but this time it will be quasi orders, not on the set of items, but on the collection of subsets of items. Mathematically, the outcome is another generalization of Birkhoff's result (Theorem 3.1). It may in fact be considered as a more direct generalization than Theorem 4.2; it is, for instance, like Birkhoff's Theorem, also valid in the infinite case.

While the representation of Chapter 7 is indeed well suited for questioning experts, a straightforward application would again be impracticable. As indicated, we are dealing here with quasi orders on the subsets of items, so the corresponding tables are of the order 2^n if n is the number of items. This is prohibitive, even for moderately large n . Fortunately, the quasi orders in question enjoy a number of

additional properties and this implies that the above mentioned table contains a lot of redundant information. We can exploit the redundancy of this representation by making inferences on the basis of these extra properties. This saves us many information requests to the expert and generally it will appear that all the information necessary for the construction of the knowledge space that the expert is implicitly consulting is contained in a subtable that is only a tiny fraction of the full table. This search for a most efficient implementation of the new representation of Chapter 7 is the subject of Chapter 8. It culminates in the specification of an explicit algorithm for questioning an expert in a very systematic way and deriving the implied knowledge space from the obtained answers.

The notion of well-gradedness was introduced in Section 2 and it appeared to be a very reasonable extra assumption for a knowledge structure in practice (at least for a structure defined on the notions, i.e. after forming the equivalence classes of indistinguishable items). In the previous section we saw that one of the developed assessment routines even assumed that the underlying knowledge structure was well graded and this assumption is also crucial in the probabilistic model of Falmagne (1989), mentioned above. This prompts the question of how the special case of a well graded knowledge space appears in the two alternative characterizations of knowledge spaces. For the representation of Chapter 7 this is, as yet, unknown (at least, no easily checkable criterion is available). For the representation by surmise mappings, however, we know which extra condition we have to impose to make it equivalent to the concept of a well graded knowledge space. This is discussed in Chapter 9, where it is put in the context of a number of more and more restrictive extra conditions for a surmise mapping, corresponding to smaller and smaller classes of knowledge spaces.

References

- Birkhoff, G. (1937). Rings of sets. *Duke Mathematical Journal*, **3**, 443-454.
- Doignon, J.-P. & Falmagne, J.-C. (1985). Spaces for the assessment of knowledge. *International Journal of Man-Machine Studies*, **23**, 175-196.
- Falmagne, J.-C. (1989). A latent trait theory via a stochastic learning theory for a knowledge space. *Psychometrika*, in press.
- Falmagne, J.-C. & Doignon, J.-P. (1988a). A class of stochastic procedures for the assessment of knowledge. *British Journal of Statistical and Mathematical Psychology*, **41**, 1-23.
- Falmagne, J.-C. & Doignon, J.-P. (1988b). A Markovian procedure for assessing the state of a system. *Journal of Mathematical Psychology*, **32**, 232-258.
- Falmagne, J.-C., Koppen, M., Villano, M., Johannesen, L. & Doignon, J.-P. (1989). Introduction to knowledge spaces: How to build, test and search them. *Psychological Review*, to appear.

- Monjardet, B. (1970). Tresses, fuseaux, préordres et topologies. *Mathématiques et Sciences Humaines*, 30, 11-22.
- Nilsson, N. J. (1971). *Problem-Solving Methods in Artificial Intelligence*. New York: McGraw-Hill.

CHAPTER 7

HOW TO BUILD A KNOWLEDGE SPACE BY QUERYING AN EXPERT

(To appear in *Journal of Mathematical Psychology*)

How to build a knowledge space by querying an expert

Mathieu Koppen

and

Jean-Paul Doignon

New York University

Université Libre de Bruxelles

A particular field of knowledge is conceptualized as a set of problems (or questions). A person's *knowledge state* in this domain is formalized as the subset of problems this person is capable of solving. When the family of all knowledge states is closed under union, it is called a *knowledge space*. Doignon and Falmagne (1985) established a 1-1 correspondence between knowledge spaces and a class of *surmise systems*, a slight variant of AND/OR graphs. Here we rather obtain a 1-1 correspondence with a well defined class of quasi orders on the collection of all subsets of problems. The resulting approach to knowledge spaces helps to build such spaces for particular domains. We describe a procedure which relies on the answers of an expert to a carefully chosen sequence of information requests.

1. Introduction

A particular field of knowledge can be conceptualized as being comprised of a possibly large, but specified set of *notions*. The *knowledge state* of an individual in this domain can then be formalized as the subset of notions she/he has mastered. Here, a notion can be identified with a question or problem, or, rather, an equivalence class of questions or problems, testing just that notion. Doignon and Falmagne (1985) described the motivation and investigated the algebraic foundation of this approach in some detail. A number of knowledge assessment procedures based on this formalization have been developed (Falmagne and Doignon, 1988a; Falmagne and Doignon, 1988b; Falmagne, 1989; Degreef, Doignon, Ducamp and Falmagne, 1986). These procedures all start from a fixed knowledge structure of the domain, where a *knowledge structure* is defined as the collection of all possible

This work was supported by DOD grant MDA903-87-K-0002 to Jean-Claude Falmagne at New York University, and also by a Fulbright travel grant and a NATO scientific grant to Jean-Paul Doignon. The authors thank Jean-Claude Falmagne for numerous advices on previous drafts of the manuscript. Address comments and requests for reprints to M. Koppen, Dept. of Psychology NYU, 6 Washington Place, New York, NY 10003.

knowledge states. (While any knowledge state is a subset of notions, in general not any such subset is a possible state. For instance, in the field of arithmetic, mastering long division implies mastering of subtraction, so a subset of notions (problems) containing long division but not containing subtraction is not a possible knowledge state. It is indeed this kind of restrictions that lends “structure” to a field.)

In this paper we are concerned with the problem of how to build such knowledge structures for particular domains. For this, we would not want to rely on our own, restricted knowledge, but we would rather consult an expert in the field. A straightforward approach results in presenting to the expert each of all subsets of problems and ask whether it constitutes a possible state. This seems a rather demanding task for the expert. Presented with a subset of problems he has to decide how plausible it is that a subject has mastered this subset *and no other problem*. This means that the expert always has to consider a complete response pattern on the total set of problems, even if the presented subset has only a few elements. Besides, this approach is quite unmanageable in practice because of the sheer number of subsets.

Without any a priori assumption on the family of possible knowledge states, the above approach is the only one conceivable. If we assume a knowledge structure is such that any intersection and any union of knowledge states are again knowledge states – we call such a structure *closed* under union and intersection –, then, by a theorem of Birkhoff (1937, see Theorem 3.1 here), it can be equivalently well specified by a quasi order on the set of problems. Such a representation of knowledge structures by quasi orders would be well suited for eliciting the relevant information from experts (as we will argue in the next section). However, Doignon and Falmagne (1985), who recalled Birkhoff’s result in this context, argue convincingly that the assumption of closure under intersection is not a realistic one for knowledge structures in practice. Assuming only closure under union, they showed that a representation by surmise systems (a variant of AND/OR graphs), instead of quasi orders, is possible. This result of Doignon and Falmagne, together with some other motivation they present, explains the central role played in the theory by knowledge structures closed under union; these are called *knowledge spaces*.

Accepting the assumption of closure under union, we can now reformulate our problem as that of designing a procedure for building a knowledge space for a particular domain by querying experts. Unfortunately, the representation by surmise systems is not very promising in this respect. The point of this paper, then, is to derive an alternative representation for knowledge spaces that is fitting our purpose. Again, quasi orders will enter the picture, but this time they will be relations on the power set of the set of problems. This representation is at the basis of a procedure

that translates the responses of an expert to a set of queries of a specified form into a corresponding knowledge space. The principles of such a procedure are illustrated on a small example. A real-life, large scale application is to be reported elsewhere (Kambouri, Koppen, Villano and Falmagne, 1989); it is based on the algorithm presented in Koppen (1989) as a result of a deeper elaboration of these principles. We also notice here that similar theoretical work was done independently by Mueller (1989).

The paper is organized as follows. In the next section we illustrate the quasi order representation (Birkhoff's theorem), and its limitation, and we describe more precisely the kind of procedure we are going to develop. In Section 3, we formally state the classical result of Birkhoff (1937) and we introduce the key concept of a *Galois connection*. Following Monjardet (1970), Birkhoff's result is derived from a Galois connection between the collection of knowledge structures and the collection of binary relations on the set of problems. Section 4 contains our main theoretical results. We establish a more general Galois connection between knowledge structures and relations on the power set of the set of problems. Theorem 4.5 can be seen as another extension of Birkhoff's result. This leads to the characterization of knowledge spaces as a well defined kind of quasi orders on this power set (Corollary 4.6). In Section 5 the surmise systems of Doignon and Falmagne (1985) are put in this context. The final section describes how the theoretical results of Section 4 give rise to an algorithm for obtaining a knowledge space from the expert's answers and the principles of such an algorithm are shown by running the procedure on the small example presented in Section 2.

2. Background

Let us recall here a small academic example with five problems a, b, c, d, e , presented in Doignon and Falmagne (1985). The content of the five problems (drawn from the field of elementary probability) was examined and in a first analysis the following family K' of possible knowledge states was obtained (\emptyset denotes the empty set):

$$K' = \{ \emptyset, \{c\}, \{e\}, \{b, e\}, \{c, e\}, \{a, b, e\}, \{b, c, e\}, \\ \{c, d, e\}, \{a, b, c, e\}, \{b, c, d, e\}, \{a, b, c, d, e\} \}.$$

Notice that this family K' is closed under both union and intersection. By an old theorem of Birkhoff (1937, Theorem 3.1 here), this means that the family K' can be represented by a quasi order (a reflexive, transitive relation) P' on the set of problems. The correspondence between P' and K' is such that a pair of problems

(x, y) is in P' if and only if any state of K' containing y also contains x . Thus P' can be thought of as a *surmise relation*: if a student solves y correctly and the pair (x, y) is in P' , we can surmise that this student has also mastered x . In our case P' is the partial order sketched in Figure 1. Note that the above construction is reversible; that is, the family K' can be fully recovered from P' .

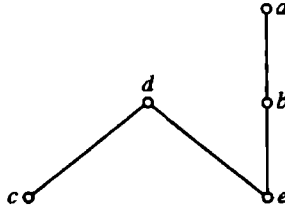


Figure 1. Hasse diagram of the surmise relation P' corresponding to the knowledge structure K' .

This representation by quasi orders suggests an alternative approach to the problem of obtaining a knowledge structure by querying an expert. Instead of asking the expert directly for the knowledge structure, we ask him to specify the quasi order. That is, the expert is presented with a pair of problems and is asked whether solving the first one would imply solving the second one. This task of the expert seems much simpler than deciding directly about the “state-hood” of a subset: now he has at each instance only two problems to consider and he can ignore the rest. The maximum number of questions he will have to answer (i.e., the number of pairs) is only quadratic in the number of problems. Besides, using the transitivity of the quasi order we can make inferences and thus save information requests. Finally, quasi orders and the corresponding knowledge structures are related in such a way that inclusions are reversed. Thus, when we stop the process halfway and end up with a quasi order that is a part of the true quasi order, the corresponding knowledge structure will include the true knowledge structure. There may be some “nuisance” states left (states that will never be assigned to any students and that only act to slow down the assessment procedure), but at least no vital states are missing. So, this procedure of obtaining a knowledge structure by asking an expert about a quasi order is clearly an attractive option. It is not available, however, for building a general knowledge space, since representability by a quasi order requires the extra assumption of closure under intersection.

To continue our example, in a second analysis of the five problems Doignon and Falmagne (1985) argued that a case could be made for including one more subset of problems, the set $\{a, c, d, e\}$, as a possible state. Thus we obtain the knowledge structure

$$\mathbf{K} = \mathbf{K}' \cup \{ \{a, c, d, e\} \}.$$

It can easily be checked that \mathbf{K} is still closed under union, but it is no longer closed under intersection (since $\{a, b, c, e\} \cap \{a, c, d, e\}$ is not a state). This means that \mathbf{K} is a knowledge space which cannot be represented by a quasi order on the problems. For instance, we have to remove from \mathbf{P}' the arrow from b to a . (It is no longer true that knowing a implies knowing b .) This leads to the partial order \mathbf{P} of Figure 2. Now the state $\{a, c, d, e\}$ has indeed been added, but some additional states have been introduced ($\{a, e\}$ and $\{a, c, e\}$) that do *not* appear in the knowledge structure \mathbf{K} .

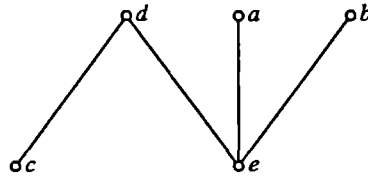


Figure 2. Hasse diagram of the surmise relation \mathbf{P} that “contains” the knowledge structure \mathbf{K} .

Generalizing Birkhoff’s theorem, Doignon and Falmagne (1985) showed that knowledge structures closed under union are in 1-1 correspondence with so-called *surmise systems*. The precise statement of this correspondence will be recalled in Section 5, but the essential point here is that this characterization does not seem to lend itself to an efficient procedure for obtaining a knowledge space from an expert.

In this paper an alternative representation for knowledge spaces is obtained by means of well-defined relations on the power set of the set of notions; in particular, subsets of notions are quasi ordered.[†] The derived representation leads to a procedure for uncovering a knowledge space that is very much like the quasi order procedure for a knowledge structure closed under union and intersection. In fact, it is an extension of this procedure. We start with presenting to the expert pairs of problems, querying whether a student that fails the first one would – in principle – also fail the second problem. This is in effect the same as we did in the procedure sketched above, and if we were willing to assume closure under intersection we could stop the process after this block of questions. Assuming only closure under union, however, we need more information and we continue with a second block of questions. Here the expert is presented a pair of problems together with a single problem and he is to indicate whether a student failing both problems of the pair

[†] We owe this basic idea to Jean-Claude Falmagne.

would also fail the single problem. Next, in a third block, the query is whether failing a triple of problems would imply failing another, single problem. And so on. In general, the information asked from the expert is whether failing a particular subset of problems implies failing another specified problem.

This procedure saves – to a great extent – the advantages of the previously described procedure for the case where we assume closure under intersection. Of course we have lost the quadratic bound on the maximal number of questions, but because the encoding of the knowledge space is again in terms of quasi orders, we can use transitivity to draw inferences, both from positive and negative answers of the expert, thus reducing the information to be requested explicitly. (It will appear that we have in fact more inferences than just from transitivity.) Also, the queries to the expert are still as “local” as possible. In the first block he has, at each instance, to consider two problems only, in the second block three problems (a pair and a single problem), in the third block four, etc. Clearly things get bad toward the end, but because the inferences work all the way down, the number of questions to be asked in the later blocks can be greatly reduced. Still, if it happens that for practical reasons we have to stop the procedure after, say, the second or third block of questions, and the encoding quasi order is only partly known, then, again, we have the advantage that, going from the quasi order to the encoded knowledge space, the inclusions are reversed. The obtained quasi order that is a part of the true quasi order translates into a knowledge space that includes the true knowledge space. So again we cannot lose vital states by interrupting the procedure.

3. Preliminary definitions and results

Let X be some fixed set; in our interpretation it is the collection of problems in some field of knowledge. In the present and next sections we do *not* assume X to be finite, although in practical applications this will be the case. Any knowledge structure \mathbf{K} on X , that is, any family \mathbf{K} of subsets of X , induces a relation Q on X by the definition

$$xQy \quad \text{iff} \quad (\text{for all } K \in \mathbf{K} : y \in K \text{ implies } x \in K).$$

Our interpretation of xQy is that from the mastery of problem y it can be surmised that the same student also masters problem x . Denoting by \mathbf{K}_x the subfamily of \mathbf{K} consisting of the sets that contain the element x , we can write this equivalence more compactly as

$$xQy \quad \text{iff} \quad \mathbf{K}_x \supseteq \mathbf{K}_y. \quad (3.1)$$

It is immediate that the resulting relation Q is a quasi order (that is, it is reflexive and transitive). Conversely, any relation Q on X yields a knowledge structure K on X by setting

$$K \in \mathbf{K} \quad \text{iff} \quad (\text{for all } (x, y) \in Q : y \in K \text{ implies } x \in K). \quad (3.2)$$

It is easy to check that the resulting knowledge structure K is closed both under union and intersection. The above discussion is summarized in the following classical result.

3.1. Theorem. (Birkhoff, 1937.) *The formulae (3.1) and (3.2) define a 1-1 correspondence between knowledge structures K on X that are closed under both union and intersection, and quasi orders Q on X .*

As shown in Monjardet (1970), this result can be considered as a corollary to the fact that the mappings (3.1) and (3.2) define a *Galois connection* between knowledge structures on X and relations on X (both collections ordered by inclusion).

3.2. Definition. (Birkhoff, 1967.) Let (Y, \leq) and (Z, \leq) be two partially ordered sets. A pair of mappings $f : Y \rightarrow Z$ and $g : Z \rightarrow Y$ define a *Galois connection* between (Y, \leq) and (Z, \leq) iff, for $y, y' \in Y$ and $z, z' \in Z$,

- (i) $y \leq y'$ implies $f(y') \leq f(y)$ and $z \leq z'$ implies $g(z') \leq g(z)$;
- (ii) $y \leq (g \circ f)(y)$ and $z \leq (f \circ g)(z)$.

To state results in this context, the notion of a closure on a partially ordered set is useful.

3.3. Definition. (Birkhoff, 1967.) Let (Y, \leq) be a partially ordered set. A mapping $h : Y \rightarrow Y$ is a *closure* on (Y, \leq) when, for $y, y' \in Y$,

- (i) $y \leq h(y)$;
- (ii) $y \leq y'$ implies $h(y) \leq h(y')$;
- (iii) $(h \circ h)(y) = h(y)$.

Any $y \in Y$ such that $y = h(y)$ is called a *closed element* (under h).

3.4. Theorem. (Birkhoff, 1967.) *If the mappings $f : Y \rightarrow Z$ and $g : Z \rightarrow Y$ define a Galois connection between (Y, \leq) and (Z, \leq) , then they induce a 1-1 correspondence between $g(Z)$ and $f(Y)$. More specifically, $f \circ g$ and $g \circ f$ are closures on (Y, \leq) and (Z, \leq) with the closed elements collected in $g(Z)$ and $f(Y)$, respectively; the restrictions of f and g to $g(Z)$ and $f(Y)$, respectively, are inverse order-reversing isomorphisms.*

As an example, (3.1) and (3.2) define a Galois connection leading to the 1-1 correspondence in Theorem 3.1. Here, the closure of a knowledge structure \mathbf{K} on X is the smallest family of subsets of X containing \mathbf{K} that is closed under union and intersection. For the relations on X , “closure” here means the reflexive transitive closure, mapping any relation \mathbf{R} to the smallest quasi order containing \mathbf{R} .

Knowledge spaces, being closed only under union, cannot be captured by quasi orders on X . In the next section we will show, by establishing another Galois connection, that they are still in a 1-1 correspondence with a class of quasi orders, but we have to move one level up and consider quasi orders on the power set of X .

4. A Galois connection for knowledge spaces

Let us first introduce some more notation. We denote by Ω the collection of all knowledge structures on X , so

$$\Omega = 2^{(2^X)},$$

the power set of the power set of X . Clearly, Ω is ordered by inclusion:

$$\mathbf{K} \subseteq \mathbf{K}' \quad \text{iff} \quad (\text{for all } K \in 2^X : K \in \mathbf{K} \text{ implies } K \in \mathbf{K}').$$

We denote by Ψ the collection of all binary relations on the power set of X :

$$\Psi = 2^{(2^X \times 2^X)}.$$

This collection is again ordered by inclusion:

$$\mathbf{R} \subseteq \mathbf{R}' \quad \text{iff} \quad (\text{for all } A, B \in 2^X : A \mathbf{R} B \text{ implies } A \mathbf{R}' B).$$

We proceed to construct a Galois connection (r, k) between the two partially ordered sets (Ω, \subseteq) and (Ψ, \subseteq) . With any knowledge structure \mathbf{K} on X we can associate a relation $r(\mathbf{K})$ on the power set of X by the definition

$$A r(\mathbf{K}) B \quad \text{iff} \quad (\text{for all } K \in \mathbf{K} : B \cap K \neq \emptyset \text{ implies } A \cap K \neq \emptyset).$$

Introducing the notation \mathbf{K}_A for the subcollection of sets in \mathbf{K} that “meet” the subset A of X :

$$\mathbf{K}_A = \{ K \in \mathbf{K} : A \cap K \neq \emptyset \},$$

this can be written as

$$A r(\mathbf{K}) B \quad \text{iff} \quad \mathbf{K}_A \supseteq \mathbf{K}_B.$$

The empirical interpretation of $A r(\mathbf{K}) B$ is that if a student masters some question in B , he also masters some question in A ; or, equivalently, if he does not master any

question in A , he also does not master any question in B . Notice the similarity with the corresponding definition (3.1) in the Birkhoff case. The next proposition follows easily from the definition of K_A .

4.1. Proposition. *For the mapping $r : \Omega \rightarrow \Psi$ defined, for $A, B \in 2^X$, by*

$$Ar(K)B \quad \text{iff} \quad K_A \supseteq K_B$$

we have, for $K, K' \in \Omega$:

- (i) $r(K)$ extends \supseteq on 2^X (that is, $A \supseteq B$ implies $Ar(K)B$);
- (ii) $r(K)$ is transitive (and thus, by (i), a quasi order);
- (iii) if $Ar(K)B_i$, for all i in some index set I , then $Ar(K)(\cup_{i \in I} B_i)$;
- (iv) if $K \subseteq K'$ then $r(K) \supseteq r(K')$;
- (v) if K^c is the closure under union of K , then $r(K) = r(K^c)$.

Relations on 2^X that have the properties (i) to (iii) of Proposition 4.1 will play an important role in the sequel, so we investigate these more closely:

4.2. Lemma. *For a transitive extension P of \supseteq on 2^X , the following conditions are equivalent:*

- (i) APB , for all i in some index set I implies $AP(\cup_{i \in I} B_i)$;
- (ii) Any $AP = \{Z \in 2^X : APZ\}$ has a maximum (for the inclusion) element A^* ;
- (iii) APB iff (for all $x \in B : AP\{x\}$).

These conditions imply

- (iv) APB implies $AP(B \cup A)$

and in case X is finite, (iv) is equivalent to (i) - (iii).

A transitive extension of \supseteq on 2^X for which these conditions hold is called an *entail relation* for X . (Note that an entail relation for X is a relation on 2^X .) One has APB iff $A^* \supseteq B$, with A^* as in (ii).

Proof. That (iii) implies (i) is immediate, as is the implication from (i) to (ii) and (iv) if we notice that APA . The “only if” part of (iii) is given (P extends \supseteq) and the “if” part follows from (ii): If $AP\{x\}$ for all $x \in B$, then, using the definition in (ii), $A^* \supseteq \{x\}$ for all $x \in B$, which implies $A^* \supseteq B$. So we have $APA^* \supseteq B$, thus APB . To see that (iv) implies (i) in the finite case it suffices to prove that APB and APC imply $AP(B \cup C)$. From $(A \cup B) \supseteq A$ we have $(A \cup B)PA$. By (iv) and APC , we also have $AP(A \cup C)$. There follows $(A \cup B)P(A \cup C)$. Applying (iv), we get $(A \cup B)P(A \cup B \cup C)$. Since $(A \cup B \cup C) \supseteq (B \cup C)$ and, by (iv), $AP(A \cup B)$, we finally derive $AP(B \cup C)$. ■

Notice that (iv) is not equivalent to (iii) in general. For a counterexample, take X

infinite and define \mathbf{P} by $A \mathbf{P} B$ iff $B - A$ is finite.

Having defined the mapping $r : \Omega \rightarrow \Psi$, we now introduce a mapping that goes the other way. Again, the definition will be reminiscent of the corresponding mapping (3.2) in the Birkhoff case:

4.3. Proposition. *For the mapping $k : \Psi \rightarrow \Omega$ defined by*

$$K \in k(\mathbf{R}) \quad \text{iff} \quad (\text{for all } (A, B) \in \mathbf{R} : B \cap K \neq \emptyset \text{ implies } A \cap K \neq \emptyset)$$

we have, for $\mathbf{R}, \mathbf{R}' \in \Psi$:

- (i) $k(\mathbf{R})$ is closed under union;
- (ii) if $\mathbf{R} \subseteq \mathbf{R}'$ then $k(\mathbf{R}) \supseteq k(\mathbf{R}')$.

The following lemma, giving an alternative characterization of the mapping k for the case of an entail relation, will be very useful.

4.4 Lemma. *If $\mathbf{P} \in \Psi$ is an entail relation for X , then, with k as in Proposition 4.3,*

$$K \in k(\mathbf{P}) \quad \text{iff} \quad (\text{for all } Z \in 2^X : (X - K) \mathbf{P} Z \text{ implies } (X - K) \supseteq Z).$$

Proof. If $K \in k(\mathbf{P})$, then for any $A, B \in 2^X$ such that $A \mathbf{P} B$ and $K \cap A = \emptyset$ we have $K \cap B = \emptyset$. Putting $A = X - K$ it follows that $(X - K) \mathbf{P} B$ implies $(X - K) \supseteq B$ for all $B \in 2^X$. Conversely, let K be such that $(X - K) \mathbf{P} Z$ implies $(X - K) \supseteq Z$ for all $Z \in 2^X$ and suppose for some $A, B \in 2^X$ we have $A \mathbf{P} B$ and $K \cap A = \emptyset$. This means $(X - K) \supseteq A \mathbf{P} B$, thus $(X - K) \mathbf{P} B$ and by our assumption $(X - K) \supseteq B$. So, $K \cap B = \emptyset$ and $K \in k(\mathbf{P})$. ■

Now we are ready to formulate our main result:

4.5. Theorem. *The pair (r, k) , where r is as in Proposition 4.1 and k as in Proposition 4.3, is a Galois connection between the partially ordered sets (Ω, \subseteq) and (Ψ, \subseteq) . The closed elements in Ω are the knowledge spaces and the closed elements in Ψ are the entail relations for X .*

Proof. The inclusion $\mathbf{R} \subseteq (r \circ k)(\mathbf{R})$ can be checked easily from the definitions of the mappings r and k . So in view of Propositions 4.1(iv) and 4.3(ii) the Galois connection is established if we can show that $\mathbf{K} \subseteq (k \circ r)(\mathbf{K})$ for $\mathbf{K} \in \Omega$. Since by Proposition 4.1(i) to (iii) any $r(\mathbf{K})$ is an entail relation, we can use Lemma 4.4 to see that

$$K \in (k \circ r)(\mathbf{K}) \quad \text{iff} \quad (\text{for all } Z \in 2^X : \mathbf{K}_{X-K} \supseteq \mathbf{K}_Z \text{ implies } X - K \supseteq Z).$$

For any $K, Z \in 2^X$, the inclusion $K_{X-K} \supseteq K_Z$ implies $K \notin K_Z$, so clearly $K \in K$ together with this inclusion implies $X-K \supseteq Z$. Consequently $K \subseteq (k \circ r)(K)$.

Now let us determine the closed elements. If K is a space and $K \in 2^X$ we can, above, take $Z = X - K_0$, where K_0 is the unique maximal element of K included in K . Then clearly $K_{X-K} \supseteq K_{X-K_0}$, and thus $K \in (k \circ r)(K)$ implies $K \subseteq K_0$, that is, $K \in K$. So if K is a space we have $K = (k \circ r)(K)$. This shows that any space is in the image of k and since Proposition 4.3(i) says that any element of $k(\Psi)$ is a space, we see that the knowledge spaces constitute the closed elements in Ω .

The closed elements in Ψ are the images by r and by Proposition 4.1(i) to (iii) any element of $r(\Omega)$ is an entail relation for X . We will show the converse by establishing $R = (r \circ k)(R)$ for such a relation R . One inclusion being checked above it remains to show that $R \supseteq (r \circ k)(R)$. So, suppose $A (r \circ k)(R) B$ and let A^* be the maximum element of $A R$ (Lemma 4.2(ii)). Clearly, for all $Z \in 2^X$, $A^* R Z$ implies $A R A^* R Z$ and thus $A^* \supseteq Z$. By Lemma 4.4, this means $X - A^* \in k(R)$. From $A^* \supseteq A$ it follows $X - A^* \notin k(R)_A$, which, in view of $A (r \circ k)(R) B$, implies $(X - A^*) \cap B = \emptyset$. This gives $A R A^* \supseteq B$, thus $A R B$. ■

4.6. Corollary. *The mappings r and k induce a 1-1 correspondence between the collections of knowledge spaces on X and entail relations for X .*

4.7. Remark. We point out here that, independently, a very similar 1-1 correspondence was recently obtained by Burigana (1988). He did not derive it from a Galois connection, but formulated it directly in terms of closures (see below for this approach). Curiously enough, Burigana's motivation differs completely from ours: his paper is devoted to the study of regularity in sequences of stimuli.

4.8. An alternative formulation. The 1-1 correspondence induced by the Galois connection (r, k) can be described alternatively in terms of closures (cf. Definition 3.3). Observe first that any knowledge space K defines a closure h_K on $(2^X, \subseteq)$ by letting $h_K(A)$ denote the maximum element $Z \in 2^X$ such that $K_A \supseteq K_Z$. (For a space, this is well-defined and it follows that $X - h_K(A)$ must be a state of K , in fact the largest state included in $X - A$.) With our interpretation, $h_K(A)$ collects the problems that we can infer a person will fail if we know this person fails all problems in A . On the other hand, for any entail relation P for X the mapping $A \rightarrow A^* = \cup A P$ (see Lemma 4.2(ii)) yields a closure h_P on $(2^X, \subseteq)$. Now a space K and an entail relation P are paired in the 1-1 correspondence of Corollary 4.6 ($r(K) = P$ and $k(P) = K$) if and only if the induced closures are the same (iff $h_K = h_P$).

The mappings $K \rightarrow h_K$ and $P \rightarrow h_P$ are bijections. More specifically, for an arbitrary closure h on $(2^X, \subseteq)$ we find the knowledge space K such that $h = h_K$ as

$$K = \{X - h(A) : A \in 2^X\},$$

and the entail relation P such that $h = h_P$ is, for $A, B \in 2^X$, defined by

$$A P B \quad \text{iff} \quad h(A) \supseteq B.$$

It may be checked that the defining properties of a closure ensure that the thus defined family K is indeed a space and the relation P indeed an entail relation for X .

4.9. Birkhoff revisited. By Lemma 4.2(iii) we know that for an entail relation P for X we do not lose any information by restricting its range to the singleton elements of 2^X . If, in addition, P is such that for $A \in 2^X$ and $y \in X$

$$A P \{y\} \quad \text{iff} \quad (\text{for some } x \in A : \{x\} P \{y\}),$$

we can also restrict the domain of P to the singleton sets and we have for $A, B \in 2^X$:

$$A P B \quad \text{iff} \quad (\text{for any } y \in B \text{ there is } x \in A \text{ such that } \{x\} P \{y\}). \quad (*)$$

It is easy to check (with Lemma 4.4, for instance) that for such a P the knowledge space $k(P)$ is closed under intersection. Conversely, if a knowledge space K is closed under intersection, then the relation $P = r(K)$ satisfies (*). For, if $A r(K) B$ and $y \in B$, we have $A r(K) \{y\}$, meaning that any state containing y contains an element of A . In particular, the state $\bigcap K_{\{y\}}$ contains some $x \in A$; consequently this x appears in any element of $K_{\{y\}}$ and we obtain $K_{\{x\}} \supseteq K_{\{y\}}$. Thus the "only if" part of (*) is established; the "if" part poses no problems.

Since r and k are inverse mappings when restricted to knowledge spaces and entail relations for X , respectively, we obtain in this way a 1-1 correspondence between knowledge structures closed under union as well as intersection and quasi orders on 2^X satisfying (*). By identifying singleton sets with their element, the latter collection is in a natural 1-1 correspondence with the collection of quasi orders on X , which gives us Birkhoff's result.

5. The relation with surmise systems

In this section we will assume that the set X is finite. Under that restriction Doignon and Falmagne (1985) obtained a different generalization of Birkhoff's result. They noted that a knowledge structure's being closed under intersection forces a notion to have one unique set of prerequisites, which is too restrictive. So, dropping the requirement of closure under intersection, Doignon and Falmagne (1985) were led to

the definition of a *surmise system*, in which to each notion is assigned a collection of subsets of notions, representing the possible sets of prerequisites for that notion. Again by establishing a Galois connection, they derived a 1-1 correspondence between the collection of knowledge spaces and the collection of *space-like* surmise systems on X :

5.1. Definition. A *surmise system* on X is a mapping $\sigma: X \rightarrow \Omega$. The elements of $\sigma(x)$ are called the *clauses* for x . The *states* of σ are the sets $Z \in 2^X$ that contain a clause for each element:

$$\text{for all } x \in Z \text{ there is } C \in \sigma(x) \text{ such that } C \subseteq Z.$$

The surmise system σ is called *space-like* if each clause for x is a state containing x and the clauses for x are pairwise incomparable (with respect to inclusion).

5.2. Theorem. (Doignon and Falmagne, 1985.) *The collection of knowledge spaces on X is in 1-1 correspondence with the collection of space-like surmise systems on X . In this correspondence, the clauses for $x \in X$ in the space-like surmise system constitute the minimal states in the knowledge space that contain x .*

Since knowledge spaces are in 1-1 correspondence with both entail relations for X and space-like surmise systems, obviously the last two collections must be in 1-1 correspondence. We will investigate here how each entail relation for X induces a space-like surmise system.

If P is a relation on 2^X , we denote by \bar{P} the complement of P :

$$\bar{P} = (2^X \times 2^X) - P.$$

From Lemma 4.2(iii) we know that in case P is an entail relation, it is uniquely determined by the sets $\bar{P}\{x\}$, $x \in X$. To simplify notation, we will in the sequel write $A\bar{P}x$ for $A\bar{P}\{x\}$ and also APx for $AP\{x\}$. So

$$\bar{P}x = \{Z \in 2^X : \text{not } ZP\{x\}\}.$$

We have for APB the interpretation: if a subject has mastered a notion in B , he must also have mastered a notion in A . Thus $A\bar{P}x$ means: it is still possible to know x without knowing anything in A ; or, in terms of surmise systems, there is still a clause for x contained in $X - A$. With this interpretation, the following theorem is not really surprising.

5.3. Theorem. *For an entail relation \mathbf{P} for X , define a mapping $s(\mathbf{P}): X \rightarrow \Omega$ by assigning to each $x \in X$ the collection of complements of maximal elements of $\bar{\mathbf{P}}x$. That is, $C \in s(\mathbf{P})(x)$ iff $(X - C) \bar{\mathbf{P}}x$ and no strict superset of $X - C$ has this property. Then $s(\mathbf{P})$ is a space-like surmise system and the states of $s(\mathbf{P})$ coincide with the states of the knowledge space $k(\mathbf{P})$.*

Proof. The clauses for any $x \in X$ are pairwise incomparable by definition and each clause C for x contains x since $(X - C) \bar{\mathbf{P}}x$ implies $(X - C) \not\bar{\mathbf{P}}\{x\}$. We have proved that $s(\mathbf{P})$ is space-like if we can show that any clause for x is a state. So, suppose $C \in s(\mathbf{P})(x)$ and $x' \in C$. Since $X - C$ is maximal in $\bar{\mathbf{P}}x$, it must be the case that $((X - C) \cup \{x'\}) \mathbf{P}x$. In view of this, the assumption $(X - C) \mathbf{P}x'$ would lead to $(X - C) \mathbf{P}((X - C) \cup \{x'\}) \mathbf{P}x$, contradicting $(X - C) \bar{\mathbf{P}}x$. Thus, $(X - C) \bar{\mathbf{P}}x'$ and, using finiteness of X , there is some maximal $A \supseteq (X - C)$ such that $A \bar{\mathbf{P}}x'$. In other words, $(X - A) \in s(\mathbf{P})(x')$ and $(X - A) \subseteq C$.

For the states of $k(\mathbf{P})$ we use the characterization of Lemma 4.4. To show that any state of $s(\mathbf{P})$ is a state of $k(\mathbf{P})$ we take $K, Z \in 2^X$ such that $(X - K) \mathbf{P}Z$ and for all $x \in K$ there is $C_x \subseteq K$ with $(X - C_x) \bar{\mathbf{P}}x$. The proof that then $(X - K) \supseteq Z$ follows by contradiction: if $x \in K \cap Z$ we have $(X - K) \cup Z \supseteq \{x\}$ and, since $(X - K) \mathbf{P}Z$ implies $(X - K) \mathbf{P}((X - K) \cup Z)$, we derive

$$(X - C_x) \supseteq (X - K) \mathbf{P}((X - K) \cup Z) \supseteq \{x\},$$

contradicting $(X - C_x) \bar{\mathbf{P}}x$. On the other hand, if K is a state of $k(\mathbf{P})$, Lemma 4.4 yields $(X - K) \bar{\mathbf{P}}x$ for each $x \in K$. Using finiteness of X we may conclude that then for any $x \in K$ there must be a minimal $C_x \subseteq K$ such that $(X - C_x) \bar{\mathbf{P}}x$. This means that K is a state of $s(\mathbf{P})$. ■

6. The construction of the knowledge space

The theoretical results of Section 4 find a useful practical application in the problem of obtaining the knowledge structure of a particular domain from querying an expert in that field. We will sketch here a straightforward way of doing this that derives directly from the results of Section 4, and we illustrate the procedure on the 5 problem example of Section 2. For applications of practical importance (i.e., for a larger problem set) we need a more sophisticated version. Such a practicable procedure, based on an elaboration of the same principles, is described elsewhere (Koppen, 1989).

The basic idea is to not ask the expert directly for the states in the knowledge structure, but rather for the corresponding entail relation for the set of notions. So, in principle, we present the expert with two subsets A and B of notions, and we ask

whether or not it is safe to conclude that if a subject fails all problems in A (has not mastered any notion in A), then the same subject will fail all problems in B (has not mastered any notion in B). When this conclusion is valid, the pair (A, B) is in the relation, otherwise it is not.

Knowing that the relation we are looking for is an entail relation for X , we certainly need not offer the expert all pairs of subsets of notions. (We put ourselves in the comfortable position that our expert is perfectly reliable.) For one thing, from Lemma 4.2(iii) we know that we need only consider singleton sets in the range of such an entail relation. That is, we have to ask the expert only questions of the kind: if a subject fails all problems in A , implies this that he will fail problem x ? This would be a silly question to ask when x is an element of A , which reflects the fact that an entail relation extends the superset relation. In order to deal only with singletons in the range of our relation, we need the following easy Lemma, which characterizes such restrictions of entail relations. (The statement of this Lemma is in terms of X as the restricted range, identifying the singleton $\{x\}$ with the element x .)

6.1. Lemma. *A relation \hat{P} included in $2^X \times X$ is the restriction of an entail relation for X iff it satisfies the following two conditions:*

- (i) \hat{P} includes the reverse membership relation rmo ;
- (ii) $[(B \hat{P} z \text{ for all } z \in Z) \ \& \ Z \hat{P} y]$ implies $B \hat{P} y$.

The Lemma follows, since P defined by

$$A P B \quad \text{iff} \quad A \hat{P} b \text{ for all } b \in B$$

defines an entail relation for X if and only if \hat{P} satisfies the conditions in Lemma 6.1.

As a consequence of Lemma 6.1, the expert will be queried only about pairs of the form A, x . Moreover, not all such pairs need to be presented because inferences can be drawn, both from positive responses (validating the implication) and from negative responses (denying it). Let us, for notational convenience, introduce the shorthand A_x for the set $A \cup \{x\}$ with $A \in 2^X$ and $x \in X$. According to condition (ii) of Lemma 6.1, then, any positive inference that *directly* involves a new observation $A P x$ must be of one of the two forms (a) or (b):

- (a) We had already established $B P a$ for all $a \in A$, and then from $A P x$ we infer $B P x$ ($y = x$ and $Z = A$ in 6.1(ii)).
- (b) We had already established $A_x P y$, and then $A P x$ leads to $A P y$ ($B = A$ and $Z = A_x$ in 6.1(ii)).

The two cases (a) and (b) can be combined in one equivalent inference rule:

$$A P x \text{ implies } B P y \quad \text{whenever} \quad B P A \ \& \ A_x P y. \quad (6.1)$$

Indeed, (a) and (b) correspond to the special cases $y = x$ and $B = A$, respectively,

while conversely any inference BP_y obtained by (6.1) also follows by first inferring AP_y (case (b)) and next BP_y (case (a)). Note that in this discussion we only considered *direct* inferences from AP_x , that is, inferences in which A and/or x appear in the conditions of the rule. We find *all* positive inferences by applying rule (6.1) iteratively to these direct inferences, and so on. In general, repeated application can yield new inferences.

A positive response can also lead to *negative* inferences according to similar rules, for instance,

$$AP_x \text{ implies } B\bar{P}_y \text{ whenever } A_xPB \ \& \ A\bar{P}_y \quad (6.2)$$

and

$$AP_x \text{ implies } B\bar{P}_y \text{ whenever } B\bar{P}_x \ \& \ B_yPA. \quad (6.3)$$

Both rules follow because the opposite assumption, BP_y would lead to a contradiction. Note that, by Lemma 4.2(iv), the observation AP_x is equivalent to APA_x . Thus, in the case of (6.2), BP_y would give us APA_xBP_y , contradicting $A\bar{P}_y$. Similarly, with (6.3) we would have BP_yPA , contradicting $B\bar{P}_x$. When we get a negative response $A\bar{P}_x$ from the expert, we can only obtain negative inferences. We can formulate the rule

$$A\bar{P}_x \text{ implies } B\bar{P}_y \text{ whenever } APB \ \& \ B_yP_x, \quad (6.4)$$

since BP_y under the conditions of (6.4) would imply AP_x .

It must be noted that rules (6.2) to (6.4) do not necessarily find all possible negative inferences, even when rule (6.4) is applied iteratively to the obtained negative inferences. Rule (6.4), for instance, is a special case of the more general rule

$$A\bar{P}_x \text{ implies } B\bar{P}_C \text{ whenever } APB \ \& \ (B \cup C)P_x.$$

But the inference $B\bar{P}_C$ cannot be processed in the restricted range version of P . It means that for some $y \in C$ we must have $B\bar{P}_y$, but unless for every but one element z of C we had already established BP_z , we do not know which y to pick. (Of course, if we had already $B\bar{P}_y$ for some $y \in C$, the inference tells us nothing new.) So in general we would have to save all such inferences $B\bar{P}_C$ and after each new inference we would have to check whether any saved inference can now be consummated. This appears to be rather heavy. In practice we will use only the rules (6.2) to (6.4) that apply directly to the singleton set range of P . Below will be indicated how we deal with their being incomplete.

Using the implications (6.1) really means that we get to our ultimate relation P via a number of intermediate relations, all of which are entail relations. In the absence of any information we start with a relation P_0 which is set to the superset

So at each moment t of the construction process we have a knowledge space \mathbf{K}_t , that is a conservative approximation of the final knowledge space \mathbf{K} ; conservative in the sense that any \mathbf{K}_t contains all states of \mathbf{K} . As mentioned in Section 2, this fact is of some practical significance. It means that when we want the knowledge space for future use in knowledge assessment procedures, we may interrupt the construction process at any time practical considerations lead us to do so, and we will end up with a knowledge space that possibly contains some “nuisance” states, but in which at least no vital states are missing.

In particular, we can start the construction process with the singleton subsets (“if a subject fails x , is it safe to conclude that he will also fail y ?”) and stop when all such questions have been asked. If P_t is the obtained relation at this point, then, by construction, P_t satisfies condition (*) of 4.9 and clearly any addition to P_t will invalidate the “only if” part of this equivalence. So the corresponding knowledge space \mathbf{K}_t is closed under intersection (as it should be, since until now we have, in effect, asked questions regarding the quasi order on X , so we are still in the Birkhoff case) and it is the smallest such space. Thus, stopping the process at this point leaves us with the closure under intersection of the final knowledge space \mathbf{K} .

6.2. Example. Let us illustrate this procedure with the example from Section 2, where $X = \{a, b, c, d, e\}$. To simplify notation we will in the sequel denote subsets of X as strings without surrounding braces and separators between the elements. So abc denotes the subset $\{a, b, c\}$. In this way we lose the distinction between a singleton set and its element, but the correct interpretation will be clear from context. Furthermore, pairs of subsets are given in a dot notation: $abc.de$ represents the pair of subsets $(\{a, b, c\}, \{d, e\})$.

Now suppose we want to recover the knowledge space

$$\mathbf{K} = \{\emptyset, c, e, be, ce, abe, bce, cde, abce, acde, bcde, X\},$$

given in Section 2, by gradually constructing the corresponding entail relation for X from the expert’s responses. We have to ask the expert questions of the kind: does failing all in Z imply failing x , where Z runs through the subsets of X and x through X . The order of these questions is in principle arbitrary, but it makes sense to start with the simpler ones. We will adopt here a very straightforward rule for choosing the next question to ask. We order the subsets by increasing cardinality and within classes of equal cardinality we choose, arbitrarily, the lexicographic ordering; this ordering is also used for the elements of X . Thus we obtain an ordering of the pairs (Z, x) by letting (Z, x) precede (Z', x') iff Z precedes Z' or $Z = Z'$ and x precedes x' . We choose for the next question the first undecided pair in this ordering.

At Start :		$N_0 = \{\emptyset x : x \in X\}$	$P_0 = \supseteq$	$K_0 = 2^X$
<i>t</i>	Observed	Adding to N	Adding to P	Deleting from K
1 to 10	$aNx, x \neq a$ $bNx, x \neq b$ cNa, cNb	$a.b a.c a.d a.e$ $b.a b.c b.d b.e$ $c.a c.b$		
11	cPd	$d.a d.b cd.a$ $cd.b$	$c.d ac.d bc.d$ $ce.d abc.d ace.d$ $bce.d abce.d$	$abde bde ade$ $abd de bd$ $ad d$
12	cNe	$c.e d.e cde$		
13	dNc	$d.c$		
14	ePa		$e.a be.a ce.a$ $de.a bce.a bde.a$ $cde.a bcde.a$	$abcd acd abc$ $ac ab a$
15	ePb		$e.b ae.b ce.b$ $de.b ace.b ade.b$ $cde.b acde.b$	$bcd bc b$
16	eNc	$e.c ab.c ae.c$ $be.c abe.c$		
17	ePd	$ad.c bd.c de.c$ $abd.c ade.c$ $bde.c abde.c$	$e.d ae.d be.d$ $abe.d$	cd
18	$abNd$	$ab.d ab.e$		
19	$acNb$	$ac.b ac.e ad.b$ $ad.e acd.b acd.e$		
20	$bcPa$		$bc.a bcd.a$	ae
21	$bcNe$	$bc.e bd.e abc.e$ $abd.e bcd.e abcd.e$		
22	$bdPa$		$bd.a$	ace

Table 1. Observed responses of the expert, based on the space K of Example 6.1, and inferences, yielding successive approximations of the entail relation P and the corresponding knowledge space K .

Using this design in questioning a perfectly reliable expert we would get the results gathered in Table 1. Here and in the sequel we show results for the relations P_t and N_t only restricted to the singleton set range. We start with a relation P_0 set equal to the superset relation on the power set of X and the corresponding knowledge space K_0 equal to the power set. In principle $N_0 = \emptyset$ (any pair of subsets can be in the relation P) and we would start inquiring about the pairs $\emptyset x$, for $x \in X$.

	t = 11		t = 14		t = 15		t = 17		t = 20		t = 22	
	abcde	K	abcde	K	abcde	K	abcde	K	abcde	K	abcde	K
∅	nnnnn	↓	nnnnn	↓	nnnnn	↓	nnnnn	↓	nnnnn	↓	nnnnn	↓
a	nnttt	↓	nnttt	↓	nnttt	↓	nnttt	↓	nnttt	↓	nnttt	↓
b	ntttt	↓	ntttt	↓	ntttt	↓	ntttt	↓	ntttt	↓	ntttt	↓
c	nttp.		nttp.		nttp.		nttp.		nttp.		nttp.	
d	nt.n	↓	nt.n	↓	nt.n	↓	nt.n	↓	nt.n	↓	nt.n	↓
e	↓	p		pp . . .		ppn . .		ppn . .		ppn . .	
ab	nb . . .	↓	nb . . .	↓	nb . . .	↓	nb . . .	↓	nb . . .	↓	nb . . .	↓
ac	n . . . p		n . . . p		n . . . p		n . . . p		n . . . p		n . . . p	
ad	n	↓	n	↓	n	↓	n	↓	n	↓	n	↓
ae	n	↓	n	↓	np . . .		npn . .		npn . .		npn . .	
bc p	 p	 p	 p	 p	 p	
bd	↓	↓	↓	↓	↓	↓
be	↓	p		pn . . .		pn . . .		pn . . .		pn . . .	
cd	nn . . .	↓	nn . . .	↓	nn . . .	↓	nn . . .	↓	nn . . .	↓	nn . . .	↓
ce p		p		pp . . .		pp . . .		pp . . .		pp . . .	
de	↓	p		pp . . .		pp . . .		pp . . .		pp . . .	
abc	nb . . . p		nb . . . p		nb . . . p		nb . . . p		nb . . . p		nb . . . p	
abd	nb	↓	nb	↓	nb	↓	nb	↓	nb	↓	nb	↓
abe	nb	↓	nb	↓	nb	↓	nb	↓	nb	↓	nb	↓
acd	n	↓	n	↓	n	↓	n	↓	n	↓	n	↓
ace	n p		n p		np . . .		np . . .		np . . .		np . . .	
ade	n	↓	n	↓	np . . .		np . . .		np . . .		np . . .	
bcd	↓	↓	↓	↓	↓	↓
bce p		p		pn . . .		pn . . .		pn . . .		pn . . .	
bde	↓	p		pn . . .		pn . . .		pn . . .		pn . . .	
cde	↓	p		pp . . .		pp . . .		pp . . .		pp . . .	
abcd	nb p		nb p		nb p		nb p		nb p		nb p	
abce	nb p		nb p		nb p		nb p		nb p		nb p	
abde	nb	↓	nb	↓	nb	↓	nb	↓	nb	↓	nb	↓
acde	n	↓	n	↓	np . . .		np . . .		np . . .		np . . .	
bcde	↓	p		pp . . .		pp . . .		pp . . .		pp . . .	
X	nb p	↓	nb p	↓	nb p	↓	nb p	↓	nb p	↓	nb p	↓

Table 2. Situations after each positive response of the expert (see Table 1). Indicated are positive and negative inferences from the latest response (p and n), earlier inferences (p and n) and earlier observed responses (p and n). The “.” means that the row subset contains the column element, a dot that the corresponding pair is still undecided. The complement of a subset indexing a checked row is in the current knowledge space; if the check is in boldface, it is known to be in the final space.

But a moment's thought will reveal that $\emptyset P x$ means precisely that x does not appear at all in the knowledge space. We will save us five questions here by assuming we know the specification of the set of problems X is right in the sense that $X = \cup K$. On this assumption we can then "infer" $\emptyset N x$ for all $x \in X$, giving us our initial N_0 . So our real first question to the expert is: does failing a imply failing b ? As we can check from the above specification, K has states containing b and not containing a (be , for instance), so we will get a negative response from our expert. The same is true for the next 9 questions, up to and including the question: does failing c imply failing b ? Since we can not infer anything from negative responses in the absence of any positive response, we can only add these pairs of subsets to the relation N_0 to obtain N_{10} , while still $P_{10} = P_0$ and $K_{10} = K_0$.

At $t = 11$ we get our first success: since any state of K containing d also contains c , we get a positive response from our expert for the pair $c.d$. As we can see in Table 1, this gives us some extra inferences. The negative inference dNa , for instance, follows, by transitivity of P , from cNa and cPd , and since cPd and the known cPc imply $cPcd$, we similarly obtain $cdNb$ from cNb , and so on. These are all simple instances of rule (6.2). The positive inferences all follow from a special case of rule (6.1). With cPd we have $C \supseteq cPd$ for any set C containing c and this adds to P_{11} all pairs $C.d$ where C contains c but not d . Adding pairs to P_{11} means deleting elements from K_{11} and the correspondence is via Lemma 4.4. Adding $c.d$ to P , for instance, means that it is no longer true that cPZ implies $c \supseteq Z$ for all $Z \in 2^X$, so according to Lemma 4.4 we have to drop $X - c = abde$ from K_{11} . In the same way, $acePd$ forces us to discard bd from K_{11} and so on.

The situation after the first positive response at $t = 11$ is depicted in Table 2 in the column with that heading. Here the problem of finding the relation P is represented as filling out a matrix, the rows of which are indexed by the subsets of X and the columns by the elements of X . In this and in the next columns of Table 2 boldface p and n entries represent positive, respectively negative inferences made from the expert's positive response to the latest question asked; italic p (n) denotes earlier positive (negative) inferences and roman p (n) denotes previously observed positive (negative) responses from the expert. A "⊃" entry corresponds to row-column pairs that are in P because the subset of the row contains the column element, while a dot indicates that the corresponding pair is still undecided. In a subcolumn is checked which rows contribute states to the current knowledge space; according to Lemma 4.4 the complement of a row indexing subset is a state as long as there is no positive (p , p or p) entry in that row. The checkmark is in boldface when the corresponding row has been completed, meaning that the complementary set is bound to be a state of the final knowledge space.

Returning to Table 1 we see that at times $t = 12$ and $t = 13$ we get negative responses, giving some negative inferences, and the next positive response is obtained at $t = 14$. From this we get a number of positive inferences and consequently a number of states can be discarded, all of this in the same way as we have seen above. In this way the process continues, until at $t = 22$ the relation \mathbf{P} is completely known. With \mathbf{K}_{22} we have indeed reconstructed the knowledge space \mathbf{K} we used for deducing the expert's responses. In Table 2 we show the situations after each newly obtained positive response and consequent change in \mathbf{K}_t . The inquisitive reader may want to check that at each moment t the knowledge structure \mathbf{K}_t is closed under union, that for $t = 0, \dots, 17$ \mathbf{K}_t is also closed under intersection while this no longer holds from $t = 20$, and finally that indeed \mathbf{K}_{17} is the closure under intersection of \mathbf{K}_{22} . The corresponding relation \mathbf{P}_{17} , considered as a partial order on X , equals exactly the partial order \mathbf{P} we found in Section 2 (see Fig. 2) as our best try for a quasi order representation of \mathbf{K} .

This example illustrates the correspondence between knowledge spaces and entail relations, and it shows in principle how we can obtain the space by questioning an expert about the entail relation. However, it will also be clear from this example that, as such, the procedure would not be practicable, even for a very moderate number of problems in X . For instance, for a 20 problem set the equivalent of Table 2 would consist of over one million rows. The size of this table simply doubles with each additional problem. On closer inspection, however, it appears that, generally, many rows are redundant: the complete information on entail relation and knowledge space is contained in a subtable of considerably smaller size. An essential part of the algorithm presented in Koppen (1989) deals precisely with the issue of constructing just this minimal subtable, dynamically, in the course of questioning the expert. (The minimal subtable depends on the obtained responses.) The algorithm of Koppen (1989) cannot avoid the theoretical – but in practice uninteresting – worst case where all subsets are states (the expert will give only negative responses and the minimal subtable is the complete table), but it has proved to be applicable to the real-life situation of a set of 50 problems in U.S. high school mathematics. The results of this application will be reported elsewhere (Kambouri et al., 1989). Let us here just mention that actually constructed minimal subtables in this case were in the order of 2000 rows, while the naive “Table 2 version” would contain well over 10^{15} rows. This gives an idea of the reduction that can be obtained once we go beyond the straightforward procedure described in this section for illustrative purposes.

References

- Birkhoff, G. (1937). Rings of sets. *Duke Mathematical Journal*, **3**, 443-454.
- Birkhoff, G. (1967). *Lattice Theory*. Providence, R.I.: American Mathematical Society Colloquium Publications, Vol. XXV.
- Burigana, L. (1988). Organization by rules in finite sequences. University of Padua. *Submitted*.
- Degreef, E., Doignon, J.-P., Ducamp, A. & Falmagne, J.-C. (1986). Languages for the assessment of knowledge. *Journal of Mathematical Psychology*, **30**, 243-256.
- Doignon, J.-P. & Falmagne, J.-C. (1985). Spaces for the assessment of knowledge. *International Journal of Man-Machine Studies*, **23**, 175-196.
- Falmagne, J.-C. (1989). A latent trait theory via a stochastic learning theory for a knowledge space. *Psychometrika*, in press.
- Falmagne, J.-C. & Doignon, J.-P. (1988a). A class of stochastic procedures for the assessment of knowledge. *British Journal of Statistical and Mathematical Psychology*, **41**, 1-23.
- Falmagne, J.-C. & Doignon, J.-P. (1988b). A Markovian procedure for assessing the state of a system. *Journal of Mathematical Psychology*, **32**, 232-258.
- Kambouri, M., Koppen, M., Villano, M. & Falmagne, J.-C. (1989). Knowledge assessment: Tapping human expertise. Dept. of Psychology, New York University. *In progress*.
- Koppen, M. (1989). Extracting human expertise for constructing knowledge spaces: An algorithm. *Mathematical Studies in Perception and Cognition* 89-1, Dept. of Psychology, New York University. *Submitted*. (Chapter 8.)
- Monjardet, B. (1970). Tresses, fuseaux, préordres et topologies. *Mathématiques et Sciences Humaines*, **30**, 11-22.
- Mueller, C. (1989). A procedure to facilitate an expert's judgements on a set of rules. In Degreef, E. & Roskam, E. (Eds.), *Progress in Mathematical Psychology*, vol. 2. To appear.

CHAPTER 8

EXTRACTING HUMAN EXPERTISE FOR CONSTRUCTING KNOWLEDGE SPACES : AN ALGORITHM

(Submitted)

Extracting human expertise for constructing knowledge spaces : An algorithm

Mathieu Koppen

New York University

In a theory for the efficient assessment of knowledge – introduced by Doignon and Falmagne (1985, *International Journal of Man-Machine Studies*, 23, 175-196) and elaborated by them and co-workers in a number of subsequent publications – the cognitive organization of a field of information is represented by a *knowledge space*, the collection of all possible *knowledge states*. The construction of such spaces for particular domains is a problem of some practical importance. We present a method for achieving this by confronting experts in the field with a carefully chosen sequence of questions about specific relationships between the various *items* of the domain.

1. Introduction

Any more advanced computerized instruction system must include a component for assessing the initial knowledge of a student and updating this assessment as the student progresses through the administered course material. In designing such procedures we have to specify very concretely how we represent and measure “the knowledge” of a student. A framework for doing this has been developed by Doignon and Falmagne (1985). They conceptualize a particular domain of knowledge as the collection of *problems* (or *items*) in that field and the *knowledge state* of an individual in this domain is the subcollection of these problems that this individual is capable of solving. In any particular population of students, only some of the subsets of problems will constitute possible knowledge states; this family of feasible states is called the *knowledge structure* of the domain for this population.

This formalization is at the basis of a knowledge assessment project, a comprehensive description of which is given by Falmagne, Koppen, Villano,

This work was supported by DOD grant MDA903-87-K-0002 to Jean-Claude Falmagne at New York University. The author wants to thank Jean-Claude Falmagne for his comments on a previous draft. Address comments and requests for reprints to M. Koppen, Dept. of Psychology NYU, 6 Washington Place, New York, NY 10003.

Johannesen and Doignon (1989). In this context, Falmagne (1989) has developed a procedure for testing a given knowledge structure against empirical data and possibly refining it by a sequence of likelihood ratio tests. Also, a number of procedures have been designed that search such structures in order to locate or at least approximate the state of particular students (Falmagne and Doignon, 1988a; Falmagne and Doignon, 1988b; Degreef, Doignon, Ducamp and Falmagne, 1986). As this description suggests, all these procedures start from a fixed, predetermined knowledge structure. How to arrive at such a structure in the first place is another problem; it is the one we will be concerned with in this paper. We will present here a method for eliciting the necessary information from experts in the field (in our case, for instance, experienced teachers or tutors).

It must be realized that it is quite unrealistic to ask an expert directly for the list of possible knowledge states. For one thing, this list may a priori be very large. In addition, it is more than likely that an expert does not have recourse to an explicit representation of this list; rather, it is something she uses implicitly. Accordingly, we resort to an indirect approach. The theoretical, mathematical basis for this method is provided in Koppen and Doignon (1989). The central concept is the special kind of knowledge structure where the union of any number of knowledge states is again a state. Doignon and Falmagne (1985) argue that this assumption of closure under union is a reasonable one for knowledge structures in practice and they reserve a special name for it: such a structure is called a *knowledge space*.

In Koppen and Doignon (1989) an alternative characterization of knowledge spaces is derived. It is shown that they are in 1-1 correspondence with a class of relations between the different items in our domain of information, and by asking an expert about these relations between items we can in fact recover the knowledge space that, from the viewpoint of this expert, represents the domain. In the next section, we discuss the kind of questions we have to ask the expert and we show how the intended knowledge space can be derived from the expert's answers to these questions. It turns out that we can substantially reduce the number of questions to be asked by making appropriate inferences. In the following sections, we will describe these inferences in some detail and investigate how they can be exploited to turn an impractical straightforward approach into a querying procedure that has proved to be workable in practice. In the last two sections, these discussions of the various aspects converge in a description of the resulting algorithm and we add some remarks about the applicability of this method.

2. Deriving the knowledge space from the expert's answers

We consider a field of knowledge, represented by a (finite) set X of problems. We assume that, implicitly, experts in a field have some knowledge space of that field. We want to make this knowledge space explicit, but we cannot do that by simply asking an expert to list all the possible knowledge states. Instead, we use an indirect method and ask her questions of the following kind:

[Q] *“Suppose that a student under examination has just failed all the problems a_1, a_2, \dots, a_n . Is it then practically certain that this student will also fail problem b ?”*

The expert is told to assume optimal examination conditions, meaning that there is no “noise” in the data in the form of lucky guesses or careless errors by the student. We suppose that, in this situation, the response of the expert reflects a particular knowledge space consulted implicitly. Defining $A = \{a_1, a_2, \dots, a_n\}$ as the subset of X that is known to be failed by the hypothetical student, our interpretation of the above question is:

[I] *Does it hold for the knowledge space in question that there is no state that is disjoint from A and contains the element b ?*

Indeed, if there were some such state, a student in that state would fail all problems a_1 to a_n , yet give a correct answer to b . It is clear now that the answers to the questions [Q] tell us something about the knowledge space. In Koppen and Doignon (1989) it is formally shown that a knowledge space is in fact completely characterized by the answers to all possible questions of the form [I] (i.e., [Q]). That is, from these answers we can fully recover the knowledge space that was operative in producing them and we want to describe here how this can be done.

2.1. A straightforward approach.

In a sense, our interpretation [I] tells us exactly what to do. A priori, before obtaining any response from the expert, there are no restrictions on the knowledge space; that is, we consider every subset of X as a possible knowledge state. So, we draw up the list of all these subsets and start querying the expert. Whenever the expert indicates that, indeed, failing some set of problems A would imply failing some problem b , we apply [I]: we go through this list and remove as a possible knowledge state every set that contains b but is disjoint from A .

Notice that after applying this operation the resulting, trimmed collection of states is still a space, i.e., it is still closed under union. (Obviously, the power set that we start with is a space.) Indeed, if $K = K_1 \cup K_2$ is removed because of the above observation, that is, if $b \in K$ and $K \cap A = \emptyset$, then both $K_1 \cap A = \emptyset$ and $K_2 \cap A = \emptyset$ while also $b \in K_1$ or $b \in K_2$. Thus, K_1 or K_2 has also been removed: the

operation cannot create a counterexample for the closure under union. We will come back to this issue in the final section.

In several respects, however, the above straightforward procedure is too simplistic. For one thing, it does not take into account that we can save us a lot of queries from the expert by drawing appropriate inferences from her answers. Secondly, “drawing up the list of all these subsets” is not much of a problem when we are dealing with, say, 5 items ($2^5 = 32$), but it becomes rather prohibiting when the number of items equals 50 ($2^{50} > 10^{15}$). (The algorithm that we are going to describe has been applied with this number of items.) Fortunately, by fully using all inferences from the expert’s answers, we can, in general, greatly reduce the collection of subsets that require consideration. Rather than starting with the full tableau of subsets of X , the algorithm will generate subsets that are of interest dynamically, in the process. These two related issues, the question of drawing inferences and the question of using these inferences to generate a minimal subtable of subsets of X will be dealt with in some detail in the next two sections.

2.2. The relation P .

Below we give examples of some basic inferences and show how these are actually implied by the above mentioned straightforward procedure. Next we describe how the knowledge space can be recovered from the collection of all inferences. This means that the knowledge space can remain implicit; we need only maintain the table of inferences and how this can be done most economically is the subject of the next sections. We denote the situation where failing A entails failing b by $A P b$. This defines P as a binary relation between subsets of items and items. The expert’s task is it to tell us, in principle for each pair (A, b) , whether or not $A P b$ holds.

2.3. Example.

With this interpretation, then, it is clear that $A P b$ implies $A' P b$ for any set $A' \supseteq A$. This inference is actually already put into practice when we implement $A P b$ according to the interpretation [I]: after all sets that contain b and are disjoint from A have been removed from the list of possible states, there certainly are no sets left containing b and disjoint from A' . Thus, a “new” observation $A' P b$, with $A \subseteq A'$, would be totally uninformative and we would not want to ask the expert the corresponding question.

2.4. Example.

Suppose we have observed $a P b$ and $b P c$. (If the left argument to the relation P is a singleton set $\{a\}$, we simply write $a P b$ for $\{a\} P b$, etc.) According to the interpretation of P , it is then tempting to conclude that we must also have $a P c$.

(That is, restricted to single items, the relation \mathbf{P} is transitive.) Again, this is really a foregone conclusion. Consider a subset of X that contains c but not a . This set either contains b , in which case it has been removed by the observation $a\mathbf{P}b$, or it does not contain b , in which case it has been removed by $b\mathbf{P}c$.

2.5. Introduction of the complementary relation \mathbf{N} .

This transitivity property of \mathbf{P} makes it clear that there may also be inferences from negative answers. We denote by \mathbf{N} the complement of \mathbf{P} , that is, we write $A\mathbf{N}b$ when it is not the case that failing all items in A implies failing the item b . In terms of the knowledge space, $A\mathbf{N}b$ indicates that there is indeed a state that is disjoint from A and that contains b .

2.6. Example.

Suppose now that we have observed $a\mathbf{P}b$ and $a\mathbf{N}c$. Then we must conclude $b\mathbf{N}c$, since the alternative $b\mathbf{P}c$, together with the observation $a\mathbf{P}b$, would produce the inference $a\mathbf{P}c$, contradicting the observation $a\mathbf{N}c$.

Inferences can get more complicated. The transitivity of \mathbf{P} in the domain of singleton sets, for instance, generalizes in the following way.

2.7. Example.

Suppose we have $A\mathbf{P}b_i$ for some subset A and items $b_i, i = 1, \dots, k$, and suppose that, for the subset $B = \{b_1, \dots, b_k\}$ and item c , we have also observed $B\mathbf{P}c$. As in the single item case above, the interpretation of \mathbf{P} given in 2.2 would lead us to conclude $A\mathbf{P}c$ and we note that this conclusion has already been implemented: a subset containing c and disjoint from A has been removed by $B\mathbf{P}c$ if it is disjoint from B and by $A\mathbf{P}b_i$ if it contains b_i . So we do not want to bother the expert with a question involving a set A and an item c , whenever there is some set B for which we have observed (or inferred!) $B\mathbf{P}c$ and for each element b of which we know $A\mathbf{P}b$ to hold. There are corresponding generalizations for negative inferences.

2.8. Extended versions of the relations.

This example suggests a natural way of extending the use of \mathbf{P} to the case where both arguments are subsets. If A, B are subsets of X , this extension is defined by

$$A\mathbf{P}B \text{ iff } (A\mathbf{P}b \text{ for all } b \in B). \quad (1)$$

This defines effectively \mathbf{P} as a relation on the power set of X and the above example shows that, as such, \mathbf{P} is transitive. We can make the corresponding extension of \mathbf{N} to the relation on the power set that is the negation of \mathbf{P} . Accordingly,

$$A \cap B \text{ iff } (A \cap b \text{ for some } b \in B). \quad (2)$$

2.9. Recovering the knowledge space from the relation.

All of this shows that it is not a good idea to simply translate a response of the expert into its consequences for the knowledge space under construction, and, after that, forget about it. We will, in fact, take quite the opposite approach: we collect all observations and inferences in a subsets-by-items table and leave the implied knowledge space implicit until the end. It turns out that the final knowledge space can be easily read off from the completed table. To see this, suppose that for any subset $A \subseteq X$ and any $x \in X$ we have established whether $A \text{ P } x$ or $A \text{ N } x$ holds. Consider now, for some $A \subseteq X$, the collection A^* of items that any student will fail who fails all items in A :

$$A^* = \{x \in X : A \text{ P } x\}. \quad (3)$$

This mapping, associating with each subset A the above defined A^* , has the following properties:

$$A^* \supseteq A, \quad (4)$$

$$A \supseteq B \text{ implies } A^* \supseteq B^*, \quad (5)$$

$$(A^*)^* = A. \quad (6)$$

These properties are immediate from the interpretation of P and definition (3). Specifically, (6) represents the fact that, observing a student who fails all of A , we may conclude that he will fail all of A^* , *but nothing more than that*: it is possible to master exactly those problems that are not in A^* . In other words, the set of these problems, $X - A^*$, is a knowledge state. We denote this complement of A^* by A^\perp ; it is the collection of items that can still be solved by a student who is known to fail all of A :

$$A^\perp = \{x \in X : A \text{ N } x\}. \quad (7)$$

We have seen that any set A^\perp is a knowledge state and the argument can easily be reversed. If $K \subseteq X$ is a knowledge state, then there may exist a student who solves correctly all of K , but fails all of $X - K$. Consequently, from failing $X - K$ we can surmise no more than just that: $(X - K)^* = X - K$, and, thus, $K = (X - K)^\perp$. This shows how the knowledge space, let us call it \mathbf{K} , can be extracted from the table in which the expert's answers and resulting inferences are stored:

$$\mathbf{K} = \{A^\perp : A \subseteq X\}. \quad (8)$$

Accordingly, the problem of determining, in an efficient way, the knowledge space

of an expert amounts to the problem of determining, for any subset A and any item b , whether the expert thinks APb or ANb is the case. In the following sections we will investigate how the latter task can be handled most efficiently, using inferences like the ones shown above.

3. Collecting the inferences

Inferences can be made by exploiting certain properties of the relation P . In order to develop the possible inferences in full generality, we will consider the definitions (1) and (2) of P and N , respectively, as relations between subsets of X . We have to keep in mind, however, that, for practical reasons, we can maintain only the original version of these relations in a subsets-by-items table. Accordingly, we will at the end of this section give the appropriate version of the inferences. As we have indicated, maintaining a full subsets-by-items table is impractical even with a moderate number of items and in the next section we will investigate how to generate a minimal part of this table that still contains all the information.

3.1. The basic properties.

Our interpretation of P shows that APx must hold for any $x \in A$. Combining this with the extension definition (1), the following properties follow easily: for any $A, B \subseteq X$,

$$A \supseteq B \text{ implies } APB \quad (9)$$

and

$$APB \text{ iff } AP(B \cup A). \quad (10)$$

We have already observed the transitivity of P : for any $A, B, C \subseteq X$,

$$APB \ \& \ BPC \text{ implies } APC. \quad (11)$$

In Koppen and Doignon (1989) it was proved that (9), (10) and (11) fully characterize the relations that are generated by the positive answers to the questions $[Q]$ (II). So these are the properties we can use in making inferences.

The situation in which we want to generate inferences is as follows. For each pair (A, B) of subsets of X we want to establish whether APB or ANB holds. By collecting information (responses of the expert) and making inferences from it, we have arrived at a decision for some of these pairs. Now for some pair, (A, A') say, information is obtained, resulting in either APA' or ANA' , and the question is: for which of the pairs of sets (B, C) that are still undecided does this new information

provide the missing link that makes an inference, either BPC or BNC , possible ?

3.2. Inferences from transitivity.

Let us start with the more familiar case of inferences just due to the transitivity condition (11). Consider first the case where the new information is positive: we add APA' to our list. This allows a positive inference for pairs (B, C) for which this new entry completes a path in \mathbf{P} from B to C , that is, for which there were already such paths from B to A and from A' to C :

$$APA' \text{ implies } BPC \text{ whenever } BPA \ \& \ A'PC. \quad (12a)$$

The new observation APA' can also lead to negative inferences for pairs (B, C) . One possibility is that, in \mathbf{P} , there is a path from A' to B , while it is known that there is no such path from A to C :

$$APA' \text{ implies } BNC \text{ whenever } A'PB \ \& \ ANC. \quad (12b)$$

A negative inference is also obtained when, in \mathbf{P} , A can be reached from C , while it is known that A' cannot be reached from B :

$$APA' \text{ implies } BNC \text{ whenever } CPA \ \& \ BNA'. \quad (12c)$$

Finally we look at the case where the new information is negative: ANA' is obtained. This allows only negative inferences BNC and these are obtained whenever the alternative BPC would complete a path in \mathbf{P} from A to A' :

$$ANA' \text{ implies } BNC \text{ whenever } APB \ \& \ CPA'. \quad (12d)$$

These four inference rules are represented in Figure 1 in the form of diagrams. Note that, in this whole discussion, paths may have zero length; in the last case, for instance, we might have $A = B$ or $C = A'$, and similarly in (12a,b,c).

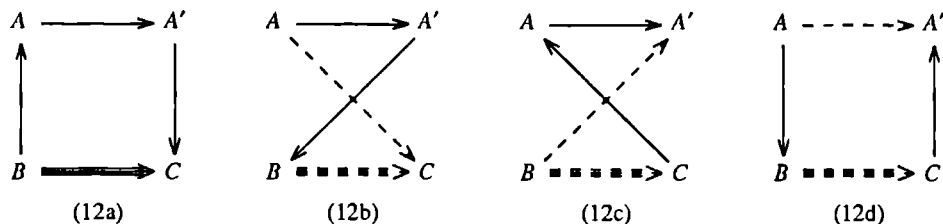


Figure 1. Diagrams of the inferences from transitivity given in rules (12). Solid arrows connect pairs of sets known to be in \mathbf{P} , dashed arrows represent pairs known to be in \mathbf{N} . In each diagram, the new observation is in the top row and the thick arrow in the bottom row is the inference implied by the other arrows.

3.3. The general case.

The inferences in (12) are solely based on the transitivity (11). They may be strengthened by incorporating the properties (9) and (10). The latter tells us that the observation APA' is equivalent to $AP(A' \cup A)$ (thus, ANA' equivalent to $AN(A' \cup A)$). Similarly, BPC and BNC are equivalent to $BP(C \cup B)$ and $BN(C \cup B)$, respectively. This allows us to replace, in any condition in (12a) to (12d), A' by the equivalent $A' \cup A$ and C by the equivalent $C \cup B$. In each rule we may choose the alternative that makes the condition the weakest (and thus the corresponding inference rule the most powerful). In order to apply (12a), for instance, we have to find a subset B such that BPA and a subset C such that either $A'PC$, or $(A' \cup A)PC$, or $A'P(C \cup B)$, or $(A' \cup A)P(C \cup B)$. The weakest of these conditions is $(A' \cup A)PC$, in view of the facts that the relation P contains the relation \supseteq (property (9)) and that P is transitive (11). The alternative condition $A'P(C \cup B)$, for instance, leads to $(A' \cup A) \supseteq A'P(C \cup B) \supseteq C$, which, by (9) and (11), implies $(A' \cup A)PC$. Similarly for the other two alternatives. In general, whenever there appears in a condition a pair of sets that is in the relation P , we get the weakest form of this condition by making the left member of this pair as big and the right member as small as possible. Accordingly, we choose $A' \cup A$ in (12a,b), A' in (12d); C in (12a), $C \cup B$ in (12c,d).

It looks as though the same kind of reasoning with respect to the pairs in N appearing in the conditions of rules (12b,c) indicates that we should make the right member of such a pair as big as possible. In (12b), for instance, it is clear that, by (9), ANC implies $AN(C \cup B)$, so the latter condition is weaker. On closer inspection, however, the two versions appear to be completely equivalent. Any B we are considering here satisfies the other condition of (12b): $A'PB$. Together with the new observation APA' this implies APB . So whenever ANx for some $x \in C \cup B$, we cannot have $x \in B$: $AN(C \cup B)$ implies ANC . We do not gain, then, by replacing C by $C \cup B$ in (12b) and in practice we only lose. We deal with the relations P and N here in terms of their extended definition, as relations between subsets, to describe the inference rules in full generality, but in practice these relations are only available in their restricted form, as relations between a subset and an item. That is, the condition ANC is established by verifying that indeed ANx for some $x \in C$ and, obviously, it does not make any sense to extend the domain of search to $C \cup B$, if we know that we will not find such an x in B . Therefore, we do not replace C by $C \cup B$ in (12b) and, for the same reason, we choose A' over $A' \cup A$ in (12c).

All in all, by taking the properties (9) and (10) into account we have arrived at the following transformation of the inference rules (12):

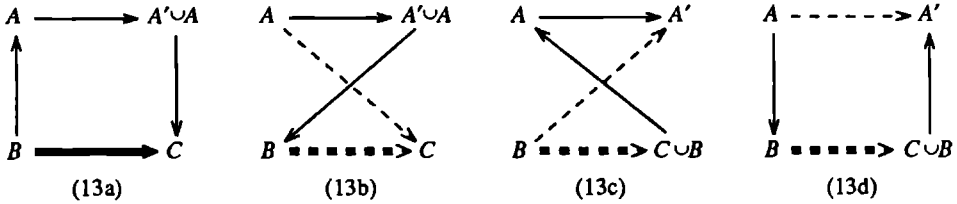


Figure 2. Diagrams of the general inference rules (13). Same conventions as in Figure 1.

$$APA' \text{ implies } BPC \text{ whenever } BPA \text{ \& } (A' \cup A)PC, \quad (13a)$$

$$APA' \text{ implies } BNC \text{ whenever } (A' \cup A)PB \text{ \& } ANC, \quad (13b)$$

$$APA' \text{ implies } BNC \text{ whenever } (C \cup B)PA \text{ \& } BNA', \quad (13c)$$

$$ANA' \text{ implies } BNC \text{ whenever } APB \text{ \& } (C \cup B)PA'. \quad (13d)$$

See Figure 2 for the corresponding diagrams. These are the general inference rules for a relation P between subsets of X that satisfies (9), (10) and (11), and where N denotes the complement (negation) of P .

3.4. The practical implementation.

Now we want to apply these rules in the actual situation where we deal with P and N only in their restricted version, as relations between subsets and items. Note that the conditions in (13) can all be checked in this restricted version. According to definition (1), the condition APB is fulfilled if and only if for all $x \in B$ it has been established that APx . Similarly, by (2), ANB follows if and only if there is some $x \in B$ for which ANx has been established. The restricted version forces only changes in the input and output of the rules (13), the new observation and the derived inference.

For one thing, our expert provides us with answers to questions of the form $[Q]$; that is, the new information we get does not concern a pair of subsets A and A' , but a subset A and an item x . Accordingly, in the above rules (13) we may substitute x for A' everywhere. As for the new inferences, we can, in the restricted version of the relations P and N , only handle conclusions of the form $BP y$ or $BN y$. This is no real restriction in (13a), where the conclusion is positive. Indeed, the conclusion BPC is equivalent to the collection of conclusions $BP y$ for all $y \in C$, and these can all be implemented.

The situation is more complicated, however, for the negative conclusions BNC in (13b,c,d). This means that for some $y \in C$ it must be true that $BN y$, but unless we have already established $BP z$ for every but one $z \in C$, we would not know which

$y \in C$ to pick. (If there is already some established BNz , $z \in C$, then the conclusion BNC does not tell us anything new, of course.) We must conclude that the negative inference BNC cannot be implemented in full generality in the restricted version of P and N . Since it is utterly impractical to work with the extended version of these relations, we are content with just collecting all inferences of the form BNy . This restriction amounts to substituting y for C everywhere in the rules (13).

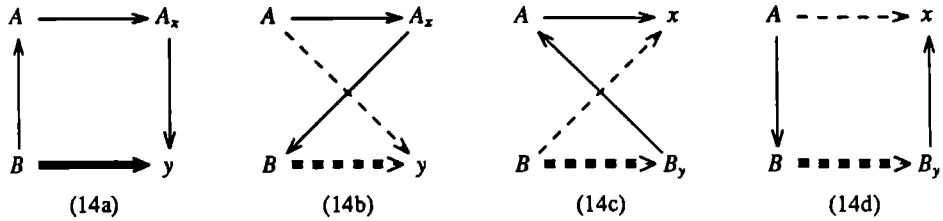


Figure 3. Diagrams of the implemented rules (14). Same conventions as in Figures 1 and 2.

In sum, then, the following specialization of (13) is the collection of inference rules that are actually applied in our algorithm. We use the shorthands A_x for $A \cup \{x\}$ and B_y for $B \cup \{y\}$; the corresponding diagrams are given in Figure 3:

$$APx \text{ implies } BP_y \text{ whenever } BPA \ \& \ A_xPy, \quad (14a)$$

$$APx \text{ implies } BNy \text{ whenever } A_xPB \ \& \ ANy, \quad (14b)$$

$$APx \text{ implies } BN_y \text{ whenever } B_yPA \ \& \ BN_x, \quad (14c)$$

$$ANx \text{ implies } BN_y \text{ whenever } APB \ \& \ B_yPx. \quad (14d)$$

So, we must be aware that, in our tables, we do collect all positive, but not necessarily all negative inferences. It appears, in practice, that rules (14b,c,d) still find most negative inferences. However, they are in principle incomplete and how we deal with this situation will be discussed in Section 5.

4. Generating a minimal subtable of inferences

As we already indicated in Section 2, even for a moderately large number of items it becomes soon impractical to generate a full subsets-by-items table. Fortunately, it is, in general, not necessary to construct this full table. The inferences we described in the previous section may render complete rows of such a table uninformative. That is, there may be a number of subsets A of X , such that for any $b \in X$, it can be inferred whether APb or ANb holds by inspection of some other rows of the table.

It would be nice, of course, if we could avoid including all these redundant subsets in the table. Below we will describe a method for generating a minimal subcollection of the power set of X , such that the corresponding subtables of \mathbf{P} and \mathbf{N} still contain all necessary information.

Since we are dealing here with rows of the table, it appears that the properties of the relation \mathbf{P} can most conveniently be used in the form of Eqs. (4) to (6). In the previous section, the inferences were based on the properties (9) to (11), but it is not difficult to establish that the two sets are equivalent: in Koppen and Doignon (1989) it is shown that, under definition (3), the properties (4) to (6) and the properties (9) to (11) characterize the same collection of relations between subsets of a finite set X . Using (4), (5) and (6), we can find some important examples of redundant subsets.

4.1. Example. Suppose we have $A \subseteq B \subseteq A^*$ for subsets A, B of X . By (5), this implies $A^* \subseteq B^* \subseteq (A^*)^*$ and by (6) it follows $A^* = B^*$. Thus, whenever we have determined A^* for some subset A , we have in fact determined B^* for all subsets B that are between A and A^* . It is of no use to have such subsets in the table: since $B^* = A^*$, it is also the case that $B^\perp = A^\perp$, so B does not add an additional state to the knowledge space defined by (8). Also, any inference involving the row corresponding to B can be made on the basis of the row corresponding to A , so the table of inferences is still complete.

4.2. The case of equivalent subsets.

We call two subsets A and B of X *equivalent* whenever $A^* = B^*$. (This means they have identical rows in the subsets-by-items table.) From (7) and (8) it is clear that equivalent subsets contribute one and the same state to the knowledge space, so for recovering the complete space only one representative of a collection of equivalent subsets is needed in the subtable. However, the primary purpose of the subtable is to collect all inferences and this implies that two equivalent subsets must both be in the subtable whenever their equivalence is only established empirically, on the basis of the expert's answers to questions regarding these subsets. In this case, the two rows contain independent information on which inferences can be based. Only if we can *infer*, ahead of time, that some subset is equivalent to another subset, then that subset can be discarded from the subtable without throwing away inferences. The above example is an important case of this situation; below is another one.

4.3. Example. Let A, B be two subsets of X such that $A^* = B^*$ and consider an arbitrary subset C . Since $(A \cup C) \supseteq A$, we can use (5) and (4) to obtain $(A \cup C)^* \supseteq A^* = B^* \supseteq B$; also $(A \cup C)^* \supseteq C$. Thus $(A \cup C)^* \supseteq (B \cup C)$, and, by (5) and (6), $(A \cup C)^* \supseteq (B \cup C)^*$. We have the same derivation with the roles of A

and B interchanged, so we conclude that $(A \cup C)^* = (B \cup C)^*$. Consequently, whenever we have found two subsets A and B for which $A^* = B^*$, we need to consider in our table only supersets of one of the two. Any $B' \supseteq B$ is known to be equivalent to the set $A \cup (B' - B)$, which is a superset of A .

4.4. Partitioning the table into blocks.

Once it is clear that all information may be contained in a subtable that is often of considerably smaller size, the problem presents itself of generating just this subtable instead of the full table. The approach we have chosen here is to proceed by *blocks*, where a block is defined by the cardinality of the set A in [I], or, the number n in [Q]. First we ask the expert for any single item a and for any b whether failing a would imply failing b . This collection of questions constitutes block 2. Next comes block 2, in which the generic question is whether failing a pair a_1, a_2 would imply failing b . Then, in block 3, we consider a triple a_1, a_2, a_3 , etc.. There are good reasons for choosing this ordering of the questions. It would seem that, for the expert, questions are easier when there are less items involved, so in this way she starts with the easiest questions. By using the inferences, we can then try to minimize the number of more difficult questions that have to be asked later on. And, in fact, we get more inferences from a positive answer $A \text{ P } b$ the smaller the set A is. This is clear from Example 2.1: a smaller set simply has more supersets. It can also be seen in the Examples 4.1 and 4.2 above. If A is small, there may be more "room" between A and A^* in 4.1 and if A and B are small in 4.2, there are more supersets of B that can be represented by corresponding supersets of A .

4.5. The construction of a new block.

The idea is not only to proceed by block, but to actually construct block k of the table only after blocks 1 to $k-1$ have been completed. These previous blocks are then consulted to decide which subcollection of the k -subsets of X needs to be considered in block k . That is, the table is constructed dynamically, depending on the answers and inferences obtained so far. At the start, we do not know anything about equivalent subsets, so block 1 is constructed to consist of all singleton subsets. Block k is constructed in such a way that

- (i) it includes only those k -subsets of X that are not known to be equivalent to some subset in the previous blocks 1 to $k-1$, and
- (ii) it includes only one representative of any number of k -sets that are known to be mutually equivalent.

Condition (i) provides us with a termination condition: we know that blocks 1 to k contain all information regarding \mathbf{P} – and, thus, all information regarding the knowledge space – if and only if the next block to be constructed, block $k+1$,

appears to be empty: any subset not in blocks 1 to k has a representative there to which it is equivalent.

How is it actually decided which subsets are to be in block k ? First, because of the above conditions, we need to consider for “admission” to block k only k -subsets that are of the form $A \cup \{x\}$ where A is a subset appearing in block $k-1$ and where $x \notin A^*$. Indeed, if A is a $(k-1)$ -subset that is not present in block $k-1$, then, by construction, there must be a subset B , somewhere in blocks 1 to $k-1$, to which A is equivalent. But then the k -subset $A \cup \{x\}$ is equivalent to $B \cup \{x\}$ (cf. Example 4.2) and thus this subset needs no separate consideration. And if $x \in A^*$, then, by (4), $(A \cup \{x\}) \subseteq A^*$, and, by (5) and (6), $A^* \subseteq (A \cup \{x\})^* \subseteq (A^*)^* = A^*$, so in this case $A \cup \{x\}$ is equivalent to (and can be represented by) the subset A in the previous block.

Secondly, if a number of subsets A_1, \dots, A_m in block $k-1$ have turned out to be equivalent, then we pick one $i \in \{1, \dots, m\}$, say $i = 1$, and we do not include in block k any extension of the sets A_2, \dots, A_m . The only extensions we consider are the sets $A_1 \cup \{x\}$ with $x \notin A_1^* = \dots = A_m^*$. This is a direct application of Example 4.2.

Finally, any remaining candidate must be subjected to the test of Example 4.1. That is, if B is such a set, then we have to go through the blocks 1 to $k-1$ in order to check whether there is some subset A such that $A \subseteq B \subseteq A^*$. Only if this is not the case, is the subset B finally accepted as a member of block k .

4.6. The transport of previous inferences to a new block.

Until now we have bypassed one crucial problem. In Section 3 we discussed what inferences could be drawn from observed responses of the expert and in that discussion it was tacitly understood that we were in a position to actually implement any of these inferences. We acted as if the full subsets-by-items table were available, while we have seen above that in fact only a subtable is created, dynamically – block by block. Clearly, if we are dealing with block k , we can, on the one hand, not implement any inferences pertaining to following blocks – they are not there, yet – and on the other hand we do not need to implement any inferences in the previous blocks – they have already been completed. The conclusion must be that in our algorithm the inference rules (14) are only applied within the current block. This conclusion is correct, but it poses the problem of how to recover, after the construction of a new block, the inferences for that block that we “forgot” to implement when we were dealing with a previous block. In fact, a full recovery is possible, again by switching from the properties (9), (10) and (11) to the equivalent set (4), (5) and (6).

Let A be a subset in a newly constructed block k . Then collecting all positive inferences APx amounts to finding the maximal subset A^+ of X such that from the data in blocks 1 to $k-1$ it can be inferred that $A^+ \subseteq A^*$. Without looking at any data, but just at (4), we know $A \subseteq A^+$. We also know, by (5) and (6), that whenever $B \subseteq A^*$, then in fact $B^* \subseteq A^*$. This means that the same holds with respect to A^+ : for any subset B we may from $B \subseteq A^+$ infer $B^* \subseteq A^+$. These observations lead to the following procedure for determining A^+ :

- (i) Set initially A^+ to A .
- (ii) While there is a subset B in blocks 1 to $k-1$ such that $B \subseteq A^+$ and $B^* \not\subseteq A^+$, add B^* to A^+ .

Note that the loop (ii) will not be entered when $k=1$, and indeed in this case the initialization step (i) gives the correct result. Note also that the while-loop cannot be implemented by one pass through blocks 1 to $k-1$: each time that A^+ is adjusted, it grows and therefore previously considered subsets B that were rejected because $B \not\subseteq A^+$ must be reconsidered.

Similarly, collecting the negative inferences ANx amounts to finding the maximal subset A^- of X such that $A^- \subseteq A^\perp$. Since $C \subseteq A^\perp$ is equivalent to $A^* \subseteq (X-C)$, we have to use the properties (4), (5) and (6) to derive bounds on A^* . Well, from (5) and (6) it follows that whenever $A \subseteq B^*$, then also $A^* \subseteq B^*$, for any subset B . Taking complements we see that $B^\perp \subseteq A^\perp$ for any B such that $A \subseteq B^*$ and we obtain the following procedure for finding A^- :

- (i) Set initially A^- to the empty set.
- (ii) For every subset B in blocks 1 to $k-1$ such that $B^* \supseteq A$, add B^\perp to A^- .

Again, in block 1 only the initialization step (i) is executed, with the correct result. This time, the loop (ii) is just one pass through the completed blocks.

4.7. Summary.

This completes the description of how the relevant subtable is constructed. To recapitulate, the construction proceeds by blocks of subsets of equal cardinality, starting with the singleton subsets. After a block has been completed, the next block is constructed according to the specifications in 4.5. Next, the procedures of 4.6 are applied in order to collect the positive and negative inferences for this new block that can be drawn from the data in the previous, completed blocks. Only then do we start asking the expert questions from the collection in this block for which no inference was obtained. In this process, we derive new inferences within this block by applying the inference rules (14) of the previous section. The whole procedure terminates when the newly constructed block is in fact empty (and this is surely the case after block $|X|$).

5. Choosing the next question

In the preceding two sections we have dealt with the main design issues. Before we move to a presentation of the resulting algorithm in the next section, we discuss here one remaining question that is of a more practical nature. Nevertheless, the general solution to the problem of deciding which of the open questions to ask the expert next will also give us a way to deal with the theoretical problem of the negative inferences being incomplete (see Section 3.4).

On a general level, the order of the questions follows from the above design decisions: the questions are presented in blocks. In block 1, the expert has to consider whether (negative) information regarding a single item is sufficient to arrive at a (negative) conclusion for another item; in following blocks, information regarding pairs, triples, etc. of items is offered. Within one block, however, we are completely free to pick one or the other out of the pool of open questions in that block as the next one to be asked. It will be practical considerations that guide us in this choice. Our objective is to minimize the work load of the expert; we would like to optimally use possible inferences so as to minimize the total number of questions that have to be asked. But this criterion is not practicable in its generality: it would require us to compute, for each of the questions under consideration, all possible continuations.

To become practical, we drastically limit our scope: we try to find a best question to ask next by considering the number of inferences each candidate question would yield when asked at this moment. Since the table of inferences is created dynamically (see previous section) and thus the later blocks are not there, “number of inferences” has to be interpreted with the qualification “in the current block”. This attempt at optimization by looking at the immediate consequences of a question may be compared with trying to win a chess game by looking just one move ahead: it does not guarantee success, but it certainly beats doing moves at random.

5.1. Criteria for a “best” question.

It is not obvious what constitutes a “best” question, since which inferences are made, and thus also the “number of inferences” alluded to above, will depend on the expert’s answer. Several approaches are conceivable, all based on different ways of combining two values: the numbers of inferences in case of a positive and a negative answer, respectively. If we think the expert is – overall – equally likely to say “yes” or “no”, the mean – or, equivalently, the sum – of these two numbers is a good measure: we pick a question for which this sum is maximal. If we think there is some other fixed probability for a “yes” answer, we maximize instead a weighted sum, the weights being determined by the probability of the corresponding event.

We can try to be more sophisticated and estimate dynamically the probability of a “yes” answer (based on observed relative frequencies so far). It seems doubtful, however, that we will obtain consistent estimates in this way; the probability in question may change drastically over the time course of the query procedure (for instance, going from one block to the next), so the estimates based on prior information may be misleading.

Effects may also vary widely for different experts and different sets of items and to get some hold on what criterion is best in what situation would require extensive empirical investigations. It must be noted that the choice of a criterion is indeed fully a matter of practical considerations: one criterion may be substituted for another without any consequence for the rest of the algorithm.

5.2. The adopted selection rule.

The selection rule that we implemented satisfies a kind of “maximin” criterion, which may be regarded as more conservative in that it not so much tries to maximize the immediate gain, as it tries to minimize the immediate cost resulting from a very poor question. For each candidate question we consider the number of inferences for a positive answer and that for a negative answer, and we select the questions for which the minimum of these two numbers is maximal over all questions. From these, a best question is picked as one for which the total number of inferences – or, equivalently, the “other” number – is maximal.

In practice it is not always feasible to consider all the open questions in this process. Whenever the number of candidate questions is too large, we just select a sample and pick the best from this sample. The idea is again that – especially when there are many open questions – there may be many questions that are “best” or approximately so and that the important thing is to avoid picking a particularly bad question.

5.3. Dealing with the incomplete negative inferences.

Regardless of what criterion is used for choosing the next question, the point that the possible inferences of a question are computed before the question is actually posed to the expert is of some significance in another respect. It allows us to discover negative inferences that were missed by the incomplete inference rules (14b,c,d). For any question that we consider in the selection procedure, we compute the list of inferences for a negative answer and the corresponding list for a positive answer. Whenever this last list contains a contradiction with the inferences established thus far, we must conclude that a positive answer to this question is excluded by the data collected earlier. In other words, we have a negative inference for this question; since the list of inferences for a negative answer is still available, we can add all

these inferences to the table. After processing this pseudo-observation, we have to restart the selection procedure, since the collection of open questions and, more generally, the available data have changed. By picking up missed negative inferences in this way, we make sure that any question posed to the expert is indeed an open question: both answers are compatible with the existing body of data and, thus, the obtained answer will really be an additional piece of information.

6. The algorithm.

Now we are ready to put the pieces together and present the integrated algorithm for obtaining a knowledge space from an expert's answers to questions of the form [Q]. The algorithm consists of two parts; first it is established – directly or indirectly – for any subset A of X and any $x \in X$ whether APx or ANx holds. This is the main part; the second part consists just in translating the constructed table into the corresponding knowledge space. In the description of the construction of the table, we will fall back on notation used in Section 4.6: for any subset A , A^+ is the variable denoting the maximal subset for which *at the moment* $A^+ \subseteq A^*$ has been established. Similarly, A^- is the variable corresponding to the maximal subset for which currently $A^- \subseteq A^+$ has been verified. We have complete information regarding a subset A whenever $A^+ \cup A^- = X$; then necessarily $A^+ = A^*$ and $A^- = A^\perp$.

In the following, the algorithm is presented in a top-down fashion. Expressions in italics are names of procedures whose definition follows; plain text is supposed to be self-explanatory and words in bold face denote key words of our “programming language”.

Main program:

```

Initialize first block;
while new block is non-empty
do
  Fill block;
  Construct next block
od;
Output space.

```

Initialize first block:

```

Generate singleton sets  $\{x\}$ , for all  $x \in X$ ;
Initialize  $\{x\}^+ := \{x\}$ ,  $\{x\}^- := \emptyset$ , for all  $x \in X$ .

```

(*Initialize first block* is just a special case of *Construct next block*.)

Fill block:

```

while open question left
do
  Choose next question;
  Obtain answer;
  Add appropriate list of inferences to table
od.

```

Choose next question:

```

decide on sample size;
for sample size number of questions
do
  Collect inferences for negative answer;
  Collect inferences for positive answer;
  if this is currently best question
  then save the two lists of inferences
fi
od.

```

Collect inferences for negative answer:

apply rule (14d) of Section 3.4.

Collect inferences for positive answer:

apply rules (14a,b,c) of Section 3.4;

if contradiction with existing data

then

add list of inferences for negative answer to table;

goto "open question left" test in *Fill block*

fi.

Construct next block:

find all subsets in new block by applying Section 4.5;

initialize A^+ and A^- for each subset A by applying Section 4.6.

Output space:

output all sets $A^- = A^\perp$ collected in the table.

(According to (8) in Section 2.9 these sets constitute the knowledge space.)

This completes the description of the algorithm and in the final section we will turn to some questions regarding its practical applicability.

7. Applying the algorithm

Any evaluation of the performance of the algorithm presented in the previous section must start with the observation that there is a terrible worst case. If, actually, any subset of X is a knowledge state, then we will not observe any positive answer from the expert, the inference rules will not give us any additional inferences and the minimum subtable containing all the information is the full subsets-by-items table. To fix ideas, for a 50 item set this would mean that we would have to complete a 2^{50} by 50 table by asking an expert $2^{50} \cdot 50/2 \approx 2.8 \cdot 10^{16}$ questions. This could take some time.

There are several reasons why this worst case, in itself, need not bother us too much. First, cases where a domain of knowledge has hardly any structure at all (the constituting notions can be acquired in almost any order) are rare and generally not very interesting. Second, such fields, if they appear, are usually recognizable as such from the outset. Consequently, we need not question many experts – and certainly not in the format of the above procedure – to find out that “almost anything goes”. Third, even if we would start the procedure in such a case, we would find out about the lack of structure after completing the first block. It is hard to come up with fields of information where there are no – or very few – implications from a single item to another item, but where there suddenly are many from pairs or triples (etc.) of items. That is, if we get hardly any positive responses of the expert in block 1, it seems a safe bet that the situation will not get much better later on and we may decide to stop the procedure. Completing the first block will, in the worst case, take a number of questions that is quadratic in the number of items.

What remains as a lesson from consideration of the worst case is the fact that there are no convenient theoretical bounds on time and space requirements. Performance will depend strongly on the amount of “structure” present in the chosen domain of application and to find out how well the algorithm is suited to its task, we simply have to apply it to domains of knowledge that are interesting in practice. Such a practical test has actually been performed: a number of experts went through the presented procedure where the set X consisted of 50 problems in high school mathematics. Indeed, the algorithm fared pretty well here. More extensive analysis of the results of this application will be the subject of a forthcoming paper; we will drop here just some numbers that give an indication of what may be expected in practical cases.

With one exception (that we will go into below) the experts finished the procedure by completing no more than five blocks. (Note that theoretically 50 blocks might be necessary.) Thus, after no more than 5 blocks, *Construct next block* produced an empty new block, indicating that all information had been collected.

For these experts, the number of questions they had to answer varied between 1000 and 2500, a far cry from the theoretical maximum of $\pm 2.8 \cdot 10^{16}$.

Indeed, there was one expert who did not finish the procedure; after the completion of block 3, the constructed fourth block appeared still to be very large and it was decided that it was not feasible to let her continue. This shows the dependence of efficiency on the amount of perceived structure (this expert was more conservative than the others), but it also permits us to highlight an additional advantage of our algorithm. The fact is that this procedure produces a knowledge space, even when it is interrupted at an arbitrary moment, and this "half-way" knowledge space is guaranteed to include the final space that would have been obtained, had the procedure terminated regularly.

This is discussed more fully in Koppen and Doignon (1989), but a direct justification for this assertion is available, if we accept that our algorithm is just a more efficient implementation of the straightforward algorithm described in Section 2.1. There we started with the full power set of X ; every positive answer of the expert led us to remove a number of states and we noted that every intermediate knowledge structure had to be a space. In particular, when we interrupt the procedure after block 1, we obtain a knowledge space that is known to be the closure under intersection of the final space. That is, this space contains all the states of the final space, plus possibly a number of intersections of these states that do not appear in the final space. Such a space that is closed under intersection corresponds to the interesting special case where the relationships between the items can be fully described by a partial order. Again, for more details see Koppen and Doignon (1989).

To return to our stopped expert, we were indeed able to compute her knowledge space at the end of the third block. (In terms of the algorithm of the previous section, this requires an addition to the procedure *Output space* which consists in establishing explicitly inferences for later, not considered subsets from the data in the completed blocks – just as this happens in the procedure *Construct next block*.) As expected, the resulting knowledge space was bigger than those of the other experts, but it was in the same order of magnitude and nowhere near the theoretical maximum. In general, by far the most reduction in the size of the knowledge space takes place in block 1 (i.e., is caused by positive answers in block 1) and the rate of reduction drops continually in later blocks. This was one of the reasons to treat the blocks in this order. (Again referring to the straightforward procedure of Section 2.1, it is for instance clear that the first positive response in block 1 results in a reduction of 25 per cent, removing all subsets containing one, but not the other of a pair of items.) So, it appears that even when the query procedure cannot be completed, but has to be stopped after, say, three or four blocks, it still produces a

knowledge space and generally one that is close to the one that would have been obtained without interruption.

References

- Degreef, E., Doignon, J.-P., Ducamp, A. & Falmagne, J.-C. (1986). Languages for the assessment of knowledge. *Journal of Mathematical Psychology*, **30**, 243-256.
- Doignon, J.-P. & Falmagne, J.-C. (1985). Spaces for the assessment of knowledge. *International Journal of Man-Machine Studies*, **23**, 175-196.
- Falmagne, J.-C. (1989). A latent trait theory via a stochastic learning theory for a knowledge space. *Psychometrika*, in press.
- Falmagne, J.-C. & Doignon, J.-P. (1988a). A class of stochastic procedures for the assessment of knowledge. *British Journal of Statistical and Mathematical Psychology*, **41**, 1-23.
- Falmagne, J.-C. & Doignon, J.-P. (1988b). A Markovian procedure for assessing the state of a system. *Journal of Mathematical Psychology*, **32**, 232-258.
- Falmagne, J.-C., Koppen, M., Villano, M., Johannesen, L. & Doignon, J.-P. (1989). Introduction to knowledge spaces: How to build, test and search them. *Psychological Review*, to appear.
- Koppen, M. & Doignon, J.-P. (1989). How to build a knowledge space by querying an expert. *Journal of Mathematical Psychology*, to appear. (Chapter 7.)

CHAPTER 9

SURMISE MAPPINGS AND WELL GRADED KNOWLEDGE SPACES

(Submitted)

Surmise mappings and well graded knowledge spaces

Mathieu Koppen

New York University

We discuss the 1-1 correspondence (introduced by Doignon and Falmagne, 1985, in the context of knowledge assessment) between *surmise mappings*, a generalization of quasi orders, and *knowledge spaces*, families of sets satisfying the axiom of closure under union. Possible additional conditions for surmise mappings are presented, with their consequences for the corresponding spaces. In particular, the condition corresponding to *well graded* knowledge spaces (Falmagne and Doignon, 1988; Falmagne, 1989) is detected. Results are related to the mathematical theory of convex geometries.

1. Representing domains of knowledge by quasi orders

In a model for the assessment of knowledge, introduced and motivated in Doignon and Falmagne (1985), a body of knowledge is formalized as a finite set X , consisting of all the *questions* (or *problems*) in that domain. Assessing a person's knowledge means finding out what she knows and does not know. Accordingly, in this model, an individual's *knowledge state* in a particular domain is defined as the subset of questions this individual is capable of solving. The critical point here, to avoid trivialities, is that not every subset of problems is in fact a possible knowledge state. The mastery of some (more difficult) problem may imply the mastery of some (easier) problems. In the field of elementary arithmetic, for instance, it seems safe to assume that a pupil will not be capable of multiplying 3-digit numbers unless he is capable of adding such numbers. Thus, any subset of problems containing multiplication of 3-digit numbers, but not containing addition of 3-digit numbers, is excluded as a possible knowledge state. A domain X , then, is characterized by the collection of all possible knowledge states, a particular subset of the power set of X

This work was supported by DOD grant MDA903-87-K-0002 to Jean-Claude Falmagne at New York University. The author wants to thank Jean-Claude Falmagne and Jean-Paul Doignon for their useful comments on a previous draft. Address comments and requests for reprints to M. Koppen, Dept. of Psychology NYU, 6 Washington Place, New York, NY 10003.

which is called the *knowledge structure* of that domain.

To get a handle on the kind of restrictions a knowledge structure might impose, we can generalize the above example and consider a relation S on X , where xSy is interpreted as: from a correct answer to problem y it may be surmised that problem x will be answered correctly. A knowledge structure is said to be *compatible* with a relation S if it allows for the above interpretation of S , that is, if xSy implies that any knowledge state containing y contains also x . There are two viewpoints possible now.

On the one hand, we may, for a fixed knowledge structure K , consider the collection of relations that K is compatible with. Regarding this collection we may observe the following facts: (i) it is non-empty, since any knowledge structure is compatible with the empty relation; (ii) it has a unique maximal element, since if K is compatible with a number of relations, it is also compatible with their union; (iii) this maximal element is a reflexive and transitive relation (i.e., a quasi order), since if K is compatible with a relation, it is also compatible with the reflexive transitive closure of that relation. This maximal quasi order compatible with K will be called the *surmise relation* corresponding to K .

We may, on the other hand, start with a fixed relation S on X and consider the collection of knowledge structures on X that are compatible with S . Again, some properties of this collection follow: (i) it is non-empty, since the empty knowledge structure is compatible with any relation; (ii) it has a unique maximal element, since if a number of knowledge structures are compatible with S , then so is their union; (iii) this maximal element is closed under union and intersection (i.e., any union or intersection of states is again a state), since if K is compatible with S , then so is the structure consisting of all unions and intersections of states in K . In this context the empty set \emptyset and the full set X are considered to be the union and intersection, respectively, of zero states. Thus, with any relation we can associate the maximal knowledge structure compatible with it, and this structure will be closed under union and intersection.

A classical result by Birkhoff (1937) tells us that the above correspondences between relations and knowledge structures are in fact 1-1 when restricted to quasi orders and knowledge structures closed under union and intersection. More precisely, we have:

1.1. Theorem. (Birkhoff, 1937.) *The formula*

$$xSy \text{ iff } (\text{for all } K \in \mathbf{K} : y \in K \text{ implies } x \in K) \quad (1.1)$$

defines an order reversing isomorphism between quasi orders S and knowledge structures K that are closed under union and intersection. ■

(Both collections are ordered by inclusion. See Monjardet (1970) for a formulation of this result in terms of a *Galois connection* between knowledge structures and binary relations.) Note that (1.1) defines for any knowledge structure \mathbf{K} the surmise relation S and for any relation S the maximal compatible knowledge structure \mathbf{K} . The impact of the theorem is that the compositions of these two mappings are identities on the collection of quasi orders, respectively the collection of knowledge structures closed under union and intersection. A knowledge structure is fully characterized by its surmise relation if and only if it is closed under union and intersection.

In the next section we discuss Doignon and Falmagne's (1985) generalization of this result to *surmise mappings*, describing a less restricted class of knowledge structures. (An alternative way of generalizing Theorem 1.1 is obtained by Koppen and Doignon, 1989.) Section 3 introduces a number of possible additional conditions for surmise mappings. The consequences of these conditions for the corresponding knowledge structures are explored in the next section with special attention to the important case of a *well graded knowledge space* (Falmagne and Doignon, 1988; Falmagne, 1989). The final section relates the obtained results to the mathematical theory of convex geometries.

2. The generalization to surmise mappings

Attractive as it may seem, representing a knowledge structure by its surmise relation is in general too strict a model. The reason for this is that this representation cannot deal with the by no means exceptional case where there are various ways of solving a problem. Suppose we observe a correct response of a pupil to problem x , and suppose we know that there are two ways of solving x : either by the mastery of problem a , or by the mastery of both problems b and c . From the correct answer to x , then, we cannot surmise a correct answer to a : the student might have solved x via b and c . Nor can we surmise a correct answer to either b or c , since the student may have taken the route via a . In short, if S is the surmise relation of the knowledge structure, none of aSx , bSx or cSx will be valid. If the knowledge structure were given by (1.1), the implication would be that there is a state containing x , but not containing a , one containing x and not b , and one containing x and not c . Since each knowledge structure defined by (1.1) is closed under intersection, this would imply a knowledge state containing x and none of a , b or c . This, however, is contradictory with our starting assumption that any state containing x contains either a or both b and c . Thus we see that a knowledge structure completely defined by its surmise relation cannot deal with alternative solutions to a

problem, and in the process we have got a hint of which of the implied properties of the knowledge structure causes this trouble: the closure under intersection seems to be the culprit.

Doignon and Falmagne (1985) remedied this situation by defining a generalization of the notion of a surmise relation. This generalization can best be understood by taking a slightly different viewpoint on the surmise relation. The usual way of looking at a (binary) relation on X is as a subset of the Cartesian product $X \times X$. Another, equally valid way, however, is as a mapping from X into the power set of X . This can be done in two different ways. Here, we will identify a relation S on X with the mapping that associates with any $x \in X$ the subset

$$Sx = \{y \in X : ySx\}.$$

A relation's being a quasi order translates then into the following properties of this mapping: if $x, y \in X$ then

$$x \in Sx, \tag{2.1}$$

$$y \in Sx \text{ implies } Sy \subseteq Sx. \tag{2.2}$$

The former property represents the reflexivity, the latter the transitivity of a quasi order. In terms of a surmise relation, the set Sx is the collection of problems that can be surmised from problem x . We could call it the set of prerequisites, or, another suggestive term, antecedents for x , since it consists of the problems that must have been acquired before mastery of problem x can take place. A surmise relation S associates with any problem x a unique set of prerequisites, Sx . From this perspective, the solution for the case of multiple paths for a problem is quite natural: generalize the mapping S to a mapping, call it σ , that associates with any problem not just one subset of X , but rather a family of subsets. So, σ maps X into the power set of the power set of X ; the elements of $\sigma(x)$ are called the *clauses* for x and they constitute the possible sets of antecedents for x . The idea is that knowing that a person has mastered problem x , we can infer that he must in fact have mastered at least one clause for x in its totality. Doignon and Falmagne (1985) show that such mappings do indeed describe an interesting class of knowledge structures, provided that some axioms are imposed. The first of these is that any $x \in X$ has at least one clause:

$$\sigma(x) \neq \emptyset. \tag{2.3}$$

The next two axioms generalize the reflexive and transitive properties (2.1) and (2.2) of a quasi order. First, it is required that any problem be contained in any of its clauses:

$$C \in \sigma(x) \text{ implies } x \in C. \tag{2.4}$$

Next, we demand that any clause contain a clause for any of its elements:

$$y \in C \in \sigma(x) \text{ implies } C' \subseteq C \text{ for some } C' \in \sigma(y). \quad (2.5)$$

This requirement is reasonable in view of our interpretation of a clause for x as a possible set of antecedents for x : if y appears in such a set, then y itself must be "reachable", that is, it must have a set of antecedents within this set. With the interpretation in terms of antecedents it is also clear that it does not make sense to have two clauses for x where one is a subset of the other. Thus, as our final axiom, we want the clauses of any problem to be incomparable with respect to inclusion:

$$C, C' \in \sigma(x) \ \& \ C \subseteq C' \text{ implies } C = C'. \quad (2.6)$$

Any mapping σ from X into the power set of the power set of X satisfying (2.3) to (2.6) is called a *surmise mapping* on X .[†] It can be checked easily that a surmise relation (interpreted as a mapping) is a surmise mapping; it is the special case where any problem has just one clause. Note that with surmise mappings we have no trouble describing the situation that was at the start of this section, where a correct answer to x implied that either a , or both b and c were also mastered. This is represented by having (at least) two clauses for x : one of them being $\{a, x\}$, the other one $\{b, c, x\}$.

With these surmise mappings, we have the following generalization of Birkhoff's Theorem 1.1.

2.1. Theorem. (Doignon and Falmagne, 1985.) *The collection of surmise mappings on X is in 1-1 correspondence with the collection of knowledge structures on X that are closed under union. ■*

We remark here, as an aside, that Theorem 2.1 is not the only way of generalizing surmise relations to obtain a correspondence with knowledge structures closed under union. Koppen and Doignon (1989) describe how this can be achieved in a different way, replacing surmise relations on X not with surmise mappings on X , but rather with a class of quasi orders on the power set of X . The ensuing characterization appears particularly suitable for extracting the knowledge structure governing some field from experts in that field; it has already been put into practice for this purpose.

Comparing Theorems 1.1 and 2.1, we see that the generalization from surmise relations to surmise mappings corresponds, in the domain of the knowledge structures, to dropping the axiom of closure under intersection. (Something like this

[†] We call here simply surmise mapping what in Doignon and Falmagne (1985) was termed a *space-like* surmise mapping.

was expected in the discussion of the example in the first paragraph of this section.) Doignon and Falmagne (1985) reserve a special name for the resulting concept: a *knowledge space* on X is a family of states (subsets of X) such that any union of states is again a state.

The correspondence between knowledge spaces and surmise mappings can be made fully explicit. For a surmise mapping we consider the collection of all of its clauses and close this collection under the union. This yields the corresponding knowledge space. For a knowledge space, we call a state K *minimal for x* if it contains x and no state properly included in K does. The surmise mapping corresponding to a knowledge space maps any $x \in X$ to the collection of minimal states for x in the space.

There is an alternative description, using the notion of the *basis* of a knowledge space. A subcollection \mathbf{B} of a knowledge space \mathbf{K} is called a *basis for \mathbf{K}* if (i) \mathbf{B} is independent in the sense that no element of \mathbf{B} can be written as the union of a number of other elements of \mathbf{B} , and (ii) \mathbf{B} is complete in the sense that any state in \mathbf{K} can be obtained as the union of some elements in \mathbf{B} . Any knowledge space \mathbf{K} on a finite set X has a unique basis, consisting of the states of \mathbf{K} that are not unions of other states of \mathbf{K} (the sup-irreducible elements of \mathbf{K} , in lattice-theoretic terms). The empty set is considered as the union of zero states and, thus, is never a basis element. In terms of bases of knowledge spaces, the correspondence with surmise mappings can be stated very simply. The clauses of a surmise mapping constitute the basis of the corresponding knowledge space and a basis element is a clause for those of its members that are not contained in any properly included basis element. (By definition, any basis element must contain at least one such member.)

We see that a knowledge space is fully characterized by its surmise mapping. In the previous section, the surmise relation was described as a partial characterization of any knowledge structure, so we might well wonder what the connection is between these two notions.

2.2. Proposition. *Let \mathbf{K} be a knowledge space on X , σ its surmise mapping and S its surmise relation. Then, for any $x \in X$,*

$$Sx = \bigcap \sigma(x).$$

Proof. By definition (1.1), ySx implies that any state of \mathbf{K} containing x contains also y . In particular, any state in \mathbf{K} that is minimal for x , that is, any clause for x in σ , contains y . For the reverse inclusion, let y be contained in any clause for x in σ . This means that y is contained in any state of \mathbf{K} that is minimal for x . But since any state of \mathbf{K} containing x must include a state that is minimal for x , it follows that y must be present in any such state and (1.1) yields ySx . ■

Intuitively, this proposition is clear. The problems that one can surmise from a correct answer to problem x are exactly the problems that are common to all sets of prerequisites of x .

The next proposition deals with the question of equivalent problems in a knowledge space, that is, problems that are *indistinguishable* in the space in the sense that any state containing one of the problems also contains the other. It appears that regarding equivalence all information is already contained in the surmise relation.

2.3. Proposition. *Let \mathbf{K} be a knowledge space on X with surmise mapping σ and surmise relation S . Two problems are equivalent with respect to the space \mathbf{K} and the surmise mapping σ iff they are equivalent with respect to the surmise relation S . More formally: for any $x, y \in X$,*

$$(for\ all\ K \in \mathbf{K},\ x \in K\ iff\ y \in K)\ iff\ (\sigma(x) = \sigma(y))\ iff\ (xSy\ \&\ ySx).$$

Proof. That equivalence in σ means the same as equivalence in \mathbf{K} is an immediate consequence of Theorem 2.1 and either equivalence implies equivalence in S since S is defined in terms of \mathbf{K} (by (1.1)) or σ (by the above Proposition 2.2). The only thing to prove is that $\sigma(x) = \sigma(y)$ whenever both xSy and ySx . By symmetry it suffices to show that in this case any clause for x is also a clause for y . Take $C \in \sigma(x)$. Since ySx we have, by Proposition 2.2, $y \in C$ and thus, by Condition (2.5), $C' \subseteq C$ for some $C' \in \sigma(y)$. Now xSy and (2.5) imply the existence of $C'' \in \sigma(x)$ such that $C'' \subseteq C' \subseteq C$. By Condition (2.6), then, $C'' = C$, and thus $C = C' \in \sigma(y)$. ■

Thus, two problems that are equivalent in the surmise relation are indistinguishable in the corresponding knowledge space: a subject has mastered one problem whenever he has mastered the other one. A tempting interpretation of this situation is that two such problems test in fact one and the same *notion*. In basic arithmetic, for instance, there may be a number of problems that are all equivalent instances of the notion "addition of two 2-digit numbers without carry". The above proposition shows that it is not a real restriction to assume that a surmise relation is in fact a partial order (i.e., it is antisymmetric). This only assumes that all elements of X are distinguishable and this is by definition true if we look at X as the collection of notions, where each notion may consist in an equivalence class of problems. Whenever convenient, we will in the sequel assume that such a reduction has been carried out and that all elements of X are distinguishable.

3. Extra conditions for a surmise mapping

Here we want to discuss a number of additional conditions that may be imposed on a surmise mapping. A need for this may arise in view of the interpretation of surmise mappings. A clause for a problem x is interpreted as a possible collection of antecedents or prerequisites for x . This suggests that there should not too much "cyclicity" be present in the system of clauses. After all, subjects are supposed to be able to move in some reasonable way through the corresponding knowledge space, from the null state to the full set X . In this spirit, Falmagne and Doignon (1988) and Falmagne (1989) introduced the important case of a *well graded* knowledge space. A space is called well graded when any non null state K has a member x such that $K - \{x\}$ is a state. This means that the space consists of a number of learning paths or *gradations* along which subjects can move, learning one problem at a time, from the empty state to the state X . Well-gradedness appears to be a very reasonable assumption for knowledge spaces in practice.

Below, we consider a number of ways in which the idea of cycles in a surmise mapping can be defined. In the next section we investigate what the impact is of these various conditions on the knowledge spaces they describe via the correspondence of Theorem 2.1. We will there especially be interested in which condition characterizes surmise mappings corresponding to well graded knowledge spaces.

In describing the conditions, the following two relations P_σ and R_σ on X , indexed by a surmise mapping σ on X , will be useful: for $x, y \in X$,

$$x P_\sigma y \quad \text{iff} \quad x \in \bigcap \sigma(y)$$

and

$$x R_\sigma y \quad \text{iff} \quad x \in \bigcup \sigma(y).$$

In words, $x P_\sigma y$ iff x is a member of all clauses for y , and $x R_\sigma y$ iff x appears in some clause for y . It is clear that $P_\sigma \subseteq R_\sigma$ (any pair (x, y) in P_σ is also in R_σ) and by (2.4) P_σ and R_σ are reflexive. By (2.5), the relation P_σ is also transitive; thus, P_σ is a quasi order. Indeed, Proposition 2.2 shows that P_σ is precisely the surmise relation of the knowledge space corresponding to σ . As indicated in the discussion after Proposition 2.3, it is not a real restriction to assume, as we will do, that P_σ is in fact a partial order (i.e., it is antisymmetric).

Proposition 2.3 tells us that the worst case of cyclic behavior cannot materialize in surmise mappings: if the mastery of problem x always implies the mastery of y , then we cannot have the same situation the other way around, without x and y being equivalent throughout. The following example shows, however, that a next worse

case is in fact possible. (In examples, we denote elements of X by a, b, c , etc. To simplify notation we write subsets of X without separators and braces; that is, we write simply abc for the subset $\{a, b, c\}$, etc.)

3.1. Example. Let $X = \{a, b, c, d\}$ and σ be defined by:

$$\begin{aligned}\sigma(a) &= \{abc, abd\} & \sigma(c) &= \{c\} \\ \sigma(b) &= \{abc, bd\} & \sigma(d) &= \{d\}\end{aligned}$$

It can be checked easily that the conditions for a surmise mapping are satisfied. We see that b appears in every clause for a . And while it is not true that a appears in every clause for b , this is the case for some clause for b (namely, abc). In this situation, any clause for b containing a may be considered as not a very realistic one. We could require that it be not possible for one problem to appear in a clause for another problem when the latter appears in every clause of the former. This condition can be more formally expressed as

$$x \mathbf{P}_\sigma y \ \& \ y \mathbf{R}_\sigma x \ \text{implies} \ x = y. \quad (3.1)$$

In the next example this condition is satisfied.

3.2. Example. Let $X = \{a, b, c\}$ and σ be defined by:

$$\sigma(a) = \{ab, ac\} \quad \sigma(b) = \{ab, bc\} \quad \sigma(c) = \{c\}$$

Although (3.1) is satisfied, there is still a problem with the interpretation of this surmise mapping. The set ab is a clause for a , and as such gives rise to the interpretation: it is possible to arrive at a via b . In order to get at b , then, a clause for b must be fulfilled *that is contained in this clause for a* . The only such clause, however, is ab itself, which now must be interpreted as a way of getting at b via a . We observe again a cyclic pattern: it is not clear which comes first, the chicken or the egg. To avoid this situation we must require that the surmise mapping be *exclusive* in the sense that a clause cannot be shared by two distinct problems:

$$x \neq y \ \text{implies} \ \sigma(x) \cap \sigma(y) = \emptyset. \quad (3.2)$$

The following surmise mapping is exclusive.

3.3. Example. Let $X = \{a, b, c, d\}$ and σ be defined by:

$$\begin{aligned}\sigma(a) &= \{abc, ad\} & \sigma(c) &= \{c\} \\ \sigma(b) &= \{abd, bc\} & \sigma(d) &= \{d\}\end{aligned}$$

In the above example we can, however, still observe some form of cycles. On the one hand there is a clause for a containing b ("we can reach a via b ") while on the

other hand a appears in a clause for b ("b can be reached via a"). If we want to avoid this situation where with two distinct problems each appears in some clause for the other, we demand in effect that the relation R_σ be *antisymmetric*:

$$x R_\sigma y \ \& \ y R_\sigma x \ \text{implies} \ x = y. \quad (3.3)$$

By extension, a surmise mapping satisfying (3.3) will be called *antisymmetric*.

3.4. Example. Let $X = \{a, b, c, d, e, f\}$ and σ be defined by:

$$\begin{array}{lll} \sigma(a) = \{abd, ae\} & \sigma(c) = \{ace, cf\} & \sigma(e) = \{e\} \\ \sigma(b) = \{bcf, bd\} & \sigma(d) = \{d\} & \sigma(f) = \{f\} \end{array}$$

Now all "direct" cycles have indeed disappeared, but there are in the above example still cycles of the same kind left, only with length exceeding 2. Here, there is a clause for a containing b , a clause for b containing c and, finally, a clause for c containing a . It looks like any objections against a surmise mapping not being antisymmetric are also valid in this situation. That is, we might then just as well require that R_σ , at least its irreflexive part, be *acyclic* instead of just antisymmetric:

$$(x_i R_\sigma x_{i+1}, i = 1 \cdots n, x_{n+1} = x_1) \ \text{implies} \ (x_i = x_1, i = 1 \cdots n). \quad (3.4)$$

Such a surmise mapping is also called *acyclic*.[†] An example is given below.

3.5. Example. Let $X = \{a, b, c, d, e, f\}$ and σ be defined by:

$$\begin{array}{lll} \sigma(a) = \{abd, ae\} & \sigma(c) = \{cf\} & \sigma(e) = \{e\} \\ \sigma(b) = \{bcf, bd\} & \sigma(d) = \{d\} & \sigma(f) = \{f\} \end{array}$$

Finally we might consider the very special situation where the extension $P_\sigma \subseteq R_\sigma$ is trivial and the two relations are in fact identical:

$$R_\sigma = P_\sigma. \quad (3.5)$$

From the definitions it is clear that this is the case if and only if each $x \in X$ has just one clause: σ essentially coincides with the surmise relation P_σ , as in the following example.

3.6. Example. Let $X = \{a, b, c, d, e, f\}$ and σ be defined by:

$$\begin{array}{lll} \sigma(a) = \{abcfe\} & \sigma(c) = \{cf\} & \sigma(e) = \{e\} \\ \sigma(b) = \{be\} & \sigma(d) = \{def\} & \sigma(f) = \{f\} \end{array}$$

[†] This notion of an acyclic surmise mapping was already introduced in Doignon and Falmagne (1985).

The foregoing discussion strongly suggests the following lemma:

3.7. Lemma. *Regarding the Conditions (3.1) to (3.5) we have following chain of implications:*

$$(3.5) \Rightarrow (3.4) \Rightarrow (3.3) \Rightarrow (3.2) \Rightarrow (3.1)$$

and none of the reverse implications hold in general.

Proof. The only implication that is not completely obvious is $(3.2) \Rightarrow (3.1)$, but this follows by the same argument we used in Proposition 2.3. Suppose (3.2) and $x \mathbf{P}_\sigma y$ and $y \mathbf{R}_\sigma x$ for $x, y \in X$. The last condition means we have $y \in C$ for some $C \in \sigma(x)$. By (2.5) this implies $C' \subseteq C$ for some $C' \in \sigma(y)$. Using $x \mathbf{P}_\sigma y$ and again (2.5), we obtain $C'' \subseteq C' \subseteq C$, where $C, C'' \in \sigma(x)$ and $C' \in \sigma(y)$. Consequently $C'' = C' = C$ and (3.2) yields $x = y$. The examples show that none of the implications can be reversed (and that (3.1) is not implied by the axioms for a surmise mapping). ■

4. The connection with the knowledge spaces

In this section we explore the impact of Conditions (3.1) to (3.5) on corresponding knowledge spaces. Actually, the situation is clear for the most restrictive Condition (3.5). Here, surmise mapping and surmise relation coincide, which means that Theorem 1.1 applies: the corresponding spaces are the ones that are closed under intersection.

In order to describe the effect of the other conditions, we consider for any knowledge space \mathbf{K} a collection of relations $\{\leq_K\}_{K \in \mathbf{K}}$. For any $K \in \mathbf{K}$, \leq_K is the following relation on K :

$$x \leq_K y \text{ iff } y \in K' \in \mathbf{K} \ \& \ K' \subseteq K \text{ implies } x \in K'.$$

It is easy to check that any such \leq_K is a quasi order. In fact, this definition coincides with (1.1) if we restrict our universe of problems to K , so \leq_K is the usual quasi order (surmise relation) of the knowledge space on K that is induced by \mathbf{K} by considering only states that are subsets of K . Such a subcollection is indeed still closed under union and, on the other hand, any subspace of \mathbf{K} is obtained in this way.

The following lemma gives an alternative definition in terms of surmise mappings:

4.1. Lemma. *Let \mathbf{K}_σ be the knowledge space on X corresponding to the surmise mapping σ and let $K \in \mathbf{K}_\sigma$. Then*

$$x \leq_K y \text{ iff } C \in \sigma(y) \text{ \& } C \subseteq K \text{ implies } x \in C.$$

Proof. Follows easily from the fact that a state of \mathbf{K}_σ contains an element x if and only if it includes some $C \in \sigma(x)$. ■

The next lemma collects some immediate consequences.

4.2. Lemma. *Let \mathbf{K}_σ be the knowledge space on X corresponding to the surmise mapping σ .*

- (i) *If $K_1, K_2 \in \mathbf{K}_\sigma$ and $x, y \in X$ such that $x, y \in K_1 \subseteq K_2$, then $x \leq_{K_2} y$ implies $x \leq_{K_1} y$.*
- (ii) $\mathbf{P}_\sigma = \leq_X$.
- (iii) $\mathbf{R}_\sigma = \cup_{K \in \mathbf{K}_\sigma} \leq_K$. (That is, $x \mathbf{R}_\sigma y$ iff $x \leq_K y$ for some $K \in \mathbf{K}_\sigma$.) ■

Notice that, as a consequence of (i) and (ii), $x \mathbf{P}_\sigma y$ implies $x \leq_K y$ for any $K \in \mathbf{K}_\sigma$ containing both x and y .

Lemma 4.2(ii) and (iii) give a direct translation of Conditions (3.1), (3.3) and (3.4) in terms of the knowledge space \mathbf{K}_σ :

4.3. Lemma. *Let \mathbf{K}_σ be the knowledge space corresponding to the surmise mapping σ .*

- (i) *Condition (3.1) on σ is equivalent to the implication $x = y$ whenever we have on the one hand $x \leq_X y$, while on the other hand there is $K \in \mathbf{K}_\sigma$ such that $y \leq_K x$.*
- (ii) *σ is antisymmetric (Condition (3.3)) iff $x = y$ whenever there are $K_1, K_2 \in \mathbf{K}_\sigma$ such that $x \leq_{K_1} y$ and $y \leq_{K_2} x$.*
- (iii) *σ is acyclic (Condition (3.4)) iff $x_1 = x_2 = \dots = x_n$ whenever there are $K_1, \dots, K_n \in \mathbf{K}_\sigma$ such that $x_i \leq_{K_i} x_{i+1}$ for $i = 1, \dots, n$ and $x_{n+1} = x_1$.* ■

The translation of Condition (3.2) in terms of the relations \leq_K is more interesting.

4.4. Theorem. *Let \mathbf{K}_σ be the knowledge space on X generated by the surmise mapping σ . Then the following conditions are equivalent:*

- (i) σ is exclusive;
- (ii) for any $K \in \mathbf{K}_\sigma$, \leq_K is a partial order;
- (iii) \mathbf{K}_σ is well graded.

Proof. ((i) \Rightarrow (ii)) The only thing to prove is the antisymmetry of \leq_K . So, let $K \in \mathbf{K}_\sigma$ and suppose $x \leq_K y$ and $y \leq_K x$. Then, obviously, $y \in K$ and thus there is some clause C for y contained in K . As before, we can, from the assumptions, construct the situation $C'' \subseteq C' \subseteq C$, with $C, C'' \in \sigma(y)$ and $C' \in \sigma(x)$. Consequently, $C = C' \in \sigma(x) \cap \sigma(y)$ and since σ is exclusive we obtain $x = y$.

((ii) \Rightarrow (iii)) Let K be a non null state of \mathbf{K}_σ . By assumption, \leq_K is a partial order on K and since K is finite and non null, we can find $x \in K$ that is maximal in this partial order, i.e., such that $x \leq_K y$ implies $x = y$. Thus, whenever $y \neq x$ we cannot have $x \leq_K y$. By definition of \leq_K , this means that for any $y \in K, y \neq x$, there is a state K_y included in K that contains y , but does not contain x . But then, since \mathbf{K}_σ is a space, $\cup_{y \neq x} K_y = K - \{x\}$ is a state of \mathbf{K}_σ . We have proved that \mathbf{K}_σ is well graded.

((iii) \Rightarrow (i)) This implication follows by contradiction. Suppose $C \in \sigma(x) \cap \sigma(y)$ for some $x, y \in X$. Then, since C is a clause in $\sigma, C \in \mathbf{K}_\sigma$ and $C \neq \emptyset$. And C being a clause for both x and y , there can be no state strictly included in C containing x or y . In other words, any such state (and \emptyset is one) differs from C by at least two elements. Consequently, \mathbf{K}_σ is not well graded. ■

We see that the exclusive surmise mappings characterize the collection of well graded knowledge spaces. Alternatively, a well graded knowledge space is characterized by the fact that the surmise relations of all its subspaces are partial orders.

In Section 2 we described the relation between a surmise mapping and the basis of the corresponding space. The collection of families $\sigma(x), x \in X$, appeared to cover the basis corresponding to σ : any set in the basis is a clause for some element x . A surmise mapping is exclusive if and only if this covering is in fact a partitioning of the basis. For an exclusive surmise mapping we have a kind of inverse mapping from the basis onto X , assigning any basis element to the unique element of X for which it is a clause. The partitioning of the basis corresponds to the equivalence with respect to this mapping. The above theorem, then, provides a characterization of the basis of a well graded knowledge space: it is partitioned by the subcollections consisting of the sets that are minimal for the various $x \in X$.

Notice, by the way, that the restriction to exclusive surmise mappings can very easily be built into the definition of surmise mappings: it corresponds to the replacement of

$$y \in C \in \sigma(x) \text{ implies } C' \subseteq C \text{ for some } C' \in \sigma(y) \tag{2.5}$$

by

$$y \in C \in \sigma(x) \text{ implies } C' \subseteq C - \{x\} \text{ for some } C' \in \sigma(y).$$

5. The theory of convex geometries

The results of the preceding section appear to be related to the mathematical theory of convex geometries. This connection, and, indeed, the very existence of this subfield of mathematics were pointed out to me by Jean-Paul Doignon. Convex geometries appear in an abstract, combinatorial approach to the notion of convexity; they were introduced independently by Paul Edelman and Robert Jamison. Equivalent structures have been described by other authors under the names of “shelling structures” and “selectors”. The following brief sketch of some basic concepts of this theory is based on a joint paper by Edelman and Jamison (1985), where further references can be found.

Edelman and Jamison (1985) consider a finite set X and *alignments* of X , that is, families of subsets of X that are closed under intersection. The subsets of X in such an alignment are called *convex sets*. A *convex geometry* on X , then, is an alignment on X that is such that for every convex set $C \neq X$ there is an element x not contained in C for which $C \cup \{x\}$ is convex. (The synonym *antimatroid* is also used in the literature.) A *copoint attached at x* is a maximal convex set *not* containing x . For every convex set C a *C -factor relation* is defined between elements that are *not* in C : a pair (x, y) of such elements is in this relation if and only if x is contained in any convex set containing $C \cup \{y\}$.

The connection between this theory and that of the knowledge spaces is made, once we realize that whenever a family of subsets is closed under intersection, the family of complementary sets is closed under union, and vice versa. Via this complementation mapping $C \rightarrow X - C$, we get indeed a direct translation of the above notions into our terminology. An alignment corresponds to a knowledge space, a convex set to a knowledge state and a convex geometry to a well graded knowledge space. A copoint at x refers in the same way to a minimal state containing x , or, in the language of surmise mappings, to a clause for x . The C -factor relation on $X - C$, finally, is exactly the relation \leq_K defined in the previous section, where $K = X - C$.

Having made this translation, it becomes clear that our Theorem 4.4 is a rediscovery and combination of Edelman and Jamison’s (1985) Theorems 2.3 and 2.4 which consider conditions under which an alignment is a convex geometry. Their Theorem 2.3 states that this is the case if and only if all C -factor relations are partial orders and their Theorem 2.4 states that an equivalent condition is that every copoint is attached at a unique point (in our language: the surmise mapping is exclusive).

References

- Birkhoff, G. (1937). Rings of sets. *Duke Mathematical Journal*, **3**, 443-454.
- Doignon, J.-P. & Falmagne, J.-C. (1985). Spaces for the assessment of knowledge. *International Journal of Man-Machine Studies*, **23**, 175-196.
- Edelman, P. H. & Jamison, R. E. (1985). The theory of convex geometries. *Geometriae Dedicata*, **19**, 247-270.
- Falmagne, J.-C. (1989). A latent trait theory via a stochastic learning theory for a knowledge space. *Psychometrika*, in press.
- Falmagne, J.-C. & Doignon, J.-P. (1988). A Markovian procedure for assessing the state of a system. *Journal of Mathematical Psychology*, **32**, 232-258.
- Koppen, M. & Doignon, J.-P. (1989). How to build a knowledge space by querying an expert. *Journal of Mathematical Psychology*, to appear. (Chapter 7.)
- Monjardet, B. (1970). Tresses, fuseaux, préordres et topologies. *Mathématiques et Sciences Humaines*, **30**, 11-22.

DISCUSSION OF PART II

1. A survey of the various representations

Let us first review the results of the preceding chapters in terms of the correspondences that have been established between (mathematical) objects of different kind. At the basis of the knowledge assessment project is the conceptualization of a domain of knowledge as a finite set X of all problems or items in that domain. The knowledge state of a student can then be formalized as the subset of problems this student is capable of solving and this leads to the central concept of a knowledge structure as the collection of all possible such knowledge states. Thus, a knowledge structure is a subset of the power set of X .

In order to find alternative characterizations some restriction has to be imposed and the knowledge space was defined as a structure closed under union: any $K \subseteq 2^X$ is a knowledge space if

$$K_1 \in K \ \& \ K_2 \in K \ \Rightarrow \ K_1 \cup K_2 \in K. \quad [S1]$$

(Throughout this discussion we will only consider the finite case.) In Chapter 6, Section 4, we described a class of mappings from X into the power set of the power set of X , the surmise mappings introduced by Doignon and Falmagne (1985). Any $\sigma: X \rightarrow 2^{2^X}$ is a surmise mapping, provided

$$\begin{aligned} \sigma(x) &\neq \emptyset \\ C \in \sigma(x) &\Rightarrow x \in C \\ y \in C \in \sigma(x) &\Rightarrow C' \subseteq C \text{ for some } C' \in \sigma(y) \\ C \in \sigma(x) \ \& \ C' \in \sigma(x) &\Rightarrow C \not\subseteq C'. \end{aligned} \quad [M1]$$

These surmise mappings turn out to be in 1-1 correspondence with the knowledge spaces on X and the nature of this correspondence is discussed in Chapter 6 and again in Chapter 9.

The original contribution of Chapter 7 in this respect is to add a third alternative. This time we consider binary relations on the power set of X , in particular the class of what were called entail relations. Any $P \in 2^{(2^X \times 2^X)}$ is an entail relation when

$$\begin{aligned}
 A \supseteq B &\Rightarrow APB \\
 APB \ \& \ BPC &\Rightarrow APC \\
 APB_1 \ \& \ APB_2 &\Rightarrow AP(B_1 \cup B_2).
 \end{aligned}
 \tag{R1}$$

In Chapter 7, these entail relations were shown to be in 1-1 correspondence with knowledge spaces, and thus also with surmise mappings. In Figure 1, these three concepts, knowledge spaces, surmise mappings and entail relations are represented by the large boxes and the equivalences by the two-sided arrows connecting these boxes.

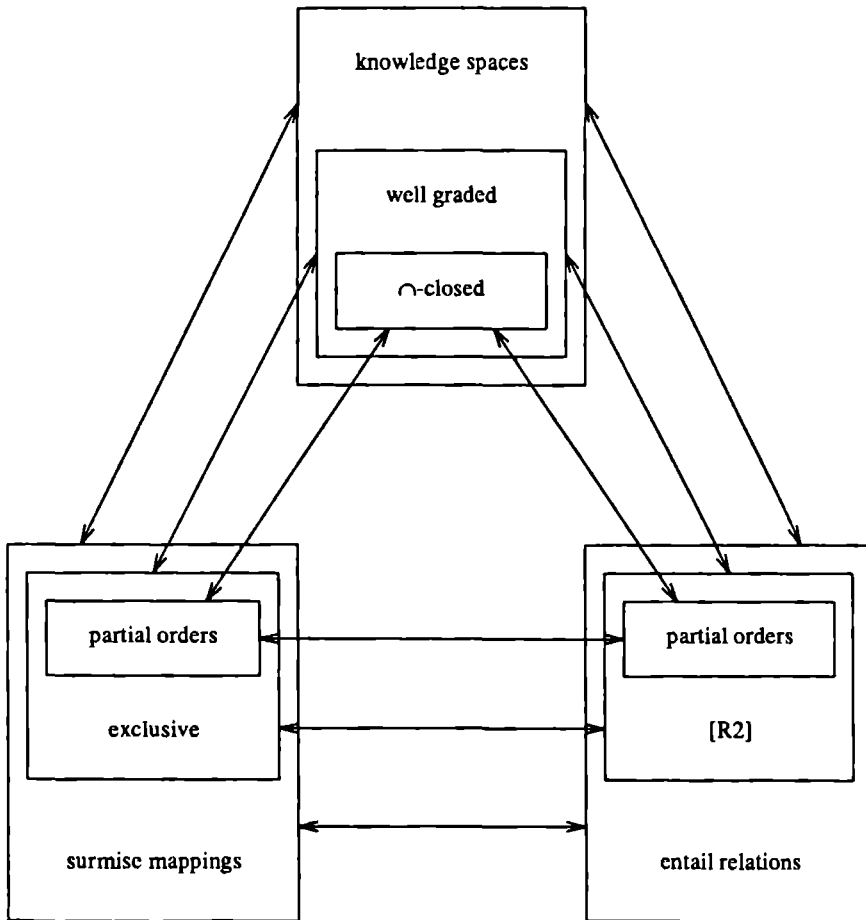


Figure 1. The three equivalent concepts with equivalent special cases.

In the domain of knowledge spaces we can consider interesting special cases. To avoid trivial formal complications, we will assume throughout that X is the set of notions of the field under investigation, that is, there are no indistinguishable elements in X . In Chapter 6, two interesting classes of knowledge spaces were introduced (and they reappeared in the following chapters). First, a knowledge space may be well graded, which according to Chapter 6, Section 4, is obtained by adding the following axiom to [S1]:

$$\emptyset \neq K \in \mathbf{K} \Rightarrow K - \{x\} \in \mathbf{K} \text{ for some } x \in K. \quad [\text{S2}]$$

Or, a knowledge space may be closed under intersection, which means adding to [S1]

$$K_1 \in \mathbf{K} \ \& \ K_2 \in \mathbf{K} \Rightarrow K_1 \cap K_2 \in \mathbf{K}. \quad [\text{S3}]$$

It turns out that, under [S1], [S2] is implied by [S3]: any (discriminating) knowledge space that is closed under intersection is well graded (see, e.g., Chapter 6, Section 3). Accordingly, the middle size box in the top box of Fig. 1 denotes the subcollection of well graded knowledge spaces and the smallest box the family of knowledge spaces closed under intersection.

Now we consider in the domain of the surmise mappings the extra conditions that correspond to the subclasses of knowledge spaces defined by [S2] and [S3]. In Chapter 9 it was established, among other things, that the well graded knowledge spaces are in 1-1 correspondence with the exclusive surmise mappings, which are obtained by supplementing the axioms [M1] with

$$x \neq y \Rightarrow \sigma(x) \cap \sigma(y) = \emptyset. \quad [\text{M2}]$$

The surmise mappings corresponding to knowledge spaces closed under intersection were known from the beginning: these are the mappings that essentially coincide with a partial order (quasi order if we want to allow for indistinguishable elements). They are distinguished by having only one clause for each element:

$$C \in \sigma(x) \ \& \ C' \in \sigma(x) \Rightarrow C = C'. \quad [\text{M3}]$$

The lower left box in Fig. 1 contains boxes representing the collection of surmise mappings defined by [M2] and [M3], with the appropriate arrows to the top boxes.

For the entail relations we have not found at this moment an easy, direct axiom defining the subclass that corresponds to well graded knowledge spaces (exclusive surmise mappings). The only conditions we can formulate are the ones that go via the correspondence with knowledge spaces or surmise mappings; these conditions are indirect and as such not very satisfying. For instance, the correspondence between entail relations and knowledge spaces is such that the complement (with respect to X) of a subset A is a state in the space \mathbf{K} if and only if in the

corresponding entail relation \mathbf{P} we only have $A \mathbf{P} Z$ for $Z \subseteq A$ (see Lemma 4.4 of Chapter 7). Via this substitution, [S2] translates into

$$(A \neq X \ \& \ (A \mathbf{P} B \Leftrightarrow A \supseteq B) \text{ for all } B \subseteq X) \Rightarrow \quad [R2]$$

$$(A + \{x\}) \mathbf{P} B \Leftrightarrow (A + \{x\}) \supseteq B \text{ for all } B \subseteq X, \text{ for some } x \notin A).$$

This is indeed rather cumbersome and too close to [S2] to be enlightening. The subclass of entail relations corresponding to the knowledge spaces closed under intersection, on the other hand, poses no problems. As with surmise mappings, these are the entail relations that are essentially partial orders on X (quasi orders for indistinguishable elements). They are obtained by adding to [R1] the axiom

$$A \mathbf{P} \{y\} \Leftrightarrow \{x\} \mathbf{P} \{y\} \text{ for some } x \in A \quad [R3]$$

(see Chapter 7, Section 4.9). In Fig. 1 we have again drawn the appropriate boxes (in the lower right box) and the connecting arrows. This completes the description of Fig. 1; it presents a concise picture of the various theoretical concepts that played a role in the preceding chapters, and their interrelationships.

2. Another view on surmise mappings and entail relations

Both surmise mappings and entail relations were developed as alternative characterizations of knowledge spaces. It was clear that this implies a 1-1 correspondence between surmise mappings and entail relations, but the specification of this correspondence always involved the knowledge space as intermediary concept. It is possible to make the equivalence of these two concepts almost immediate and that is what we set out to do in this section. It will prove helpful for this purpose to give a translation of both surmise mapping and entail relation into yet another domain, that of propositional logic. These translations will again be guided by the interpretations of surmise mappings and entail relations in terms of the corresponding knowledge space, but once the translation has been made, the equivalence of surmise mapping and entail relation will be obvious “syntactically”, without recourse to the knowledge space interpretation.

Let the set of items be $X = \{x_1, x_2, \dots, x_n\}$. For any $x_i \in X$ we define a logical variable \bar{x}_i , that is, a variable that can take one of the two values *TRUE* or *FALSE*. Any subset A of X may now be interpreted as a mapping of these variables into $\{FALSE, TRUE\}$. That is, A defines a *valuation* or truth assignment v_A by the rule

$$v_A(\bar{x}_i) = TRUE \Leftrightarrow x_i \in A. \quad (1)$$

According to this direct 1-1 correspondence between subsets of X and valuations of

$\bar{x}_1, \dots, \bar{x}_n$, a knowledge structure \mathbf{K} on X amounts to the particular collection of valuations $\{v_K : K \in \mathbf{K}\}$. It is clear that a valuation of the variables $\bar{x}_1, \dots, \bar{x}_n$ determines the truth value of any (syntactically well formed) logical formula in $\bar{x}_1, \dots, \bar{x}_n$. Such a formula may involve various logical operations, like the conjunction (“ \wedge ”), disjunction (“ \vee ”), implication (“ \rightarrow ”), equivalence (“ \leftrightarrow ”) and negation (“ \neg ”). Its truth value is computed from the truth tables for these operations; for instance, for any valuation v , $v(\bar{x}_1 \wedge \bar{x}_2) = TRUE$ if and only if $v(\bar{x}_1) = TRUE$ and $v(\bar{x}_2) = TRUE$.

Surmise mappings and entail relations were devised as descriptions of what inferences can be made in the corresponding knowledge structure. In other words, they describe formulae in $\bar{x}_1, \dots, \bar{x}_n$ that are *TRUE* under all of the valuations v_K with $K \in \mathbf{K}$. In both cases, the formulae can be presented in the form of implications $\bar{x} \rightarrow \phi$, where the subformula ϕ represents the inferences that can be drawn from the presence of x in a state of \mathbf{K} . This is pretty obvious in the case of the surmise mapping. The very idea of a surmise mapping σ corresponding to \mathbf{K} was to collect in $\sigma(x)$ all possible prerequisites of x . Concretely, $\sigma(x) = \{C_1, C_2, \dots, C_m\}$ represents the fact that, for any $K \in \mathbf{K}$,

$$x \in K \implies C_1 \subseteq K \text{ or } C_2 \subseteq K \text{ or } \dots \text{ or } C_m \subseteq K. \quad (2)$$

Let the clause C_1 consist of k elements denoted by x_1, \dots, x_k and define the formula $\bar{\gamma}_1$ by

$$\bar{\gamma}_1 = \bar{x}_1 \wedge \bar{x}_2 \wedge \dots \wedge \bar{x}_k. \quad (3)$$

For the other clauses we construct similar formulae $\bar{\gamma}_2, \dots, \bar{\gamma}_m$. Then it is clear that the statement “ $C_j \subseteq K$ ” can be expressed in terms of the valuation v_K as $v_K(\bar{\gamma}_j) = TRUE$. Defining now the formula $\bar{\sigma}_x$ by

$$\bar{\sigma}_x = \bar{\gamma}_1 \vee \bar{\gamma}_2 \vee \dots \vee \bar{\gamma}_m, \quad (4)$$

we see that (2) is tantamount to requiring that

$$v_K(\bar{x} \rightarrow \bar{\sigma}_x) = TRUE \quad (5)$$

for any $K \in \mathbf{K}$. In this way a surmise mapping σ can be identified with the collection of formulae $\bar{\sigma}_x$ defined by (4) and (3). In case \mathbf{K} is a space, σ determines \mathbf{K} completely; intuitively this means that the formula $\bar{\sigma}_x$ represents *all* that can be inferred from the presence of x in a state of \mathbf{K} . Formally this is expressed by the statement that any (other) formula ψ such that $v_K(\bar{x} \rightarrow \psi) = TRUE$ for any $K \in \mathbf{K}$ must be logically implied by $\bar{\sigma}_x$: $v(\bar{\sigma}_x \rightarrow \psi) = TRUE$ for any valuation v .

It appears that we can proceed similarly with the entail relation \mathbf{P} corresponding to a knowledge structure \mathbf{K} . The interpretation of $A \mathbf{P} x$ for $A \subseteq X$ and $x \in X$ is that

there is no state of \mathbf{K} disjoint with A and containing x ; or, equivalently, if $K \in \mathbf{K}$ and $x \in K$, then $A \cap K \neq \emptyset$. This holds for any A such that APx ; thus if $\{A_1, A_2, \dots, A_l\}$ is the collection of all such subsets, we have

$$x \in K \implies A_1 \cap K \neq \emptyset \text{ and } A_2 \cap K \neq \emptyset \text{ and } \dots \text{ and } A_l \cap K \neq \emptyset. \quad (6)$$

We remark here that we may assume that in (6) we have collected only the *minimal* subsets A such that APx . These determine \mathbf{P} completely, since an entail relation contains the superset relation (if $A' \supseteq A$ and APx , then also $A'Px$), and it is clear that the truth value of the right hand side only depends on the minimal A_j appearing there. If $A_1 = \{x_1, \dots, x_h\}$, we define the formula $\tilde{\delta}_1$ by

$$\tilde{\delta}_1 = \bar{x}_1 \vee \bar{x}_2 \vee \dots \vee \bar{x}_h, \quad (7)$$

and similar formulae $\tilde{\delta}_2, \dots, \tilde{\delta}_l$ for the other sets A_j . In this way " $A_j \cap K \neq \emptyset$ " is equivalent to $v_K(\tilde{\delta}_j) = TRUE$. If we now define the formula $\tilde{\rho}_x$ by

$$\tilde{\rho}_x = \tilde{\delta}_1 \wedge \tilde{\delta}_2 \wedge \dots \wedge \tilde{\delta}_l, \quad (8)$$

we can reformulate (6) as

$$v_K(\bar{x} \rightarrow \tilde{\rho}_x) = TRUE \quad (9)$$

for any $K \in \mathbf{K}$. Again, if \mathbf{K} is a space, \mathbf{P} determines \mathbf{K} completely and $\tilde{\rho}_x$ describes all inferences from the presence of x in a state of \mathbf{K} . Formally: any formula ϕ for which $v_K(\bar{x} \rightarrow \phi) = TRUE$ for any $K \in \mathbf{K}$ is logically implied by $\tilde{\rho}_x$: $v(\tilde{\rho}_x \rightarrow \phi) = TRUE$ for any valuation v .

In sum, then, we have in a knowledge space \mathbf{K} on X two formulae, $\tilde{\sigma}_x$ and $\tilde{\rho}_x$, describing the inferences from x and one must be implied by the other: in the above statement we can take $\phi = \tilde{\sigma}_x$ and in the corresponding statement after Eq. (5) in the preceding paragraph we can take $\psi = \tilde{\rho}_x$. Thus, $\tilde{\sigma}_x$ and $\tilde{\rho}_x$ are logically equivalent formulae: $v(\tilde{\sigma}_x \leftrightarrow \tilde{\rho}_x) = TRUE$ for any valuation v . Looking more closely at $\tilde{\sigma}_x$ we see that this formula is of a special form: it is a disjunction (cf. (4)) of conjunctions of variables (cf. (3)); such a formula is said to be in *disjunctive normal* form. Similarly, $\tilde{\rho}_x$ has the special form of a conjunction (cf. (8)) of disjunctions of variables (cf. (7)); it is in *conjunctive normal* form. By a well known result of propositional logic, any formula has logically equivalent versions in conjunctive and disjunctive normal form. (In the general definition of these forms, the variables that appear may or may not be negated. The reason why, in a space, all inferences from the presence of x in a state can be collected in conjunctive or disjunctive forms without negated variables – as testified by (3) and (7) – is that any space on X contains X as a state. In other words, negated variables can be avoided since the presence of an element x can never lead to negative conclusions regarding any

element y .)

We may conclude that surmise mappings and entail relations are in a sense very similar: both collect for any $x \in X$ the inferences that can be made from the presence of x in a state of the corresponding space. The distinction is only in the form in which these inferences are presented: the choice of the disjunctive normal form leads to a representation of the space by a surmise mapping; the choice of the conjunctive normal form amounts to the characterization of the space by an entail relation. It is easy to switch from one form to the other by using the distributive laws of logic; for instance,

$$(\bar{x}_1 \wedge \bar{x}_2) \vee (\bar{x}_3 \wedge \bar{x}_4) \Leftrightarrow (\bar{x}_1 \vee \bar{x}_3) \wedge (\bar{x}_1 \vee \bar{x}_4) \wedge (\bar{x}_2 \vee \bar{x}_3) \wedge (\bar{x}_2 \vee \bar{x}_4),$$

with a disjunctive normal form on the left and an equivalent conjunctive normal form on the right (“ \Leftrightarrow ” means “is logically equivalent to”). In this way, we can make the transition from the surmise mapping representation to the entail relation representation and vice versa on a “low”, syntactical level; i.e., without referring in any way to the represented knowledge space.

The above discussion might seem to be somewhat abstract, but the translation of surmise mapping and entail relation in logical formulae is in fact rather simple and the ensuing correspondence is very direct. Let us finally illustrate these assertions by way of a small example. Consider the following knowledge space K on $X = \{a, b, c, d\}$:

$$K = \{ \emptyset, \{a\}, \{b\}, \{a, b\}, \{a, c\}, \{b, c\}, \{a, b, c\}, \{a, b, d\}, \{b, c, d\}, X \}.$$

The corresponding surmise mapping σ can easily be found by checking, for the various $x \in X$, the minimal states containing x . As for the entail relation P , the collections $\rho(x)$ of minimal subsets A such that $A P x$ – which, as we have seen, are sufficient in (6) – can also be found from inspection of K . For instance, any state containing d contains also b and an element of $\{a, c\}$ and it is easily checked that these are all the minimal such sets. The relevant data are collected in the following table, where also the corresponding formulae $\bar{\sigma}_x$ and $\bar{\rho}_x$ have been computed.

x	$\sigma(x)$	$\bar{\sigma}_x$	$\bar{\rho}_x$	$\rho(x)$
a	$\{a\}$	\bar{a}	\bar{a}	$\{a\}$
b	$\{b\}$	\bar{b}	\bar{b}	$\{b\}$
c	$\{a, c\}, \{b, c\}$	$(\bar{a} \wedge \bar{c}) \vee (\bar{b} \wedge \bar{c})$	$(\bar{a} \vee \bar{b}) \wedge \bar{c}$	$\{a, b\}, \{c\}$
d	$\{a, b, d\}, \{b, c, d\}$	$(\bar{a} \wedge \bar{b} \wedge \bar{d}) \vee (\bar{b} \wedge \bar{c} \wedge \bar{d})$	$(\bar{a} \vee \bar{c}) \wedge \bar{b} \wedge \bar{d}$	$\{a, c\}, \{b\}, \{d\}$

Indeed, all the $\bar{\rho}_x$ can directly be computed from $\bar{\sigma}_x$ and vice versa. It is clear that $\bar{\sigma}_c$ and $\bar{\sigma}_d$ can be obtained from $\bar{\rho}_c$ and $\bar{\rho}_d$, respectively, by a single application of a

distributive law. But even if we do not recognize the special form of $\bar{\sigma}_c$ (the common “factor” \bar{c}), there is still a simple computation from $\bar{\sigma}_c$ to $\bar{\rho}_c$, by mechanically applying a distributive law and next simplifying:

$$\begin{aligned}\bar{\sigma}_c &= (\bar{a} \wedge \bar{c}) \vee (\bar{b} \wedge \bar{c}) && \text{(distributive law)} \\ \Leftrightarrow & (\bar{a} \vee \bar{b}) \wedge (\bar{a} \vee \bar{c}) \wedge (\bar{c} \vee \bar{b}) \wedge (\bar{c} \vee \bar{c}) && \text{(simplifying)} \\ \Leftrightarrow & (\bar{a} \vee \bar{b}) \wedge (\bar{a} \vee \bar{c}) \wedge (\bar{c} \vee \bar{b}) \wedge \bar{c} && \text{(simplifying)} \\ \Leftrightarrow & (\bar{a} \vee \bar{b}) \wedge \bar{c} = \bar{\rho}_c.\end{aligned}$$

We can go from $\bar{\sigma}_d$ to $\bar{\rho}_d$ by a similar, somewhat longer computation, or simply by noting the common “factor” $\bar{b} \wedge \bar{d}$. Anyway, it will be clear that we can move between surmise mapping and entail relation without even being aware of an implied knowledge space.

3. Conclusion

The theoretical concepts of the previous chapters, whose relationships we reviewed in the preceding sections, are of a very different character from those used in traditional approaches to psychometric testing. In the typical psychometric model, mental test results are analyzed in terms of the concept of ability. For this concept often a unidimensional representation is sought and if multidimensionality is allowed, the emphasis is on identifying unidimensional and preferably “independent” components. The abilities are usually represented by way of numerical scales. This is a well established approach to the assessment of knowledge, which is certainly sensible when one is interested in broad, long-term predictions concerning an individual's performance.

Here, the situation is quite different: our objective is rather to build procedures capable of assessing very accurately the current knowledge of an individual in a specific domain. For such purposes, the search for a kind of abstract, unidimensional “abilities” does not seem appropriate. We are willing to accept the multidimensional character of the situation, with all its interdependencies, and we are not trying to describe a person's performance by some summary statistic like the total number of items solved. We want to capture very concretely what is known and what is not known at a particular moment in time; that is, we want our model to deal with the full response pattern. This leads naturally to a combinatorial rather than numerical approach and to the basic definitions of knowledge state and knowledge structure.

These notions are in some sense very concrete, as indicated above. In another sense they are sufficiently abstract to make them applicable in superficially very different contexts, like that of expert systems and pattern recognition. The case of computerized medical diagnosis, which has received some attention in the area of expert systems (see, e.g., Shortliffe, 1976), is an example. Such a computerized system deals with a finite number of symptoms; the presence of some combinations of symptoms indicates specific diseases. The analogy is clear: the symptoms correspond to the items of the domain and the diseases, defined as particular subsets of symptoms, are the states. The structure of the domain consists of the collection of diseases that can be diagnosed by this system. The stages involved in building such a system are also very similar to those encountered in the knowledge assessment project. First, the structure has to be determined; that is, it has to be established which collections of symptoms correspond to diseases. This usually involves an extensive consultation of experts in the field. Next, efficient assessment procedures have to be developed. These are to determine the disease (if any) of a patient by way of a carefully designed sequence of verifications of which symptoms are present. (Note that the analogy does not break down when some symptoms are not binary ("present" vs. "absent"), but have a multicategory response (e.g., present to some specified degrees). Through appropriate dichotomizations of the multiple responses such a symptom can always be turned into a collection of binary symptoms.)

Above it was mentioned that the knowledge space approach is primarily combinatorial instead of numerical; in view of the stated objectives of efficiency and practicability this signals imminent danger. Combinatorial algorithms tend to "explode" at some point, both in terms of space and time. Indeed, knowledge structures are objects whose size grows exponentially with the number of items in the domain and the danger alluded to is very real. Chapter 8 was in fact concerned with this issue. It showed how to improve dramatically on a straightforward approach that would be completely impracticable, even for a very moderate number of items, by fully using all obtained information and by representing only the essential part of the collected data. Since there was still a horrendous worst case conceivable for the resulting algorithm, it had to prove its value in practice, on some domain of interest. As was already indicated in Chapter 8, such a test has in fact taken place and the algorithm appeared to do quite well: the practical performance was very, very far from the theoretical worst case. A number of experts went through the procedure; the results of this experiment will be the subject of a forthcoming paper (Kambouri, Koppen, Villano and Falmagne, 1989, in progress).

This brings us to the prospects for further research in this project. Obtaining knowledge spaces from experts by the procedure of Chapter 8 is only the first step in

building the knowledge structure that is ultimately going to be used in the assessment routines. Of course, different experts give us different spaces and we are confronted with the task of devising reasonable procedures for integrating all these spaces in one combination space. Such a space should feature aspects common to different experts and discard the idiosyncrasies of the individual spaces. When concrete rules have been specified to satisfy this general requirement, and a combined space is obtained, then this is still not the end of the process. At this point, the space is tested against empirical data, using the learning model developed by Falmagne (1989). Via successive likelihood ratio tests the space may possibly be pruned down, until a space results that presents a most economical, yet satisfactorily fitting model for the domain under investigation. This space, then, will finally be the representation of the domain in terms of which the computerized routines will perform the knowledge assessment.

References

- Doignon, J.-P. & Falmagne, J.-C. (1985). Spaces for the assessment of knowledge. *International Journal of Man-Machine Studies*, 23, 175-196.
- Falmagne, J.-C. (1989). A latent trait theory via a stochastic learning theory for a knowledge space. *Psychometrika*, in press.
- Kambouri, M., Koppen, M., Villano, M. & Falmagne, J.-C. (1989). Knowledge assessment: Tapping human expertise. Department of Psychology, New York University. *In progress*.
- Shortliffe, E. H. (1976). *Computer-based Medical Consultation: Mycin*. New York: American Elsevier.

SUMMARY

This thesis reports on theoretical investigations in two areas where the data are binary (e.g., "correct - incorrect") and reflect a dominance relation between two sets of entities (generically called "subjects" and "items"). The Guttman scale is the classical notion by which to represent data of this kind and both cases we consider here involve multidimensional versions of this basic idea. In the first half, we investigate the problem of representing the binary data as an intersection or union of (a minimal number of) biorders. A *biorder* is the equivalent in terms of binary relations for the concept of a Guttman scale, and representations by an intersection or union of such biorders give a similar translation of the conjunctive and disjunctive models, extensions of the Guttman scale that were proposed by Coombs and Kao, back in 1955. In the second half of the thesis we also consider data that can be modeled by multiple Guttman scales, but now we have a more specific application in mind. The "subjects" are students, the "items" notions in some field of knowledge and a collection of Guttman scales represents the different orders in which the notions can be mastered by these students. We consider this representation in the context of computerized instruction systems, more precisely, as the basis for the knowledge assessment component of such systems.

In the first chapter we discuss the common aspects of the two parts and introduce the basic concepts and terminology. Because of the ordinal character of the binary data (a Guttman scale amounts to a joint ordering of subjects and items), the mathematics involved is mainly the algebra of sets and binary relations, in particular order relations.

The chapters 2 to 5 constitute the first part of the thesis, dealing with biorder representation. Chapter 2 provides the background: the notion of a biorder is introduced and the translation is made from the conjunctive / disjunctive model of Coombs and Kao to the problem of representing a binary relation as the intersection / union of such biorders. We present part of the mathematical theory of biorder representation that was developed by Doignon, Ducamp and Falmagne in 1984, and concentrate here on the problem of finding, for any relation, the minimal number of biorders needed for such a representation. This number is called the biorder dimension, or *bidimension*, of the relation. Of central importance in this respect is Doignon *et al.*'s characterization of this bidimension as the chromatic number of some hypergraph that may be associated with any binary relation.

This equivalence is used in chapter 3 as the basis for a procedure for determining the bidimension of an arbitrary relation. We prove a theorem on how hypergraphs may be reduced to subhypergraphs without changing the chromatic number and it is shown how this result for hypergraphs in general can be effectively applied to the special kind of hypergraph defined by Doignon *et al.* Finally, we present a recurrence relation for the chromatic number of a hypergraph; the proposed reduction mechanism can be applied to the subhypergraphs generated at each level of the recursion.

Chapter 4 deals with the problem of constructing actual biorder representations for a relation. The treatment leans again heavily on the connection with the hypergraph. Algorithms for the stated purpose are derived, one version of which generates some representations in the minimum dimensionality as a by-product of the chapter 3 procedure for computing the bidimension. Another version produces exactly the collection of subhypergraphs that is needed in the recurrence relation used in chapter 3. As a result, we have a completely specified procedure for computing the bidimension.

In chapter 5 we review the results of the preceding chapters. We discuss the close connection between biorders and partial orders (the former are a generalization of the latter) and relate the success of the hypergraph approach to this connection. We also consider the prospects for application of the biorder representation to empirical data and signal two problems in this respect: the lack of uniqueness of solutions and the completely deterministic character of the model. First attempts at escaping the latter problem are sketched.

In the second part of the thesis, chapters 6 to 10, multiple Guttman scales appear in the context of knowledge assessment. Again their ordinal character is fully respected: a Guttman scale corresponds here to a possible order of mastering the various items in some specified domain of knowledge. In chapter 6 we show how a restricted collection of such orderings gives rise to a restricted collection of *knowledge states*, that is, possible subsets of items that a student in this field may have mastered. Such a family of knowledge states is called the *knowledge structure* of the domain. If we impose some extra conditions, alternative representations for such knowledge structures are possible and we discuss in chapter 6 the important case of *knowledge spaces*, i.e., structures in which the union of any two states is again a state. These can be represented by *surmise mappings*, a variant of AND/OR graphs defined by Doignon and Falmagne in 1985.

In order to efficiently assess the knowledge state of a student, the (computerized) assessment procedures must be provided with an accurate description of the set of possible states, that is, the structure of the domain under consideration. Chapters 7 and 8 address the question how we can arrive at such a description by systematically

interviewing experts in the field. It is not feasible to ask the experts simply for the list of states. Therefore, in chapter 7 an alternative representation for knowledge spaces is derived that is better suited for querying experts. This is the *entail relation*, which encodes dependencies of a certain kind between the various items of the domain; we suppose that an expert, drawing on his knowledge and experience in the field, can give us reliable answers when questioned about these dependencies. It follows from the theory in chapter 7 how the responses to these questions about the entail relation can be translated into a knowledge space, consulted implicitly by the expert.

A straightforward way of querying the expert about the entail relation would again be impracticable: even for a moderate number of items, far too many questions would need to be asked. However, since an entail relation is a relation with certain well defined properties, not all entries in such a relation are independent. In other words, in obtaining responses from the expert, we can make inferences regarding other questions, thereby reducing the number of questions that have to be asked. How to find and exploit the possible inferences is the subject of chapter 8. Here, an algorithm is specified for deriving the corresponding knowledge space from the responses of an expert to queries about the entail relation, and in this procedure at each instant possible inferences are derived in order to minimize the number of queries needed to obtain the space. This algorithm has been applied to a 50 item set, for which a straightforward approach would be unthinkable.

In chapter 9, we come back to the equivalence between knowledge spaces and surmise mappings. We present possible additional conditions on surmise mappings and investigate the consequences of these extra conditions for the corresponding classes of knowledge spaces. In particular, the condition on surmise mappings is detected that corresponds to the restriction to *well graded* knowledge spaces, an important subclass defined by Falmagne and Doignon, 1988, and Falmagne, 1989.

In chapter 10, we review the relationships between the various alternative representations derived and used in the preceding chapters. A direct equivalence between surmise mappings and entail relations (both equivalent to knowledge spaces) is established. We conclude with a few remarks on the character of and the prospects for the knowledge assessment project that is based on these theoretical concepts.

SAMENVATTING

In dit proefschrift wordt verslag gedaan van theoretisch onderzoek op twee gebieden met binaire data (b.v. "juist - onjuist") die een dominantierelatie weergeven tussen elementen van twee verschillende verzamelingen (in het algemeen aangeduid als "proefpersonen" en "items"). Voor de representatie van dit soort data bestaat de klassieke notie van de Guttmanschaal en in de beide gevallen die we hier bekijken hebben we te doen met multidimensionale versies van dit fundamentele idee. In de eerste helft stellen we ons het probleem om de relatie gegeven door een binaire data-matrix te schrijven als de doorsnede of vereniging van (een minimaal aantal) biordes. Een *biorde* is het equivalent in termen van binaire relaties voor de notie van een Guttmanschaal, en de biorde-representatie correspondeert met het conjunctieve / disjunctieve model dat al in 1955 werd geïntroduceerd door Coombs en Kao. Het tweede gedeelte van het proefschrift is ook gewijd aan data die gemodelleerd kunnen worden door meerdere Guttmanschalen, maar we hebben nu een meer specifieke toepassing op het oog. De "proefpersonen" zijn leerlingen, de "items" noties in een of ander kennisdomein en een verzameling Guttmanschalen geeft de verschillende volgordes weer waarin leerlingen zich de noties kunnen eigen maken. Deze representatie is ontwikkeld als de basis voor procedures om de kennis van leerlingen te peilen in de context van geautomatiseerde onderwijssystemen.

In het eerste hoofdstuk worden de gemeenschappelijke aspecten van de twee gedeeltes besproken en worden de basisbegrippen en -terminologie gegeven. Vanwege het ordinale karakter van de binaire data (een Guttmanschaal komt neer op een gezamenlijke ordening van proefpersonen en items) gebruiken we vooral verzamelingstheoretische wiskunde en de wiskunde van binaire relaties, met name orde relaties.

Het eerste gedeelte van het proefschrift, bestaande uit de hoofdstukken 2 tot en met 5, gaat over biorde-representatie. In hoofdstuk 2 wordt de achtergrond gegeven: het begrip biorde wordt geïntroduceerd en het conjunctieve / disjunctieve model van Coombs en Kao wordt vertaald als het probleem om een binaire relatie te schrijven als de doorsnede / vereniging van zulke biordes. Een gedeelte van de wiskundige theorie van biorde-representatie, ontwikkeld door Doignon, Ducamp en Falmagne in 1984, komt aan de orde, waarbij we ons concentreren op het probleem om het minimaal aantal biordes te vinden dat is vereist voor zo'n representatie. Dit aantal heet de biorde-dimensie of *bidimensie* van de relatie. We gaan met name in op de karakterisering door Doignon *et al.* van deze bidimensie als het kleurgetal van een

voor elke binaire relatie gedefinieerde hypergraph.

Op deze equivalentie wordt in hoofdstuk 3 een procedure gebaseerd voor de bepaling van de bidimensie van een willekeurige relatie. We geven aan hoe een hypergraph kan worden gereduceerd tot een subhypergraph zonder zijn kleurgetal te veranderen en laten zien hoe dit resultaat voor algemene hypergraphen op doeltreffende wijze kan worden toegepast op de hypergraph gedefinieerd door Doignon *et al.* Tenslotte wordt een recurrente betrekking gegeven voor het kleurgetal van een hypergraph; het voorgestelde reductie-mechanisme kan worden toegepast op de subhypergraphen die op elk niveau van de recursie worden gegenereerd.

In hoofdstuk 4 komt het probleem aan de orde om daadwerkelijk biorde-representaties van een relatie te construeren. De behandeling steunt weer nadrukkelijk op het verband met de hypergraph. We leiden algoritmes af voor het gestelde doel; een van de versies produceert een aantal representaties in de minimum dimensionaliteit als een bijproduct van de procedure voor het berekenen van de bidimensie volgens hoofdstuk 3. Een andere versie genereert precies de verzameling subhypergraphen vereist in de recurrente betrekking van hoofdstuk 3; hiermee hebben we een volledig expliciete procedure voor de berekening van de bidimensie.

Hoofdstuk 5 vormt een terugblik op de resultaten in de voorafgaande hoofdstukken. We gaan in op de nauwe samenhang tussen biorde en partiële ordes (de eerste vormen een generalisatie van de laatste) en brengen het succes van de hypergraph-benadering in verband met deze samenhang. We beschouwen ook de vooruitzichten voor toepassing van biorde-representatie op empirische data en wijzen in dit verband op twee problemen: oplossingen zijn verre van eenduidig en het model is volledig deterministisch. We schetsen enkele eerste pogingen tot een oplossing van het laatstgenoemde probleem.

In de tweede helft van het proefschrift beschouwen we collecties Guttmanschalen die corresponderen met de mogelijke volgordes waarin de verschillende noties in een bepaald kennisgebied kunnen worden verworven. In hoofdstuk 6 laten we zien hoe een welbepaalde verzameling van zulke volgordes aanleiding geeft tot een welbepaalde collectie *kennistoestanden*, d.w.z. mogelijke deelverzamelingen van geleerde noties in een vakgebied. De verzameling van alle mogelijke kennistoestanden heet de *kennisstructuur* van het gebied. Als we extra condities opleggen kunnen we alternatieve representaties voor zulke kennisstructuren vinden en we gaan in hoofdstuk 6 met name in op het belangrijke speciale geval van *kennisruimten*, gedefinieerd als structuren waarin de vereniging van elk tweetal kennistoestanden weer een kennistoestand is. Voor deze kennisruimten is er een representatie in de vorm van "surmise" functies, Doignon en Falmagne's (1985) variant op het idee van AND/OR graphs.

Voor een efficiënte meting van de kennistoestand van een leerling moet de (geautomatiseerde) procedure beschikken over een precieze beschrijving van de verzameling mogelijke toestanden, d.w.z. de kennisstructuur van het gebied in kwestie. In de hoofdstukken 7 en 8 onderzoeken we hoe zo'n beschrijving kan worden verkregen via systematische vragen aan experts in het gebied. We kunnen hen niet eenvoudigweg vragen naar de volledige lijst van mogelijke kennistoestanden. Daarom leiden we in hoofdstuk 7 een alternatieve representatie voor kennisruimten af, de zgn. "entail" relatie. Deze geeft bepaalde afhankelijkheidsrelaties tussen de verschillende noties weer en we nemen aan dat een expert, op basis van zijn kennis en ervaring in het gebied, ons betrouwbare antwoorden kan verschaffen op specifieke vragen betreffende deze afhankelijkheden. De theorie in hoofdstuk 7 vertelt ons hoe we op basis van deze antwoorden de kennisruimte kunnen construeren die de expert impliciet raadpleegt.

Als we de expert zonder meer alle vragen zouden stellen betreffende de "entail" relatie zouden we, zelfs met een zeer beperkt aantal noties, al spoedig voor een ondoenlijke taak staan: het aantal vereiste vragen zou veel te groot worden. Een "entail" relatie bezit echter, formeel, bepaalde eigenschappen die we kunnen gebruiken om gevolgtrekkingen te maken op basis van verkregen antwoorden. Op deze manier wordt een groot aantal vragen aan de expert overbodig. Hoofdstuk 8 beschrijft hoe we de mogelijke gevolgtrekkingen kunnen vinden en gebruiken. Dit resulteert in een uitgewerkt algoritme voor het afleiden van de kennisruimte van een gebied via een zorgvuldig gekozen serie vragen aan een expert. Dit algoritme is toegepast op een verzameling van 50 items, en dit aantal zou met de naïeve aanpak ondenkbaar zijn geweest.

In hoofdstuk 9 komen we terug op de equivalentie tussen kennisruimten en de zgn. "surmise" functies. We beschrijven een aantal extra condities die aan zulke functies kunnen worden opgelegd en gaan na welke consequenties deze hebben voor de overeenkomstige kennisruimten. Met name wordt vastgesteld welke klasse van "surmise" functies correspondeert met de zgn. "well graded" kennisruimten gedefinieerd door Falmagne en Doignon, 1988, en Falmagne, 1989.

In het laatste hoofdstuk wordt een kort overzicht gegeven van de diverse alternatieve wiskundige representaties die we in de vorige hoofdstukken zijn tegengekomen. We presenteren ook nog een directe equivalentie tussen "surmise" functies en "entail" relaties, die immers beide equivalent zijn met kennisruimten. We eindigen met enkele opmerkingen van algemene aard over het project dat is gebaseerd op deze theoretische begrippen.

CURRICULUM VITAE

Mathieu Koppen werd geboren op 16 augustus 1955 in Ospel. Na het eindexamen Gymnasium β aan het Bisschoppelijk College te Weert werd in 1973 begonnen met de studie wiskunde aan de K.U. te Nijmegen. In 1976 werd het bijvak psychologie de officiële studierichting en in 1978 werd, aan de K.U. te Nijmegen, het kandidaatsexamen in de psychologie behaald. Hierna moest hij in militaire dienst, en in september 1980 begon hij de doctoraalstudie psychologie aan de Rijksuniversiteit Utrecht. Via zelfstudie werd toch ook de wiskunde nog bijgehouden; dit resulteerde in januari 1983 in het M.O.-A examen en in juli 1983 in het kandidaatsexamen in de wiskunde, beide aan de K.U. Nijmegen. In mei 1984 werd te Utrecht het doctoraalexamen in de psychologie behaald (cum laude), met een gecombineerde hoofdrichting mathematische psychologie en psychologische functieleer. Onmiddellijk na zijn afstuderen werd hij als onderzoeksmedewerker aangesteld in een Z.W.O.-project onder leiding van Norman Verhelst betreffende multidimensionale uitbreidingen van de Guttmanschaal. In dit onderzoek ontstond een samenwerking met Jean-Paul Doignon van de Université Libre te Brussel en kwam hij in contact met Jean-Claude Falmagne, toen aan New York University. Deze laatste nodigde hem uit te komen werken in zijn onderzoeksgroep en sinds 1987 is hij werkzaam te New York, waar dit proefschrift tot stand kwam.

1. In onze verwachtingen betreffende multidimensionale 'analyses' van een binaire data-matrix moeten we niet uit het oog verliezen dat ieder data-punt in zo'n matrix inderdaad slechts 1 bit informatie bevat.
2. Een formulering in termen van het kleurgetal van een hypergraph is mogelijk voor diverse definities van de 'dimensie' van een binaire relatie. Deze herformulering hoeft op zich echter niet van nut te zijn bij het probleem van de berekening van zo'n dimensie. (*Dit proefschrift*).
3. Dat gevonden oplossingen niet uniek zijn is niet een speciaal probleem van de biorde-representatie, maar een kenmerk van vele multidimensionale modellen, zoals b.v. de factor-analyse. Echter, terwijl in de factor-analyse het gebrek aan uniciteit welomschreven is en, gegeven één oplossing, de volledige klasse van oplossingen is bepaald, wordt in biorde-representatie de klasse van oplossingen slechts gegeven door een volledige opsomming van haar elementen. (*Dit proefschrift*).
4. Bij toepassing in minder gestructureerde gebieden bestaat het gevaar dat geen hanteerbare *knowledge spaces* geconstrueerd kunnen worden zonder dat er leerlingen in de ruimte verloren gaan (d.w.z. zonder essentiële *knowledge states* weg te laten).
5. Hoewel om theoretische redenen te restrictief bevonden, zal een representatie van een *knowledge structure* door middel van een partiële orde in de praktijk vaak een goede benadering geven tegen aanzienlijk minder kosten.
6. Het is een groot misverstand, met name voorkomend onder zgn. 'holistisch' ingestelde psychologen, dat een eventuele volledig materialistische, fysiologische verklaring van psychische verschijnselen een aanval zou betekenen op de 'waardigheid' van de menselijke geest, zelfs het bestaan ervan zou ontkennen. Het tegendeel is veeleer het geval, net zoals we tot de tegengestelde conclusie zouden komen wanneer in het bovenstaande 'fysiologisch' wordt vervangen door 'quantum-mechanisch', 'psychisch' door 'atomair' en 'menselijke geest' door 'molecuul'.
7. Het feit dat we zelf zo slecht zijn in strict logisch redeneren en deze vaardigheid op zijn best zeer gedeeltelijk, laat in de ontwikkeling en onder kunstmatige omstandigheden leren, pleit niet voor een analyse van ons cognitief functioneren in termen van logische schakel-elementen.

8. Dat in de uitoefening van bepaalde functies, b.v. op wetenschappelijk gebied, kwaliteit niet altijd perfect en volledig objectief kan worden gemeten, is geen reden om geheel van dergelijke beoordelingen af te zien.

9. Hoewel het woord 'psycholoog' bij het algemene publiek nog altijd bepaalde vaste associaties oproept, heeft van alle wetenschapsgebieden de psychologie de vogels van de meest diverse pluimage onder haar hoede.

10. Zolang we niet het equivalent hebben van de Amerikaanse honkbalverslaggever die ons feilloos, in drie decimalen, het honkslagpercentage weet te melden van *déze* slagman, in *déze* situatie: eerste en tweede honk bezet en twee man uit, tegen linkshandige werpers, over de laatste vijf seizoenen, zal het nooit echt wat worden met de Nederlandse sportverslaggeving.

11. De meeste promovendi hebben de neiging om één stelling teveel op te nemen.

Stellingen behorende bij het proefschrift van Mathieu Koppen,
Ordinal data analysis: biorder representation and knowledge spaces,
Nijmegen, 21 augustus 1989.

