



HAL
open science

MCP: Learning Propositional Formulas from Binarized Data

Miki Hermann, Gernot Salzer

► **To cite this version:**

Miki Hermann, Gernot Salzer. MCP: Learning Propositional Formulas from Binarized Data. ADEMAL, Miki Hermann; Nina Narodytska; Corina S. Pasareanu, Jul 2023, Rome, France. 10.4230/OA-SIcs.ADEMAL.2023 . hal-04307792

HAL Id: hal-04307792

<https://hal.science/hal-04307792>

Submitted on 26 Nov 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

MCP: Learning Propositional Formulas from Binarized Data

Miki Hermann  

LIX, CNRS, École Polytechnique, France

Gernot Salzer  

Technische Universität Wien, Austria

Abstract

Experimental data is often given as bit vectors, with vectors corresponding to observations, and coordinates to attributes, with a bit being true if the corresponding attribute was observed. Observations are usually grouped, e.g. into positive and negative samples. Among the essential tasks on such data, we have compression, the construction of classifiers for assigning new data, and information extraction.

Our system, MCP, approaches these tasks by propositional logic. For each group of observations, MCP constructs a (usually small) conjunctive formula that is true for the observations of the group, and false for the others. Depending on the settings, the formula consists of Horn, dual-Horn, bijnunctive or general clauses. To reduce its size, only relevant subsets of the attributes are considered. The formula is a (lossy) representation of the original data and generalizes the observations, as it is usually satisfied by more bit vectors than just the observations. It thus may serve as a classifier for new data. Moreover, (dual-)Horn clauses, when read as if-then rules, make dependencies between attributes explicit. They can be regarded as an explanation for classification decisions.

2012 ACM Subject Classification Computing methodologies → Machine learning

Keywords and phrases machine learning, binary data, propositional logic, rules, Horn formulas, bijnunctive formulas

Digital Object Identifier 10.4230/OASlcs.ADEMAL.2023.

1 Introduction and Related Work

Since several years, computer science applications are challenged by large quantities of data, commonly referred to as *Big Data*, that needs to be interpreted, captured, treated, and transformed. There exist several approaches to cope with this challenge, mainly from the field of Artificial Intelligence. One of these approaches is the *Logical Analysis of Data*. This document presents a tool called MCP, performing logical analysis of big data, producing a propositional formula. The basic idea behind this tool programmed in C++ is to describe a very large data set by a propositional formula.

Logical Analysis of Data is a part of Machine Learning, which has been developed by Hammer and his colleagues [5, 10]. There also exists another approach through mechanized hypothesis formation, the GUHA Project developed in Prague by Hájek and his colleagues [13, 15].

2 Preliminaries

We recall the main structures of Boolean algebra. A *literal* is either a variable, called positive literal, or its negation, called negative literal. A *clause* is a disjunction of literals. A *formula* in *conjunctive normal form* is a conjunction of clauses. A *Horn* clause is a clause with at most one positive literal. A *dual Horn* clause is a clause with at most one negative literal. A *bijnunctive* clause is a clause consisting of at most two literals. An *affine* clause is a linear equation of the form $x_1 + \dots + x_k = b$, where x_i are variables, $+$ is the exclusive-or operator,

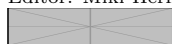


© M. Hermann and G. Salzer;

licensed under Creative Commons License CC-BY 4.0

1st Workshop on Automated Deduction for Machine Learning (ADEMAL 2023).

Editor: Miki Hermann and Nina Narodytska and Corina S. Păsăreanu; Article No. ; pp. :1–:12



Open Access Series in Informatics

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

XX:2 MCP: Learning Propositional Formulas from Binarized Data

44 and $b \in \{0, 1\}$ is a Boolean value. A Horn, dual Horn, bijnunctive, or affine formula is a
45 conjunction of only Horn, dual Horn, bijnunctive, or affine clauses, respectively.

46 We will work with vectors, also called tuples, of finite arity over a domain D . This domain
47 is either Boolean, i.e., $D = \{0, 1\}$, or finite, i.e., $|D| = n$ for some natural number $n \geq 2$.
48 Vectors (a_1, \dots, a_k) of arity k will be shortened to $a_1 \cdots a_k$ when the elements a_i are clear.

49 Let $\mathbf{a} = a_1 \cdots a_k$, $\mathbf{b} = b_1 \cdots b_k$, and $\mathbf{c} = c_1 \cdots c_k$ be Boolean vectors of the same arity k .
50 There exist different closures of these Boolean vectors.

- 51 ■ *Horn closure* of \mathbf{a} and \mathbf{b} is the vector $\mathbf{d} = d_1 \cdots d_k$, such that $d_i = a_i \wedge b_i$;
- 52 ■ *Dual Horn closure* of \mathbf{a} and \mathbf{b} is the vector $\mathbf{d} = d_1 \cdots d_k$, such that $d_i = a_i \vee b_i$;
- 53 ■ *Bijnunctive closure* of \mathbf{a} , \mathbf{b} , and \mathbf{c} is the vector $\mathbf{d} = d_1 \cdots d_k$, such that $c_i = \text{maj}(a_i, b_i, c_i)$,
54 where maj is the associative-commutative majority operator;
- 55 ■ *Affine closure* of \mathbf{a} , \mathbf{b} , and \mathbf{c} is the vector $\mathbf{d} = d_1 \cdots d_k$, such that $d_i = a_i + b_i + c_i$, where $+$
56 is the exclusive-or operator in the Boolean ring \mathbb{Z}_2 ;

57 all for each $i = 1, \dots, k$. Given a set of Boolean vectors S of arity k , we denote by $\langle S \rangle_C$
58 the C -closure of S for C being Horn, dual Horn, bijnunctive, or affine. A basic result from
59 universal algebra states that for an arbitrary set of Boolean vectors S of the same arity k ,
60 the C -closure is the set of satisfying assignments for some C -formula φ [3, 4].

61 3 Running Example: Mushrooms

62 To illustrate the MCP system, we will use the Mushroom Data Set from the UCI Machine
63 Learning Repository as a running example.¹ It contains 8124 records with 22 attributes each.
64 Each record describes an instance of one of 23 species of mushrooms, categorized as *edible* or
65 *poisonous*. Among the attributes we find e.g. *cap-shape*, *odor* or *habitat*. Each attribute may
66 take a value from a finite set. The odor, for instance, may be one of *almond*, *fishy*, *foul* and
67 six further odors. A record looks like

68 `edible,convex,smooth,white,yes,almond,...,purple,several,woods`

69 where the first field specifies the category and the subsequent ones the values of the 22
70 attributes in a fixed order; e.g., the sixth field specifies the odor as *almond*.

71 To process data like the mushrooms with the MCP system, it needs to be *binarized*. MCP
72 offers a utility, `mcp-trans`, for transforming the data. The particular encoding has to be
73 specified in a file that describes the mapping for every attribute and value. To ease the
74 burden on the user, the utility `mcp-guess` takes the original data and generates a draft of
75 this specification. So the typical workflow consists of first running `mcp-guess`, then checking
76 and manually adjusting the generated specification, then running `mcp-trans` on the data,
77 and finally feeding its binary output to the core tools of the MCP system. The next section
78 takes a closer look at the functionality of the MCP core, while the utility programs are
79 discussed in the subsequent section.

80 4 Core of the MCP System

81 MCP consists of several modules. The core modules generate a propositional formula from
82 given sets of binary tuples, according to the following specification.

¹ <https://archive.ics.uci.edu/ml/datasets/mushroom>

83 ► **Problem** (MCP Problem). Given two sets of Boolean vectors (tuples) of arity k over the
 84 Boolean domain $D = \{0, 1\}^k$, representing positive samples $T \subseteq D$ and negative samples
 85 $F \subseteq D$, **compute** a Horn, dual Horn, bijunctive, or general CNF formula φ , respectively,
 86 such that (1) $T \models \varphi$ and (2) for each $f \in F$, $f \not\models \varphi$.

87 There are several reasons why we focus on the aforementioned four subclasses of
 88 propositional formulas. Horn, dual Horn, bijunctive, and affine formulas are the four
 89 families of Boolean formulas, whose satisfiability problem can be decided in polynomial
 90 time. Moreover, Horn formulas are the foundation of logic-oriented programming, with Horn
 91 clauses being naturally interpreted as rules.

92 ► **Example.** Suppose we aim for a formula that characterizes the *edible* mushrooms in our
 93 running example. Then the records categorized as *edible* become the basis for the positive
 94 samples, T , and the records labeled as *poisonous* result in the negative samples, F . Since the
 95 attributes of the example are mostly non-binary, the dimension of the tuples will not equal
 96 the number of attributes 22, but will be larger and will depend on the chosen transformation.
 97 The one described later results in 111 Boolean attributes.

98 4.1 Strategies for Computing the Closure

99 An instance of the MCP problem is not solvable if some tuple occurs in the set of positive
 100 and negative samples at the same time, hence we require $T \cap F = \emptyset$. For a C -formula (with
 101 C being Horn, dual Horn, bijunctive or affine), the condition has to be strengthened to
 102 $\langle T \rangle_C \cap F = \emptyset$, as each tuple in the C -closure of T necessarily satisfies any C -formula for T .

103 With this constraint, instances of the MCP problem are solvable. In general, there is a
 104 range of solutions. While every formula solving the instance is satisfied by the tuples in T
 105 and falsified by those in F , its behavior for the remaining tuples is undetermined. MCP
 106 offers two **strategies**. The **large** strategy (the default) computes a formula satisfied by
 107 a maximal number of tuples, while the **exact** strategy minimizes the number of satisfying
 108 tuples.

109 4.2 Minimal Section

110 For a set of tuples, S , let $S|_A$ denote the restriction of the tuples to the coordinates in A .
 111 Clearly, if two sets S_1 and S_2 have an empty intersection for a subset of the coordinates,
 112 then the original sets have an empty intersection, too: $S_1|_A \cap S_2|_A = \emptyset$ implies $S_1 \cap S_2 = \emptyset$.
 113 In general, each coordinate contributes a variable to the formula, hence minimizing A will
 114 reduce the number of different variables in the formula.

115 Finding an A of minimal cardinality such that $T|_A \cap F|_A = \emptyset$ (or $\langle T \rangle_C|_A \cap F|_A = \emptyset$) is an
 116 NP-complete problem. MCP implements several approximations differing in the **direction**
 117 the coordinates are tried, skipping coordinates whose removal would render the problem
 118 unsolvable. The following directions are available:

119 **begin:** Prefer coordinates to the left (at the begin) of the tuples by removing coordinates
 120 from the right. This direction is the default.

121 **end:** Prefer coordinates to the right (at the end) of the tuples by removing coordinates from
 122 the left.

123 **lowcard:** Prefer coordinates with a lower Hamming weight, by removing coordinates with
 124 high Hamming weight.

125 **highcard:** Prefer coordinates with a higher Hamming weight, by removing coordinates with
 126 small Hamming weight.

127 **random**: Remove coordinates in random order.

128 **nosect**: Use all coordinates, do not remove any.

129 4.3 Effective Learning of Formulas

130 The MCP system learns *Horn* formulas by the following procedure. For each $f \in F$ it
 131 determines if $f \in \langle T \rangle_{\text{Horn}}$ efficiently, without computing the Horn closure. Then it computes
 132 the minimal section of $\langle T \rangle_{\text{Horn}}$ and F , followed by the computation of the corresponding
 133 Horn formula according to the chosen direction and strategy on the (approximate) minimal
 134 section of $\langle T \rangle_{\text{Horn}}$ and F . It uses different algorithms for the strategies: that of Angluin *et*
 135 *al* [1] for the large strategy and another of Hébrard and Zanuttini [14] for the exact strategy.

136 Learning of *dual Horn* formulas is done very easily. MCP system first swaps the polarity
 137 of the Boolean vectors in T and F , producing the new sets T' and F' , respectively. Then it
 138 computes the Horn formula φ' for T' and F' , followed by swapping the polarity of literals
 139 in φ' , producing the dual Horn formula φ .

140 There is no known possibility to determine if $f \in \langle T \rangle_{\text{bijunctive}}$ for each $f \in F$ without
 141 computing the bijunctive closure $\langle T \rangle_{\text{bijunctive}}$. Moreover, the bijunctive closure $\langle T \rangle_{\text{bijunctive}}$
 142 can be (and usually also is) very much time and space consuming. We adopted the following
 143 solution to produce *bijunctive* formulas by MCP system: It computes the minimal section
 144 using an intersection test, followed by application of the *Baker-Pixley Theorem* [2] (projection
 145 on every pair of coordinates), which implicitly guarantees the bijunctive closure.

146 Learning a *general CNF* formula presents several challenges. Its advantage is that We
 147 get a propositional formula in any case, provided that $T \cap F = \emptyset$. Its drawback is that the
 148 produced formula is usually very big. We adopted two different approaches in the MCP
 149 system, depending on the applied strategy. In case of large strategy, for each false element
 150 $f \in F$ the MCP system produces the unique clause c_f which falsifies f . The resulting
 151 formula φ is the conjunction of all falsification clauses c_f . In case of exact strategy, the MCP
 152 system uses an algorithm producing a CNF formula in time $O(|T|k^2)$, where k is the arity
 153 of vectors in T , using a Boolean restriction of a larger algorithm from [12].

154 Learning *affine* formulas reveals more from linear computer algebra than from logic,
 155 therefore we did not implement it in the MCP system for the time being. We may implement
 156 it in a further version if there is demand.

157 4.4 First Postprocessing: Redundancy Elimination

158 The inferred formula φ can contain redundant literals and clauses, which can and must be
 159 eliminated to produce the smallest possible formula. There are several stages, which can
 160 be applied for *redundancy elimination*, called **cooking** inside the MCP system, with the
 161 following options: **raw** performs no redundancy elimination, **bleu** performs unit resolution,
 162 **medium** performs unit resolution and clause subsumption, and finally **well done**, which
 163 is the default, performs unit resolution, clause subsumption, and implied clause removal.
 164 Moreover, the *exact* strategy includes a **primality** step, reducing the clauses by elimination
 165 of unnecessary literals, using an algorithm from [12].

166 4.5 Second Postprocessing: Set Cover

167 In case of the *large* strategy, we are mainly interested in producing a formula φ falsified by
 168 each tuple $f \in F$. However, the inferred formula φ may contain more clauses than necessary,
 169 even after full redundancy elimination. Our task is to keep the smallest number of clauses

170 in φ which are necessary to guarantee falsification by all tuples $f \in F$. For this purpose in
 171 the MCP system, we use *Set Cover* where a clause $c \in \varphi$ covers a vector $f \in F$ if f falsifies c .
 172 Set Cover is a well-known NP-complete problem, therefore we use Johnson's approximation
 173 algorithm (see e.g. [11]), where the measure of a clause is the number of covered tuples. Of
 174 course, this approach is inapplicable for the *exact* strategy.

175 4.6 Input Format and Action Possibilities

176 The input file of the MCP system core, is a Boolean matrix, one Boolean vector per row.
 177 Each vector is prefixed by a string g , identifying a group to which the vector belongs. The
 178 MCP system core collects first the vectors from the input matrix and distributes them into
 179 the identified groups. Each input file starts with an indication line, containing two boolean
 180 values. If both values are equal to 0, the following lines are the rows of the Boolean matrix
 181 with leading group indicators. If the first value is equal to 1, the following line contains
 182 the variable names ordered by coordinates. If the second value is equal to 1, there is one
 183 more line of supplementary information before the matrix. However, this supplementary
 184 information is unused by the MCP system, but it is still maintained for compatibility reasons
 185 with data sets used in [8, 9].

186 Let G be the set of identified groups. The actual computation is determined by the
 187 **action**, which determines how the sets of positive samples T and negative samples F are
 188 constituted. There are two options, **one** and **all**.

189 The option *one* consecutively selects two groups $g, g' \in G$, determines the vectors
 190 belonging to the group g as the positive samples T and the vectors belonging to the group g'
 191 as the negative samples F , then starts the computation of the corresponding formula with
 192 minimal section. If there are n groups in the set G , this action proceeds with the computation
 193 of $n(n - 1)$ formulas.

194 The option *all*, which is the default, consecutively selects a group $g \in G$, determines
 195 the vectors belonging to the group g as the positive samples T and all vectors belonging
 196 to any group from $G \setminus \{g\}$ as the negative samples F , then starts the computation of the
 197 corresponding formula with minimal section. For n groups in the set G , this action proceeds
 198 with the computation of n formulas.

199 4.7 Parallelization

200 For a set of n groups, the MCP system computes either n or $n(n - 1)$ formulas. These
 201 computations are independent, therefore they can be performed in parallel. This is called
 202 *outer parallelism* in the MCP core.

203 In case of Horn closure of the positive samples T , the MCP core needs to determine
 204 if a given vector $f \in F$ from negative samples belongs to $\langle T \rangle_{\text{Horn}}$, without computing the
 205 closure itself. This procedure is quite time consuming when the set T is quite large. It can
 206 be computed in parallel, each time taking only a determined chunk of T . This is called *inner*
 207 *parallelism* in the MCP core.

208 We adopted three types of parallelization within the MCP core: the **Message Passing**
 209 **Interface** (MPI) [16], the **POSIX threads** (pthreads) [7], and a **hybrid** version combining
 210 both. These parallelizations are effective only on very large input data sets. The MPI version
 211 is applied only for outer parallelism, the pthreads version to both, and in the hybrid version
 212 MPI is applied for outer parallelism and pthreads for inner parallelism.

213 **4.8 Invocation**

214 MCP core is called by one of the following commands and options:

sequential version:	mcp-seq	} -i <i>input-file</i> -o <i>output-file</i>
MPI version:	mcp-mpi	
POSIX threads version:	mcp-pthread	
hybrid version:	mcp-hybrid	
		-l <i>formula-prefix</i> -c <i>closure</i>
		-d <i>direction</i> -s <i>strategy</i>
		--cook <i>cooking</i> --setcover <i>y/n</i>

216 Each of these core modules produces files *formula-prefix_g.log* containing the learned
 217 formula for each group *g* inside *input-file*. Consult the manual pages for more detailed
 218 information.

219 ► **Example.** Suppose we have transformed the data of our running example to 15 858 binary
 220 tuples of length 111, stored in a file `mushroom.bin`, as described at the end of section 5.1.
 221 Running the command

```
222 mcp-seq -i mushroom.bin -o mushroom.frm
```

223 we obtain, within a minute or so, Horn formulas for the edible as well as for the poisonous
 224 mushrooms. E.g., the edible mushrooms are characterized by 17 rules like the following ones:

```
225 cap-color ≠ yellow ∨ odor ≠ foul
226 cap-color ≠ gray ∨ odor ≠ foul
227 spore-print-color = white → odor = none
```

228 The first two formulas specify that for an edible mushroom, its *odor* is either not *foul*, or its
 229 *cap-color* is different from *yellow* and *gray*. Moreover, if the *spore-print-color* is *white*, then
 230 an edible mushroom has no *odor* at all.

231 **5 Prequel and Sequel Modules**232 **5.1 Data Binarization**

233 The core of the MCP system accepts only Boolean vectors. However, data are usually
 234 spanning much larger domains: finite, or infinite but countable, or uncountable. In the latter
 235 two cases, every very large finite data set contains only a finite subset of the domain, but it
 236 can be intractable due to the amount of data to be treated. The MCP system copes with
 237 this situation by *binarization*.

238 **Binarization** is the process of transforming data of any domain into binary vectors to
 239 make classifier algorithms, in our case the MCP system core, more efficient. Its advantage is
 240 that we obtain the possibility to treat any data by propositional formulas. Its drawback is a
 241 possible exponential explosion. Binarization concerns both, particular values, especially for
 242 finite domains, as well as intervals, usually used for infinite ones. MCP system adopts both
 243 approaches.

244 Binarization in the MCP system is a two-step procedure. The first step consists of
 245 scanning of the CSV file and generating a meta-file template. This step is performed by the
 246 command

```
247 mcp-guess -i csv-file -o meta-template
```

248 where it is implicitly assumed that the *csv-file* contains one data vector per line, the vector
 249 elements are separated by commas or semicolons or space or tabs, vector element can be
 250 quoted, missing elements are denoted by a question mark. The template generated by
 251 *mcp-guess* cannot be used directly by the next module, but it must be manually adapted to
 252 a proper meta-file. This command just creates indications if the values of a given coordinate
 253 are Boolean, enumerated strings, enumerated integers, integers in a range, or floats in a
 254 range.

255 The second step of the binarization process is performed by the command

256 **mcp-trans** -i *data-file* -m *meta-file* -o *binarized-file* [--pvt *pivot-file*] [-r y/n]

257 which generates a *binarized-file*, ready to be treated by the MCP system core, from the
 258 original *data-file* using a *meta-file*. This meta file consists of transformation commands. Each
 259 transformation command describes the treatment of one attribute and has the following
 260 format:

261 *identifier* = *coordinate* : *indicator* ; {# *comment*}

262 where # starts an optional comment stretching until end of line, the symbols = and : and ;
 263 are syntactic sugar, *identifier* will become the name of the variable for the given attribute
 264 *coordinate* and the *indicator* has one of the following forms:

concept		identifier of the concept to be learned
pivot		identifiers for pivot values in prediction
bool	[' <i>elem₀</i> <i>elem₁</i> ']	boolean 2-element set
enum	[' <i>elem₀</i> ... <i>elem_ℓ</i> ']	enumerated set of $\ell + 1$ elements
up	[' <i>elem₀</i> ... <i>elem_ℓ</i> ']	enumerated set of increasing $\ell + 1$ elements
down	[' <i>elem₀</i> ... <i>elem_ℓ</i> ']	enumerated set of decreasing $\ell + 1$ elements
int	<i>min max</i>	integers in the range between <i>min</i> and <i>max</i>
dj	<i>n min max</i>	interval [<i>min</i> , <i>max</i>) cut in <i>n</i> disjoint chunks
cp	[^] <i>number₀</i> ... <i>number_ℓ</i> [\$]	intervals determined by checkpoints (The numbers must build an increasing sequence. Consecutive numbers <i>a</i> and <i>b</i> determine the interval [<i>a</i> , <i>b</i>). The optional symbols ^ and \$ stand for minimal and maximal infinity, respectively)
over	<i>n min max ℓ</i>	[<i>min</i> , <i>max</i>) cut in <i>n</i> chunks with overlaps of length ℓ
span	ℓ <i>min max</i>	[<i>min</i> , <i>max</i>) cut in disjoint chunks, each of length ℓ
warp	ℓ_0 <i>min max</i> ℓ_1	[<i>min</i> , <i>max</i>) cut in chunks of length ℓ_0 , overlaps of length ℓ_1

266 The square brackets '[' and ']', written in quotation marks to distinguish them from optional
 267 parameter indications, are just syntactic sugar for a better orientation of the parser.

268 Some data vectors can contain missing values indicated by a quation mark. The default
 269 treatment by **mcp-trans** is their elimination. If however we wish to include these data rows
 270 with missing values, we need to generate for them *robust* extensions [6], where the question
 271 marks are replaced by all other data for this coordinate gathered from all other data rows.
 272 This is achieved by the optional -r flag. The *robust* option is incompatible with the *pivot*
 273 option.

274 The optional *pivot-file* will contain data identifiers used for prediction (see the command
 275 **mcp-predict**). The **concept** and **pivot** indicators are exclusive. The **concept** indicator is used
 276 during the learning process on training data, whereas the **pivot** indicator is used during the
 277 checking process on testing data, when a prediction for these data has to be done. Usually,

XX:8 MCP: Learning Propositional Formulas from Binarized Data

278 the pivot coordinate contains the identifiers, one per data item, for which a prediction is
279 done by means of the `mcp-predict` command.

280 After the transformation and binarization, the produced binary vectors do not need to
281 be unique. To avoid this situation, we can use the command

```
282 mcp-uniq -i input-file -o output-file
```

283 which eliminates row doublets. Essentially, this command acts as the Linux command `uniq`,
284 but the rows do not need to be sorted and duplicate rows need not be consecutive.

285 ► **Example.** Our mushroom example uses non-binary attributes. Therefore, we run `mcp-guess`
286 on the original data to draft a specification for the binarization process. The command

```
287 mcp-guess -i mushroom.data -o mushroom.spec
```

288 analyzes the values occurring for every attribute and generates a file starting as follows.

```
289 id0 = 0: bool [edible poisonous];  
290 id1 = 1: enum string [bell conical convex flat knobbed sunken];  
291 id2 = 2: enum string [fibrous grooves scaly smooth];  
292 id3 = 3: enum string [brown buff cinnamon gray green pink ...  
293 id4 = 4: bool [no yes];  
294 id5 = 5: enum string [almond anise creosote fishy foul musty ...
```

295 We adapt this file by marking `id0` as the column that defines the category (the ‘concept’),
296 and by replacing the generic identifiers `idx` by mnemonic labels, to improve readability. If we
297 can also store the names of the attributes in the file `mushroom.names` and run the command

```
298 mcp-guess -i mushroom.data -o mushroom.spec -n mushroom.names
```

299 with the `-n` option. The identifiers of the attributes are then read from that file and replace
300 the identifiers `id0`, ..., `id5`. Moreover, we check the encodings proposed by `mcp-guess`.
301 Attributes taking more than two different, unordered values are tagged as `enum`, which tells
302 the binarization utility to use a separate propositional variable for every value. Attributes
303 with just two values are marked as `bool`, which results in a single variable for both values.
304 While saving on variables reduces the size of the problem, it may make the information
305 harder to access and prevent MCP from constructing, for instance, a Horn formula. In our
306 final specification, the lines above take the following form.

```
307 class = 0: concept;  
308 cap-shape = 1: enum [bell conical convex flat knobbed sunken];  
309 cap-surface = 2: enum [fibrous grooves scaly smooth];  
310 cap-color = 3: enum [brown buff cinnamon gray green pink purple ...  
311 bruises = 4: bool [no yes];  
312 odor = 5: enum [almond anise creosote fishy foul musty ...
```

313 Now we binarize the original data with the command

```
314 mcp-trans -i mushroom.data -m mushroom.spec -o mushroom.bin
```

315 The binarized data, `mushroom.bin`, starts with the following lines.

```
316 1 0  
317 cap-shape_5:cap-shape==sunken:cap-shape!=sunken ...  
318 edible 1 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 1 ...  
319 poisonous 0 0 0 1 0 0 1 0 0 0 0 1 0 0 0 0 0 0 ...
```

320 The first line with 1 0 indicates that line 2 is a header and that the actual data starts in line 3.
 321 The quite verbose header consists of space-separated labels, one for each binary attribute.
 322 Each label starts with the name of the attribute (`cap-shape_5` in the first label above),
 323 followed by two expressions separated by colons. The first expression, `cap-shape==sunken`,
 324 specifies that a value 1 for the attribute `cap-shape_5` means that the original attribute
 325 `cap-shape` had the value *sunken*, whereas the second expression, `cap-shape!=sunken`
 326 reminds us that `cap-shape_5` being 0 means that `cap-shape` had a value different from
 327 *sunken*. This information is used by MCP to enhance the readability of the generated
 328 formulas. From line 3 to the end of the file, each line contains a tuple of 111 binary values,
 329 which is prefixed by one of the strings `edible` or `poisonous` giving the category of the tuple.
 330 The total number of lines is actually almost twice as large as in the original data file. The
 331 increase is caused by unknown values (appearing as a question mark, `?`) in the data file.
 332 Records with such values can be either dropped, decreasing the number of tuples, or can be
 333 replaced by all possible values of the respective attribute. The latter option may lead to a
 334 significant increase of the data set.

335 5.2 Formula Evaluation and Classification Prediction

336 If we are interested only in the produced formula, then the output file generated by the
 337 MCP core contains the satisfied formulas for each group of Boolean vectors. However, if we
 338 want to evaluate the accuracy of the produced formula, we must proceed further. The first
 339 prerequisite for a possibility to check the accuracy of a formula, is to have two sets of vectors:
 340 one for learning the formula, the other for checking its accuracy. Either we have these two
 341 sets of vectors already from the beginning or we need to split the original set of Boolean
 342 vectors into the learning part and the checking part before running the MCP core on the
 343 learning part. The latter is performed by the command

```
344 mcp-split -i input-file -l learn-file -c check-file -r ratio
```

345 that splits uniformly at random the *input-file* into a *learn-file* and *check-file*, where *ratio* is
 346 the percentage of vectors from the *input-file* populating the *check-file*. If the options `-l` or
 347 `-c` are not explicitly stated, the software deduces the file identifiers from the base name of
 348 the *input-file* and adding the suffix `.lrn` or `.chk` to it, respectively. The *ratio* default is 10.
 349

The accuracy of the formula for a given group *g* is checked by the command

```
350 mcp-check -i check-file -l formula-file -o output-file
```

351 where *formula-file* is the file *formula-prefix*_g.log produced by the MCP core. Its *output-file*
 352 reproduces the formula and reports the following statistical entities, measured on the vectors
 353 from *check-file*: true positives (*tp*), true negatives (*tn*), false positives (*fp*), false negatives
 354 (*fn*), sensitivity ($tp/(tp+fn)$), miss rate ($fn/(fn+tp)$), specificity ($tn/(tn+fp)$), and precision
 355 ($tp/(tp+fp)$). The optimal situation would be to have neither false positives nor false
 356 negatives. If, however, these values are non-zero, it can be either due to an insufficient
 357 cardinality of learning data, or a wrong binarization, or else the data itself are not precise.

358 Once a learning process has been done and corresponding formulas have been generated,
 359 a prediction for testing data is performed by the following command:

```
360 mcp-predict -i input -o output -l formula-prefix -pdx prediction [-pdt pivot]
```

361 where *input* is the testing file containing Boolean vectors without group identifiers for which
 362 the prediction will be done; *output* is the file which will contain the report of the prediction

XX:10 MCP: Learning Propositional Formulas from Binarized Data

363 run; *formula-prefix* is the prefix for files containing formulas produced by MCP core and
364 where the corresponding *formula-file* is supposed to have the name *formula-prefix_g.log* for
365 some group G; the file *prediction* will contain the prediction results in the form

366 *pivot-value, group identifiers*

367 where *group-identifiers* is a list of groups (only one in the best case) separated by the sign +,
368 for which the corresponding formula from *formula-prefix_g* is satisfied by that *pivot-value*;
369 and finally the optional file *pivot* contains the corresponding pivot identifiers to identify the
370 prediction results, one per line.

371 The difference between check files and test files is only the presence (in check files) or
372 absence (for test files) of group identifiers for attribute data. This is also visible in the
373 difference in the semantics of the last two commands. The command **mcp-check** checks the
374 accuracy of the solution with respect to existing knowledge contained in the data, whereas the
375 command **mcp-predict** synthesises the missing attribute from data by means of a previously
376 learned formula.

377 We can transform a check file into a test file by means of the following command:

378 **mcp-chk2tst** -i *input-file* -o *output-file*

379 which essentially discards the group information from data.

380 6 System Distribution and Examples

381 The MCP system is available at the github.com/miki-hermann/mcp. Please, follow the
382 instructions in `README.md` file at the root. It is indispensable to run the installation
383 instructions described in that file to be able to run the MCP system properly.

384 The overall performance of the MCP system is very competitive, both in terms of time,
385 as well as in terms of quality of the produced formulas. The performance of the system has
386 been measured on a DELL computer with an Intel Core™ i7-9700 CPU @ 3.00GHz × 8 with
387 16GB of memory, running under Linux Fedora 38. All examples from [8, 9] run under one
388 second.

389 We have been testing the MCP system on several examples from the UCI Machine
390 Learning Repository (archive.ics.uci.edu/ml). All examples in the subdirectories are
391 equipped with a `Makefile` simplifying the application of the MCP system on them. The
392 directory *uci* contains the following treated examples:

<i>abalone</i>	identifying abalone with 27 rings;
<i>accent</i>	identifying several accents in spoken English language;
<i>balance-scale</i>	identifying psychological experiments balancing a scale;
<i>balloons</i>	a toy example, where specific formulas are required to be produced;
<i>banknote</i>	identifying forged and genuine banknotes;
<i>breast-cancer-wisconsin</i>	identifying benign and malignant breast cancer cases in Wisconsin;
<i>car</i>	identifying very good cars;
<i>divorce</i>	predicting if a marriage will end up in a divorce according to an analysis of responses to psychological investigation;
<i>forest-fire</i>	predicting forest fires in July, August, and September;
<i>iris</i>	identifying three types of iris flowers;
<i>monk</i>	the well-know monk examples;

<i>mushroom</i>	identifying edible and poisonous mushrooms;
<i>nursery</i>	for admission of children into a nursery;
<i>optdigits</i>	for determining digits from optical reading;
<i>shuttle</i>	the shuttle example;
<i>vote</i>	identifying democrats and republicans in the House of Representatives according to the 1984 US Congressional Voting Records.

393

394 We would especially drive the readers attention to the *mushroom* example, which identifies
 395 the edible and poisonous mushrooms always with 100% accuracy. This illustrates very well
 396 the strength of the MCP system.

397 **7** Concluding Remarks

398 The MCP system consists of more than 7000 lines of C++ code, using only the standard
 399 library. Parallel execution requires installation of the MPI software. Future versions of MCP
 400 will include a web GUI to enhance usability, as well as support for finite domains [12] to
 401 obviate the need for data binarization.

402 ——— References ———

- 403 1 D. Angluin, M. Frazier, and L. Pitt. Learning conjunctions of Horn clauses. *Machine Learning*,
 404 9(2-3):147–164, 1992.
- 405 2 K. A. Baker and A. F. Pixley. Polynomial interpolation and the Chinese Remainder Theorem
 406 for algebraic systems. *Mathematische Zeitschrift*, 143(2):165–174, 1975.
- 407 3 E. Böehler, N. Creignou, S. Reith, and H. Vollmer. Playing with Boolean blocks, part I: Post’s
 408 lattice with applications to complexity theory. *SIGACT News*, 34(4):38–52, 2003.
- 409 4 E. Böehler, N. Creignou, S. Reith, and H. Vollmer. Playing with Boolean blocks, part II:
 410 Constraint satisfaction problems. *SIGACT News*, 35(1):22–35, 2004.
- 411 5 Endre Boros, Yves Crama, Peter L. Hammer, Toshihide Ibaraki, Alexander Kogan, and
 412 Kazuhisa Makino. Logical analysis of data: classification with justification. *Annals of*
 413 *Operations Research*, 188(1):33–61, 2011.
- 414 6 Endre Boros, Toshihide Ibaraki, and Kazuhisa Makino. Monotone extensions of Boolean data
 415 sets. In Ming Li and Akira Maruoka, editors, *Proceedings 8th International Conference on*
 416 *Algorithmic Learning Theory, (ALT ’97), Sendai (Japan)*, volume 1316 of *Lecture Notes in*
 417 *Computer Science*, pages 161–175. Springer, 1997.
- 418 7 David R. Butenhof. *Programming with POSIX threads*. Addison-Wesley, 1997.
- 419 8 Arthur Chambon, Tristan Boureau, Frédéric Lardeux, and Frédéric Saubion. Logical
 420 characterization of groups of data: a comparative study. *Applied Intelligence*, 48(8):2284–2303,
 421 2018.
- 422 9 Arthur Chambon, Frédéric Lardeux, Frédéric Saubion, and Tristan Boureau. Computing sets
 423 of patterns for logical analysis of data. Technical report, Université d’Angers, 2017.
- 424 10 Yves Crama and Peter L. Hammer. *Boolean Functions - Theory, Algorithms, and Applications*,
 425 volume 142 of *Encyclopedia of Mathematics and its Applications*. Cambridge University Press,
 426 2011.
- 427 11 M. R. Garey and D. S. Johnson. *Computers and intractability: A guide to the theory of*
 428 *NP-completeness*. W.H. Freeman and Co, 1979.
- 429 12 A. Gil, M. Hermann, G. Salzer, and B. Zanuttini. Efficient algorithms for constraint description
 430 problems over finite totally ordered domains. *SIAM Journal on Computing*, 38(3):922–945,
 431 2008.

XX:12 MCP: Learning Propositional Formulas from Binarized Data

- 432 **13** Petr Hájek, Martin Holeňa, and Jan Rauch. The GUHA method and its meaning for data
433 mining. *Journal of Computer and System Sciences*, 76(1):34–48, 2010.
- 434 **14** J.-J. Hébrard and B. Zanuttini. An efficient algorithm for Horn description. *Information*
435 *Processing Letters*, 88(4):177–182, 2003.
- 436 **15** Petr Hájek and Tomáš Havránek. *Mechanizing Hypothesis Formation*. Springer, 1978.
- 437 **16** Marc Snir, Steve W. Otto, Steven Huss-Lederman, David W. Walker, and Jack Dongarra.
438 *MPI: The complete reference*. MIT Press, 1995.