

From-Below Boolean Matrix Factorization Algorithm Based on MDL

Tatiana Makhalova^{1,2} and Martin Trnečka³

National Research University Higher School of Economics, Moscow, Russia, LORIA,
(CNRS – Inria – University of Lorraine), Vandœuvre-lès-Nancy, France
Dept. Computer Science, Palacký University Olomouc, Olomouc, Czech Republic
tpmakhalova@hse.ru, martin.trnecka@gmail.com

Abstract. During the past few years Boolean matrix factorization (BMF) has become an important direction in data analysis. The minimum description length principle (MDL) was successfully adapted in BMF for the model order selection. Nevertheless, an BMF algorithm performing good results from the standpoint of standard measures in BMF is missing. In this paper, we propose a novel from-below Boolean matrix factorization algorithm based on formal concept analysis. The algorithm utilizes the MDL principle as a criterion for the factor selection. On various experiments we show that the proposed algorithm outperforms—from different standpoints—existing state-of-the-art BMF algorithms.

1 Introduction

Boolean matrix factorization (BMF), also known as Boolean matrix decomposition, is a powerful and widely used data mining tool. Like a classical matrix factorization methods, e.g. non-negative matrix factorization (NNMF) or singular value decomposition (SVD), BMF provides a different description (see Section 3.2) of Boolean data, via new, more fundamental variables called factors.

In BMF a given input data matrix is approximated by a product of so-called object-factor and factor-attribute matrices. All matrices contain zeros and ones only. The quality of the factorization—i.e. the quality of factors themselves—is usually measured by standard measures in BMF, namely by the number of factors and by the coverage (how large is the portion of data is described by factors, see Section 3). Both can be easily implemented—in fact each successful BMF algorithm already utilized them—in an arbitrary BMF algorithm. Moreover, both are very important in the evaluation of the factorization quality [1]. On the other hand, other aspects of the quality of factors, e.g. the interpretability, that are often neglected in the factor evaluation, are also an important parts of the matrix factorization.

By now, various approaches to assessment of the quality of factors were developed [1, 12]. One of the most fundamental—but surprisingly not often used—is based on the well-known minimum description length principle (MDL). In terms of MDL, the best factorization is the factorization with the minimal description.

Due to the MDL principle, such factorization is useful and easily interpretable. Nevertheless, it was many times shown (see, e.g. [1, 2]) that mixing MDL and BMF produces a poor results with respect to the BMF standard error measures (the number of factors and the coverage). More details will be provided in Section 3.

Recent results [13] in the field of formal concept analysis (FCA)—which is related to the BMF (see Section 3.3)—involving the minimum description length (MDL) motivate us to revise the use of MDL in BMF.

We propose a new heuristic BMF algorithm for from-below matrix factorization that outperforms existing state-of-the-art algorithms and produces very good result w.r.t. the standard BMF measures. The algorithm utilizes formal concept analysis and the MDL principle. Additionally, we present an extensive experimental evaluation of factors delivered by the proposed algorithm and its comparison with some already existing algorithms.

The rest of the paper is organized as follows. In the following Section 2 we provide a brief overview of the related work. Then, in Section 3, a notation used in the paper, a short introduction to BMF and MDL, and a background of the paper are presented. Section 4 describes a design of our algorithm. The algorithm is experimentally evaluated in Section 5. Section 6 draws a conclusion and future research directions.

2 Related Work

In the last decade, many BMF methods were developed [12, 2, 14, 3, 11, 19]. It was shown [18] that applying existing non Boolean methods (e.g. NNMF, SVD) on Boolean data is inappropriate, especially from the interpretation standpoint.

A good overview of BMF and related topics can be found e.g. in [1, 2, 14]. In general, BMF and BMF algorithms are addressed in various papers involving formal concept analysis [3, 8], role mining [4], binary databases [6] or bipartite graphs [16].

In many application of BMF, instead of a general Boolean factorization—which can be computed for instance by well-known ASSO algorithm—only a certain class of factorization, so-called from-below matrix factorization [2], is considered (see Section 3).

In the recent years, the minimum description length principle [7] has been applied in BMF. It was used mostly to solve the model order selection problem [15]—i.e. separation of global structure from noise—or as a factor selection criteria in BMF algorithms, e.g. in the state-of-the-art algorithm PANDA⁺ [12] (an improvement and generalized version of PANDA algorithm [11]). As a special case of application of MDL in BMF HYPER [19] algorithm can be considered, its objective is to minimize the description of factors instead of the minimization of the description length (for more details see [12]).

Another related work is [13], where a set of formal concepts with MDL is considered for the classification task. Our algorithm can be used for similar tasks. Instead of [13] our algorithm does not require computing the whole set

of formal concepts, that makes it applicable in practice. Moreover we used a different approach to MDL measuring.

This paper is, to the best of the author’s knowledge, the first to address the from-below decomposition based on the MDL.

3 Background and Basic Definitions

3.1 Notation

Through the paper we use a matrix terminology and in some convenient places a relational terminology. Matrices are denoted by upper-case bold letters (\mathbf{I}). \mathbf{I}_{ij} denotes the entry corresponding to the row i and the column j of \mathbf{I} . The set of all $m \times n$ Boolean (binary) matrices is denoted by $\{0, 1\}^{m \times n}$. The number of 1s in Boolean matrix \mathbf{I} is denoted by $\|\mathbf{I}\|$, i.e. $\|\mathbf{I}\| = \sum_{i,j} \mathbf{I}_{ij}$.

We interpret input data $\mathbf{I} \in \{0, 1\}^{m \times n}$ primarily as an object-attribute incidence matrix, i.e. a relation between the set of objects and the set of attributes. That is, the entry \mathbf{I}_{ij} is either 1 or 0, indicating that the object i does or does not have the attribute j .

If $\mathbf{A} \in \{0, 1\}^{m \times n}$ and $\mathbf{B} \in \{0, 1\}^{m \times n}$, we have the following element-wise matrix operations. The *Boolean sum* $\mathbf{A} \oplus \mathbf{B}$ which is the normal matrix sum where $1 + 1 = 1$. The *Boolean subtraction* $\mathbf{A} \ominus \mathbf{B}$ which is the normal matrix subtraction, where $0 - 1 = 0$.

3.2 Boolean Matrix Factorization

A general aim in BMF is for a given Boolean matrix $\mathbf{I} \in \{0, 1\}^{m \times n}$ to find matrices $\mathbf{A} \in \{0, 1\}^{m \times k}$ and $\mathbf{B} \in \{0, 1\}^{k \times n}$ for which

$$\mathbf{I} \approx \mathbf{A} \circ \mathbf{B} \tag{1}$$

where \circ is Boolean matrix multiplication, i.e. $(\mathbf{A} \circ \mathbf{B})_{ij} = \max_{l=1}^k \min(\mathbf{A}_{il}, \mathbf{B}_{lj})$, and \approx represents approximate equality assessed by $\|\cdot\|$. The corresponding metric E is defined for matrices $\mathbf{I} \in \{0, 1\}^{m \times n}$, $\mathbf{A} \in \{0, 1\}^{m \times k}$ and $\mathbf{B} \in \{0, 1\}^{k \times n}$ by

$$E(\mathbf{I}, \mathbf{A} \circ \mathbf{B}) = \|\mathbf{I} \ominus (\mathbf{A} \circ \mathbf{B})\|. \tag{2}$$

A decomposition of \mathbf{I} into $\mathbf{A} \circ \mathbf{B}$ may be interpreted as a discovery of k factors that exactly or approximately explain the data: interpreting \mathbf{I} , \mathbf{A} , and \mathbf{B} as the object–attribute, object–factor, and factor–attribute matrices, the model (1) has the following interpretation: the object i has the attribute j , i.e. $\mathbf{I}_{ij} = 1$, if and only if there exists factor l such that l applies to i and j is one of the particular manifestations of l .

Note also an important geometric view of BMF: a decomposition $\mathbf{I} \approx \mathbf{A} \circ \mathbf{B}$ with k factors represents a coverage of the 1s in \mathbf{I} by k rectangular areas in \mathbf{I} full of 1s, the l th rectangle is the Boolean sum of the l th column in \mathbf{A} and the l th row in \mathbf{B} . For more details see, e.g. [9].

If the rectangular areas cover only non zero elements in the matrix \mathbf{I} , the $\mathbf{A} \circ \mathbf{B}$ is called the *from-below matrix decomposition* [2]. An example of the from-below BMF follows.

Example 1. Let us consider Boolean matrix with rows $1, \dots, 8$ and columns a, \dots, h depicted in Figure 1. The Boolean matrix is given in the shape of table, where nonzero entries are marked by crosses. Two different factorizations of the data are shown in Figure 2.

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>
1	×	×	×				×	×
2	×	×	×				×	×
3	×	×	×	×				
4	×	×	×	×		×		
5		×	×	×				×
6		×	×	×	×	×		
7		×	×	×	×	×	×	×
8					×	×	×	×

Fig. 1. Example data.

3.3 BMF with Help of Formal Concept Analysis

Formal concept analysis (FCA) [5] provides a basic framework for dealing with factors. The main notion of FCA is *formal context*, which is usually represented as a Boolean matrix, it is defined as a triple $\langle \mathcal{X}, \mathcal{Y}, \mathcal{I} \rangle$, where \mathcal{X} is a nonempty set of objects, \mathcal{Y} is a nonempty set of attributes and \mathcal{I} is a binary relation between \mathcal{X} and \mathcal{Y} . Hence the formal context $\langle \mathcal{X}, \mathcal{Y}, \mathcal{I} \rangle$ with m objects and n attributes is a Boolean matrix $\mathbf{I} \in \{0, 1\}^{m \times n}$.

To every Boolean matrix $\mathbf{I} \in \{0, 1\}^{n \times m}$, one might associate the pair $\langle \uparrow, \downarrow \rangle$ of operators (in FCA well known as the arrow operators) assigning to sets $C \subseteq \mathcal{X} = \{1, \dots, m\}$ and $D \subseteq \mathcal{Y} = \{1, \dots, n\}$ the sets $C^\uparrow \subseteq \mathcal{Y}$ and $D^\downarrow \subseteq \mathcal{X}$ defined by

$$C^\uparrow = \{j \in \mathcal{Y} \mid \forall i \in C : \mathbf{I}_{ij} = 1\},$$

$$D^\downarrow = \{i \in \mathcal{X} \mid \forall j \in D : \mathbf{I}_{ij} = 1\},$$

where C^\uparrow is the set of all attributes (columns) shared by all objects (rows) in C and D^\downarrow is the set of all objects sharing all attributes in D .

The pair $\langle C, D \rangle$ for which $C^\uparrow = D$ and $D^\downarrow = C$ is called the *formal concept*. C and D are called the *extent* and the *intent* of formal concept $\langle C, D \rangle$, respectively. The concepts are partially ordered as follows: $\langle A, B \rangle \leq \langle C, D \rangle$ iff $A \subseteq C$ (or

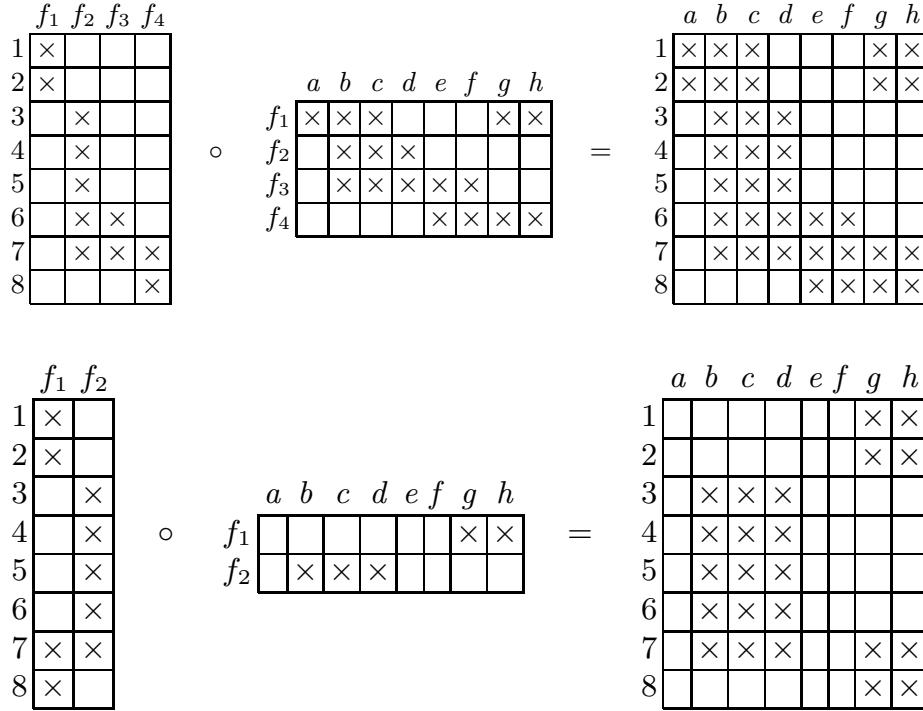


Fig. 2. Two examples of data factorization.

$D \subseteq B$), a pair $\langle A, B \rangle$ is a subconcept of $\langle C, D \rangle$, while $\langle C, D \rangle$ is a superconcept of $\langle A, B \rangle$. The set of all formal concepts we denote by

$$\mathcal{B}(\mathbf{I}) = \{\langle C, D \rangle \mid C \subseteq \mathcal{X}, D \subseteq \mathcal{Y}, C^\uparrow = D, D^\downarrow = C\}.$$

The whole set of partially ordered formal concepts is called the *concept lattice* of \mathbf{I} .

Given a set $\mathcal{F} = \{\langle C_1, D_1 \rangle, \dots, \langle C_k, D_k \rangle\} \subseteq \mathcal{B}(\mathbf{I})$ (with a fixed indexing of the formal concepts $\langle C_l, D_l \rangle$), induces the $m \times k$ and $k \times n$ Boolean matrices $\mathbf{A}_{\mathcal{F}}$ and $\mathbf{B}_{\mathcal{F}}$ by

$$(\mathbf{A}_{\mathcal{F}})_{il} = \begin{cases} 1, & \text{if } i \in C_l, \\ 0, & \text{if } i \notin C_l, \end{cases} \quad (3)$$

and

$$(\mathbf{B}_{\mathcal{F}})_{lj} = \begin{cases} 1, & \text{if } j \in D_l, \\ 0, & \text{if } j \notin D_l, \end{cases} \quad (4)$$

for $l = 1, \dots, k$. That is, the l th column and l th row of $\mathbf{A}_{\mathcal{F}}$ and $\mathbf{B}_{\mathcal{F}}$ are the characteristic vectors of C_l and D_l , respectively. The set \mathcal{F} is also called a set of *factor concepts*. Clearly, $\mathbf{A}_{\mathcal{F}} \circ \mathbf{B}_{\mathcal{F}}$ is the from-below matrix decomposition.

Example 2. Let us consider two factorizations depicted in Figure 2. The first one corresponds to the set

$$\mathcal{F} = \{\langle\{1, 2\}, \{a, b, c, g, h\}\rangle, \langle\{3, 4, 5, 6, 7\}, \{b, c, d\}\rangle, \langle\{6, 7\}, \{b, c, d, e, f\}\rangle, \langle\{7, 8\}, \{e, f, g, h\}\rangle\}.$$

The second one corresponds to the set

$$\mathcal{F} = \{\langle\{1, 2, 7, 8\}, \{g, h\}\rangle, \langle\{3, 4, 5, 6, 7\}, \{b, c, d\}\rangle\}.$$

For more details how formal concept analysis is utilized in BMF and the advantages of such approach see the pioneer work [3].

3.4 A Brief Introduction to MDL

The minimum description length (MDL) principle, which is a computable version of Kolmogorov complexity [7], is a formalization of the law of parsimony, well known as Occam’s razor. In terms of MDL, it is formulated as follows: the best model is the model that ensures the best compression of the given data.

More formally, for a given set of models \mathcal{M} and data (in our case represented via Boolean matrix \mathbf{I}) the best model $M \in \mathcal{M}$ is the one that minimizes the following cost function:

$$L(M) + L(\mathbf{I} | M), \tag{5}$$

where $L(M)$ is the encoding length of M in bits and $L(\mathbf{I} | M)$ is the encoding length in bits of the data \mathbf{I} encoded with M .

In general, we are only interested in the length of the encoding, and not in the coding itself, i.e. we do not have to materialize the codes themselves.

Note that MDL requires the compression to be lossless in order to allow for a fair comparison between different models.

3.5 The Quality of Factorization

The quality of the obtained factorization (1) is usually evaluated via some variants of metric (2). From the BMF perspective there are two basic viewpoints, emphasizing the role of the first k factors and the need to account for a prescribed portion of data, respectively. They are known as the discrete basis problem (DBP) and the approximate factorization problem (AFP), see [14] and [3, 2]. Both of them emphasize the coverage of data, i.e. the geometric view of BMF.

In many applications of BMF, the interpretation of factors plays a crucial role. It is reasonable instead of the coverage of the obtained factorization emphasize a different quality measures that access the interpretability of factors, e.g. the MDL.

On the other hand, the geometric view of BMF is very important and an interpretable factorization should reflect it.

In the next section, we propose a novel BMF algorithm which is based on well-known GRECOND algorithm [3]. The algorithm computes from-below factorization via minimization of the cost function (5). The results of experiments show that it preserves a lot of information from the original data w.r.t. the error measure (2).

4 Design of Algorithm

4.1 MDL in From-below Matrix Factorization

For matrices $\mathbf{A}_{\mathcal{F}} \in \{0, 1\}^{m \times k}$, $\mathbf{B}_{\mathcal{F}} \in \{0, 1\}^{k \times m}$, and $\mathbf{I} \in \{0, 1\}^{m \times n}$ where $\mathbf{I} \approx (\mathbf{A}_{\mathcal{F}} \circ \mathbf{B}_{\mathcal{F}})$ we define an error matrix \mathbf{E} as follows:

$$\mathbf{I} = (\mathbf{A}_{\mathcal{F}} \circ \mathbf{B}_{\mathcal{F}}) \oplus \mathbf{E}.$$

One may observe that matrix \mathbf{E} can be easily computed via metric (2), i.e. $\mathbf{E} = E(\mathbf{I}, \mathbf{A}_{\mathcal{F}} \circ \mathbf{B}_{\mathcal{F}})$. Hence, to provide a lossless compression of \mathbf{I} it is sufficient to encode the matrices $\mathbf{A}_{\mathcal{F}}$, $\mathbf{B}_{\mathcal{F}}$ and \mathbf{E} , i.e. the MDL cost function (5) has the following form

$$L(\mathbf{A}_{\mathcal{F}} \circ \mathbf{B}_{\mathcal{F}}) + L(\mathbf{E}). \quad (6)$$

According to the MDL principle, the best factorization of \mathbf{I} minimizes function (6). In the following we explain how to compute the length of the encoding of matrices $\mathbf{A}_{\mathcal{F}}$, $\mathbf{B}_{\mathcal{F}}$ and \mathbf{E} in bits. We use a similar approach as in [15] and we modify it for the from-below matrix factorization.

More precisely, to use optimal prefix codes we need to encode the dimensions of the matrices and the matrices themselves, i.e.

$$\begin{aligned} L(\mathbf{A}_{\mathcal{F}} \circ \mathbf{B}_{\mathcal{F}}) + L(\mathbf{E}) &= L(m) + L(n) + L(k) + \\ &+ L(\mathbf{A}_{\mathcal{F}}) + L(\mathbf{B}_{\mathcal{F}}) + L(\mathbf{E}). \end{aligned}$$

For the sake of simplicity we may encode the dimensions m, n, k with block-encoding, which give us $L(m) = L(n) = L(k) = \log(\max(m, n, k))$.

To not introduce some influencing between factors, these are encoded per factor, i.e. we encode $\mathbf{A}_{\mathcal{F}}$ per column and $\mathbf{B}_{\mathcal{F}}$ per row.

In order to use optimal prefix code, we need to first encode the probability of encountering 1 in a particular column or row respectively, i.e. we need $\log m$ bits for each extent in set \mathcal{F} and $\log n$ bits for each \mathcal{F} intent in set \mathcal{F} , respectively.

For simplicity, extent C and intent D of factor concept $\langle C, D \rangle$ can be seen as characteristic vectors, i.e. $C \in \{0, 1\}^{m \times 1}$ and $D \in \{0, 1\}^{1 \times n}$. We need to encode all ones and zeros. The length of optimal code is determined by Shannon entropy. This gives us the number of bits required for the encoding of matrices $\mathbf{A}_{\mathcal{F}}$ and $\mathbf{B}_{\mathcal{F}}$:

$$L(\mathbf{A}_{\mathcal{F}}) = \sum_{\langle C, D \rangle \in \mathcal{F}} \log m - (\|C\| \cdot \log \frac{\|C\|}{m} + (m - \|C\|) \cdot \log \frac{m - \|C\|}{m}),$$

$$L(\mathbf{B}_{\mathcal{F}}) = \sum_{\langle C, D \rangle \in \mathcal{F}} \log n - (\|D\| \cdot \log \frac{\|D\|}{n} + (n - \|D\|) \cdot \log \frac{n - \|D\|}{n}).$$

In a similar way we can compute the number of bits required for the encoding of matrix \mathbf{E} :

$$L(\mathbf{E}) = \log mn - (\|\mathbf{E}\| \cdot \log \frac{\|\mathbf{E}\|}{mn} + (mn - \|\mathbf{E}\|) \cdot \log \frac{mn - \|\mathbf{E}\|}{mn}).$$

Note, we can encode matrix \mathbf{E} element-by-element without any influence, because these elements are clearly independent.

4.2 Algorithm

In this section we propose a BMF algorithm, called MDLGRECOND¹, that uses the above described MDL cost function. The algorithm is a modified version—it utilizes a similar search strategy—of the GRECOND² algorithm [3], which is one of the most successful from-below matrix decomposition algorithms (see e.g. [1]).

Pseudocode of MDLGRECOND is depicted in Algorithm 1. The algorithm works as follows.

The algorithm computes a candidate $\langle C, D \rangle$ to a factor concept that minimizes the cost function (6) stored in variable *total_cost*. This is done via searching of a promising column j that is not included in D (lines 8–21). Note that the adding of j to D is realized via \uparrow and \downarrow operators mentioned in Section 3.3. Only the best column j is considered (lines 16–20). If a new column is added to $\langle C, D \rangle$, i.e. the $\langle C, D \rangle$ is changed, the modified $\langle C, D \rangle$ is used as a new candidate and another promising column is searched for. If there is no column that reduce the cost function (line 6), already computed candidate is added to the output set \mathcal{F} of factor concepts. The algorithm ends if there is no candidate that allows for reduction of the cost function.

¹ MDLGRECOND is an abbreviation of Minimum Description Length Greedy Concept on Demand.

² GRECOND is an abbreviation of Greedy Concept on Demand.

Input: Boolean matrix \mathbf{I} .
Output: Set \mathcal{F} of factor concepts.

```

1  $\mathcal{F} \leftarrow \emptyset$ 
2  $total\_cost \leftarrow \infty$ 
3  $\mathbf{E} \leftarrow \mathbf{I} \ominus (\mathbf{A}_{\mathcal{F}} \circ \mathbf{B}_{\mathcal{F}})$ 
4 while  $L(\mathbf{A}_{\mathcal{F}} \circ \mathbf{B}_{\mathcal{F}}) + L(\mathbf{E})$  is decreasing do
5    $\langle C, D \rangle \leftarrow \langle \emptyset, \emptyset \rangle$ 
6   while  $\langle C, D \rangle$  is changing do
7      $total\_cost' \leftarrow total\_cost$ 
8     foreach  $j \notin D$  do
9        $D' \leftarrow (D \cup \{j\})^{\downarrow\uparrow}$ 
10       $C' \leftarrow D'^{\downarrow}$ 
11      if  $\langle C', D' \rangle \in \mathcal{F}$  then
12        | continue with next  $j$ 
13      end
14       $\mathcal{F}' \leftarrow \mathcal{F} \cup \langle C', D' \rangle$ 
15       $cost \leftarrow L(\mathbf{A}_{\mathcal{F}'} \circ \mathbf{B}_{\mathcal{F}'}) + L(\mathbf{I} \ominus (\mathbf{A}_{\mathcal{F}'} \circ \mathbf{B}_{\mathcal{F}'}))$ 
16      if  $cost < total\_cost'$  then
17        |  $total\_cost' \leftarrow cost$ 
18        |  $C'' \leftarrow C'$ 
19        |  $D'' \leftarrow D'$ 
20      end
21    end
22     $total\_cost \leftarrow total\_cost'$ 
23     $C \leftarrow C''$ 
24     $D \leftarrow D''$ 
25  end
26   $\mathcal{F} \leftarrow \mathcal{F} \cup \langle C, D \rangle$ 
27 end
28 return  $\mathcal{F}$ 

```

Algorithm 1: MDLGRECOND algorithm

4.3 Computational Complexity

The Boolean matrix factorization problem is NP-hard [17] as well as the computation of factorization that minimizes the cost function (6). The proposed algorithm is heuristic. One may easily derive an exact algorithm with an exponential time complexity. Such algorithm is inapplicable in practice.

We do not provide the time complexity analysis, since the time complexity is not a main concern of Boolean matrix factorization. The presented algorithm is only slightly slower than GRECOND algorithm, which is, probably, the fastest BMF algorithm (see e.g. [2]). Both of them are able to factorize, in order of second, on ordinary PC, all the data presented in Section 5.

5 Experimental Evaluation

In this section, the results of an experimental comparison of BMF algorithms with MDLGRECOND are presented.

5.1 Datasets

We use 6 different real-world datasets, namely **Breast**, **Ecoli**, **Iris** and **Mushroom** from UCI repository [10], and **Domino** and **Emea** from [4]. The characteristics of the datasets are shown in Table 1. All of them are well known and widely used as benchmark datasets in BMF.

Table 1. Datasets and their characteristics.

dataset	size	dens. \mathbf{I}	$\ \mathcal{B}(\mathbf{I})\ $
Breast	699×20	0.499	642
Domino	79×231	0.400	73
Ecoli	336×34	0.235	813
Emea	3046×35	0.068	780
Iris	150×19	0.263	164
Mushroom	8124×90	0.252	186332

5.2 Algorithms

GRECOND [3] algorithm is based on the “on demand” greedy search for formal concepts of \mathbf{I} . It is designed to compute an exact from-below factorization. Instead of going through all formal concepts, which are the candidates for factor concepts, it constructs the factor concepts by adding sequentially “promising columns” to candidate $\langle C, D \rangle$ to factor concept. More formally, a new column j that minimizes the error

$$E(\mathbf{I}, \mathbf{A}_{\mathcal{F} \cup \langle (D \cup j)^\downarrow, (D \cup j)^\uparrow \rangle} \circ \mathbf{B}_{\mathcal{F} \cup \langle (D \cup j)^\downarrow, (D \cup j)^\uparrow \rangle})$$

is added to $\langle C, D \rangle$. This is repeated until no such columns exist. If there is no such column, the $\langle C, D \rangle$ is added to the set \mathcal{F} . The algorithm ends if $E(\mathbf{I}, \mathbf{A}_{\mathcal{F}} \circ \mathbf{B}_{\mathcal{F}})$ is smaller than the prescribed parameter ϵ or the prescribed number of factors is reached. For more details see [2]. Note, that usually $\epsilon = 0$, i.e. the whole matrix \mathbf{I} is covered by factors. Such setting was adopted in our experiments.

PANDA⁺ [12] is an algorithmic framework based on PANDA [11] algorithm. The algorithm aims to extract a set \mathcal{F} of pairs $\langle C, D \rangle$ that minimizes the cost function:

$$\sum_{\langle C, D \rangle \in \mathcal{F}} (|C| + |D|) + E(\mathbf{I}, \mathbf{A}_{\mathcal{F}} \circ \mathbf{B}_{\mathcal{F}}).$$

Every $\langle C, D \rangle$ in \mathcal{F} is computed in two stages. On the first stage the core of $\langle C, D \rangle$ is computed, on the second stage the core is extended. A core is a rectangle, not necessarily a formal concept, contained in \mathbf{I} and it is computed by adding columns from a sorted list. Extension to $\langle C, D \rangle$ is performed by adding columns and rows to a core while such an addition allows for reducing the cost. Note, that PANDA⁺ does not produce the from-below factorization. The computation of PANDA⁺ is driven by several parameters (see [12]). All of them are tuned for each dataset. The best obtained results are reported.

HYPER [19] algorithm aims to extract a set \mathcal{F} of pairs $\langle C, D \rangle$ that minimize the cost function which is defined as follows:

$$\sum_{\langle C, D \rangle \in \mathcal{F}} (|C| + |D|) / E(\mathbf{I}, \mathbf{A}_{\mathcal{F}} \circ \mathbf{B}_{\mathcal{F}}).$$

As candidates to factors the set of all formal concepts $\mathcal{B}(\mathbf{I})$ together with all single attribute rectangles in data are considered. Each candidate is divided into a set of single row rectangles that are sorted according to the number of uncovered elements in \mathbf{I} . Then the algorithm tries to add the single row rectangles back to the candidate, until the above mentioned cost function decreases. After this, the algorithm in each iteration selects the concept $\langle C, D \rangle$ from the modified set of candidates that minimizes the cost function. HYPER algorithm produces the from-below factorization. The size of $\mathcal{B}(\mathbf{I})$ can be exponentially large. In such case HYPER has the exponential time complexity. To reduce computational cost authors of [19] propose to use only frequent formal concepts (the frequency is an additional parameter of the algorithm). Our experiments show that the frequency affects highly the performance of the algorithm. In our experiments we use the whole set of formal concepts $\mathcal{B}(\mathbf{I})$, (for the set sizes see the last column of Table 1).

5.3 Evaluation

In our experiments we compare MDLGRECOND algorithm with GRECOND, HYPER and PANDA⁺. We study factors themselves and how well they cover the analyzed datasets.

The number of factors One of the main characteristic of BMF algorithms is the number of factors they produce. We measure not only the total number of factors, but also how many non-trivial factors are computed. Under trivial factors we mean the single-attribute ones. The results are shown in Table 2.

As it can be seen from the table, PANDA⁺ tends to produce only few factors (w.r.t. the number of attributes, see Table 1).

HYPER returns the number of factors which is close to the number of attributes. Moreover, more than a half of them are trivial. This is true on all datasets with an exception of **Breast** and **Mushroom** data.

On average (see Figure 3), the number of non-trivial factors of GRECOND is better than the number in case of HYPER algorithm. MDLGRECOND generates

Table 2. The number of factors.

dataset	algorithm	no. of factors	
		non-trivial	trivial
Breast	GRECOND	15	4
	PANDA ⁺	4	0
	HYPER	36	0
	MDLGRECOND	6	1
Domino	GRECOND	13	8
	PANDA ⁺	3	0
	HYPER	10	132
	MDLGRECOND	7	3
Ecoli	GRECOND	38	3
	PANDA ⁺	6	0
	HYPER	35	30
	MDLGRECOND	8	1
Emea	GRECOND	9	33
	PANDA ⁺	3	0
	HYPER	3	35
	MDLGRECOND	7	2
Iris	GRECOND	8	12
	PANDA ⁺	8	0
	HYPER	13	15
	MDLGRECOND	7	0
Mushroom	GRECOND	98	3
	PANDA ⁺	8	0
	HYPER	89	2
	MDLGRECOND	50	0

a small set of factors, most of them are non-trivial. PANDA⁺ tends to produce the smallest number of factors. All of them are non-trivial.

However, considering only the number of factors might be insufficient, since usually one wants to find not just the smallest number of factors, but the set of factors that capture (coverage) a large part of data. Further we will show how the algorithms capture the analyzed data.

Data coverage Another important characteristic of factors is how much information from the analyzed dataset they retain. We measure it by coverage rate. We differentiate *data coverage* and *object coverage*. Data coverage measures the rate of “crosses” covered by factors in the dataset—this is a standard measure in BMF, see e.g. [1]. However, data coverage might be an inappropriate measure in cases where a dataset contains a lot of redundant attributes. Taking into consideration these cases, we measure the object coverage rate, i.e. how many objects are covered at least by one factor. The following example explains how the coverage measures are computed.

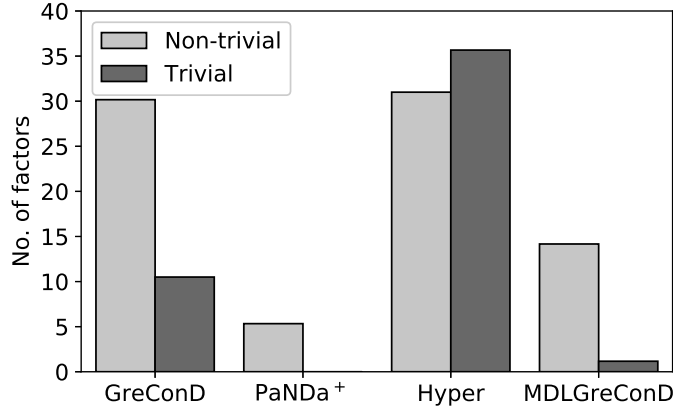


Fig. 3. The average number of factors.

Example 3. The factor set of the first factorization (Figure 2) covers almost all crosses in data, while the second set covers around a half of crosses. The coverings for both of them are given below. The crosses covered by one factor are light gray, the crosses covered by more factors are colored with darker gray.

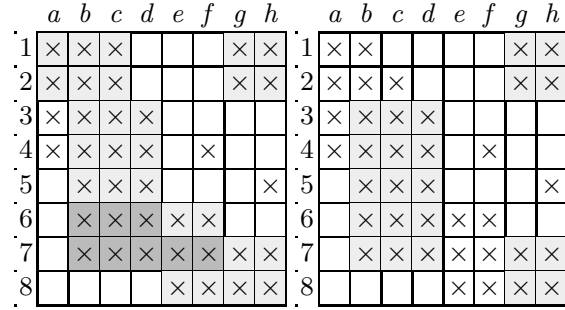


Fig. 4. The covering with factors from the running examples.

Note, both factor sets cover all objects, i.e. every row in the dataset has at least one colored cross, thus the object coverage rates is equal to 1 for both factorizations.

For the first factorization, the cross coverage rate is $35/39 = 0.897$. In the case of the second factorization, the cross coverage rate is $23/39 = 0.589$. Obviously, the bigger value is better.

Average values of data coverage and object coverage rates over all datasets as well as the minimal, maximal values and quantiles are shown in Figures 5 and 6 respectively. The average data coverage rate of non-trivial factors of MDLGRE-

COND is slightly lower than the analogous measure for GRECOND and HYPER. It is important to note that MDLGRECOND provides more stable results, in other words, the data coverage rate does not depend a lot on datasets, while for HYPER algorithm, the data coverage rate changes from 0.2 to 1.0. PANDA⁺ covers slightly more than a half of data by a small set of factors. Moreover, if we take into account results regarding the number of factors from Section 5.3 MDLGRECOND outperforms all remaining algorithms. Namely, it provides a large coverage by a smaller number of factors.

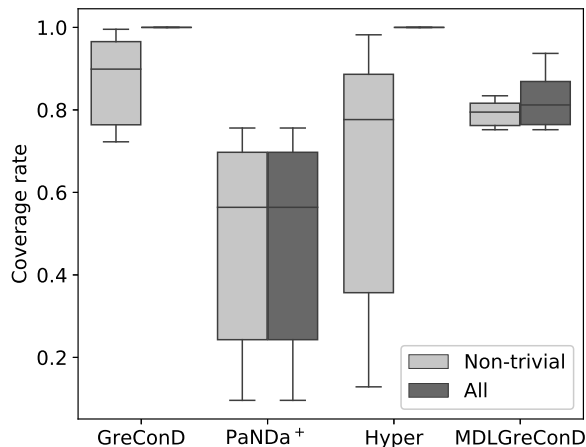


Fig. 5. The average data coverage rate.

Regarding the object coverage rate, all the algorithms have similar performance, however a large number of non-trivial factors in HYPER ensures its high coverage rate for all chosen datasets.

Redundancy of factors An important characteristic of a factor set is redundancy. The factor set is redundant if it contains repetitive information, i.e. if it contains some overlaps between factors. We measure redundancy by *overlapping rate* (see Example 4), i.e. how many times the covered crosses are covered by several factors.

Example 4. For the factor sets from Figure 2 the average overlapping rate is computed as follows. We count the total area of factors $area(\langle C, D \rangle) = \|C\| \cdot \|D\|$. In the case of the first factorization we obtain $area(f_1) = 10$, $area(f_2) = 15$, $area(f_3) = 10$ and $area(f_4) = 8$. The total area is 43, the number of covered crosses is 35, thus, the average overlapping rate is $43/35$. The second factorization is without overlapped crosses, thus its average overlapping rate is 1.

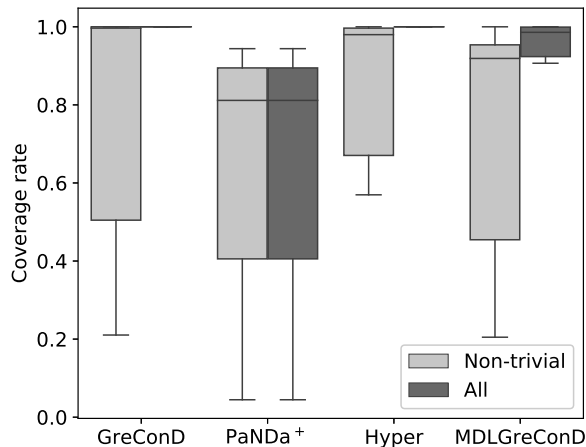


Fig. 6. The average object coverage rate.

Averages values of overlapping rate are shown in Figure 7. Our experiments show that factor sets with minimal redundancy are produced by HYPER algorithm. It can be explained regarding the previous experiments (see Section 5.3), where it was shown that HYPER algorithm tends to produce a large number of trivial factors. PANDA⁺ tends to produce a very small number of factors with low coverage rate. As one may clearly observe, GRECOND produces factorizations with the largest overlapping rate. MDLGRECOND generates a non-redundant set.

5.4 Discussion

Let us summarize the experimental evaluation. GRECOND and HYPER are both able to explain the whole data. However, the quality of factorizations they produce is lower than the quality of MDLGRECOND. More precisely, HYPER produces a large number of trivial factors. GRECOND produce a less number of trivial factors, but with a lot of overlappings between them.

The quality of factorization obtained via PANDA⁺ algorithm is low as well. The factors delivered by PANDA⁺ cover only a small part of input data.

According to the experimental evaluation, MDLGRECOND algorithm provide a factor set with well-balanced characteristics. The number of factors is reasonably small, factors themselves explain a large portion of data and are not redundant.

6 Conclusions

In this paper an MDL-based from-below factorization algorithm, which utilizes formal concept analysis, has been proposed. It produces a small subset of formal concepts having a low information loss rate.

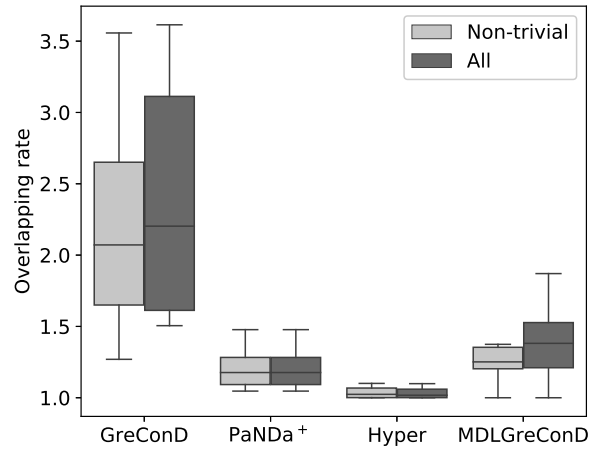


Fig. 7. The average overlapping rate.

The proposed algorithm does not require computing the whole set of formal concepts, that makes it applicable in practice. More than that, it computes factor sets that have better overall characteristics than factor sets computed by the existing BMF algorithms. The MDLGRECOND-generated factor sets are small, contain few single-attribute factors and have a high coverage with low overlapping rate.

An important direction of future work is application of the proposed method under supervised settings, i.e. for dealing with classification tasks.

Bibliography

- [1] Radim Belohlavek, Jan Outrata, and Martin Trnecka. Toward quality assessment of boolean matrix factorizations. *Inf. Sci.*, 459:71–85, 2018.
- [2] Radim Belohlavek and Martin Trnecka. From-below approximations in boolean matrix factorization: Geometry and new algorithm. *J. Comput. Syst. Sci.*, 81(8):1678–1697, 2015.
- [3] Radim Belohlavek and Vilem Vychodil. Discovery of optimal factors in binary data via a novel method of matrix decomposition. *J. Comput. Syst. Sci.*, 76(1):3–20, 2010.
- [4] Alina Ene, William G. Horne, Nikola Milosavljevic, Prasad Rao, Robert Schreiber, and Robert Endre Tarjan. Fast exact and heuristic methods for role minimization problems. In Indrakshi Ray and Ninghui Li, editors, *13th ACM Symposium on Access Control Models and Technologies, SACMAT 2008, Estes Park, CO, USA, June 11-13, 2008, Proceedings*, pages 1–10. ACM, 2008.
- [5] B. Ganter and R. Wille. *Formal Concept Analysis Mathematical Foundations*. Springer-Verlag, Berlin, Heidelberg, 1999.
- [6] Floris Geerts, Bart Goethals, and Taneli Mielikäinen. Tiling databases. In Einoshin Suzuki and Setsuo Arikawa, editors, *Discovery Science, 7th International Conference, DS 2004, Padova, Italy, October 2-5, 2004, Proceedings*, volume 3245 of *Lecture Notes in Computer Science*, pages 278–289. Springer, 2004.
- [7] Peter D. Grünwald. *The Minimum Description Length Principle (Adaptive Computation and Machine Learning)*. The MIT Press, 2007.
- [8] Dmitry I. Ignatov, Elena Nenova, Natalia Konstantinova, and Andrey V. Konstantinov. Boolean matrix factorisation for collaborative filtering: An fca-based approach. In Gennady Agre, Pascal Hitzler, Adila Alfa Krishnadh, and Sergei O. Kuznetsov, editors, *Artificial Intelligence: Methodology, Systems, and Applications - 16th International Conference, AIMS 2014, Varna, Bulgaria, September 11-13, 2014. Proceedings*, volume 8722 of *Lecture Notes in Computer Science*, pages 47–58. Springer, 2014.
- [9] Ki Hang Kim. *Boolean matrix theory and applications*, volume 70. Dekker, 1982.
- [10] M. Lichman. UCI machine learning repository, 2013.
- [11] Claudio Lucchese, Salvatore Orlando, and Raffaele Perego. Mining top-k patterns from binary datasets in presence of noise. In *Proceedings of the SIAM International Conference on Data Mining, SDM 2010, April 29 - May 1, 2010, Columbus, Ohio, USA*, pages 165–176. SIAM, 2010.
- [12] Claudio Lucchese, Salvatore Orlando, and Raffaele Perego. A unifying framework for mining approximate top-k binary patterns. *IEEE Trans. Knowl. Data Eng.*, 26(12):2900–2913, 2014.
- [13] Tatiana P. Makhhalova, Sergei O. Kuznetsov, and Amedeo Napoli. A first study on what MDL can do for FCA. In Dmitry I. Ignatov and Lhouari

- Nourine, editors, *Proceedings of the Fourteenth International Conference on Concept Lattices and Their Applications, CLA 2018, Olomouc, Czech Republic, June 12-14, 2018.*, volume 2123 of *CEUR Workshop Proceedings*, pages 25–36. CEUR-WS.org, 2018.
- [14] Pauli Miettinen, Taneli Mielikäinen, Aristides Gionis, Gautam Das, and Heikki Mannila. The discrete basis problem. *IEEE Trans. Knowl. Data Eng.*, 20(10):1348–1362, 2008.
 - [15] Pauli Miettinen and Jilles Vreeken. Model order selection for boolean matrix factorization. In Chid Apté, Joydeep Ghosh, and Padhraic Smyth, editors, *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Diego, CA, USA, August 21-24, 2011*, pages 51–59. ACM, 2011.
 - [16] S. D. Monson, S. Pullman, and R. Rees. A survey of clique and biclique coverings and factorizations of $(0,1)$ -matrices. In *Bulletin of the ICA*, 14, pages 17–86, 1995.
 - [17] L. J. Stockmeyer. *The Set Basis Problem is NP-complete*. Research reports. IBM Thomas J. Watson Research Division, 1975.
 - [18] Nikolaž Tatti, Taneli Mielikäinen, Aristides Gionis, and Heikki Mannila. What is the dimension of your binary data? In *Proceedings of the 6th IEEE International Conference on Data Mining (ICDM 2006), 18-22 December 2006, Hong Kong, China*, pages 603–612. IEEE Computer Society, 2006.
 - [19] Yang Xiang, Ruoming Jin, David Fuhry, and Feodor F. Dragan. Summarizing transactional databases with overlapped hyperrectangles. *Data Min. Knowl. Discov.*, 23(2):215–251, 2011.