



**HAL**  
open science

## A generalized approach for Boolean matrix factorization

Rodrigo Cabral Farias, Sebastian Miron

► **To cite this version:**

Rodrigo Cabral Farias, Sebastian Miron. A generalized approach for Boolean matrix factorization. *Signal Processing*, 2023, 206, pp.108887. <10.1016/j.sigpro.2022.108887>. <hal-03929288>

**HAL Id: hal-03929288**

**<https://hal.science/hal-03929288v1>**

Submitted on 8 Jan 2025

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.



Distributed under a Creative Commons CC BY-NC 4.0 - Attribution - Non-commercial use - International License

# A generalized approach for Boolean matrix factorization

Rodrigo Cabral Farias<sup>a,\*</sup>, Sebastian Miron<sup>b</sup>

<sup>a</sup> *Université Côte d'Azur, CNRS, I3S Laboratory, 06900 Sophia-Antipolis, France*

<sup>b</sup> *Université de Lorraine, CNRS, CRAN, 54000 Nancy, France*

---

## Abstract

In this paper, we propose a generalized framework for fitting Boolean matrix factorization models to binary data. In this generalized setting, the binary rank-1 components of the underlying model can be combined by any Boolean function, thus extending the standard Boolean matrix factorization model, where the combination is restricted to the logical 'OR' function. We introduce two algorithms relying on a relaxation of the binary constraints on the factors of the model and on a polynomial representation of the Boolean function that combines the rank-1 components. One of the algorithms is based on the gradient descent optimization method, while the other is based on block coordinate descent. A detailed presentation of the algorithms is given, along with numerical experiments both on simulated and real datasets. A comparison with other algorithms from the literature is presented in the standard Boolean matrix setting allowing to assess the advantages and shortcomings of the proposed methods in terms of factor retrieval and data denoising performance, convergence behavior and time complexity.

*Keywords:* Binary data, Binary matrix factorizations, Boolean matrix factorizations, Data mining.

---

## 1. Introduction

Binary data matrices are one of the most natural ways of numerically encoding data in many applications. They can encode *yes/no* answers to surveys, voting records, tables indicating the presence or absence of traits of a group of individuals or indicating the proximity between the individuals, implicit feedback or thresholded explicit feedback (grades) in audio/video streaming platforms, input/output relationships in digital circuits, among many others. For this reason, a large number of data mining techniques specially tailored for binary data matrices have been developed

---

\*Corresponding author

*Email address:* [Rodrigo.CABRAL-FARIAS@univ-cotedazur.fr](mailto:Rodrigo.CABRAL-FARIAS@univ-cotedazur.fr) (Rodrigo Cabral Farias)

*URL:* <https://www.i3s.unice.fr/rcabral/en/home> (Rodrigo Cabral Farias)

8 recently. Among these techniques, a diverse number of models and algorithms relying on matrix  
9 factorizations have appeared in the last two decades. They have been successfully used in a great  
10 number of applications such as text mining [1], recommender systems [2–4], genetics [5, 6], protein  
11 complex prediction [7], role mining [8] and telecommunications [9, 10].

12 Different factorization models for binary matrices have been proposed in the literature. One  
13 group of models [2, 11–14] relies on a modification of logistic regression where a logistic function is  
14 applied to a constrained matrix factorization model such as principal component analysis (PCA).  
15 Most models of this group are named logistic PCA or binary PCA. Although the matrix factors  
16 are not constrained to be binary, the use of the logistic function allows the model to constrain the  
17 elements of its output matrix to lie in the interval  $[0, 1]$ . A generalized version of logistic PCA,  
18 based on the family of mean-parameterized Bernoulli models, was recently studied in [15].

19 Another family of models aims at factorizing the data matrix into two binary matrices, which is  
20 equivalent to decomposing it as a sum of rank-1 binary matrices. For the needs of the presentation, in  
21 this paper we will distinguish between several sub-families of this factorization model, depending on  
22 the definition of the “sum”. We call binary matrix factorization (bMF) the decomposition using the  
23 classical (algebraic) sum between the rank-1 terms [16]. This is also the term used in the literature to  
24 refer to this generic family of models. The limitations of bMF appear whenever its expected rank-1  
25 components have overlapping supports. In this case, the classical sum of the rank-1 components  
26 does not lead to a binary matrix. One way to counter this issue is to assume that the presence of a  
27 ‘1’ in the data matrix is due to the contributions of several ‘1’ in the rank-1 terms. Note that this  
28 corresponds to simply replacing the arithmetic sums in the decomposition model by logical ‘OR’  
29 operations. This leads to a particular matrix factorization model, that we will call Boolean matrix  
30 factorization (BMF). There is also an interesting link between BMF and the so called *Subtropical*  
31 *Matrix Factorization (SMF)* model [17–19]. In the *subtropical (or max-times)* algebra framework,  
32 BMF can be seen as a special case of SMF where the entries of the matrices are restricted to the  
33 binary set  $\{0, 1\}$ .

34 A major issue for fitting optimally a bMF or BMF model to data is the discrete nature of the  
35 model parameters, which makes the underlying optimization problem difficult to solve. It has been  
36 shown that fitting a single component model is already a NP-hard problem [20]. Due to its hardness,  
37 algorithms do not aim to solve exactly the original fitting problem. A class of methods, such as the  
38 association rules algorithm (ASSO) [21] or formal concept analysis (FC) [22], rely on low complexity  
39 iterative rules that extract in a greedy-like manner binary rank-1 components that are expected to

40 approximate the optimal ones. A revisited version of the FC-based algorithms of [22], significantly  
41 faster, was recently introduced in [23]; another variant, that uses the minimum description length  
42 principle (MDL) for factor selection was proposed in [24]. A different approach, called the penalty  
43 function algorithm (PF) [16], solves a relaxed version of the underlying fitting problem. In PF, the  
44 bMF factor elements are relaxed to the nonnegative orthant, while a penalization term forcing the  
45 factors to be close to binary is added to the original data fitting objective function. Such a relaxation  
46 allows to use multiplicative gradient algorithms to retrieve the factors in a similar manner as for  
47 nonnegative matrix factorization (NMF) [25].

48 A modification of the PF algorithm explicitly tailored for BMF has been proposed in [26]. The  
49 authors propose to apply a threshold function to the bMF model, such that the output matrix  
50 elements are either 0 or 1. Since the threshold function is not differentiable, to be able to use a  
51 multiplicative gradient algorithm as in the PF algorithm, a smooth approximation of the threshold  
52 function is used. The resulting BMF method is called post-nonlinear PF algorithm (PNL-PF).  
53 A heuristic model selection algorithm for estimating the number of binary sources in the BMF  
54 setting was also proposed in [27], based on stability criteria. This method constructs an ensemble of  
55 random matrices that are slight perturbations of the initial matrix to test the stability of the Boolean  
56 decomposition. Other algorithms have also been developed under a stochastic setting, where the  
57 elements of the factors are supposed to be random [28, 29].

58 In this paper, we propose an approximate factorization approach for binary valued matrices that  
59 generalizes BMF to arbitrary Boolean “*sum*” functions. Instead of considering combinations of the  
60 rank-1 components with logical ‘OR’, we assume that an arbitrary Boolean function with known  
61 truth table is used. Our approach is based on the relaxation of the binary constraints, as in PF and  
62 PNL-PF, but instead of representing the behavior of the logical combiner with a threshold function,  
63 we represent it as a multivariate polynomial of the elements of each component. Since a multivariate  
64 polynomial is a differentiable function, such a representation allows developing a gradient algorithm  
65 for fitting the generalized BMF model, without the need to resort to smooth approximations, as in  
66 PNL-PF. We also propose a generalized BMF approximation algorithm based on block coordinate  
67 descent. The algorithm alternatively updates the columns of the factors to be retrieved in a similar  
68 manner as in hierarchical alternating least squares (HALS) for NMF [30]. Since the multivariate  
69 polynomial representing the logical combiner is multilinear in the elements of the components, the  
70 updates required in the block coordinate descent algorithm can be obtained in closed-form.

71 We present implementation details of our approach in the specific case of of BMF approximation,

72 and under this setting we compare the performance of the two resulting algorithms with state-of-the  
73 art BMF methods. The performance of the methods are evaluated in terms of denoising and factor  
74 retrieval capabilities, but also in terms of convergence behavior, time complexity and sensitivity to  
75 initializations. To illustrate our approach in a more practical setting, we apply one of the proposed  
76 methods to retrieve the BMF of 4 real datasets. Finally, we also show simulation results concerning  
77 the application of the general version of the proposed algorithms to retrieve factorizations where  
78 the component combining functions are the logical exclusive ‘OR’ (‘XOR’) and the 3-term majority  
79 function.

### 80 1.1. Outline

81 In Section 2, we present the binary and Boolean factorization models along with the polynomial  
82 representation of the general Boolean factorization. We introduce two algorithms for the generalized  
83 Boolean factorization in Section 3 and illustrate their implementation in the specific case of BMF.  
84 Section 4 shows the results of the conducted numerical experiments to compare the performance of  
85 the proposed algorithms with state-of-the art methods. Results on 4 real datasets are also given.  
86 We also provide numerical simulation results in a more general setting, where the data do not follow  
87 the standard BMF model. Finally, we conclude the paper in Section 5.

### 88 1.2. Notations

89 Scalars are represented by lower case letters  $x$ , while vectors are represented by bold-face lower  
90 case letters  $\mathbf{x}$ . Bold-face upper case letters  $\mathbf{X}$  are used to represent matrices. A single subscript  $x_i$   
91 is used to represent the  $i$ -th element of a vector or the  $i$ -th column of a matrix  $\mathbf{x}_i$  ( $i$ -th column of  
92  $\mathbf{X}$ ). A double subscript  $x_{ij}$  denotes the  $(i, j)$ -th element of a matrix. Superscripts (or subscripts)  
93 of the form  $\mathbf{x}^{1:n}$  denote the tuple  $(\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n)$ .

94 Matrix transpose is denoted  $\mathbf{X}^\top$ , while the Frobenius norm of a matrix is symbolized by  $\|\mathbf{X}\|_F$ .  
95 To denote a matrix of size  $I \times J$  with all elements equal to 1, we use  $\mathbf{1}_{I \times J}$ . The symbol  $\square$  denotes  
96 the Hadamard (entry-wise) matrix product and  $\text{Diag}(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n)$  denotes a block diagonal  
97 matrix with matrices  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$  in its diagonal blocks. The column-major vectorization of  
98 a matrix is denoted  $\text{vec}(\mathbf{X})$ .

99 Logical ‘OR’ is symbolized by  $\vee$  and the same symbol is used for its entry-wise matrix version.  
100 Logical ‘XOR’ is denoted by  $\oplus$ .

101 **2. General binary and Boolean factorizations**

102 We are interested in exactly or approximately decomposing a  $I \times J$  data matrix  $\mathbf{Y}$  with binary  
 103 elements  $y_{ij} \in \{0, 1\}$ , for  $(i, j) \in \{1, 2, \dots, I\} \times \{1, 2, \dots, J\}$ , into  $R \geq 2$  binary rank-1 matrices  
 104  $\mathbf{X}^{1:R} = \{\mathbf{X}^1, \mathbf{X}^2, \dots, \mathbf{X}^R\}$ . Each binary  $\mathbf{X}^r$  with  $r \in \{1, 2, \dots, R\}$  is written as

$$\mathbf{X}^r = \mathbf{a}_r \mathbf{b}_r^\top, \quad (1)$$

105 where  $\mathbf{a}_r$  and  $\mathbf{b}_r$  are vectors of sizes  $I$  and  $J$  respectively and with their elements constrained  
 106 to be binary  $[\mathbf{a}_r]_i \in \{0, 1\}$ , for  $(i, r) \in \{1, 2, \dots, I\} \times \{1, 2, \dots, R\}$ ,  $[\mathbf{b}_r]_j \in \{0, 1\}$ , for  $(j, r) \in$   
 107  $\{1, 2, \dots, J\} \times \{1, 2, \dots, R\}$ . The vectors  $\mathbf{a}_r$  and  $\mathbf{b}_r$  can be stored in matrices  $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_R]$   
 108 and  $\mathbf{B} = [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_R]$ . As presented in [26] and [27], different decompositions can be considered  
 109 depending on how one precisely defines the way that  $\mathbf{X}_r$  are combined to approximate  $\mathbf{Y}$ . If we  
 110 specify a function  $f : \{0, 1\}^R \rightarrow \mathcal{Y} \subset \mathbb{R}$  defined on  $R$  binary inputs and resulting in a value on a  
 111 finite subset  $\mathcal{R}$  of the integers, general binary factorization corresponds to approximate the elements  
 112 of  $\mathbf{Y}$  as follows:

$$y_{ij} \approx f(x_{ij}^1, x_{ij}^2, \dots, x_{ij}^R) = f(x_{ij}^{1:R}) = f(a_{i,1:R} b_{j,1:R}), \quad (2)$$

113 where  $a_{i,1:R} b_{j,1:R} = \{a_{i1} b_{j1}, a_{i2} b_{j2}, \dots, a_{iR} b_{jR}\}$ . Denoting the matrix resulting of the element-wise  
 114 application of  $f(\cdot)$  to the rank-one matrices  $\mathbf{X}^{1:R}$  simply by  $f(\mathbf{X}^{1:R})$ , the approximation problem  
 115 (2) can be cast as the following minimization problem:

$$\begin{aligned} & \text{minimize} && \mathcal{F}(\mathbf{A}, \mathbf{B}) = \frac{1}{2} \|\mathbf{Y} - f(\mathbf{X}^{1:R})\|_F^2 \\ & \text{(where} && \mathbf{X}^{1:R} = \{\mathbf{X}^1 = \mathbf{a}_1 \mathbf{b}_1^\top, \dots, \mathbf{X}^R = \mathbf{a}_R \mathbf{b}_R^\top\}, \\ & \text{with respect to} && \mathbf{A} \in \{0, 1\}^{I \times R}, \mathbf{B} \in \{0, 1\}^{J \times R}. \end{aligned} \quad (3)$$

116 A solution for this problem is guaranteed to exist since the cardinal of the feasible set is finite.  
 117 If a solution  $\mathbf{A}^*, \mathbf{B}^*$  of (3) achieves  $\mathcal{F}(\mathbf{A}, \mathbf{B}) = 0$ , we say that it is an exact factorization of  $\mathbf{Y}$ . In  
 118 the data analysis literature, 3 common types of factorization problems which are special cases of the  
 119 general form above are the following:

120 *Binary matrix factorization (bMF)*. When  $f(x_{ij}^{1:R}) = \sum_{r=1}^R x_{ij}^{1:R}$  is the usual sum on  $\mathbb{R}$ . See the left  
 121 column of Tab. 1 for an example when  $R = 2$ . We say that  $\mathbf{A}$  and  $\mathbf{B}$  are approximate factors of  $\mathbf{Y}$ ,  
 122 since in this case  $\mathbf{Y} \approx \mathbf{A} \mathbf{B}^\top$ .

123 *Boolean matrix factorization (BMF)*. When  $f(x_{ij}^{1:R}) = \bigvee_{r=1}^R x_{ij}^{1:R}$  is the  $R$ -term logical ‘OR’. See  
 124 the middle column of Tab. 1 for an example. Similarly to the previous case, we can say that  $\mathbf{A}$  and  
 125  $\mathbf{B}$  are approximate Boolean factors of  $\mathbf{Y}$ , since  $\mathbf{Y} \approx \mathbf{A} \wedge \mathbf{B}^\top$  where  $(\cdot \wedge \cdot)$  is the matrix product  
 126 defined in the Boolean semi-ring (sums are replaced by logical ‘OR’).

Inputs		Output		
$x^{(1)}, x^{(2)}$	bMF (+)	BMF ( $\vee$ )	F2MF ( $\oplus$ )	
0, 0	0	0	0	
0, 1	1	1	1	
1, 0	1	1	1	
1, 1	2	1	0	

Table 1: Results for different  $f(\cdot)$  with  $R = 2$  inputs used in different factorizations.

127  $\mathbb{F}_2$  matrix factorization (F2MF). When  $f(x_{ij}^{1:R}) = \bigoplus_{r=1}^R x_{ij}^{1:R}$  is the  $R$ -term modulo-2 sum, that  
128 is, a cascade of  $R$  logical ‘XOR’ operations applied sequentially to  $x_{ij}^{1:R}$ . In this case, we can write  
129  $\mathbf{Y} \approx \mathbf{A} \odot \mathbf{B}^\top$ , where  $(\cdot \odot \cdot)$  is the matrix product in the  $\mathbb{F}_2$  field (Galois field of two elements). Here  
130 the sums are replaced by logical ‘XOR’. Therefore, we call this model  $\mathbb{F}_2$  matrix factorization.

131 Observe that if one wants to factorize a binary data matrix  $\mathbf{Y}$  without errors using bMF, its  
132 factors should contain columns with disjoint supports. This is due to the presence of possible values  
133 larger than 1 in the outputs of the sum operation for bMF (see Tab. 1).

134 For a given  $\mathbf{Y}$ , the characteristics of its factorization such as rank or uniqueness may change  
135 depending on the chosen  $f(\cdot)$ . Before focusing on algorithms for general Boolean factorizations,  
136 which are the main contribution of this paper, we briefly illustrate with some toy examples, some  
137 important differences between factorizations with different  $f(\cdot)$ .

### 138 2.1. Ranks and uniqueness of binary factorizations

139 As in standard matrix factorizations, the minimal number of columns of the factors  $R$  for which  
140 exact factorizations of  $\mathbf{Y}$  exist is called the rank of  $\mathbf{Y}$  and we denote it  $\text{rank}_f(\mathbf{Y})$ :

$$\text{rank}_f(\mathbf{Y}) = \min \left\{ R \mid \mathbf{Y} = f(\mathbf{X}^{1:R}), \mathbf{A} \in \{0, 1\}^{I \times R}, \mathbf{B} \in \{0, 1\}^{J \times R} \right\}. \quad (4)$$

141 Following the denominations in [26, 31], if  $f(\cdot)$  is the standard sum, we call this rank the *binary rank*  
142 and we denote it  $\text{rank}_{\{0,1\}}(\mathbf{Y})$ . If  $f(\cdot)$  is the logical ‘OR’ then this rank is the *Boolean rank* and it  
143 is denoted  $\text{rank}_{\mathbb{B}}(\mathbf{Y})$ . We call it  $\mathbb{F}_2$  rank, when the combining function is the modulo-2 sum and we  
144 denote it  $\text{rank}_{\mathbb{F}_2}(\mathbf{Y})$ . In what follows, we give some toy examples from the literature to illustrate  
145 the fact that these ranks can be different for a given matrix. Consider the  $3 \times 3$  binary matrix [21]

$$\mathbf{Y}_1 = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 1 \end{bmatrix}. \quad (5)$$

146 Minimum rank decompositions of  $\mathbf{Y}_1$  for bMF and BMF are

$$\mathbf{Y}_1 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}^\top + \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}^\top + \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}^\top = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}^\top \vee \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}^\top$$

147 and the same factors of bMF can be used for F2MF. Thus, we have  $\text{rank}_{\mathbb{F}_2}(\mathbf{Y}_1) = 3$ , which is  
 148 larger than  $\text{rank}_{\mathbb{B}}(\mathbf{Y}_1) = 2$ . Note also that the rank of  $\mathbf{Y}_1$  on the reals is  $\text{rank}(\mathbf{Y}_1) = 3$ , since the 3  
 149 columns of  $\mathbf{Y}_1$  are linearly independent.

150 One can easily find cases where the relations between these ranks are different from the previous  
 151 example. Consider a matrix  $\mathbf{Y}_2$  which is equal to  $\mathbf{Y}_1$  except for the central element  $[\mathbf{Y}_1]_{2,2}$  which is  
 152 flipped to zero. Then the BMF factors for  $\mathbf{Y}_1$  give an exact F2MF for  $\mathbf{Y}_2$  with a minimum number  
 153 of columns. In this case,  $\text{rank}_{\mathbb{B}}(\mathbf{Y}_2) = \text{rank}(\mathbf{Y}_2) = 3$ , but  $\text{rank}_{\mathbb{F}_2}(\mathbf{Y}_2) = 2$ .

154 By increasing the size of the data matrix, one can also find cases where the  $\text{rank}(\mathbf{Y}) < \text{rank}_{\{0,1\}}(\mathbf{Y})$ .  
 155 For example, for [32]

$$\mathbf{Y}_3 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix},$$

156 we have  $\text{rank}(\mathbf{Y}_3) = 3$ , while  $\text{rank}_{\{0,1\}}(\mathbf{Y}_3) = \text{rank}_{\mathbb{B}}(\mathbf{Y}_3) = 4$ .

157 When the rank-one components  $\mathbf{X}^r$  have disjoint supports, all of the previously mentioned ranks  
 158 coincide. Thus, for [26]

$$\mathbf{Y}_4 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix}, \quad (6)$$

159  $\text{rank}_{\{0,1\}}(\mathbf{Y}_4) = \text{rank}_{\mathbb{B}}(\mathbf{Y}_4) = \text{rank}_{\mathbb{F}_2}(\mathbf{Y}_4) = \text{rank}(\mathbf{Y}_4) = 2$ , and one exact decomposition corre-  
 160 sponding to this rank has factors

$$\mathbf{A} = \mathbf{B} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix}, \quad (7)$$

161 for all factorizations presented above, but also for factorizations on  $\mathbb{R}$ .

162 In data mining applications, for interpretability reasons, it is expected that factors  $\mathbf{A}$  and  $\mathbf{B}$   
 163 could be retrieved uniquely from the data up to joint column permutations. It is well-known that,  
 164 in general, matrix factorization over the real numbers is not unique. For F2MF, the factorization  
 165 is unique only when  $R = 1$ , since for  $R \geq 2$  one can find binary  $\mathbf{T}, \mathbf{T}'$ , different from permutation  
 166 matrices, such that  $\mathbf{A} \odot \mathbf{B}^\top = (\mathbf{A} \odot \mathbf{T}) \odot (\mathbf{B} \odot \mathbf{T}')^\top$ . Regarding bMF and BMF, the binary constraints  
 167 on  $\mathbf{A}$  and  $\mathbf{B}$  allow to retrieve unique factors under some particular conditions (see [26, 31, 33–35]  
 168 for details on uniqueness conditions).

169 *2.2. Polynomial representation of a general Boolean function*

170 We focus in this paper in solving problem (3) where  $f(\cdot)$  is a general Boolean function  $f :$   
 171  $\{0, 1\}^R \rightarrow \{0, 1\}$ .

172 We consider a two-step approach: in the first step, we apply an optimization algorithm to solve  
 173 a relaxed version of (3) where the binary constraints are dropped. The elements of  $\mathbf{A}$  and  $\mathbf{B}$  are  
 174 either allowed to lie on  $\mathbb{R}$  or on the interval  $[0, 1]$ . In the second step, the resulting approximations  
 175 of  $\mathbf{A}$  and  $\mathbf{B}$ , denoted  $\hat{\mathbf{A}}$  and  $\hat{\mathbf{B}}$ , are then binarized by comparing the values of their elements with  
 176 a threshold of value 0.5. For the elements of  $\hat{\mathbf{A}}$ , this operation is

$$\left[ \mathcal{P}_{\mathbb{B}}(\hat{\mathbf{A}}) \right]_{ir} = \begin{cases} 0 & , \text{ for } \hat{a}_{ir} < 0.5, \\ 1 & , \text{ for } \hat{a}_{ir} \geq 0.5. \end{cases} \quad (8)$$

177 We have chosen a threshold value of 0.5 since in this case,  $\left[ \mathcal{P}_{\mathbb{B}}(\hat{\mathbf{A}}) \right]_{ir}$  corresponds to the projection  
 178 onto the set  $\mathbb{B} = \{0, 1\}$ , that is the solution of  $\arg \min_{z \in \{0, 1\}} (z - \hat{a}_{ir})^2$ . Other values of the threshold  
 179 could be used, but we would lose the natural interpretation of the binarization of the elements  $\hat{\mathbf{A}}$   
 180 and  $\hat{\mathbf{B}}$  as the solution of a standard projection problem.

181 To apply this approach, the Boolean function  $f(\cdot)$  must be represented by another function  $\bar{f}$   
 182 defined for real inputs.  $\bar{f}$  should be defined in such a way that both functions are equal for binary  
 183 inputs. In this work, we define  $\bar{f}(\cdot)$  relying on the fact that any Boolean function of  $R$  variables  
 184  $\mathbf{x} = [x_1, x_2, \dots, x_R]^T$  can be written as a multivariate polynomial. For any  $\mathbf{x} \in \{0, 1\}^R$ , the  
 185 following polynomial achieves the same values as  $f(\cdot)$  [36]:

$$\bar{f}(\mathbf{x}) = \sum_{\mathbf{w} \in \mathcal{W}^1} \left\{ \prod_{i|w_i=1} x_i \prod_{j|w_j=0} (1 - x_j) \right\}, \quad (9)$$

186 where  $\mathcal{W}^1 = \left\{ \mathbf{w} \in \{0, 1\}^R \mid f(\mathbf{w}) = 1 \right\}$ .

187 Examples of polynomial representation of simple Boolean functions are the following:

- 188 • Logical ‘OR’ with  $R = 2$ ,  $(x_1 \vee x_2)$ :

$$\bar{f}_{\text{OR}}(x_1, x_2) = (1 - x_1)x_2 + x_1(1 - x_2) + x_1x_2. \quad (10)$$

- 189 • Logical ‘XOR’ with  $R = 2$ ,  $(x_1 \oplus x_2)$ :

$$\bar{f}_{\text{XOR}}(x_1, x_2) = (1 - x_1)x_2 + x_1(1 - x_2). \quad (11)$$

- 190 • 3-term majority:

$$\begin{aligned} \bar{f}_{\text{MAJ}}(x_1, x_2, x_3) &= \mathbb{1}_{(\sum_i x_i) \geq 2}(x_1, x_2, x_3) = (1 - x_1)x_2x_3 + x_1(1 - x_2)x_3 \\ &\quad + x_1x_2(1 - x_3) + x_1x_2x_3. \end{aligned} \quad (12)$$

191 Observe that with this representation, the number of terms in the polynomial depends on the car-  
 192 dinal of  $\mathcal{W}^1$  (number of input combinations such that  $f(\mathbf{x}) = 1$ ). If the set  $\mathcal{W}^0 = \{\mathbf{w} \in \{0, 1\}^R \mid f(\mathbf{w}) = 0\}$   
 193 has a smaller cardinal than  $\mathcal{W}^1$ , then it may be more convenient to use another equivalent form of  
 194  $\bar{f}(\cdot)$ :

$$\bar{f}(\mathbf{x}) = 1 - \sum_{\mathbf{w} \in \mathcal{W}^0} \left\{ \prod_{i|w_i=1} x_i \prod_{j|w_j=0} (1 - x_j) \right\}. \quad (13)$$

195 Note that in the case of  $x_1 \vee x_2$ , the representation above leads to  $\bar{f}_{\text{OR}}(x_1, x_2) = 1 - (1 - x_1)(1 - x_2)$   
 196 and the corresponding  $R$ -term version of ‘OR’ has a simple expression

$$\bar{f}_{\text{OR}}(x_1, x_2, \dots, x_R) = 1 - \prod_{r=1}^R (1 - x_r) \quad (14)$$

197 since  $\mathcal{W}^0$  has only one element. In the rest of the paper, we use representation (13), since it allows  
 198 an easier presentation of the algorithms that we propose in the specific case of BMF, which is a  
 199 popular Boolean factorization in binary data analysis. Note that other Boolean functions for which  
 200 this representation may also be useful are  $x^{(1)} \vee \overline{x^{(2)}}$ , where  $\bar{x}$  denotes logical ‘NOT’,  $\overline{x^{(1)} \vee x^{(2)}}$  and  
 201  $\overline{x^{(1)}x^{(2)}}$  (NAND). In all these cases, the cardinal of  $\mathcal{W}^0$  is smaller than that of  $\mathcal{W}^1$ .

### 202 3. Algorithms

203 Two properties of  $\bar{f}(\cdot)$  are interesting from an optimization point of view: this function is differ-  
 204 entiable and it is multilinear in its inputs. Since  $\bar{f}(\cdot)$  is differentiable, gradient descent can be applied  
 205 to attempt solving the corresponding relaxed versions of (3). Multilinearity of  $\bar{f}(\cdot)$  with respect to  
 206 its inputs implies that this function is also multilinear in the columns of  $\mathbf{A}$  and  $\mathbf{B}$ . Therefore, if we  
 207 use a block-coordinate descent approach to attempt minimizing relaxed (3), and we set the blocks  
 208 of variables to be the columns of  $\mathbf{A}$  and  $\mathbf{B}$ , the block updates will be given by the solutions of  
 209 simple linear least squares problems. As a consequence, each of these properties leads to a different  
 210 algorithm for solving relaxed (3). These algorithms are detailed next.

211 *3.1. Gradient descent (GD) algorithm*

212 In the first algorithm, we consider a relaxed version of (3) where elements of the factors are  
 213 allowed to lie in  $\mathbb{R}$ . To force the solution to be close to binary, we introduce a penalty term  $\mathcal{G}(\mathbf{A}, \mathbf{B})$   
 214 in the objective function as in [16, 26]. The expression of this penalty term is

$$\mathcal{G}(\mathbf{A}, \mathbf{B}) = \frac{1}{2} \left\{ \sum_{i=1}^I \sum_{j=1}^J \sum_{r=1}^R [a_{ir}^2 (1 - a_{ir})^2 + b_{jr}^2 (1 - b_{jr})^2] \right\}. \quad (15)$$

215 Note that this penalty is minimal whenever all elements of the factors are ‘0’ or ‘1’. The new  
 216 optimization problem we have to solve is the following:

$$\begin{aligned} & \text{minimize} && \mathcal{H}(\mathbf{A}, \mathbf{B}; \lambda) = \mathcal{F}(\mathbf{A}, \mathbf{B}) + \lambda \mathcal{G}(\mathbf{A}, \mathbf{B}) \\ & \text{with respect to} && \mathbf{A} \in \mathbb{R}^{I \times R}, \mathbf{B} \in \mathbb{R}^{J \times R}, \end{aligned} \quad (16)$$

217 where  $\lambda > 0$  is a given value. Since this penalty term is differentiable,  $\mathcal{H}(\mathbf{A}, \mathbf{B}; \lambda)$  is differentiable.  
 218 Therefore, we can apply the standard gradient descent algorithm to find its critical points.

219 In standard gradient descent, the entries of the parameters vector  $\boldsymbol{\theta} = [\text{vec}(\mathbf{A})^\top \text{vec}(\mathbf{B})^\top]^\top$ ,  
 220 where  $\text{vec}(\mathbf{A}) = [\mathbf{a}_1^\top, \mathbf{a}_2^\top, \dots, \mathbf{a}_R^\top]^\top$  and  $\text{vec}(\mathbf{B}) = [\mathbf{b}_1^\top, \mathbf{b}_2^\top, \dots, \mathbf{b}_R^\top]^\top$ , are estimated jointly. The  
 221 estimate  $\hat{\boldsymbol{\theta}}_k$  of the parameter vector at iterate  $k$  is given by

$$\hat{\boldsymbol{\theta}}_k = \hat{\boldsymbol{\theta}}_{k-1} - \gamma_k \nabla_{\boldsymbol{\theta}} \mathcal{H}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_{k-1}}, \quad (17)$$

222 where  $\gamma_k$  is the step-size of the algorithm and  $\nabla_{\boldsymbol{\theta}} \mathcal{H}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_{k-1}}$  is the gradient vector of  $\mathcal{H}(\cdot)$  with  
 223 respect to all parameters  $\boldsymbol{\theta}$  evaluated at  $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}_{k-1}$ .

224 *Gradient expressions.* The full gradient vector can be written as a function of the partial gradients  
 225 with respect to  $\mathbf{A}$  and  $\mathbf{B}$  as follows

$$\nabla_{\boldsymbol{\theta}}^\top \mathcal{H}(\boldsymbol{\theta}) = \left[ \nabla_{\text{vec}(\mathbf{A})}^\top \mathcal{H}(\mathbf{A}, \mathbf{B}) \quad \nabla_{\text{vec}(\mathbf{B})}^\top \mathcal{H}(\mathbf{A}, \mathbf{B}) \right] \quad (18)$$

226 and the partial gradients are

$$\nabla_{\text{vec}(\mathbf{A})}^\top \mathcal{H}(\mathbf{A}, \mathbf{B}) = \left[ \frac{\partial \mathcal{H}(\mathbf{A}, \mathbf{B})}{\partial a_{11}} \quad \dots \quad \frac{\partial \mathcal{H}(\mathbf{A}, \mathbf{B})}{\partial a_{I1}} \quad \dots \quad \frac{\partial \mathcal{H}(\mathbf{A}, \mathbf{B})}{\partial a_{1R}} \quad \dots \quad \frac{\partial \mathcal{H}(\mathbf{A}, \mathbf{B})}{\partial a_{IR}} \right], \quad (19)$$

$$\nabla_{\text{vec}(\mathbf{B})}^\top \mathcal{H}(\mathbf{A}, \mathbf{B}) = \left[ \frac{\partial \mathcal{H}(\mathbf{A}, \mathbf{B})}{\partial b_{11}} \quad \dots \quad \frac{\partial \mathcal{H}(\mathbf{A}, \mathbf{B})}{\partial b_{J1}} \quad \dots \quad \frac{\partial \mathcal{H}(\mathbf{A}, \mathbf{B})}{\partial b_{1R}} \quad \dots \quad \frac{\partial \mathcal{H}(\mathbf{A}, \mathbf{B})}{\partial b_{JR}} \right]. \quad (20)$$

228 The elements of  $\nabla_{\text{vec}(\mathbf{A})} \mathcal{H}(\mathbf{A}, \mathbf{B})$  are given by

$$\frac{\partial \mathcal{H}(\mathbf{A}, \mathbf{B})}{\partial a_{i'r'}} = - \left\{ \sum_{j=1}^J [y_{i'j} - \bar{f}(x_{i'j}^{1:R})] \frac{\partial \bar{f}(x_{i'j}^{1:R})}{\partial x_{i'j}^{r'}} b_{jr'} \right\} + \lambda a_{i'r'} (1 - a_{i'r'}) (1 - 2a_{i'r'}) \quad (21)$$

for  $i' \in \{1, \dots, I\}$  and  $r' \in \{1, \dots, R\}$ , where the expression of the partial derivatives of  $\bar{f}(\cdot)$  are

$$\begin{aligned} \frac{\partial \bar{f}(x_{i'j}^{1:R})}{\partial x_{i'j}^{r'}} &= \sum_{\substack{\mathbf{w} \in \mathcal{W}^0, \\ w_{r'} = 0}} \left[ \prod_{s|w_s=1} x_{i'j}^s \right] \left[ \prod_{\substack{s'|w'_s=0, \\ s' \neq r'}} (1 - x_{i'j}^{s'}) \right] \\ &- \sum_{\substack{\mathbf{w} \in \mathcal{W}^0, \\ w_{r'} = 1}} \left[ \prod_{\substack{s|w_s=1, \\ s \neq r'}} x_{i'j}^s \right] \left[ \prod_{s'|w'_s=0} (1 - x_{i'j}^{s'}) \right]. \end{aligned} \quad (22)$$

229 The  $(j', r')$  element of  $\nabla_{\text{vec}(\mathbf{B})} \mathcal{H}(\mathbf{A}, \mathbf{B})$  for  $j' \in \{1, \dots, J\}$ ,  $r' \in \{1, \dots, R\}$  is

$$\frac{\partial \mathcal{H}(\mathbf{A}, \mathbf{B})}{\partial b_{j'r'}} = - \left\{ \sum_{i=1}^I [y_{ij'} - \bar{f}(x_{ij'}^{1:R})] \frac{\partial \bar{f}(x_{ij'}^{1:R})}{\partial x_{ij'}^{r'}} a_{i,r'} \right\} + \lambda b_{j'r'} (1 - b_{j'r'}) (1 - 2b_{j'r'}). \quad (23)$$

230 The partial gradients can be written in vector form as a function of  $\mathbf{A}$  and  $\mathbf{B}$  as follows

$$\begin{aligned} \nabla_{\text{vec}(\mathbf{A})} \mathcal{H}(\mathbf{A}, \mathbf{B}) &= -\text{Diag}(\mathbf{E} \boxtimes \mathbf{P}_1, \dots, \mathbf{E} \boxtimes \mathbf{P}_R) \text{vec}(\mathbf{B}) \\ &+ \lambda \text{vec}(\mathbf{A}) \boxtimes (\mathbf{1}_{IR \times 1} - \text{vec}(\mathbf{A})) \boxtimes (\mathbf{1}_{IR \times 1} - 2\text{vec}(\mathbf{A})), \\ \nabla_{\text{vec}(\mathbf{B})} \mathcal{H}(\mathbf{A}, \mathbf{B}) &= -\text{Diag}(\mathbf{E}^\top \boxtimes \mathbf{P}_1^\top, \dots, \mathbf{E}^\top \boxtimes \mathbf{P}_R^\top) \text{vec}(\mathbf{A}) \\ &+ \lambda \text{vec}(\mathbf{B}) \boxtimes (\mathbf{1}_{JR \times 1} - \text{vec}(\mathbf{B})) \boxtimes (\mathbf{1}_{JR \times 1} - 2\text{vec}(\mathbf{B})), \end{aligned} \quad (24)$$

231 where  $\mathbf{E}$  is the model error matrix

$$\mathbf{E} = \mathbf{Y} - \bar{f}(\mathbf{X}^{1:R}) = \mathbf{Y} - \mathbf{1}_{I \times J} + \sum_{\mathbf{w} \in \mathcal{W}^0} \left[ \begin{array}{c} \boxtimes \\ s|w_s=1 \end{array} \mathbf{X}^s \right] \boxtimes \left[ \begin{array}{c} \boxtimes \\ s'|w'_s=0 \end{array} (\mathbf{1}_{I \times J} - \mathbf{X}^{s'}) \right] \quad (25)$$

232 and  $\mathbf{P}_{r'}$  are  $I \times J$  matrices given by

$$\begin{aligned} \mathbf{P}_{r'} &= \sum_{\substack{\mathbf{w} \in \mathcal{W}^0, \\ w_{r'} = 0}} \left[ \begin{array}{c} \boxtimes \\ s|w_s=1 \end{array} \mathbf{X}^s \right] \boxtimes \left[ \begin{array}{c} \boxtimes \\ s'|w'_s=0, \\ s' \neq r' \end{array} (\mathbf{1}_{I \times J} - \mathbf{X}^{s'}) \right] \\ &- \sum_{\substack{\mathbf{w} \in \mathcal{W}^0, \\ w_{r'} = 1}} \left[ \begin{array}{c} \boxtimes \\ s|w_s=1, \\ s \neq r' \end{array} \mathbf{X}^s \right] \boxtimes \left[ \begin{array}{c} \boxtimes \\ s'|w'_s=0 \end{array} (\mathbf{1}_{I \times J} - \mathbf{X}^{s'}) \right]. \end{aligned} \quad (26)$$

233 *Step-size, penalty constant and initializations.* In the simplest version of the algorithm the step-size  
234  $\gamma_k$  can be set to a small constant value. The penalty factor  $\lambda$  may be chosen as variable through  
235 iterations:  $\lambda$  is set to a value close to zero in the first iterations and its increased up to a high target  
236 value.

237 Since the cost function being minimized is highly nonconvex, gradient descent may converge  
238 to spurious critical points. For this reason, it is important to test different initializations of the  
239 algorithm. For pairs of final estimates of the factors  $\hat{\mathbf{A}}$  and  $\hat{\mathbf{B}}$  and corresponding rank-1 estimates

240  $\hat{\mathbf{X}}^{1:R}$  obtained from different initializations, we select the pair which leads to the smallest data fitting  
 241 cost value  $\mathcal{F}(\hat{\mathbf{A}}, \hat{\mathbf{B}}) = \frac{1}{2} \|\mathbf{Y} - f(\hat{\mathbf{X}}^{1:R})\|_F^2$ . The algorithm can be initialized each time with random  
 242 elements for the factors. The elements of  $\hat{\mathbf{A}}$  and  $\hat{\mathbf{B}}$  can be drawn from independent and identically  
 243 distributed (iid) uniform samples:  $\hat{a}_{ir} \sim \mathcal{U}[0, 1]$ ,  $\hat{b}_{ir} \sim \mathcal{U}[0, 1]$ , for  $i \in \{1, \dots, I\}$ ,  $j \in \{1, \dots, J\}$   
 244 and  $r \in \{1, \dots, R\}$ .

### 245 3.2. Projected Hierarchical Alternating Least Squares (PHALS) algorithm

246 In the second approach, named projected hierarchical alternating least squares (PHALS), the  
 247 cost function is minimized with respect to each  $\mathbf{a}_1, \dots, \mathbf{a}_R, \mathbf{b}_1, \dots, \mathbf{b}_R$  in an alternating manner,  
 248 similar to the hierarchical alternating least squares (HALS) method [30]. The minimization with  
 249 respect to each column is performed by relaxing the binary constraints to  $\mathbb{R}^I, \mathbb{R}^J$ . After updating  
 250 all the estimates of a column of a factor, we project elements of the updated factor onto the interval  
 251  $[0, 1]$  to prevent the updates to converge to negative or large positive values.

252 Suppose that we want to update the estimate  $\hat{\mathbf{a}}_{r'}$  of  $\mathbf{a}_{r'}$ , all other columns of the factors are then  
 253 considered to be equal to  $\hat{\mathbf{a}}_r$  with  $r \neq r'$  and  $\hat{\mathbf{b}}_r$  for  $r \in \{1, \dots, R\}$ . The updated  $\hat{\mathbf{a}}_{r'}$  is then given  
 254 by the minimization of

$$\begin{aligned} \mathcal{F}_{\text{PHALS}}(\mathbf{a}_{r'}) &= \frac{1}{2} \sum_{i=1}^I \sum_{j=1}^J \left\{ y_{ij} - 1 + a_{ir'} \hat{b}_{jr'} \hat{p}_{ij}^{r'} + (1 - a_{ir'} \hat{b}_{jr'}) \hat{q}_{ij}^{r'} \right\}^2 \\ &= \frac{1}{2} \sum_{i=1}^I \sum_{j=1}^J \left[ y_{ij} - 1 + \hat{q}_{ij}^{r'} + a_{ir'} \hat{b}_{jr'} (\hat{p}_{ij}^{r'} - \hat{q}_{ij}^{r'}) \right]^2, \end{aligned} \quad (27)$$

255 where

$$\hat{p}_{ij}^{r'} = \sum_{\substack{\mathbf{w} \in \mathcal{W}^0, \\ w_{r'} = 1}} \left[ \prod_{\substack{s | w_s = 1 \\ s \neq r'}} \hat{a}_{is} \hat{b}_{js} \right] \left[ \prod_{\substack{s' | w_{s'} = 0 \\ s' \neq r'}} (1 - \hat{a}_{is} \hat{b}_{js}) \right] \quad (28)$$

256 and  $\hat{q}_{ij}^{r'}$  is similarly defined except that the summation is done through  $\mathbf{w} \in \mathcal{W}^0$  whose  $w_{r'} = 0$ .

257 The cost function  $\mathcal{F}_{\text{PHALS}}(\mathbf{a}_{r'})$  can be rewritten as  $\mathcal{F}_{\text{PHALS}}(\mathbf{a}_{r'}) = \sum_{i=1}^I \mathcal{F}_i(a_{ir'})$ , where, for a  
 258 given  $i'$ ,

$$\mathcal{F}_{i'}(a_{i'r'}) = \sum_{j=1}^J \left( y_{i'j} - 1 + \hat{q}_{i'j}^{r'} + a_{i'r'} \hat{b}_{jr'} (\hat{p}_{i'j}^{r'} - \hat{q}_{i'j}^{r'}) \right)^2. \quad (29)$$

259 Observe that each term  $\mathcal{F}_i(a_{ir'})$  of the cost function depends only on one of the  $a_{ir'}$ , thus the elements  
 260 of  $\hat{\mathbf{A}}_{r'}$  can be obtained separately by minimizing  $\mathcal{F}_{i'}(a_{i'r'})$ . The function  $\mathcal{F}_{i'}(a_{i'r'})$  is quadratic on

261  $a_{i'r'}$ , therefore its unconstrained minimum can be easily obtained. For a given  $i'$ , it is

$$\hat{a}_{i'r'} = \frac{\sum_{j=1}^J \left[ (y_{i'j} - 1 + \hat{q}_{i'j}^{r'}) (\hat{q}_{i'j}^{r'} - \hat{p}_{i'j}^{r'}) \hat{b}_{jr'} \right]}{\sum_{j=1}^J \left[ \hat{b}_{jr'} (\hat{q}_{i'j}^{r'} - \hat{p}_{i'j}^{r'}) \right]^2}. \quad (30)$$

262 Once  $\hat{\mathbf{a}}_{r'}$  has been completely updated, the arrays  $\hat{p}_{i,j}^{r'}$  and  $\hat{q}_{i,j}^{r'}$  have to be recalculated for the  
 263 update of the next  $\hat{\mathbf{a}}_{r'}$ . When all columns of  $\hat{\mathbf{A}}$  have been updated, the projection of its elements  
 264 onto  $[0, 1]$  is given by

$$\left[ \mathcal{P}_U(\hat{\mathbf{A}}) \right]_{ir} = \begin{cases} 0 & , \text{ for } \hat{a}_{ir} < 0, \\ \hat{a}_{ir} & , \text{ for } 0 \leq \hat{a}_{ir} \leq 1, \\ 1 & , \text{ for } \hat{a}_{ir} > 1. \end{cases} \quad (31)$$

265 A similar procedure is applied to the updates  $\hat{\mathbf{b}}_r$ . The updates before projection of its elements for  
 266  $j' \in \{1, 2, \dots, J\}$  are

$$\hat{b}_{j'r'} = \frac{\sum_{i=1}^I \left[ (y_{ij'} - 1 + \hat{q}_{ij'}^{r'}) (\hat{q}_{ij'}^{r'} - \hat{p}_{ij'}^{r'}) \hat{a}_{ir'} \right]}{\sum_{i=1}^I \left[ \hat{a}_{ir'} (\hat{q}_{ij'}^{r'} - \hat{p}_{ij'}^{r'}) \right]^2}. \quad (32)$$

267 After executing  $K$  updates of all columns of  $\hat{\mathbf{A}}$  and  $\hat{\mathbf{B}}$  with PHALS, the elements of  $\hat{\mathbf{A}}$  and  $\hat{\mathbf{B}}$   
 268 are projected on the set of binary values with  $\mathcal{P}_{\mathbb{B}}(\cdot)$  (8).

269 The full algorithm using PHALS to obtain an approximate generalized Boolean factorization is  
 270 given in Algorithm 1.

### 271 3.3. Algorithms for BMF

272 In the specific case of BMF, the most popular Boolean factorization used in practice, the expres-  
 273 sions of different quantities of the underlying algorithms can be easily written for any  $R$ . In this  
 274 subsection, we detail these expressions.

275 By relying on (14), BMF can be written in matrix form as follows:

$$\bar{f}(\mathbf{X}^{1:R}) = \bigvee_{r=1}^R \mathbf{X}^r = \mathbf{1}_{I \times J} - \boxed{\bullet} (\mathbf{1}_{I \times J} - \mathbf{X}^r) \quad (33)$$

276 For given  $\mathbf{X}^{1:R}$ , the reconstruction error  $\mathbf{E}$  can be written as

$$\mathbf{E} = \mathbf{Y} - \mathbf{1}_{I \times J} + \boxed{\bullet} (\mathbf{1}_{I \times J} - \mathbf{X}^r) = \mathbf{Y} - \mathbf{1}_{I \times J} + \boxed{\bullet} (\mathbf{1}_{I \times J} - \mathbf{a}_r \mathbf{b}_r^\top). \quad (34)$$

277 The  $\mathbf{P}_{r'}$  matrices (26) required in GD are

$$\mathbf{P}_{r'} = \boxed{\bullet} (\mathbf{1}_{I \times J} - \mathbf{X}^s) = \boxed{\bullet} (\mathbf{1}_{I \times J} - \mathbf{a}_s \mathbf{b}_s^\top). \quad (35)$$

$s = 1$   $s = 1$   
 $s \neq r'$   $s \neq r'$

---

**Algorithm 1** Projected hierarchical alternating least squares for general Boolean factorization (PHALS)

---

**Require:**  $Y, R, K$ .

- 1: Initialize  $\hat{\mathbf{A}}, \hat{\mathbf{B}}$  with random i.i.d. uniform elements  $\hat{a}_{ir} \sim \mathcal{U}[0, 1]$ ,  $\hat{b}_{ir} \sim \mathcal{U}[0, 1]$ , for  $i \in \{1, \dots, I\}$ ,  $j \in \{1, \dots, J\}$  and  $r \in \{1, \dots, R\}$ .
  - 2: **for**  $k \in \{1, \dots, K\}$  **do**
  - 3:   Update  $\hat{\mathbf{A}}$ :
  - 4:   **for**  $r \in \{1, \dots, R\}$  **do**
  - 5:     Update each column of  $\hat{\mathbf{A}}$ :
  - 6:     **for**  $i \in \{1, \dots, I\}$  **do**
  - 7:       Update elements of  $\hat{\mathbf{a}}_r$  (30):  $\hat{a}_{ir} = \frac{\sum_{j=1}^J [(y_{ij} - 1 + \hat{q}_{ij}^r)(\hat{q}_{ij}^r - \hat{p}_{ij}^r)\hat{b}_{jr}]}{\sum_{j=1}^J [\hat{b}_{jr}(\hat{q}_{ij}^r - \hat{p}_{ij}^r)]^2}$
  - 8:     **end for**
  - 9:     **for**  $r' \in \{1, \dots, R\}$  and  $r' \neq r$  **do**
  - 10:       Update  $\hat{p}_{ijr'}$  with (28) and  $\hat{q}_{ijr'}$  in a similar manner.
  - 11:     **end for**
  - 12:   **end for**
  - 13:   Project onto  $[0, 1]$  (31):  $\hat{\mathbf{A}} := \mathcal{P}_U(\hat{\mathbf{A}})$
  - 14:   Update  $\hat{\mathbf{B}}$ :
  - 15:   **for**  $r \in \{1, \dots, R\}$  **do**
  - 16:     Update each column of  $\hat{\mathbf{B}}$ :
  - 17:     **for**  $j \in \{1, \dots, J\}$  **do**
  - 18:       Update elements of  $\hat{\mathbf{b}}_r$  (32):  $\hat{b}_{jr} = \frac{\sum_{i=1}^I [(y_{ij} - 1 + \hat{q}_{ij}^r)(\hat{q}_{ij}^r - \hat{p}_{ij}^r)\hat{a}_{ir}]}{\sum_{i=1}^I [\hat{a}_{ir}(\hat{q}_{ij}^r - \hat{p}_{ij}^r)]^2}$
  - 19:     **end for**
  - 20:     **for**  $r' \in \{1, \dots, R\}$  and  $r' \neq r$  **do**
  - 21:       Update  $\hat{p}_{ijr'}$  and  $\hat{q}_{ijr'}$ .
  - 22:     **end for**
  - 23:   **end for**
  - 24:   Project onto  $[0, 1]$  (31):  $\hat{\mathbf{B}} := \mathcal{P}_U(\hat{\mathbf{B}})$
  - 25: **end for**
  - 26: Project onto  $\{0, 1\}$  (8):  $\hat{\mathbf{A}} := \mathcal{P}_{\mathbb{B}}(\hat{\mathbf{A}})$ ,  $\hat{\mathbf{B}} := \mathcal{P}_{\mathbb{B}}(\hat{\mathbf{B}})$
-

278 For PHALS, the quantities that vary depending on the choice of the Boolean function  $\bar{f}(\cdot)$  are  
 279  $\hat{p}_{ij}^{r'}$  and  $\hat{q}_{ij}^{r'}$ . For BMF, we have  $\hat{p}_{ij}^{r'} = 0$  for all possible tuples  $(ijr')$ , while  $\hat{q}_{ij}^{r'}$  can be written in  
 280 matrix form as  $\mathbf{P}_{r'}$  above for  $r' \in \{1, \dots, R\}$ :

$$\mathbf{Q}_{r'} = \mathbf{P}_{r'} = \begin{matrix} R \\ \boxed{\bullet} \\ s = 1 \\ s \neq r' \end{matrix} (\mathbf{1}_{I \times J} - \mathbf{a}_s \mathbf{b}_s^T). \quad (36)$$

## 281 4. Numerical experiments

282 In this section, we present the results of numerical experiments concerning the proposed algo-  
 283 rithms. We focus first on the BMF setting, that is, when the combining function  $f(\cdot)$  is the ‘OR’  
 284 function. Under this setting, we compare the performance of the algorithms with 3 other BMF  
 285 methods from the literature on simulated noisy binary data. In the first and second simulation  
 286 settings, the data are drawn from random BMF models which are then perturbed by binary flipping  
 287 noise. In the first setting, the performance of the algorithms is presented for different number of  
 288 columns  $R$  of  $\mathbf{A}$  and  $\mathbf{B}$ , and the probability of binary flipping the data (equivalent to noise intensity)  
 289 is kept constant. In the second setting,  $R$  is kept constant and results are shown for different values  
 290 of the probability of binary flipping. Simulation results will then be presented concerning the sen-  
 291 sitivity of the proposed methods to initialization, convergence behavior and time complexity. The  
 292 presentation of simulation results on BMF is then followed by its application to real datasets. We  
 293 apply PHALS to retrieve the BMF of the following datasets: the congressional voting dataset [37]<sup>1</sup>,  
 294 the zoo dataset [37]<sup>2</sup>, the New and Old Worlds (NOW) paleontological database [38] and the United  
 295 Nations voting dataset [39]<sup>3</sup>. We end the section by presenting performance results for PHALS and  
 296 GD for simulated data generated with  $f(\cdot)$  equal to the XOR with two inputs (XOR-2) and to the  
 297 3-term majority function (MAJ-3).

### 298 4.1. Performance for different $R$

299 In what follows, the performances of the two proposed algorithms PHALS and GD are compared  
 300 to 3 state-of-the-art methods for BMF discussed in the introduction: *ASSOCIATION rules* algorithm

---

<sup>1</sup><https://archive.ics.uci.edu/ml/datasets/Congressional+Voting+Records>

<sup>2</sup><https://archive.ics.uci.edu/ml/datasets/Zoo>

<sup>3</sup><https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/LEJUQZ>

301 (ASSO) [21], the *Formal Concept* (FC) analysis based algorithm (Algorithm 2 in [22]) and the *Post-*  
 302 *NonLinear Penalty Function* (PNL-PF) algorithm [26]. We have also included in the simulations a  
 303 version of GD that is initialized with PHALS, we will denote that version of GD as PHALS+GD.

304 In this set of simulations, the data matrix  $\mathbf{Y}$  is a  $20 \times 20$  matrix corresponding to a perturbed  
 305 version of a BMF with  $R$  components  $\mathbf{X} = \bigvee_{r=1}^R \mathbf{a}_r \mathbf{b}_r^\top$ . The noise matrix is binary  $\mathbf{N} \in \{0, 1\}$   
 306 and the perturbation consists in flipping the elements of  $\mathbf{X}$ . Therefore, the elements of  $\mathbf{Y}$  can be  
 307 written using the logical ‘XOR’:  $y_{ij} = x_{ij} \oplus n_{ij}$ . The elements of  $\mathbf{N}$  are drawn iid from a Bernoulli  
 308 distribution  $n_{ij} \sim \mathcal{B}(p_n)$  where  $p_n = \mathbb{P}(n_{ij} = 1)$ .

309 For a given data matrix  $\mathbf{Y}$ , PHALS and GD are initialized at random  $n_{\text{init}} = 3$  times and the  
 310 solution achieving the least reconstruction error  $\mathcal{F}(\hat{\mathbf{A}}, \hat{\mathbf{B}})$  is kept. For PHALS+GD, GD is initialized  
 311 with the best of the 3 initializations of PHALS. PNL-PF is initialized with the result of NMF [25]  
 312 applied to the data. The NMF algorithm is initialized randomly as PHALS and GD. ASSO and FC  
 313 do not require initializations.

314 GD and PNL-PF are executed with  $K_{\text{GD/PNL-PF}} = 2000$  iterations for each simulation, while  
 315 PHALS is executed  $K_{\text{PHALS}} = 1000$  iterations. The step-length of GD is set to a constant  $\gamma_k = \gamma =$   
 316  $0.1$  and its hyperparameter  $\lambda$  is increased linearly from  $0.01$  to  $10$ . The threshold value  $\tau$  for ASSO  
 317 (see [21]) is fixed to  $0.9$ .

318 Under each different simulation setting,  $N = 100$  random data matrices  $\mathbf{Y}^n$  ( $n = 1, \dots, N$ ) are  
 319 generated, each with a random pair  $(\mathbf{X}^n, \mathbf{N}^n)$ . Each  $\mathbf{X}^n$  with a given  $R$  is obtained from randomly  
 320 generated factor matrices  $\mathbf{A}^n$  and  $\mathbf{B}^n$ . The elements of  $\mathbf{A}^n$  and  $\mathbf{B}^n$  are iid samples from Bernoulli  
 321 distributions  $a_{ir} \sim \mathcal{B}(p_a)$ ,  $b_{ir} \sim \mathcal{B}(p_b)$  where  $p_a = p_b = 0.4$ . After applying the algorithms on all  
 322  $\mathbf{Y}^n$ , their performances in terms of normalized mean square errors (NMSE) of prediction of  $\mathbf{X}^n$  and  
 323 retrieval of  $\mathbf{A}^n$  and  $\mathbf{B}^n$  are evaluated. The expressions of these NMSE are the following:  $\text{NMSE}_{\mathbf{X}} =$   
 324  $\frac{1}{NIJ} \sum_{n=1}^N \left\| \mathbf{X}^n - \bigvee_{r=1}^R \hat{\mathbf{a}}_r^n (\hat{\mathbf{b}}_r^n)^\top \right\|_F^2$ ,  $\text{NMSE}_{\mathbf{A}} = \frac{1}{NIR} \sum_{n=1}^N \left\| \mathbf{A}^n - \hat{\mathbf{A}}^n \right\|_F^2$ ,  $\text{NMSE}_{\mathbf{B}} = \frac{1}{NJR} \sum_{n=1}^N \left\| \mathbf{B}^n - \hat{\mathbf{B}}^n \right\|_F^2$ .  
 325 Since the elements of all matrices are binary, these NMSE can be interpreted as error rates. Note  
 326 that due to the permutation ambiguity on the estimation of the factors, their columns should be  
 327 permuted to match in the best possible way those of the true factors before the evaluation of the  
 328 NMSE.

329 The evolution of the NMSE for the BMF methods is shown in Fig. 1a for  $p_n = 0.1$  and  $R \in$   
 330  $\{2, \dots, 6\}$ . As intuitively expected, for most methods, larger values of  $R$  lead to larger NMSE both  
 331 on prediction of  $\mathbf{X}$  and on the retrieved factors. One can also observe that there is no significant  
 332 difference in performance between the PHALS, GD and PHALS+GD. PNL-PF has a moderately

333 inferior performance compared to the proposed methods, while ASSO and FC achieve a significantly  
 334 inferior performance in terms of retrieving  $\mathbf{X}$ ,  $\mathbf{A}$  and  $\mathbf{B}$ . ASSO and FC do not seem to be adapted  
 335 to the approximation setting where noise is present, which has been also observed in previous studies  
 336 [26].

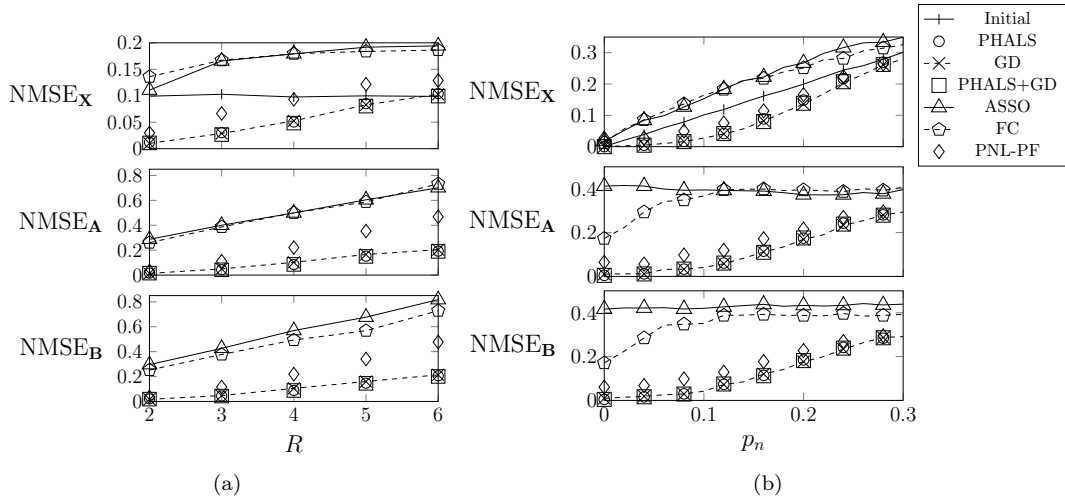


Figure 1: NMSE for the prediction of  $\mathbf{X}$  of size  $20 \times 20$  and retrieval of its BMF factors  $\mathbf{A}$  and  $\mathbf{B}$ . The factorization algorithms are executed on a noisy version  $\mathbf{Y}$  of  $\mathbf{X}$ . The noise acts by flipping the elements of  $\mathbf{Y}$  with probability  $p_n$ . Curves named “Initial” in the top subfigures indicate the NMSE of predicting  $\mathbf{X}$  using  $\mathbf{Y}$ . In (a),  $R$  varies from 2 to 6, while  $p_n = 0.1$ . In (b),  $R = 3$  and  $p_n$  is varied from 0 to 0.3 with increments of 0.02.

#### 337 4.2. Results for different $p_n$

338 The second simulation setting is very similar to the previously presented one, except that, in  
 339 this case,  $R$  is kept to a constant value,  $R = 3$ , and  $p_n$  is varied from 0 to 0.3 by increments of 0.02.  
 340 Fig. 1b displays the performances of the methods. One can see that, as it is naturally expected, the  
 341 performances degrade as  $p_n$  increases. As in the previous setting, all the proposed methods seem  
 342 to lead to similar performances and they achieve a superior performance compared to the other 3  
 343 methods from the literature. Observe also that as  $p_n$  gets close to 0.3 all NMSE of the proposed  
 344 methods are close to 0.3 and for values smaller than  $p_n < 0.3$  the NMSE seem to be smaller than  $p_n$ .  
 345 Note, on the top figure, that predicting  $\mathbf{X}$  using the noisy data is as efficient as using the proposed  
 346 methods for  $p_n = 0.3$ .

347 Finally, one can see that, when  $p_n = 0$ , the NMSE on the factors is not zero. This may be due  
 348 to convergence of the algorithms to factorizations which are not global minima of (16) or to the non  
 349 uniqueness of the approximation of some realizations  $\mathbf{Y}^n$ . However, since NMSE are very small for  
 350  $p_n = 0$ , it seems that the occurrence of such issues is very rare.

	$\mathbf{Y}_1$	$\mathbf{Y}_4$	$\mathbf{Y}_5$
PHALS	99	100	96
GD	79	100	48
PHALS+GD	99	100	96
PNL-PF	100	66	72

Table 2: Success rate  $S_{\%}$  for retrieving the exact BMF for 3 different matrices  $\mathbf{Y}_1$ ,  $\mathbf{Y}_4$  and  $\mathbf{Y}_5$  using  $n_{\text{init}} = 100$  different random initializations.

351 *4.3. Simple example with unique decomposition*

352 From the previous simulation results, it seems that ASSO and FC are not adequate in the BMF  
353 approximation setting. Therefore, in what follows, we focus on comparing only PHALS, GD and  
354 PNL-PF.

355 As previously presented, the algorithms may converge to wrong  $\hat{\mathbf{A}}$  and  $\hat{\mathbf{B}}$ , even in the exact  
356 factorization setting ( $p_n = 0$ ) when the underlying factorization is unique. Due to the non convexity  
357 of the underlying cost functions, not all initializations  $\hat{\mathbf{A}}$  and  $\hat{\mathbf{B}}$  lead to the original factors. To try  
358 to assess to which extent the algorithms are prone to this behavior, we have tested them on 3 small  
359 exact factorization problems with unique factorizations. The 3 considered data matrices are  $\mathbf{Y}_1$  (5),  
360  $\mathbf{Y}_4$  (6) and the following  $5 \times 5$  matrix from [26]:

$$\mathbf{Y}_5 = \begin{bmatrix} 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 0 & 1 \end{bmatrix}.$$

361 The matrices  $\mathbf{Y}_1$  and  $\mathbf{Y}_4$  have rank 2, while  $\mathbf{Y}_3$  has rank 3. PHALS, GD and PNL-PF are applied  
362 to these data with  $n_{\text{init}} = 100$  random initializations  $\hat{\mathbf{A}}_0^i, \hat{\mathbf{B}}_0^i, i \in \{1, \dots, n_{\text{init}}\}$ . The number of  
363 iterations for all algorithms is  $K = 2000$ . Parameters  $\gamma_k$  and  $\lambda$  of GD have been set as in the previous  
364 simulations. For each algorithm and each data matrix we have calculated the success rate  $S_{\%}$  in  
365 percent of retrieving  $\mathbf{A}$  and  $\mathbf{B}$  from data:  $S_{\%} = \left[ \text{card} \left( \left\{ i \mid \hat{\mathbf{A}}^i = \mathbf{A} \text{ and } \hat{\mathbf{B}}^i = \mathbf{B} \right\} \right) / n_{\text{init}} \right] \times 100$ ,  
366 where  $\text{card}(\cdot)$  denotes the cardinal of a set and  $\hat{\mathbf{A}}^i, \hat{\mathbf{B}}^i$  are the output factors of an algorithm when  
367 initialized with  $\hat{\mathbf{A}}_0^i$  and  $\hat{\mathbf{B}}_0^i$ . The success rates are displayed in Tab. 2. We observe that the algorithm  
368 which seem less prone to converge to spurious factors is PHALS. GD and PNL-PF may converge to  
369 spurious factors, but they do not seem to behave equally through the examples. From the results,  
370 we can also note that applying GD initialized with the resulting factors from PHALS does not lead  
371 to an improvement in the success rate.

372 *4.4. Convergence*

373 To compare the convergence behavior of PHALS, GD and PNL-PF, we generate 3 random  $\mathbf{Y}^k$   
 374 in the same manner as in Subsec. 4.1 for  $R$  equal to 2, 4 and 6. We then apply the 3 algorithms  
 375 with  $n_{\text{init}} = 100$  and  $K = 500$ . At each iteration  $k \in \{1, \dots, K\}$ , we evaluate the overall changes in  
 376 the factors using the following quantity:  $\Delta_k = \frac{\|\hat{\mathbf{A}}_k - \hat{\mathbf{A}}_{k-1}\|_F^2 + \|\hat{\mathbf{B}}_k - \hat{\mathbf{B}}_{k-1}\|_F^2}{\|\hat{\mathbf{A}}_0\|_F^2 + \|\hat{\mathbf{B}}_0\|_F^2}$ , where  $\hat{\mathbf{A}}_k$  and  $\hat{\mathbf{B}}_k$  are the  
 377  $k$ -th updates of the factors for a given algorithm. The quantity  $\Delta_k$  is small whenever the factors  
 378 do not change in consecutive iterations. Therefore, if  $\Delta_k$  reduces as  $k$  increases, the algorithm is  
 379 converging. Fig. 2 shows some statistics on  $\Delta_k$  for each algorithm. The statistics displayed are  
 380 the median, the 5-th and 95-th percentiles of  $\Delta_k$  evaluated with the  $n_{\text{init}} = 100$  available values for  
 381 each  $k$ . The overall behavior we can observe from this figure is that PHALS is the fastest method  
 382 in terms of convergence, while GD is the slowest. One can also note that the 95-th percentile of  
 383  $\Delta_k$  increases as  $R$  increases. Although PHALS is much faster than the other methods to converge,  
 384 when  $R$  increases, some initializations may lead it to converge slowly or not to converge at all.

385 *Remarks on convergence guarantees:* in their present form, we are not able to give theoretical  
 386 guarantees on convergence of the iterates of GD and PHALS.

387 Concerning GD, as presented above, it seems that in practice the algorithm converges, if the  
 388 constant step-size  $\gamma_k = \gamma$  is chosen sufficiently small. However, theoretical guarantees for convergence  
 389 of GD require that a global Lipschitz constant of the gradient of the objective function exists.  
 390 Unfortunately, this does not seem to be true for the objective in (16). A possible way to ensure  
 391 convergence is to use an adaptive step-length  $\gamma_k$  given by backtracking line-search [40]. With this  
 392 modification, since the cost function is analytic and coercive, convergence of GD is guaranteed using  
 393 the results from [41].

394 PHALS is a block coordinate descent algorithm. For this class of algorithms, convergence of the  
 395 iterates can be ensured, for example, using the results of [42]. To use the results of [42], the objective  
 396 function should be separately strongly convex in each block of variables  $\mathbf{a}_1, \dots, \mathbf{a}_R, \mathbf{b}_1, \dots, \mathbf{b}_R$ .  
 397 Unfortunately, this cannot be guaranteed, and, in practice, one can see that for difficult factorization  
 398 cases (large  $R$ ), some initializations may lead to non-converging iterates. One possibility to solve  
 399 this issue is to add proximal terms to the objective function at each update. For the updates of  $\mathbf{a}_r$ ,  
 400 one should add  $\rho \|\mathbf{a}_r - \hat{\mathbf{a}}_r\|_2^2$ , where  $\rho > 0$  is a pre-defined constant and  $\hat{\mathbf{a}}_r$  is the most recent update  
 401 of  $\mathbf{a}_r$ . Similarly, for the updates  $\mathbf{b}_r$ , the term  $\rho \|\mathbf{b}_r - \hat{\mathbf{b}}_r\|_2^2$  should be added. With modified updates  
 402 considering this additional terms, it may be possible to use the results of [42] to ensure convergence  
 403 of the iterates.

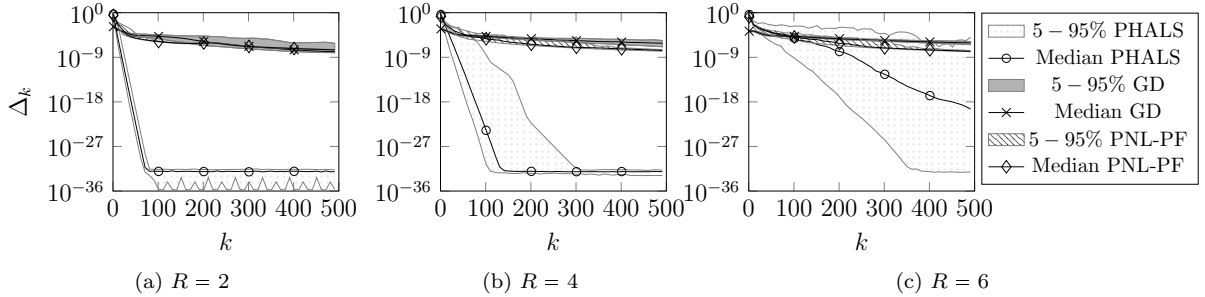


Figure 2: Statistics about the overall changes  $\Delta_k$  (4.4) in the updates  $\hat{\mathbf{A}}_k, \hat{\mathbf{B}}_k$  for GD, PHALS and PNL-PF for the approximate factorization of 3 noisy binary matrices of size  $20 \times 20$ . Each matrix is generated by randomly drawing a binary BMF model with a given  $R$ , then applying binary flipping noise. The subfigures are generated with different values of  $R$ , they are indicated in the subcaptions. The statistics are evaluated for  $n_{\text{init}} = 100$  different random initializations. The extremes of the bands (5% – 95%) around the median are the 5-th and 95-th percentiles of  $\Delta_k$ .

#### 404 4.5. Time complexity

405 Using a similar simulation setting from the previous subsection, we have also measured the  
 406 execution time for the 3 algorithms to finish 2000 iterations. The statistics on execution time<sup>4</sup> for  
 407 100 runs of the algorithms and for  $R \in \{2, \dots, 6\}$  are shown in Tab. 3. We observe that PNL-PF  
 408 takes much less time than the other methods. GD is from 10 to 40 times slower than PNL-PF and  
 409 we can clearly see a large relative increase in execution time for GD when passing from  $R = 4$  to  
 410  $R = 5$ . Such an increase can also be observed in a lesser extent in PNL-PF, while in PHALS the  
 411 relative increase is smaller than a factor of 2.

412 To test the complexity of the algorithms with larger matrices, we have also executed Monte-Carlo  
 413 simulations with matrices of sizes  $100 \times 100$ . The simulations have been carried out in a similar  
 414 setting to that of Subsec. 4.2, with the exception that we have stopped the algorithms once they  
 415 have reached  $\Delta_k = 10^{-6}$ . We have measured the total times for convergence of PNL-PF, GD and  
 416 PHALS. For the noiseless case,  $p_n = 0$ , their respective median total times in seconds are 3.750,  
 417 113.813 and 6.250. These timings remain of the same order when we increase  $p_n$  up to  $p_n = 0.2$ . For  
 418 larger  $p_n$ , we observe an overall increase in the timings. For  $p_n = 0.3$ , the timings are the following  
 419 10.375 sec. for PNL-PF, 157.375 sec. for GD and 734.828 sec. for PHALS. Clearly, for small and  
 420 moderate  $p_n$ , GD is much slower than the other algorithms, while PNL-PF and PHALS require

<sup>4</sup>These simulations are realized in *Scilab* version 6.1.0 with a processor Intel®Core™ i7-7820HQ, 2.90GHz and with 16GB of RAM.

		$R = 2$	3	4	5	6
PHALS	5-th percentile	0.844	1.547	2.188	3.0313	4.281
	Median	0.703	1.297	2.000	2.906	3.938
	95-th percentile	0.610	1.219	1.813	2.781	3.750
GD	5-th percentile	1.734	2.734	3.438	24.438	32.688
	Median	1.547	2.359	3.281	23.422	30.500
	95-th percentile	1.406	2.219	3.141	22.313	29.125
PNL-PF	5-th percentile	0.203	0.250	0.234	1.000	0.984
	Median	0.156	0.172	0.172	0.672	0.719
	95-th percentile	0.094	0.125	0.125	0.422	0.469

Table 3: Statistics on total execution times in seconds for approximate  $R$ -component BMF of  $20 \times 20$  binary matrices. The statistics are evaluated for 100 runs of the 3 algorithms with  $K = 2000$  iterations in each run.

421 about the same order of time to achieve  $\Delta_k = 10^{-6}$ . For large  $p_n$ , PHALS seems to suffer from  
422 convergence issues leading to a substantial increase in total execution time for convergence.

#### 423 4.6. Discussion on the results

424 In terms of factorization performance on simulated datasets, GD is often equivalent to PHALS.  
425 Unfortunately, GD is too slow to be applied to real datasets since their sizes may be much larger  
426 than those considered in the previously presented simulations. If we compare PNL-PF and PHALS,  
427 PHALS achieves superior factorization performance at the expense of a longer execution time per  
428 iteration and of a risk of producing slowly or non converging iterations in difficult cases (in our  
429 simulations mainly for large  $R$  or large  $p_n$ ). If we restrict the comparison to easier cases ( $R \leq 6$ ,  
430  $p_n \leq 0.2$ ), then PHALS converges much faster than PNL-PF. If a threshold on  $\Delta_k$  is used as  
431 convergence criterion to stop the algorithms, then the longer execution times of PHALS iterations  
432 are compensated by its much faster convergence. This leads to total executions times of PNL-PF  
433 and PHALS of the same order. To illustrate this behavior in a more realistic scenario, we have  
434 measured the execution times to retrieve the BMF with  $R = 4$  of the Paleontological dataset that  
435 will be presented in more details in Subsec. 4.7. This dataset is much larger than the previously  
436 simulated data ( $I = 254$  and  $J = 1375$ ). The timings for convergence to  $\Delta_k = 10^{-6}$  are of the order  
437 of minutes: 875.73 sec. for PNL-PF and 846.83 sec. for PHALS.

#### 438 4.7. Real datasets

439 In what follows, we obtain the approximate BMF of different real binary datasets. From the  
440 previous results on factorization performance obtained through simulations, factorization results  
441 are expected to be mostly similar for PHALS, GD and PNL-PF. Therefore, we have applied only

442 PHALS to factorize the real datasets. For each of the datasets,  $n_{\text{init}} = 10$  random initializations are  
443 used and the best solution in terms of data reconstruction error is selected. The maximum allowed  
444 number of iterations is set to  $K_{\text{PHALS}} = 2000$  and a convergence criterion based on  $\Delta_k$  is used as  
445 an additional stopping criterion.

446 *US Congressional voting dataset.* We first apply PHALS to a dataset containing 16 key votes of the  
447 United States congress for the year 1984 [37]<sup>5</sup>. The votes of 435 representatives are coded by binary  
448 values: ‘1’ for a vote in favor of the proposed bill and ‘0’ for a vote against it. Missing votes in a  
449 given bill have been replaced by the corresponding majority vote. This allows to fully encode the  
450 dataset into a binary matrix of size  $435 \times 16$ . The dataset also contains information on the party  
451 of each representative (democrat or republican). Since there are 2 parties, PHALS is applied to the  
452 dataset with  $R = 2$ . The dataset plot and an illustration of the results are given in Fig. 3 (a-d).  
453 One can clearly observe that the patterns related to each component have almost disjoint support,  
454 indicating an opposing voting pattern for each component. The error rate on the reconstructed  
455 data using the retrieved BMF model is of 20%. By comparing the grouping of the representatives  
456 encoded by matrix  $\hat{\mathbf{A}}$  with the information on the parties of each candidate, we found that PHALS  
457 can predict the party of the representative with an accuracy of 77%.

458 *Zoo dataset.* We also applied PHALS to a dataset containing the information on the presence or  
459 absence of a given feature, for example *hair*, *feathers*, *milk*, for different animals. The dataset [37]<sup>6</sup>  
460 contains 15 binary features and an integer feature with the number of legs. These features are given  
461 for 101 animals. The animals in the dataset are categorized in 7 classes: mammals, birds, reptiles,  
462 fishes, amphibians, insects and a class containing many different invertebrate animals (*e.g.* *crab*,  
463 *worm*, *octopus*). We have encoded the integer variable corresponding to the number of legs using  
464 one-hot encoding. We apply PHALS to the resulting  $101 \times 19$  binary matrix to try to group the  
465 animals by looking at the patterns given by the columns of  $\hat{\mathbf{A}}$ . The algorithm is applied with  $R$  in  
466 the range 2 to 7 and the data reconstruction error for these values of  $R$  are respectively 0.170, 0.116,  
467 0.0928, 0.0771, 0.0693 and 0.0620. One can clearly see that beyond  $R = 3$  the improvement on data  
468 reconstruction obtained by increasing  $R$  is mild. This result seems to be similar to what has been  
469 presented in [26] for the analysis of the same dataset with PNL-PF.

470 The dataset and an illustration of the results for  $R = 3$  are given in Fig. 3 (e-i). To simplify

---

<sup>5</sup><https://archive.ics.uci.edu/ml/datasets/Congressional+Voting+Records>

<sup>6</sup><https://archive.ics.uci.edu/ml/datasets/Zoo>

471 the interpretation of the results, the rows of the matrices, which correspond to different animals,  
472 have been reordered to correspond to continuous blocks of animals of the same category. Reordering  
473 has been carried out using the same order of the classes mentioned above, thus the first block of  
474 animals correspond to mammals, the second to birds, *etc.* One can observe that component 1 clearly  
475 corresponds to a continuous block of animals, in these case mammals. The second component mostly  
476 group together birds with two exceptions, *fruitbat* and *vampire*, which are also present in the group  
477 of mammals. The third group contains mostly fishes, but also some mammals (*e.g. dolphin*) and  
478 birds (*e.g. penguin*). Many insects and animals from the last category of invertebrate animals are  
479 not contained in any components. As the number of components is increased to  $R = 7$ , it has been  
480 observed that the retrieved BMF is not able to clearly separate the 7 underlying categories.

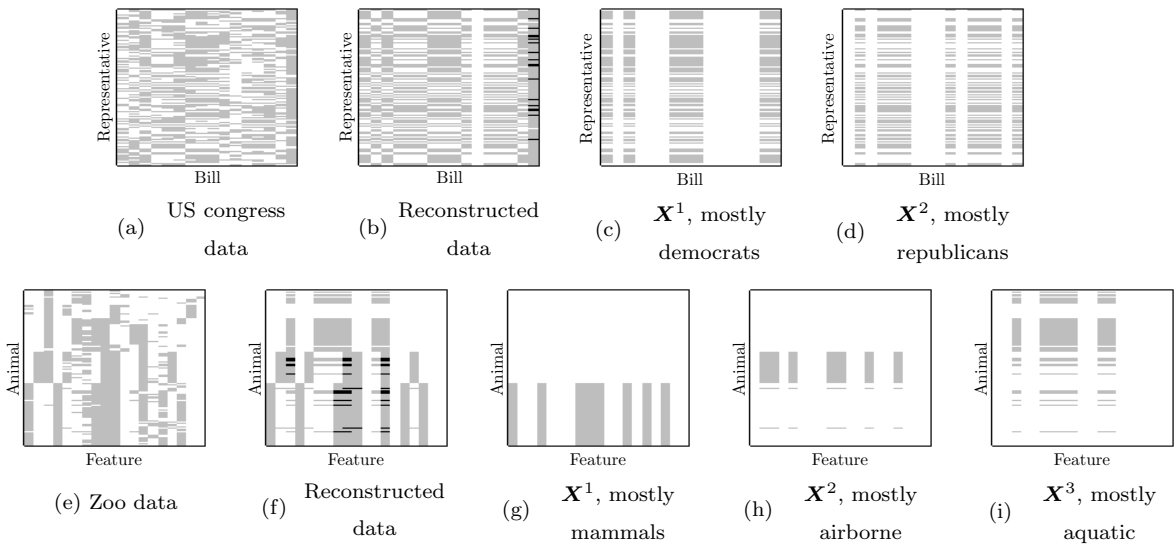


Figure 3: Real datasets and the results obtained with PHALS. In (a) and (e) the US congress voting dataset and the Zoo dataset are displayed. The gray color corresponds to a ‘1’ in the underlying dataset matrix, while white color corresponds to ‘0’. In (b) and (f), the reconstructed data using the BMF model are displayed. The black color indicates intersections between BMF components’ supports. The number of BMF components are  $R = 2$  and  $R = 3$  respectively. The rank-1 components  $\mathbf{X}^r$  retrieved with PHALS are displayed on the right of the reconstructed data in (c), (d), (g), (h) and (i).

481 *Paleontological dataset.* Following closely [15, 43], we analyze data containing information on the  
482 localization of fossil mammals [38]. The objective is to apply PHALS to factorize a binary data  
483 matrix where the rows correspond to different genera of fossil mammals and the columns correspond  
484 to the different localities where they have been found. The data obtained from [38] is preprocessed

485 in a similar manner as in [15, 43]. Fossils of small mammals are excluded from the dataset and  
 486 only those retrieved in Europe are considered. We also removed genera which are too infrequent  
 487 (less than 10 occurrences) and localities where only 1 genera has been found. As a result of this  
 488 preprocessing, a binary matrix of size 254 (genera)  $\times$  1375 (localities) is obtained, where a ‘1’ stands  
 489 for the occurrence of at least one fossil of a given genus in a given locality and a ‘0’ for its absence.  
 490 We have applied PHALS to this dataset to see if the resulting BMF allows to find communities of  
 491 mammals that appear in similar localities. The algorithm has been applied with  $R$  in the range  
 492 2 – 7. The algorithm seems to suffer from convergence issues for  $R > 4$ , generating factor matrices  
 493 with spurious empty columns. To validate the results for  $R \leq 4$ , we have followed [43] and plotted  
 494 the values of a variable related to the minimum age in millions of years of the localities where the  
 495 different genera have been found. We have observed that as  $R$  increases PHALS finds groups of  
 496 fossils with increasing minimum age. The ages of the genera in the different groups for  $R = 4$  are  
 497 displayed in Fig. 4. The genera in each of the components of this figure are the following:

- 498 • Component 1: *Amphiperatherium*, *Amphitragulus*, *Andegameryx*, *Brachyodus*, *Cainotherium*, *Cyne-*  
 499 *los*, *Diaceratherium*, *Palaeogale*, *Protaceratherium*.
- 500 • Component 2: *Amphicyon*, *Anchitherium*, *Anisodon*, *Aureliachoerus*, *Brachypotherium*, *Bunolistri-*  
 501 *odon*, *Dicrocerus*, *Dorcatherium*, *Gomphotherium*, *Hemicyon*, *Hyotherium*, *Lagomeryx*, *Lartetotherium*,  
 502 *Listriodon*, *Martes*, *Micromeryx*, *Palaeomeryx*, *Plesiaceratherium*, *Procervulus*, *Prodeinotherium*, *Prosan-*  
 503 *torhinus*, *Pseudaelurus*, *Styriofelis*, *Taucanamo*.
- 504 • Component 3: *Adcrocuta*, *Choerolophodon*, *Cremohipparion*, *Deinotherium*, *Dihoplus*, *Gazella*, *Hel-*  
 505 *ladotherium*, *Hipparion*, *Hippopotamodon*, *Hippotherium*, *Hyaenictitherium*, *Miotragocerus*, *Palaeotra-*  
 506 *gus*, *Pliodicerus*, *Tragoportax*.
- 507 • Component 4: *Bison*, *Canis*, *Cervus*, *Equus*, *Lynx*, *Mammuthus*, *Panthera*, *Stephanorhinus*, *Sus*,  
 508 *Ursus*, *Vulpes*.

509 By inspecting the median minimum ages of these groups, it seems that PHALS is able to retrieve  
 510 animal communities that have lived in different ages.

511 *UN voting dataset.* As a last example of application, we analyze the grouping of countries produced  
 512 by PHALS when used to factorize a binary matrix generated from the UN voting dataset [39]<sup>7</sup>. We  
 513 follow a similar setting as considered in [15] and we encode in a binary matrix the votes during the  
 514 cold-war period (1946 – 1990) of different countries for different UN roll-calls. We have removed

---

<sup>7</sup><https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/LEJUQZ>

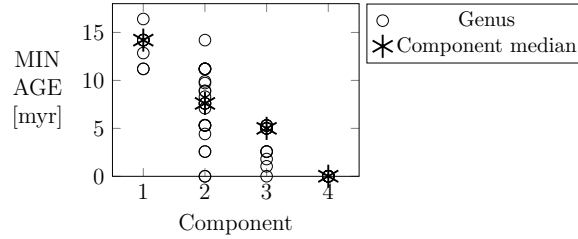


Figure 4: Minimum age in millions of years [myr] related to the genera of the different groups obtained by applying PHALS with  $R = 4$  to the paleontology dataset [38].

515 from the dataset all roll-calls whose number of unknown votes is larger than 98 (half of the listed  
516 countries) and also all the roll-calls with unanimous results. The unknown votes in the remaining  
517 roll-calls have been replaced by the majority vote. For this dataset, encoding with ‘1’ votes in favor  
518 of a roll-call leads to a very dense data matrix, whose BMF is difficult to retrieve and interpret.  
519 Therefore, to have a more sparse data matrix, we have encoded with ‘1’ votes against a roll-call and  
520 with ‘0’, votes in favor of it. PHALS has been applied to this dataset with  $R$  in the range 2 – 7. The  
521 data reconstruction error is respectively 0.0312, 0.0241, 0.0215, 0.0194, 0.0173 and 0.0158. Although  
522 a large part of the approximation improvement is observed when increasing  $R$  from 2 to 3, when we  
523 analyze the groups of countries produced for each  $R$ , interesting results seem to appear up to  $R = 6$ .

524 When  $R = 2$ , we can find a component containing the following countries: Australia, Belgium,  
525 Canada, Denmark, France, West Germany, Iceland, Israel, Italy, Japan, Luxembourg, Netherlands,  
526 New Zealand, Norway, Portugal, UK, US. The second component contains 173 countries from dif-  
527 ferent continents. If  $R$  is increased to 3, the first 2 components are similar to those obtained with  
528  $R = 2$  and a third component groups countries from the socialist block: Belarus, Bulgaria, Cuba,  
529 Czechoslovakia, East Germany, Hungary, Mongolia, Poland, Russia, Ukraine. When increasing  $R$   
530 to 4, similar results are obtained and a component with US and Israel appears. While for  $R = 5$ ,  
531 2 components with countries from the occidental block of countries are produced. For  $R = 6$ , the  
532 component containing a large number of countries seems to contain much less countries than for  
533 smaller  $R$  and a sixth component containing 36 countries from different continents appear. This  
534 last component gathers countries from the previously obtained component with a large number of  
535 countries but also countries from the occidental block (*e.g.* France) and from the socialist block (*e.g.*  
536 Cuba). Finally, for  $R = 7$ , the algorithm start finding components containing single countries (*e.g.*  
537 a component with only US).

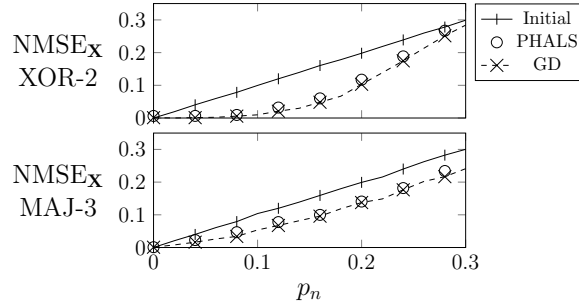


Figure 5: NMSE for the prediction of  $\mathbf{X}$  of size  $20 \times 20$  for generalized Boolean factorizations using the XOR – 2 and MAJ – 3 component combining functions. The decomposition algorithms are executed on a noisy version  $\mathbf{Y}$  of  $\mathbf{X}$ . The probability  $p_n$  of the binary noise that flips the elements of  $\mathbf{X}$  is varied from 0 to 0.3 with increments of 0.02.

#### 538 4.8. Results for XOR – 2 and MAJ – 3

539 The last experimental results concern the application of the GD and PHALS to a factorization  
540 setting different from BMF. We consider two other Boolean combining functions, the logical ‘XOR’  
541 with two inputs  $x_1 \oplus x_2$  (XOR – 2) and 3-term majority  $\mathbb{1}_{(\sum_i x_i) \geq 2}(x_1, x_2, x_3)$  (MAJ – 3). Since the  
542 uniqueness properties of these factorizations are still very little understood, we only focus on testing  
543 the methods for data denoising. We consider a simulation setting similar to the one presented in  
544 Subsection 4.2, the main differences are that the underlying (clean) data are generated with the  
545 XOR-2 and MAJ-3 functions, the maximum allowed number of iterations of the algorithms is set to  
546  $n_{\text{init}} = 5000$  and that an additional stopping criterion based  $\Delta_k$  is used for ending the iterations.  
547 The results for the NMSE of data reconstruction are displayed in Fig. 5. One can observe that  
548 the algorithms denoise the data, since their NMSE is smaller than the NMSE for the noisy data  
549 (curve named *Initial* in the plot). In both cases, GD seems to achieve a slightly superior denoising  
550 performance than PHALS. It is also possible to observe that the factorizations do not seem to have  
551 the same robustness behavior against noise. The denoising performance for F2MF (XOR-2) for small  
552 noise intensity seems much superior than for the MAJ-3 factorization. This is intuitively expected,  
553 since the MAJ-3 factorization requires the estimation of more parameters for the same amount of  
554 data. For large noise intensities ( $p_n \approx 0.3$ ), the opposite behavior is observed, with the MAJ-3  
555 factorization leading to a superior denoising performance.

#### 556 5. Conclusions and further work

557 In this paper, we have introduced a generalized framework for the Boolean factorization of binary  
558 matrices, where the “sum” between the rank-1 binary terms can be an arbitrary Boolean function.

559 We proposed two iterative algorithms for achieving this factorization, based on gradient descent  
560 (GD) and on projected hierarchical alternating least squares (PHALS) approaches, respectively.  
561 Implementation details for the algorithms have been presented for BMF and compared through  
562 numerical experiments with state-of-the art algorithms from the literature.

563 From the results of the numerical experiments, it seems that the best performing algorithm is  
564 PHALS, both in terms of performance of retrieving the factorization and of overall computation  
565 time. Although GD gives good results in terms of approximate factorization performance, its high  
566 complexity impedes its practical use on large datasets.

567 We have also tested PHALS to retrieve the BMF of real datasets. The components obtained  
568 seem to agree with those obtained in other works of the literature and with intuition on what would  
569 be possible groupings of the data. In this paper, we have not focused on the choice of the number  
570 of components  $R$  and in the presentation of the results for the real datasets, we have chosen a value  
571 of  $R$  that seemed to give stable results with components agreeing with intuition on the dataset. In  
572 practice, if no intuition on the expected components is available, a quantitative criterion for choosing  
573  $R$  may be used. Such criteria will be studied and tested in future work.

574 At the end of the experimental section, we have also presented results of applying PHALS and  
575 GD in a more general factorization setting where XOR – 2 and MAJ – 3 component combining  
576 functions are considered instead of the logical OR of BMF. The uniqueness and the interpretability  
577 of these decompositions are more problematic and further investigations are required to assess their  
578 utility in practical applications. In future work, we would like to extend our general approach to  
579 the higher-order tensor setting and to analyse the identifiability of such models.

## 580 **References**

- 581 [1] T. Li, A general model for clustering binary data, in: Proc. of the 11th ACM SIGKDD Inter-  
582 national Conference on Knowledge Discovery in Data Mining, 2005, pp. 188–197.
- 583 [2] L. Kozma, A. Ilin, T. Raiko, Binary principal component analysis in the Netflix collaborative  
584 filtering task, in: 2009 IEEE Int. Workshop Mach. Learn. Signal Process., IEEE, 2009, pp. 1–6.
- 585 [3] E. Nenova, D. I. Ignatov, A. V. Konstantinov, An FCA-based boolean matrix factorisation  
586 for collaborative filtering, FCAIR 2012 Formal Concept Analysis Meets Information Retrieval  
587 (2013) 57.

- 588 [4] M. Diop, S. Miron, A. Larue, D. Brie, Binary matrix factorization applied to Netflix dataset  
589 analysis, *IFAC-PapersOnLine* 52 (24) (2019) 13–17.
- 590 [5] E. Meeds, Z. Ghahramani, R. M. Neal, S. T. Roweis, Modeling dyadic data with binary latent  
591 factors, *Adv. Neural. Inf. Process. Syst.* 19 (2007) 977.
- 592 [6] Z.-Y. Zhang, T. Li, C. Ding, X.-W. Ren, X.-S. Zhang, Binary matrix factorization for analyzing  
593 gene expression data, *Data Min. Knowl. Discov.* 20 (1) (2010) 28–52.
- 594 [7] S. Tu, L. Xu, R. Chen, A binary matrix factorization algorithm for protein complex prediction,  
595 in: *Proc. IEEE Int. Conf. Bioinformatics Biomed. (BIBMW)*, IEEE, 2010, pp. 113–118.
- 596 [8] H. Lu, J. Vaidya, V. Atluri, Optimal Boolean matrix decomposition: Application to role engi-  
597 neering, in: *IEEE 24th Int. Conf. Data Eng.*, IEEE, 2008, pp. 297–306.
- 598 [9] S. Talwar, M. Viberg, A. Paulraj, Blind separation of synchronous co-channel digital signals  
599 using an antenna array. I. algorithms, *IEEE Trans. Signal Process.* 44 (5) (1996) 1184–1197.
- 600 [10] A.-J. Van der Veen, Analytical method for blind binary signal separation, *IEEE Trans. Signal*  
601 *Process.* 45 (4) (1997) 1078–1082.
- 602 [11] A. I. Schein, L. K. Saul, L. H. Ungar, A generalized linear model for principal component  
603 analysis of binary data, in: *International Workshop on Artificial Intelligence and Statistics*,  
604 PMLR, 2003, pp. 240–247.
- 605 [12] J. De Leeuw, Principal component analysis of binary data by iterated singular value decompo-  
606 sition, *Computational statistics & data analysis* 50 (1) (2006) 21–39.
- 607 [13] S. Lee, J. Z. Huang, J. Hu, Sparse logistic principal components analysis for binary data, *Ann.*  
608 *Appl. Stat.* 4 (3) (2010) 1579.
- 609 [14] Z. Kang, C. J. Spanos, Sequential logistic principal component analysis (SLPCA): Dimensional  
610 reduction in streaming multivariate binary-state system, in: *13th International Conference on*  
611 *Machine Learning and Applications*, IEEE, 2014, pp. 171–177.
- 612 [15] A. Lumbreras, L. Filstroff, C. Févotte, Bayesian mean-parameterized nonnegative binary matrix  
613 factorization, *Data Min. Knowl. Discov.* 34 (6) (2020) 1898–1935.
- 614 [16] Z. Zhang, T. Li, C. Ding, X. Zhang, Binary matrix factorization with applications, in: *IEEE*  
615 *Int. Conf. Data Min. (ICDM)*, IEEE, 2007, pp. 391–400.

- 616 [17] S. Karaev, P. Miettinen, Capricorn: an algorithm for subtropical matrix factorization, in: Proc.  
617 SIAM Int. Conf. Data Min., SIAM, 2016, pp. 702–710.
- 618 [18] S. Karaev, P. Miettinen, Cancer: Another algorithm for subtropical matrix factorization, in:  
619 Joint European conference on machine learning and knowledge discovery in databases, Springer,  
620 2016, pp. 576–592.
- 621 [19] S. Karaev, P. Miettinen, Algorithms for approximate subtropical matrix factorization, *Data*  
622 *Min. Knowl. Discov.* 33 (2) (2019) 526–576.
- 623 [20] N. Gillis, S. A. Vavasis, On the complexity of robust PCA and  $\ell_1$ -norm low-rank matrix ap-  
624 proximation, *Math. Oper. Res.* 43 (4) (2018) 1072–1084.
- 625 [21] P. Miettinen, T. Mielikäinen, A. Gionis, G. Das, H. Mannila, The discrete basis problem, *IEEE*  
626 *Trans. Knowl. Data Eng.* 20 (10) (2008) 1348–1362.
- 627 [22] R. Belohlavek, V. Vychodil, Discovery of optimal factors in binary data via a novel method of  
628 matrix decomposition, *J. Comput. Syst. Sci.* 76 (1) (2010) 3–20.
- 629 [23] M. Trnecka, R. Vyjidacek, Revisiting the Grecon algorithm for Boolean matrix factorization,  
630 *Knowl.-Based Syst.* (2022) 108895.
- 631 [24] T. Makhalova, M. Trnecka, From-below Boolean matrix factorization algorithm based on MDL,  
632 *Adv. Data Anal. Classif.* 15 (1) (2021) 37–56.
- 633 [25] D. D. Lee, H. S. Seung, Learning the parts of objects by non-negative matrix factorization,  
634 *Nature* 401 (6755) (1999) 788–791.
- 635 [26] S. Miron, M. Diop, A. Larue, E. Robin, D. Brie, Boolean decomposition of binary matrices  
636 using a post-nonlinear mixture approach, *Signal Process.* 178 (2021) 107809.
- 637 [27] D. DeSantis, E. Skau, D. P. Truong, B. Alexandrov, Factorization of binary matrices: Rank  
638 relations, uniqueness and model selection of boolean decomposition, *ACM Trans. Knowl. Discov.*  
639 *Data* 16 (6) (2022) 1–24.
- 640 [28] H. Nguyen, R. Zheng, Binary independent component analysis with or mixtures, *IEEE Trans.*  
641 *Signal Process.* 59 (7) (2011) 3168–3181.
- 642 [29] S. Ravanbakhsh, B. Póczos, R. Greiner, Boolean matrix factorization and noisy completion via  
643 message passing, in: *International Conference on Machine Learning*, PMLR, 2016, pp. 945–954.

- 644 [30] A. Cichocki, A.-H. Phan, Fast local algorithms for large scale nonnegative matrix and tensor  
645 factorizations, *IEICE Trans. Fundam. Electron. Comput. Sci.* 92 (3) (2009) 708–721.
- 646 [31] D. DeSantis, E. Skau, B. Alexandrov, Factorizations of binary matrices–rank relations and the  
647 uniqueness of Boolean decompositions, *arXiv preprint arXiv:2012.10496* (2020).
- 648 [32] J. E. Cohen, U. G. Rothblum, Nonnegative ranks, decompositions, and factorizations of non-  
649 negative matrices, *Linear Algebra Appl.* 190 (1993) 149–168.
- 650 [33] K. H. Kim, *Boolean matrix theory and applications*, Vol. 70, Dekker, 1982.
- 651 [34] V. L. Watts, Boolean rank of Kronecker products, *Linear Algebra Appl.* 336 (1-3) (2001) 261–  
652 264.
- 653 [35] T. Watson, Nonnegative rank vs. binary rank, *Chicago Journal of Theoretical Computer Science*  
654 2 (2016) 1–13.
- 655 [36] Y. Crama, P. L. Hammer, *Boolean functions: Theory, algorithms, and applications*, Cambridge  
656 University Press, 2011.
- 657 [37] D. Dua, C. Graff, UCI machine learning repository, <http://archive.ics.uci.edu/ml> (2017).
- 658 [38] The now community. new and old worlds database of fossil mammals (NOW)., <https://nowdatabase.org/now/database/>, accessed: 2022-06-15. (2022).
- 659
- 660 [39] E. Voeten, Data and analyses of voting in the un general assembly, in: B. Reinalda (Ed.), *Data  
661 and Analyses of Voting in the UN General Assembly*, Routledge Londres, 2013.
- 662 [40] J. Nocedal, S. Wright, *Numerical optimization*, Springer Science & Business Media, 2006.
- 663 [41] P.-A. Absil, R. Mahony, B. Andrews, Convergence of the iterates of descent methods for analytic  
664 cost functions, *SIAM J. Optim.* 16 (2) (2005) 531–547.
- 665 [42] Y. Xu, W. Yin, A block coordinate descent method for regularized multiconvex optimization  
666 with applications to nonnegative tensor factorization and completion, *SIAM Journal on imaging  
667 sciences* 6 (3) (2013) 1758–1789.
- 668 [43] E. Bingham, A. Kabán, M. Fortelius, The aspect Bernoulli model: multiple causes of presences  
669 and absences, *Pattern Anal. Appl.* 12 (1) (2009) 55–78.