



Available at

[www.ElsevierComputerScience.com](http://www.ElsevierComputerScience.com)

POWERED BY SCIENCE @ DIRECT®

INTERNATIONAL JOURNAL OF  
APPROXIMATE  
REASONING

ELSEVIER International Journal of Approximate Reasoning xxx (2003) xxx–xxx

[www.elsevier.com/locate/ijar](http://www.elsevier.com/locate/ijar)

## Hyperrelations in version space

Hui Wang <sup>a,\*</sup>, Ivo Düntsch <sup>b</sup>, Günther Gediga <sup>c</sup>,  
Andrzej Skowron <sup>d</sup>

<sup>a</sup> School of Information and Software Engineering, University of Ulster, Faculty of Informatics,  
Newtownabbey, Co. Antrim BT 37 0QB, Ireland

<sup>b</sup> Department of Computer Science, Brock University, St. Catharines, Ont., Canada L2S 3A1

<sup>c</sup> Institut für Evaluation und Marktanalysen, Brinkstr. 19, 49143 Jeggen, Germany

<sup>d</sup> Institute of Mathematics, University of Warsaw, 02-097 Warszawa, Poland

Received 1 October 2002; accepted 1 October 2003

---

### Abstract

12 A version space is a set of all hypotheses consistent with a given set of training  
13 examples, delimited by the *specific boundary* and the *general boundary*. In existing  
14 studies [Machine Learning 17(1) (1994) 5; Proc. 5th IJCAI (1977) 305; Artificial Intel-  
15 ligence 18 (1982)] a hypothesis is a conjunction of attribute-value pairs, which is shown  
16 to have limited expressive power [Machine Learning, The McGraw-Hill Companies, Inc  
17 (1997)]. In a more expressive hypothesis space, e.g., disjunction of conjunction of  
18 attribute-value pairs, a general version space becomes uninteresting unless some  
19 restriction (inductive bias) is imposed [Machine Learning, The McGraw-Hill Compa-  
20 nies, Inc (1997)].

21 In this paper we investigate version space in a hypothesis space where a hypothesis is  
22 a *hyperrelation*, which is in effect a disjunction of conjunctions of disjunctions of  
23 attribute-value pairs. Such a hypothesis space is more expressive than the conjunction of  
24 attribute-value pairs and the disjunction of conjunction of attribute-value pairs. How-  
25 ever, given a dataset, we focus our attention only on those hypotheses which are con-  
26 sistent with given data and are maximal in the sense that the elements in a hypothesis  
27 cannot be *merged* further. Such a hypothesis is called an *E-set* for the given data, and the  
28 set of all *E-sets* is the version space which is delimited by the least *E-set* (specific  
29 boundary) and the greatest *E-set* (general boundary).

---

\* Corresponding author. Tel.: +44-1232-368981; fax: 44-1232-366068.

E-mail addresses: [h.wang@ulst.ac.uk](mailto:h.wang@ulst.ac.uk) (H. Wang), [duentsch@cosc.brocku.ca](mailto:duentsch@cosc.brocku.ca) (I. Düntsch),  
[gediga@eval-institut.de](mailto:gediga@eval-institut.de) (G. Gediga), [skowron@mimuw.edu.pl](mailto:skowron@mimuw.edu.pl) (A. Skowron).

30 Based on this version space we propose three classification rules for use in different  
31 situations. The first two are based on  $E$ -sets, and the third one is based on “degraded”  
32  $E$ -sets called *weak hypotheses*, where the maximality constraint is relaxed. We present an  
33 algorithm to calculate  $E$ -sets, though it is computationally expensive in the worst case.  
34 We also present an efficient algorithm to calculate weak hypotheses. The third rule is  
35 evaluated using public datasets, and the results compare well with C5.0 decision tree  
36 classifier.

37 © 2003 Published by Elsevier Inc.

38

---

## 39 1. Introduction

40 *Version space* is a useful and powerful concept in the area of concept  
41 learning: a version space is the set of all hypotheses in a hypothesis space  
42 consistent with a dataset. A version space is delimited by the *specific boundary*  
43 and the *general boundary*—the sets of most specific and most general  
44 hypotheses respectively, to be explained below.

45 The ELIMINATE-CANDIDATE algorithm [6,8] is the flagship algorithm used to  
46 construct a version space from a dataset. The hypothesis space used in this  
47 algorithm is the conjunction of attribute-value pairs, and it is shown to have  
48 limited expressive power [9]. It is also shown by [4] that the size of the general  
49 boundary can grow exponentially in the number of training examples, even when  
50 the hypothesis space consists of simple conjunctions of attribute-value pairs.

51 Table 3 was used in [9] to show the limitation of Mitchell’s representation. It  
52 was shown that the ELIMINATE-CANDIDATE algorithm cannot come up with a  
53 consistent specific boundary or general boundary for this example due to this  
54 limitation (i.e., the chosen hypothesis space). A possible solution, by increasing  
55 the expressive power of the representation, is to extend the hypothesis space to  
56 disjunctions of conjunctions of attribute-value pairs from the restrictive con-  
57 junctions of attribute-value pairs. But it turned out that, without proper  
58 inductive bias, the specific boundary would *always* be overly specific (the dis-  
59 junction of the observed positive examples) and the general boundary would  
60 *always* be overly general (the negated disjunction of the observed negative  
61 examples) [9], so they do not carry useful information—they are “uninterest-  
62 ing”. The desired inductive bias, however, has not been explored.

63 In this paper we report a study on version space which is aimed at increasing  
64 the expressive power of the hypothesis space and, at the same time, keeping the  
65 hypotheses “interesting” through introducing an inductive bias. We explore a  
66 semilattice structure that exists in the set of all *hypertuples* of a domain, where  
67 hypertuples generalise the traditional tuples from value-based to set-based. The  
68 semilattice structure can be used as a base for a hypothesis space. We take a  
69 hypothesis to be a *hyperrelation*, i.e., a set of hypertuples. A hyperrelation can

70 be interpreted as a disjunction of hypertuples, which can be further interpreted  
 71 as a conjunction of disjunctions of attribute-value pairs. Such a hypothesis  
 72 space is much more expressive than the conjunction of attribute-value pairs  
 73 and the disjunction of conjunction of attribute-value pairs. For a dataset there  
 74 is a large number of hypertuples which are consistent with the data, some of  
 75 which can be merged (through the semilattice operation) to form a different  
 76 consistent hypertuple. We propose to focus on those hypertuples which are  
 77 consistent with the given data and cannot be merged further—they are said to  
 78 be *maximal*. A hypothesis is then a set of these hypertuples, and the version  
 79 space is the set of all these hypertuples. Clearly this version space is a subset of  
 80 the semilattice. An algorithm is presented which is able to construct the version  
 81 space.

82 A detailed analysis of Mitchell's version space shows that our version space  
 83 is indeed a generalisation of Mitchell's and that it is more expressive than  
 84 Mitchell's.

85 Classification within our version space is not trivial. We explore ways in  
 86 which data can be classified using different hypotheses in the version space. We  
 87 also examine how this whole process can be speeded up for practical benefits.  
 88 Experimental results are presented to show the effectiveness of this approach.

89 Most of the theory (and results) of version space can be expressed by the  
 90 tools offered by Boolean reasoning as investigated, for example, in [10–13,15].  
 91 Unlike some of these publications we will only be concerned with consistent  
 92 decision systems and exact classification rules.

## 93 2. Definitions and notation

94 For a set  $U$ , we let  $2^U$  be the powerset of  $U$ . If  $\leq$  is a partial ordering on  $U$ ,  
 95 and  $X \subseteq U$ , we let

$$\uparrow X = \{y \in U : (\exists x \in X)x \leq y\},$$

$$\downarrow X = \{y \in U : (\exists x \in X)y \leq x\}.$$

97 If  $X, Y \subseteq U$ , we say that  $Y$  covers  $X$ , written as  $X \preceq Y$ , if  $X \subseteq \downarrow Y$ , i.e. if for each  
 98  $x \in X$  there is some  $y \in Y$  with  $x \leq y$ . We call  $X$  dense for  $Y$ , written  $X \trianglelefteq Y$ , if for  
 99 any  $y \in Y$  there is some  $x \in X$  such that  $x \leq y$ , i.e.  $Y \subseteq \uparrow X$ . Note that  $Y$  covers  $X$   
 100 iff  $Y$  is dually dense for  $X$ .

101 A *semilattice* is an algebra  $\langle A, + \rangle$ , such that  $+$  is a binary operation on  $A$   
 102 satisfying

$$x + y = y + x,$$

$$(x + y) + z = x + (y + z),$$

$$x + x = x.$$

4

*H. Wang et al. / Internat. J. Approx. Reason. xxx (2003) xxx–xxx*

104 If  $X \subseteq A$ , we denote by  $[X]$  the sub-semilattice of  $A$  generated by  $X$ . With  
 105 some abuse of notation, we also use  $+$  for the complex sum, i.e.  
 106  $X + Y = \{x + y : x \in X, y \in Y\}$ .

107 Each semilattice has an intrinsic order defined by

$$x \leq y \iff x + y = y. \quad (2.1)$$

109  $x$  is called *minimal (maximal)* in  $X \subseteq U$ , if  $y \leq x$  ( $x \leq y$ ) implies  $x = y$ . The set  
 110 of all minimal (maximal) elements of  $X$  is denoted by  $\min X$  ( $\max X$ ).

111 A *decision system* is a tuple  $I = \langle U, \Omega, \{V_a : a \in \Omega\}, d, V_d \rangle$ , where

- 112 1.  $U = \{x_0, \dots, x_N\}$  is a nonempty finite set.
- 113 2.  $\Omega = \{a_0, \dots, a_T\}$  is a nonempty finite set of mappings  $a_i : U \rightarrow V_{a_i}$ .
- 114 3.  $d : U \rightarrow V_d$  is a mapping.

115 Set  $\mathbf{V} = \prod_{a \in \Omega} V_a$ . We let  $I : U \rightarrow \mathbf{V}$  be the mapping

$$x \mapsto \langle a_0(x), a_1(x), \dots, a_T(x) \rangle \quad (2.2)$$

117 which assigns to each object  $x$  its feature vector.  $\mathbf{V}$  is called *data space*, and the  
 118 collection  $\mathbf{D} = \{I(x) : x \in U\}$  is called the *training space*. In the sequel  $\mathbf{V}$  is  
 119 assumed finite (consequently  $\mathbf{D}$  is finite) unless otherwise stated. The mapping  
 120  $d$  defines a labelling (or partition) of  $\mathbf{D}$  into classes  $\mathbf{D}_0, \mathbf{D}_1, \dots, \mathbf{D}_K$ , where  $I(x)$   
 121 and  $I(y)$  are in the same class  $\mathbf{D}_i$  if and only if  $d(x) = d(y)$ . If  $a \in \Omega$ ,  $v \in V_a$ , then  
 122  $\langle a, v \rangle$  is called a *descriptor*.

123 We call

$$\mathcal{L} = \prod_{a \in \Omega} 2^{V_a} \quad (2.3)$$

125 the *extended data space*, its elements *hypertuples*, and its subsets *hyperrelations*.  
 126 In contrast we call the elements of  $\mathbf{V}$  *simple tuples* and its subsets *simple*  
 127 *relations*. If  $s \in \mathcal{L}$  and  $a \in \Omega$  we let  $s(a)$  be the projection of  $s$  with respect to  $a$ ,  
 128 i.e. the set appearing in the  $a$ th component.

129  $\mathcal{L}$  is a lattice under the following natural order

$$s \leq t \iff s(a) \subseteq t(a) \text{ for all } a \in \Omega,$$

131 with the least upper bound ( $+$ ) and greatest lower bound ( $\times$ ) operations given  
 132 by

$$\begin{aligned} (t + s)(a) &= t(a) \cup s(a), \\ (t \times s)(a) &= t(a) \cap s(a) \end{aligned}$$

134 for all  $a \in \Omega$ .  $s$  and  $t$  are said to be *overlapping* if there is some simple tuple  $d$   
 135 such that  $d \leq t \times s$ . By identifying a singleton with the element it contains, we  
 136 may suppose that  $\mathbf{V} \subseteq \mathcal{L}$ ; in particular, we have  $\mathbf{D} \subseteq \mathcal{L}$ .

Table 1

$a_1$	$a_2$	$y$
<i>(a) A set of simple tuples</i>		
$a$	0	$\alpha$
$a$	1	$\alpha$
$b$	1	$\alpha$
$b$	2	$\beta$
<i>(b) A set of hypertuples</i>		
$\{a, b\}$	$\{0, 1\}$	$\alpha$
$\{b\}$	$\{2\}$	$\beta$

137 Table 1(a) below shows a simple relation (i.e., a set of simple tuples) and

138 Table 1(b) a hyperrelation.

139 For unexplained notation and background reading in lattice theory, we

140 invite the reader to consult [3].

### 141 3. Hypertuples and hyperrelations

142 Suppose that  $\mathbf{D}_q$  is a labelled class. The idea of [2,19] was to collect, if  
 143 possible, across the elements of  $\mathbf{D}_q$  the values of the same attributes without  
 144 destroying the labelling information. For example, if the system generates the  
 145 rules

If  $x$  is blue, then sell,  
 if  $x$  is green, then sell,

147 we can identify the values “blue” and “green” and collect these into the single  
 148 rule

If  $x$  is blue or green, then sell.

150 Our aim is to maximise the number of attribute values on the left hand side of  
 151 the rule in order to increase the generality of the rule and also increase its base  
 152 as to guard against random influences [1].

153 This leads to the following definition [18]: We call an element  $r \in \mathcal{L}$  *equi-*  
 154 *labelled* with respect to the class  $\mathbf{D}_q$ , if

$$\downarrow r \cap \mathbf{D} \neq \emptyset \quad (3.1)$$

156 and

$$\downarrow r \cap \mathbf{D} \subseteq \mathbf{D}_q. \quad (3.2)$$

158 These two conditions express two fundamental principles in machine learning:

159 (3.1) says that  $r$  is supported by some  $d \in \mathbf{D}$ ; it is a minimality condition in the

160 sense that we do not want to consider hypertuples that have nothing to do with  
 161 the training set. Condition (3.2) tells us that  $r$  needs to be consistent with the  
 162 training set: all elements below  $r$  which are in the training set  $\mathbf{D}$  are in the  
 163 unique class  $\mathbf{D}_q$ . We denote the set of all elements equilabelled with respect to  
 164  $\mathbf{D}_q$  by  $\mathcal{E}_q$ , and let  $\mathcal{E} = \bigcup_{q \leq K} \mathcal{E}_q$  be the set of all equilabelled elements. Note that  
 165  $\mathbf{D} \subseteq \mathcal{E}$ , and that

$$q, r \leq K, q \neq r \text{ implies } \mathcal{E}_q \cap \mathcal{E}_r = \emptyset, \quad (3.3)$$

167 since  $\{\mathbf{D}_0, \dots, \mathbf{D}_K\}$  partitions  $\mathbf{D}$ . Furthermore,

$$\uparrow \mathcal{E}_q \cap \mathcal{E} = \mathcal{E}_q. \quad (3.4)$$

169 If  $\mathbf{D} \subseteq P \subseteq \mathcal{L}$ , we let  $E(P) = \{t \in \mathcal{L} : t \text{ is maximal in } [P] \cap \mathcal{E}\}$ , and set  
 170  $\mathbf{VSp} = \bigcup \{E(P) : \mathbf{D} \subseteq P \subseteq \mathbf{V}\}$ . For each  $q \leq K$ , let  $E_q(P) = E(P) \cap \mathcal{E}_q$ .

171 Each set of the form  $E(P)$  is called an  $E$ -set or a *hypothesis*. We now have

172 **Lemma 3.1.** *Let  $\mathbf{D} \subseteq P \subseteq Q \subseteq \mathcal{L}$ . Then,*

173 1.  $E(P) \preceq E(Q)$ .

174 2.  $E(P) \trianglelefteq E(Q)$ .

175 **Proof.** Both claims follow immediately from the fact that  
 176  $\emptyset \neq [P] \cap \mathcal{E} \subseteq [Q] \cap \mathcal{E}$ , and that  $[Q] \cap \mathcal{E}$  is finite.  $\square$

177 **Theorem 3.2.**  $\min \mathbf{VSp} = E(\mathbf{D})$ ,  $\max \mathbf{VSp} = E(\mathbf{V})$ .

178 **Proof.** We only prove the first part; the second part can be similarly proved.

179 “ $\subseteq$ ”: Suppose that  $t$  is minimal in  $\mathbf{VSp}$ . Then, there is some  $\mathbf{D} \subseteq P \subseteq \mathbf{V}$  such  
 180 that  $t \in E(P)$ . Since  $E(\mathbf{D}) \trianglelefteq E(P)$  by Lemma 3.1, there is some  $s \in E(\mathbf{D})$  such  
 181 that  $s \leq t$ , and the minimality of  $t$  implies  $s = t$ .

182 “ $\supseteq$ ”: If  $t \in E(\mathbf{D})$  and  $s \leq t$  is minimal in  $\mathbf{VSp}$ , then  $s \in E(\mathbf{D})$  as well. Since the  
 183 elements of  $E(\mathbf{D})$  are pairwise incomparable, we have  $s = t$ .  $\square$

184 In the sequel, we will denote  $E(\mathbf{D})$  by  $\mathbb{S}$ , and  $E(\mathbf{V})$  by  $\mathbb{G}$ , since they corre-  
 185 spond, respectively, to the specific and general boundaries of  $\mathbf{D}$  in the sense of  
 186 7; we also set  $\mathbb{S}_q = E_q(\mathbf{D})$  and  $\mathbb{G}_q = E_q(\mathbf{V})$ . By Lemma 3.1,  $\mathbb{S}$  is the least  $E$ -set  
 187 and  $\mathbb{G}$  is the greatest  $E$ -set.

188 An algorithm to find  $E(\mathbf{D})$ , called the *lattice machine* (LM), has been pre-  
 189 sented in [18], and it can easily be generalised to find  $E(P)$  for each  $\mathbf{D} \subseteq P \subseteq \mathbf{V}$ :

190 **Algorithm** (LM algorithm).

191 1.  $H_0 \stackrel{\text{def}}{=} P$ .

192 2.  $H_{k+1} \stackrel{\text{def}}{=} \text{The set of maximal elements of } [\downarrow (H_k + P)] \cap \mathcal{E}$ .

193 3. Continue until  $H_n = H_{n+1}$  for some  $n$ .

194 It was shown in [19] that  $H_n = E(P)$ . If  $|P| = m$ , then the complexity of the  
 195 algorithm is  $O(2^m)$  in the worst case.

196 To illustrate the above notions we consider the following example.

197 **Example.** A dataset is shown in Fig. 1. The small circles and crosses are simple  
 198 tuples of different classes while the rectangles are hypertuples.

199 Every hypertuples depicted here are equilabelled since they cover only  
 200 simple tuples of the same class. All the simple tuples are also equilabelled.

201 Each hypertuple here is maximal since we cannot move its boundary in any  
 202 dimension while maintaining it being equilabelled.

203 The set of all hypertuples in the figure is an  $E$ -set or a hypothesis for the  
 204 underlying concept implied by the dataset, since all hypertuples are equila-  
 205 belled and maximal, and they together cover all the simple tuples. In fact this  
 206  $E$ -set is  $E(D)$ —the specific boundary. The general boundary is not depicted  
 207 here as we need knowledge of the domain of each attribute to calculate it.  
 208 Clearly the  $E$ -set does not cover the whole data space; in other word, there are  
 209 regions that are not covered by the  $E$ -set.

210 Primary elements are those that are covered explicitly by some equilabelled  
 211 hypertuples; secondary elements are those not covered explicitly but can be  
 212 covered by extending some equilabelled hypertuple without compromising its  
 213 equilabelledness. We can easily find primary and secondary elements from this  
 214 figure.

215 Our next aim is to show that every simple tuple is covered by some equi-  
 216 labelled element:

217 **Lemma 3.3.** For each  $t \in \mathbf{V}$  there is some  $h \in \mathcal{E}$  such that  $t \leq h$ .

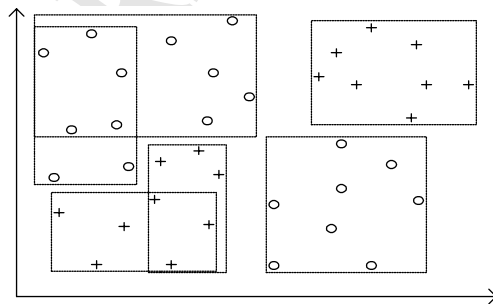


Fig. 1. A dataset and its specific boundary. Each circle or cross is a simple tuple, and each rectangle represents a hypertuple.

218 **Proof.** We show that there is some  $x \in \mathbf{D}$  such that  $\downarrow(t+x) \cap \mathbf{D} = \{x\}$ ; then,  
 219  $t \leq t+x \in \mathcal{E}$ : Let  $t \in V$ , and set  $M = \{t+x : x \in \mathbf{D}\}$ . Suppose that  $x \in \mathbf{D}$  such  
 220 that  $t+x$  is minimal in  $M$  with respect to  $\leq$ . Let  $y \in \mathbf{D}$  such that  $y \leq t+x$ , and  
 221 assume that  $y \neq x$ ; then, there is some  $a \in \Omega$  such that  $y(a) \neq x(a)$ . Since both  $y$   
 222 and  $x$  are simple tuples, we have in fact  $y(a) \cap x(a) = \emptyset$ . Now,  $y \leq t+x$  implies  
 223 that  $y(a) \subseteq t(a) \cup x(a)$ , and it follows from  $y(a) \cap x(a) = \emptyset$ —and the fact that  $t$   
 224 is simple—that  $y(a) = t(a) \subsetneq t(a) \cup x(a)$ . Hence,  $t+y \not\leq t+x$ , contradicting the  
 225 minimality of  $t+x$  in  $M$ .  $\square$

226 Since each element of  $\mathcal{E}$  is covered by some element of  $\mathbb{G}$ , and the sets  $\mathbb{G}_i$ ,  
 227  $i \leq K$ , partition  $\mathbb{G}$ , we obtain

228 **Corollary 3.4.** For each  $t \in \mathbf{V}$ , there is some  $i \leq K$  such that  $t \leq \mathbb{G}_i$ .

229 Observe that such  $\mathbb{G}_i$  need not be unique, and we need a selection procedure  
 230 to label  $t$ . Our chosen method is majority voting, combined with random  
 231 selection for tied maxima: Let

$$m'(t, n, \mathbb{G}) = |\{h \in \mathbb{G}_n : t \leq h\}| \quad \text{for each } n \leq K, \tag{3.5}$$

$$m(t, \mathbb{G}) = \max\{M(t, n, \mathbb{G}) : n \leq K\}, \tag{3.6}$$

$$M(t, \mathbb{G}) = \{i \leq K : m'(t, i, \mathbb{G}) = m(t, \mathbb{G})\}. \tag{3.7}$$

Rule 1. If  $t \in \mathbf{V}$ , then label  $t$  by a randomly chosen element of  $M(t, \mathbb{G})$ .

234 Now we consider an arbitrary hypothesis  $H$ . We can replace the  $\mathbb{G}$  by  $H$  in  
 235 Eqs. (3.5)–(3.7), and apply Rule 1 for classification. However, unlike the case  
 236 where  $H = \mathbb{G}$ , there is no guarantee that  $t \leq H$ , i.e. that  $m(t, H) \neq 0$ .

237 This motivates us to distinguish between different elements in  $\mathbf{V}$  given  
 238  $H = E(P) : t \in \mathbf{V}$  is called *primary* if there is  $H_q$  with  $t \leq H_q$ ,  $t$  is *secondary* if  
 239 there is  $h \in H_q$  such that  $t+h$  is equilabelled (i.e.,  $t+h$  does not overlap  $H_p$  for  
 240  $p \neq q$ ), and  $t$  is *tertiary* otherwise. Note that an element can be both primary  
 241 and secondary. In fact primary data is a subset of secondary data.

242 Now we let

$$m'_p(t, n, H) = |\{h \in H_n : t \leq h\}| \quad \text{for each } n \leq K, \tag{3.8}$$

$$m_p(t, H) = \max\{m'_p(t, n, H) : n \leq K\}, \tag{3.9}$$

$$M_p(t, H) = \{i \leq K : m'_p(t, i, H) = m_p(t, H)\}. \tag{3.10}$$

244 and



$$m'_s(t, n, H) = |\{h \in H_n : t + h \in \mathcal{E}\}| \quad \text{for each } n \leq K, \tag{3.11}$$

$$m_s(t, H) = \max\{m'_s(t, n, H) : n \leq K\}, \tag{3.12}$$

$$M_s(t, H) = \{i \leq K : m'_s(t, i, H) = m_s(t, H)\}. \tag{3.13}$$

*Rule 2*

- If  $t$  is primary, then label  $t$  by a randomly chosen element of  $M_p(t, H)$ .
- If  $t$  is secondary, then label  $t$  by a randomly chosen element of  $M_s(t, H)$ .
- Otherwise, label  $t$  as *unclassified*.

250 Now we explore a generalisation of both primary and secondary data in the  
 251 following sense:

252 **Lemma 3.5.** *If  $H = E(P)$ , we let  $P' = P \cup \{t\}$  and  $H' = E(P')$ . Then,*

*$t$  is secondary for  $H \Rightarrow t$  is primary for  $H'$ .*

**Proof.** Let  $h \in H$  such that  $t + h$  is equilabelled. Then,  $t + h \in [P'] \cap \mathcal{E}$ , and  
 255 thus,  $t \leq g$  for some  $g \in E(P')$ .  $\square$

256 The converse is not true: Let  $\mathbf{D} = \{a, b, c, d, e, f\}$  with

$$\begin{aligned} a &= \langle 0, 0, 0 \rangle, & b &= \langle 1, 0, 0 \rangle, & c &= \langle 1, 0, 1 \rangle, & d &= \langle 1, 1, 1 \rangle, \\ e &= \langle 0, 1, 1 \rangle, & f &= \langle 0, 1, 0 \rangle. \end{aligned}$$

258 Suppose that  $a, b, e$  are coloured blue and  $d, e, f$  are coloured red. Then,

$$E(\mathbf{D}) = \{a + b, e, c + d, f\}.$$

260 The aim is to classify  $t = \langle 0, 0, 1 \rangle$  with respect to the hypothesis  $H = E(\mathbf{D})$ .

261 Now,

$$\begin{aligned} a + t &= \langle 0, 0, 01 \rangle \text{ is equilabelled blue,} \\ e + t &= \langle 0, 01, 1 \rangle \text{ is equilabelled blue,} \\ c + t &= \langle 01, 0, 1 \rangle \text{ is equilabelled red,} \end{aligned}$$

263 while  $b + t, d + t, f + t$  are not equilabelled. Furthermore,  $q + t$  is not equi-  
 264 labelled for any  $q \in E(\mathbf{D})$ , and thus,  $t$  is not secondary with respect to  $E(\mathbf{D})$ . If  
 265 we admit that the knowledge provided by  $t$  should be admitted when we try to  
 266 classify  $t$  by  $H = E(P)$ , then we should classify by using primary data of  
 267  $E(P \cup \{t\})$ . At any rate, admitting  $t$  does not cause any inconsistencies, and  
 268 using primary data of  $H'$  is well in line with our aim of maximising consistency.

## 269 3.1. Weak hypotheses

270 In the previous discussion we have introduced two classification rules based  
 271 on  $E$ -sets. Rule 1 can be applied if  $\mathbb{G}$  can be practically constructed from  $\mathbf{V}$  by  
 272 the LM algorithm, and Rule 2 can be applied if  $\mathbb{S}$  (or  $E(P)$ ) can be practically  
 273 constructed from  $\mathbf{D}$  (for  $\mathbf{D} \subseteq P \subseteq \mathbf{V}$ ). In both cases we have to construct  $E$ -sets  
 274 for  $P$  where  $\mathbf{D} \subseteq P \subseteq \mathbf{V}$ . From the LM algorithm we know that constructing an  
 275  $E$ -set is expensive and in the worst case it is exponential. Since the most time is  
 276 spent finding maximal hypertuples, we make the following definition: A *weak*  
 277 *hypothesis* is a set of equilabelled hypertuples which covers  $\mathbf{D}$ . The following  
 278 algorithm finds a weak hypothesis  $H$  [17]:

```

279  $H \leftarrow \emptyset$ 
280 for  $q = 0$  to  $K$  do
281    $X \leftarrow D_q$ 
282   while  $X \neq \emptyset$  do
283     Order  $X$  as  $\langle g_0, \dots, g_{m(X)} \rangle$ 
284      $h \leftarrow g_0$ 
285     for  $i = 1$  to  $m(X)$  do
286       if  $h + g_i$  is equilabelled then
287          $h \leftarrow h + g_i$ 
288       end if
289     end for
290      $X \leftarrow X \setminus \downarrow h$ 
291      $H \leftarrow H \cup \{h\}$ 
292   end while
293 end for

```

294 The algorithm does not necessarily produce disjoint hypertuples: Let

$$D_0 = \{\langle 0, 1, 1 \rangle, \langle 1, 0, 1 \rangle, \langle 1, 1, 1 \rangle, \langle 1, 1, 0 \rangle, \langle 2, 1, 1 \rangle\},$$

$$D_1 = \{\langle 2, 0, 1 \rangle, \langle 0, 0, 0 \rangle\}.$$

296 Order  $X = D_0$  by

$$g_0 = \langle 0, 1, 1 \rangle,$$

$$g_1 = \langle 1, 0, 1 \rangle,$$

$$g_2 = \langle 1, 1, 1 \rangle,$$

$$g_3 = \langle 1, 1, 0 \rangle,$$

$$g_4 = \langle 2, 1, 1 \rangle.$$

298 Now,  $h_0 = g_0 + g_1 + g_2$  is equilabelled, while  $h_0 + g_3$  and  $h_0 + g_4$  are not. The  
 299 next step produces  $h_1 = g_3 + g_4$ . However,  $g_2 \leq h_0$  and  $g_2 \leq h_1$ .

300 In the worst case the time complexity for building  $H_q$  (the hypothesis for  
 301 class  $\mathbf{D}_q$ ) is  $O(|\mathbf{D}_q|^2)$ . Therefore the worst case complexity for building  $H$  (the  
 302 whole hypothesis) is  $O(K \times |\mathbf{D}_q|^2)$ , where  $K$  is the number of classes.

303 Let  $H = \bigcup_{i \leq K} H_i$  be a weak hypothesis for  $\mathbf{D}$ , where  $H_i = \{h_0, \dots, h_{t(i)}\}$  is a  
 304 weak hypothesis for class  $i$  and  $h_j$  is an equilabelled hypertuple. Let  $M_p(t, H)$   
 305 and  $M_s(t, H)$  be defined as in Eqs. (3.8) and (3.11). Then we introduce the  
 306 following rule, which is the same as Rule 2 except that the hypothesis here is  
 307 weak.

308 *Rule 3*

- 309 • If  $t$  is primary, then label  $t$  by a randomly chosen element of  $M_p(t, H)$ .
- 310 • If  $t$  is secondary, then label  $t$  by a randomly chosen element of  $M_s(t, H)$ .
- 311 • Otherwise, label  $t$  as *unclassified*.

#### 312 4. Mitchells version space

313 The situation in [9] can be described as a special decision system where  
 314  $d : U \rightarrow \{0, 1\}$ , and it is called the *target concept*. The set of positive examples  
 315 is denoted by  $D_1$ , and that of negative examples by  $D_0$ . We will describe the  
 316 concepts introduced there in our notation and we will follow the explanation in  
 317 [9, p. 22, Table 2.2].

318 A *hypothesis* is a hypertuple  $t \in \mathcal{L}$  such that for all  $a \in \Omega$

$$|t(a)| \leq 1 \text{ or } t(a) = V_a. \quad (4.1)$$

320 Thus, for each  $a \in \Omega$  we have

$$t(a) = \begin{cases} \emptyset, & \text{or} \\ m, & \text{for some } m \in V_a, \text{ or} \\ V_a. & \end{cases}$$

322 The set of all hypotheses is denoted by  $\mathcal{H}$ . Observe that  $\mathcal{H}$  is a hyperrelation,  
 323 and that  $\mathcal{H}$  is partially ordered by  $\leq$  as defined by (2.1), and that  $\mathbf{V} \subseteq \mathcal{H}$ . If  
 324  $s, t \in \mathcal{H}$ , and  $s \leq t$ , we say that  $s$  is *more specific than*  $t$  or, equivalently, that  $t$  is  
 325 *more general than*  $s$ . We say that  $s$  *satisfies hypothesis*  $t$ , if

- 326 1.  $s$  is a simple tuple, i.e.  $s \in \mathbf{V}$ ,
- 327 2.  $s \leq t$ ,

328 and denote the set of all (simple) tuples that satisfy  $t$  by  $\text{sat}(t)$ ; observe that  
 329  $\text{sat}(t) = \downarrow t \cap \mathbf{V}$ . More generally, for  $A \subseteq \mathcal{L}$  we let

$$\text{sat}(A) = \downarrow A \cap \mathbf{V}.$$

331 If  $t(a) = \emptyset$  for some  $a \in \Omega$ , then  $t$  cannot be satisfied. We interpret  $s \in \text{sat}(t)$  as  
 332 “instance  $s$  is classified by hypothesis  $t$ ”. According to [9],  $t$  is more general

333 than  $s$ , if any instance classified by  $s$  is also classified by  $t$ ; in other words,  
 334  $\text{sat}(s) \subseteq \text{sat}(t)$ . That our notion captures this concept is shown by the following  
 335 result, the easy proof of which is left to the reader.

336 **Theorem 4.1.** *Suppose that  $s, t \in \mathcal{H}$ . Then,*

$$s \leq t \iff \text{sat}(s) \subseteq \text{sat}(t).$$

338 Since we are interested in hypotheses whose satisfiable observations are  
 339 within one class of  $d$ , we say that  $t \in \mathcal{H}$  is *consistent with*  $\langle \mathbf{D}, d \rangle$ , if

$$\downarrow t \cap \mathbf{D} = D_1. \quad (4.2)$$

341 Thus, in this case, the training examples satisfying  $t$  are exactly the positive  
 342 ones. The *version space*  $\mathbf{VSp}^m$  is now the set of all hypotheses consistent with  
 343  $\langle \mathbf{D}, d \rangle$ . In other words,

$$\mathbf{VSp}^m = \{t \in \mathcal{H} : (\forall s)[s \in \text{sat}(t) \cap \mathbf{D} \iff d(s) = 1]\}.$$

345 The *general boundary*  $\mathbb{G}^m$  is the set of maximal members of  $\mathcal{H}$  consistent with  
 346  $\langle \mathbf{D}, d \rangle$ , i.e.

$$\mathbb{G}^m = \max \mathbf{VSp}^m.$$

348 The *specific boundary*  $\mathbb{S}^m$  is the set of minimal members of  $\mathcal{H}$  consistent with  
 349  $\langle \mathbf{D}, d \rangle$ , i.e.

$$\mathbb{S}^m = \min \mathbf{VSp}^m.$$

351 The two boundaries delimit the version space in the sense that for any con-  
 352 sistent hypothesis  $t$  in the version space there are  $g \in \mathbb{G}^m$  and  $s \in \mathbb{S}^m$  such that  
 353  $s \leq t \leq g$ .

354 The example in Table 2 illustrates the idea of version space. We follow [9] in  
 355 writing ? in column  $a$ , if  $a(x) = V_a$ .

## 356 5. Expressive power of version space and boolean reasoning

357 A simple tuple  $t$  can be regarded as a conjunction of descriptors

$$\langle a_0, t(a_0) \rangle \wedge \langle a_1, t(a_1) \rangle \wedge \cdots \wedge \langle a_T, t(a_T) \rangle.$$

359 If  $t$  is a satisfiable hypothesis, then no  $t(a)$  is empty, and  $t(a) = V_a$  tells us that  
 360 any value is allowed in this column. Thus, such a descriptor places no  
 361 restriction on satisfiability in column  $a$ . If  $I = \{i \leq T : t(a_i) \neq V_{a_i}\}$ , then,

$$s \in \text{sat}(t) \iff (\forall i \in I)s(a_i) = t(a_i),$$

363 so that we can interpret  $t$  as the conjunction

$$\bigwedge_{i \in I} \langle a_i, t(a_i) \rangle.$$

Table 2  
Training data, hypotheses and boundaries [9]

Sky	ATemp	Humid	Wind	Water	FCast	<i>d</i>
<i>Training data D</i>						
Sunny	Warm	Normal	Strong	Warm	Same	1
Sunny	Warm	High	Strong	Warm	Same	1
Rainy	Cold	High	Strong	Warm	Change	0
Sunny	Warm	High	Strong	Cool	Change	1
<i>Hypotheses</i>						
Sunny	Warm	?	Strong	?	?	1
Sunny	?	?	Strong	?	?	1
Sunny	Warm	?	?	?	?	1
?	Warm	?	Strong	?	?	1
Sunny	?	?	?	?	?	1
?	Warm	?	?	?	?	1
<i>Specific boundary</i>						
Sunny	Warm	?	Strong	?	?	1
<i>General boundary</i>						
Sunny	?	?	?	?	?	1
?	Warm	?	?	?	?	1

365 Such an expression is called an *exact template* in [13]. The expressive power of  
 366 this type of conjunctive hypothesis is limited. An example from [9] illustrating  
 367 this is shown in Table 3. For such a simple dataset, there is no consistent  
 368 hypothesis in the sense of (4.2).

369 A hypothesis in  $\mathbf{VSp}^m$  is a very special kind of hypertuple, and it is our aim  
 370 to extend this notion of hypothesis, so that the resulting structures are more  
 371 expressive while at the same time not so general as to carry no useful infor-  
 372 mation. As suggested by Mitchell, a possible solution is to use arbitrary dis-  
 373 junctions of conjunctions of descriptors as hypothesis representation. It is not  
 374 hard to see that this is overly general, since any positive Boolean expression can  
 375 then serve as a hypothesis. In contrast, we suggest to use a specific class of  
 376 disjunctions of conjunctions which is significantly narrower than the class of all  
 377 positive Boolean expressions. The building blocks will be the hypertuples  
 378 contained in the sub-semilattice of  $\mathcal{L} = \prod_{a \in \Omega} 2^{V_a}$  generated by  $\mathbf{V}$  with + oper-

Table 3  
Limitation of version space

	Sky	ATemp	Humid	Wind	Water	FCast	<i>d</i>
1	Sunny	Warm	Normal	Strong	Cool	Change	1
2	Cloudy	Warm	Normal	Strong	Cool	Change	1
3	Rainy	Warm	Normal	Strong	Cool	Change	0

379 ator – [V]. Since we are only interested in the finite sums of elements of **V** we  
 380 will from now on assume that each  $V_a$  is finite.

381 Suppose that  $t = \langle \{m_0^a, m_1^a, \dots, m_{t(a)}^a\}_{a \in \Omega}$  is a hypertuple. We interpret  $t$  as a  
 382 conjunction of disjunctions of descriptors

$$\bigwedge_{a \in \Omega} (\langle a, m_0^a \rangle \vee \dots \vee \langle a, m_{t(a)}^a \rangle). \tag{5.1}$$

384 By the distributivity of  $\wedge$  and  $\vee$  this can always be turned into a disjunction of  
 385 simple tuples, but not every disjunction of simple tuples (considered as a  
 386 conjunction of descriptors) is equivalent to an expression such as (5.1); con-  
 387 sider, for example,

$$(\langle a_0, t_0^0 \rangle \wedge \langle a_1, t_1^0 \rangle) \vee (\langle a_0, t_0^1 \rangle \wedge \langle a_1, t_1^1 \rangle).$$

389 A hypertuple can be viewed as a construction similar to hypercubes which  
 390 delineate solids in an appropriate space.

391 Now we compare our hypothesis space with Mitchell’s from the perspective  
 392 of expressive power. Consider a dataset **D** as defined earlier. Suppose all  
 393 attributes  $x$  are discrete, and all  $V_x$  are finite. In our hypothesis space each  
 394 attribute  $x$  takes on a subset of  $V_x$ , so there are  $2^{|V_x|}$  different subsets altogether.  
 395 As a result there are  $\prod_{x \in \Omega} 2^{|V_x|}$  different hypertuples. Since a hypothesis is a set  
 396 of hypertuples (i.e., a hyperrelation), there are  $2^{\prod_{x \in \Omega} 2^{|V_x|}}$  distinct hypotheses.

397 In Mitchell’s hypothesis space each attribute  $x$  takes on a single value in  $V_x$   
 398 plus two other special values, “?” and “ $\emptyset$ ”. Therefore there are  $|V_x| + 2$  different  
 399 values, and  $\prod_{x \in \Omega} (|V_x| + 2)$  different tuples. In his conjunctive hypothesis rep-  
 400 resentation, each hypothesis is a (simple) tuple, so there are  $\prod_{x \in \Omega} (|V_x| + 2)$   
 401 distinct hypotheses. In his disjunctive hypothesis representation, each  
 402 hypothesis is a set of (simple) tuples (i.e., simple relation), so there are  
 403  $2^{\prod_{x \in \Omega} (|V_x| + 2)}$  distinct hypotheses.

404 Clearly our hypothesis space can represent more distinct objects than  
 405 Mitchell’s can. In this sense we say our hypothesis is more expressive than  
 406 Mitchell’s.

407 Note that  $\mathcal{E}$  characterises the eligible hypotheses. So Mitchell’s version space  
 408 can be specialised from our version space in the following way:

- 409 • The  $\mathcal{E}$  is restricted to  $\mathcal{E}^m = \{\gamma(t) : t \in \mathcal{E}, \mathbf{D} \preceq t\}$ , where  $\gamma$  is an operation  
 410 to turn a hypertuple into a simple tuple in such a way that  
 411  $\gamma(t) = \langle t_0, t_1, \dots, t_T \rangle$ , where

$$t_i = \begin{cases} t(x_i), & \text{if } |t(x_i)| = 1, \\ ?, & \text{otherwise.} \end{cases}$$

413 Note that  $t(x_i)$  is the projection of tuple  $t$  onto its  $x_i$  attribute.

- 414 • Each hypothesis is an element of  $\mathcal{E}^m$ .
- 415 • There are two classes, i.e.,  $K = 2$ .
- 416 • The version space is built for only one class.

417 Given the above restrictions the specific boundary is  $\mathbb{S}^m = \{\mathbb{S}_0^m, \mathbb{S}_1^m\}$ , where  
418  $\mathbb{S}_0^m$  is the specific boundary for the negative class and  $\mathbb{S}_1^m$  is the specific  
419 boundary for the positive class. Similarly  $\mathbb{G}^m = \{\mathbb{G}_0^m, \mathbb{G}_1^m\}$ .

## 420 6. Experimental

421 We have evaluated Rule 3 using public datasets. We used 17 public datasets  
422 in our evaluation, which are available from UC Irvine Machine Learning  
423 Repository. General information about these datasets is shown in the first three  
424 columns of Table 4.

425 We used the CASEEXTRACT algorithm to construct weak hypotheses for the  
426 datasets, and applied Rule 3 to classify new data. For presentation purpose we  
427 refer to our classification procedure by GLM. The experimental results are  
428 shown in the last 5 columns of Table 4. As a comparison the C5.0 results on the  
429 same datasets are also shown. It is clear from this table that GLM performs  
430 extremely well on primary data, but it accounts for only 76.4% of all data on  
431 average. Over all data GLM compares well with C5.0.

432 Parity problems are well known to be difficult for many machine learning  
433 algorithms. We evaluated the GLM algorithm using three well known parity  
434 datasets [16]—Monk-1, Monk-2, Monk-3.<sup>1</sup> Experimental results show that  
435 GLM works well for these parity datasets.

## 436 7. Discussion and conclusion

437 Mitchell's classical work on version space has been followed by many. Most  
438 notably Hirsh and Sebag. Hirsh [5] discusses how to merge version spaces when  
439 a central idea in Mitchell's work is removed—a version space is the set of  
440 concepts *strictly consistent* with training data. This merging process can  
441 therefore accommodate noise. Sebag [14] presents what she calls a disjunctive  
442 version space approach to learning disjunctive concepts from noisy data. A

---

<sup>1</sup> Target Concepts associated to the Monk's problem: Monk-1: ( $a_1 = a_2$ ) or ( $a_5 = 1$ ); Monk-2: exactly two of  $a_1 = 1, a_2 = 1, a_3 = 1, a_4 = 1, a_5 = 1, a_6 = 1$ ; Monk-3: ( $a_5 = 3$  and  $a_4 = 1$ ) or ( $a_5 \neq 4$  and  $a_2 \neq 3$ ) (5% class noise added to the training set).

Table 4  
General information about the datasets and the classification accuracy of C5.0 on all data and of GLM on primary and secondary data

Dataset	#Attr.	#Train	#Test	Class. Accuracy (%)				%PP
				C5.0	GLM/SR	GLM/PSR	GLM/SSR	
Annealing	38	798	CV-5	96.6	96.4	98.0	62.9	92.5
Australian	14	690	CV-5	90.6	95.1	95.1	100.0	87.2
Auto	25	205	CV-5	70.7	82.4	87.0	63.0	56.1
Diabetes	8	768	CV-5	72.7	70.7	71.0	40.0	66.8
German	20	1000	CV-5	71.7	72.6	72.6	N/A	65.4
Glass	9	214	CV-5	80.4	86.6	87.6	76.5	79.4
Heart	13	270	CV-5	77.0	81.9	81.9	N/A	61.5
Hepatitis	19	155	CV-5	80.6	82.9	84.1	50.0	69.0
Horse-Colic	22	368	CV-5	85.1	82.7	82.7	N/A	67.7
Iris	4	150	CV-5	94.7	93.1	97.6	66.7	82.7
Monk-1	6	124	432	74.3	100.0	100.0	N/A	100.0
Monk-2	6	169	432	65.1	81.4	87.3	41.5	83.6
Monk-3	6	122	432	97.2	93.8	94.3	89.2	89.1
Sonar	60	208	CV-5	71.6	81.6	81.3	100.0	59.1
TTT	9	958	CV-5	86.2	96.2	96.1	100.0	94.9
Vote	18	232	CV-5	96.5	96.5	98.6	77.3	89.7
Wine	13	178	CV-5	94.3	99.0	99.0	N/A	53.9
Average				82.7	87.8	89.1	72.3	76.4

The validation method used is either 5-fold cross validation or explicit train/test validation. The acronyms are: SR—overall success ratio of primary and secondary data (i.e., the percentage of successfully classified primary and secondary data tuples over all primary and secondary data tuples), PSR—success ratio of primary data, SSR—success ratio of secondary data, PP—the percentage of primary data, and N/A—not available.



443 separate version space is learned for each positive training example, then new  
444 instances are classified by combining the votes of these different version spaces.

445 In this paper we investigate version spaces in a more expressive hypothesis  
446 space—disjunction of conjunctions of disjunctions, where each hypothesis is a  
447 set of hypertuples. Without a proper inductive bias the version space is un-  
448 interesting. We show that, with  $E$ -set as an inductive bias, this version space is a  
449 generalisation of Mitchell's original version space, which employs a different  
450 type of inductive bias.

451 For classification within the version space we proposed three classification  
452 rules for use in different situations. The first two rules are based on  $E$ -sets as  
453 hypotheses, and they can be applied when the data space ( $\mathbf{V}$ ) is finite and small.  
454 We showed that constructing  $E$ -sets is computationally expensive, so we  
455 introduced the third rule which is based on weak hypotheses. We presented an  
456 algorithm to construct weak hypotheses efficiently.

457 Experimental results show that this classification approach performs ex-  
458 tremely well on primary data, which account for over 75.0% of all data. Over  
459 all data this classification approach is comparable to C5.0.

## 460 8. Uncited reference

461 [7]

## 462 Acknowledgements

463 The work of Hui Wang was partly supported by the European Commission  
464 project ICONS, project no. IST-2001-32429.

## 465 Appendix A. Notation

$$466 \quad X \ll Y \iff (\forall x \in X)(\exists y \in Y)x \leq y$$

$$467 \quad X \triangleleft Y \iff (\forall y \in Y)(\exists x \in X)x \leq y$$

$$468 \quad U = \{x_0, \dots, x_N\}$$

$$469 \quad \Omega = \{a_0, \dots, a_T\}$$

$$470 \quad \mathbf{V} = \prod_{a \in \Omega} V_a$$

$$471 \quad \mathcal{L} = \prod_{a \in \Omega} 2^{V_a}$$

$$472 \quad I(x) = \langle a_0(x), a_1(x), \dots, a_T(x) \rangle$$

$$473 \quad \mathbf{D} = \{I(x) : x \in U\} = \mathbf{D}_0 \cup \mathbf{D}_1 \cup \dots \cup \mathbf{D}_K$$

$$474 \quad \mathcal{E}_q = \text{set of all elements equilabelled with respect to } \mathbf{D}_q$$

$$475 \quad \mathcal{E} = \bigcup_{q \leq K} \mathcal{E}_q$$

$$\begin{aligned}
476 \quad & E(P) = \{t \in \mathcal{L} : t \text{ is maximal in } [P] \cap \mathcal{E}\} \\
477 \quad & E_q(P) = E(P) \cap \mathcal{E}_q \\
478 \quad & \mathbb{S} = E(\mathbf{D}) \\
479 \quad & \mathbb{S}_q = E_q(\mathbf{D}) \\
480 \quad & \mathbb{G} = E(\mathbf{V}) \\
481 \quad & \mathbb{G}_q = E_q(\mathbf{V}) \\
482 \quad & \mathbf{VSp} = \bigcup \{E(P) : \mathbf{D} \subseteq P \subseteq \mathbf{V}\}
\end{aligned}$$

## 483 References

- 484 [1] I. Düntsch, G. Gediga, Statistical evaluation of rough set dependency analysis, *International*  
485 *Journal of Human–Computer Studies* 46 (1997) 589–604.
- 486 [2] I. Düntsch, G. Gediga, Simple data filtering in rough set systems, *International Journal of*  
487 *Approximate Reasoning* 18 (1–2) (1998) 93–106.
- 488 [3] G. Grätzer, *General Lattice Theory*, Birkhäuser, Basel, 1978.
- 489 [4] D. Haussler, Quantifying inductive bias: Ai learning algorithms and valiant’s learning  
490 framework, *Artificial Intelligence* 36 (1988) 177–221.
- 491 [5] H. Hirsh, Generalizing version spaces, *Machine Learning* 17 (1) (1994) 5–46.
- 492 [6] T.M. Mitchell, Version spaces: a candidate elimination approach to rule learning, in: *Proc. 5th*  
493 *IJCAI*, 1977, pp. 305–310.
- 494 [7] T.M. Mitchell, Version spaces: an approach to concept learning. PhD thesis, Department of  
495 *Electrical Engineering*, Stanford University, Stanford, CA, 1979.
- 496 [8] T.M. Mitchell, Generalization as search, *Artificial Intelligence* 18 (1982).
- 497 [9] T.M. Mitchell, *Machine Learning*, The McGraw-Hill Companies, Inc, 1997.
- 498 [10] H.S. Nguyen, From optimal hyperplanes to optimal decision trees, *Fundamenta Informaticae*  
499 34 (1998) 145–174.
- 500 [11] H.S. Nguyen, S.H. Nguyen, Pattern extraction from data, *Fundamenta Informaticae* 34 (1998)  
501 129–144.
- 502 [12] H.S. Nguyen, A. Skowron, Boolean reasoning for feature extraction problems, in: Z.W. Ras,  
503 A. Skowron (Eds.), *Proc. of the 10th International Symposium on Methodologies for*  
504 *Intelligent Systems (ISMIS’97)*, volume 1325 of *Lecture Notes in Artificial Intelligence*,  
505 Springer-Verlag, Berlin, 1997, pp. 117–126.
- 506 [13] S.H. Nguyen, A. Skowron, P. Synak, Discovery of data patterns with applications to  
507 decomposition and classification problems, in: L. Polkowski, A. Skowron (Eds.), *Rough sets in*  
508 *knowledge discovery*, vol. 2, Heidelberg, Physica-Verlag, 1998, pp. 55–97.
- 509 [14] M. Sebag, Delaying the choice of bias: a disjunctive version space approach, in: *Proc. of the*  
510 *13th International Conference on Machine Learning*, San Francisco, Morgan Kaufmann,  
511 1996, pp. 444–452.
- 512 [15] A. Skowron, Extracting laws from decision tables—a rough set approach, *Computational*  
513 *Intelligence* 11 (1995) 371–388.
- 514 [16] S.B. Thrun, J. Bala, E. Bloedorn, I. Bratko, B. Cestnik, J. Cheng, K. De Jong, S. Dzeroski,  
515 S.E. Fahlman, D. Fisher, R. Hamann, K. Kaufman, S. Keller, I. Kononenko, J. Kreuziger,  
516 R.S. Michalski, T. Mitchell, P. Pachowicz, Y. Reich, H. Vafaie, W. Van deWelde, W. Wenzel,  
517 J. Wnek, J. Zhang, The monk’s problems—a performance comparison of different learning  
518 algorithms. Technical Report CS-CMU-91-197, Carnegie Mellon University, 1991.
- 519 [17] H. Wang, W. Dubitzky, I. Düntsch, D. Bell, A lattice machine approach to automated  
520 casebase design: Marrying lazy and eager learning, in: *Proc. IJCAI99*, Stockholm, Sweden,  
521 1999, pp. 254–259.

- 522 [18] H. Wang, I. Düntsch, D. Bell, Data reduction based on hyper relations, in: R. Agrawal, P.  
523 Stolorz, G. Piatetsky-Shapiro (Eds.), Proceedings of KDD'98, New York, 1998, pp. 349–353.
- 524 [19] H. Wang, I. Düntsch, G. Gediga, Classificatory filtering in decision systems, International  
525 Journal of Approximate Reasoning 23 (2000) 111–136.

UNCORRECTED PROOF