

# Boolean matrix factorization meets consecutive ones property\*

Nikolaj Tatti<sup>†</sup>

Pauli Miettinen<sup>‡</sup>

## Abstract

Boolean matrix factorization is a natural and a popular technique for summarizing binary matrices. In this paper, we study a problem of Boolean matrix factorization where we additionally require that the factor matrices have consecutive ones property (OBMF). A major application of this optimization problem comes from graph visualization: standard techniques for visualizing graphs are circular or linear layout, where nodes are ordered in circle or on a line. A common problem with visualizing graphs is clutter due to too many edges. The standard approach to deal with this is to bundle edges together and represent them as ribbon. We also show that we can use OBMF for edge bundling combined with circular or linear layout techniques.

We demonstrate that not only this problem is NP-hard but we cannot have a polynomial-time algorithm that yields a multiplicative approximation guarantee (unless  $P = NP$ ). On the positive side, we develop a greedy algorithm where at each step we look for the best 1-rank factorization. Since even obtaining 1-rank factorization is NP-hard, we propose an iterative algorithm where we fix one side and find the other, reverse the roles, and repeat. We show that this step can be done in linear time using pq-trees. We also extend the problem to cyclic ones property and symmetric factorizations. Our experiments show that our algorithms find high-quality factorizations and scale well.

## 1 Introduction

Matrix factorization is an immensely popular way of summarizing data as well as discovering signal from the data. While being useful, the interpretation and visualization of discovered factor matrices may be difficult. A popular variant for factorizing binary matrices is a  $k$ -Boolean matrix factorization, which, essentially, summarizes the binary data as a union of  $k$  tiles, that is, submatrices full of 1s. However, visualizing such factorization is difficult

as the discovered rows and columns can be any sets, and there is no insightful way of visualizing them all at once.

In this paper we consider  $k$ -Boolean matrix factorization such that the resulting matrix has a certain property: we can order the columns and the rows such that the matrix consists of union of  $k$  *contiguous* tiles. We do not know the order before-hand, and we discover the order as we also discover the factorization.

Our motivation for discovering such factorization is primarily due to easy exploration of the factorization: we can draw the factorization as  $k$  tiles. While in certain cases, such a constraint may be too restrictive, there are many settings, where this constraint comes naturally. As a specific example, consider visualizing graphs. A classic technique for visualizing a graph is using linear or circular layout, where the nodes are drawn on a line or circle, and they are connected with arcs. The most common problem with visualizing graphs is clutter due to too many edges. To combat the clutter, edges are often grouped, and drawn in ribbons (see Figure 3 for an example). The problem is to discover such ribbons and the node order, while minimizing the error. We show that we can use matrix factorization on the adjacency matrix of a graph to find the order and the groups.

We show that the factorization we seek can be expressed with consecutive ones property (C1P). Namely, we will look for factor matrices  $X$  and  $Y$  whose columns can be shuffled such that each row has a form of  $[0, \dots, 0, 1, \dots, 1, 0, \dots, 0]$ . We show that the problem is NP-hard, even if  $k = 1$ , and it is inapproximable for  $k > 1$ . On the positive side, we propose a greedy algorithm that searches the factors in iterative manner. The search is done by first fixing a vector in  $X$  and finding the optimal counterpart in  $Y$ , then fixing the vector in  $Y$  and finding the optimal vector in  $X$ , and so on, until convergence. We show that we can find the optimal counterpart in linear time using pq-trees.

\*This is an extended version of the paper of the same name presented in 2019 SIAM International Conference on Data Mining.

<sup>†</sup>University of Helsinki, Helsinki, Finland,  
nikolaj.tatti@helsinki.fi

<sup>‡</sup>University of Eastern Finland, Kuopio, Finland,  
pauli.miettinen@uef.fi. Part of this work was done while the author was with MPI-INF, Saarbrücken, Germany.

We also consider 3 extensions of this factorization: the first variant, cyclic decomposition, consists of allowing factors to “wrap around the border.” the second variant is specifically designed for symmetric matrices, while the last variant combines the two. Performing cyclic and symmetric decomposition proves to be useful for cyclic layout of graphs.

The rest of the paper is organized as follows: We present preliminary notation and define the matrix factorization and the cyclic version in Section 2. We present the search algorithm in Section 3. The symmetric extensions are given in Section 4. Section 6 is dedicated to related work, and Section 5 is dedicated to experimental evaluation. Finally, we conclude the paper with remarks given in Section 7. All proofs are given in Appendix A.

## 2 Preliminary notation and problem definitions

We begin by presenting preliminary notation, and then present the two main problem definitions. Extended problems are discussed in Section 4.

### 2.1 Notation

Given an  $n$ -by- $k$  binary matrix  $A$  and a  $k$ -by- $m$  binary matrix  $B$ , the *Boolean matrix product*  $A \circ B$  is defined element-wise as

$$(1) \quad (A \circ B)_{ij} = \bigvee_{\ell=1}^k a_{i\ell} b_{\ell j} .$$

The *Boolean matrix sum* of  $A \in \{0, 1\}^{n \times m}$  and  $B \in \{0, 1\}^{n \times m}$  is defined elementwise as  $(A \vee B)_{ij} = a_{ij} \vee b_{ij}$ .

To measure the distance between two binary matrices, we use the *squared Frobenius norm* of their (normal) difference,  $\|A - B\|_F^2$ . Notice that as  $A$  and  $B$  are both binary, this is the same as calculating the number of disagreements between  $A$  and  $B$ :  $\|A - B\|_F^2 = |\{(i, j) : a_{ij} \neq b_{ij}\}|$ .

We say that a binary matrix  $X$  has a *consecutive ones property* (C1P) if its columns can be permuted such that each row has a form of  $[0, \dots, 0, 1, \dots, 1, 0, \dots, 0]$ , that is, 1s form a contiguous interval. For the sake of presentation, we will also refer these matrices as *unimodal*.

We say that a binary matrix  $X$  is *cyclic* if its columns can be permuted such that each row has a form of  $[0, \dots, 0, 1, \dots, 1, 0, \dots, 0]$  or  $[1, \dots, 1, 0, \dots, 0, 1, \dots, 1]$ .

### 2.2 Problem definitions

Next we will give our two main optimization problems.

**Problem 1** (Ordered BMF,  $\text{OBMF}$ ). Given a binary matrix  $D$  and an integer  $k \in \mathbb{N}$ , find two unimodal binary matrices  $X$  and  $Y$  that minimize the number of disagreements

$$(2) \quad \|D - (X^T \circ Y)\|_F^2 .$$

**Problem 2** (Cyclic Ordered BMF,  $\text{COBMF}$ ). Given a binary matrix  $D$  and an integer  $k \in \mathbb{N}$ , find two cyclic binary matrices  $X$  and  $Y$  that minimize the number of disagreements

$$(3) \quad \|D - (X^T \circ Y)\|_F^2 .$$

The matrix  $Z = X^T \circ Y$  given in Eq. 2 has another natural alternative characterization: the columns and the rows of  $Z$  can be permuted such that the resulting matrix is a union of  $k$  contiguous tiles of 1s. Similarly, the matrix  $Z = X^T \circ Y$  given in Eq 3 can be permuted such that the resulting matrix is a union of  $k$  contiguous tiles, but we also allow the tiles to wrap around the border.

Unsurprisingly, the problems are computationally infeasible. First, we demonstrate that  $\text{OBMF}$  is difficult even if  $k = 1$ .

**Theorem 1.** *The  $\text{OBMF}$  problem is NP-hard, even if  $k = 1$ .*

Our next result shows that not only  $\text{OBMF}$  is difficult, but it is also impossible to approximate. To show this, it is enough to demonstrate that testing for zero-error solution is expensive.

**Theorem 2.** *Deciding whether  $\text{OBMF}$  has a zero-error solution is NP-complete.*

The proofs of these and other statements are given in Appendix A.

## 3 Iterative greedy algorithm

### 3.1 Greedy algorithm

As we saw in the previous section, not only the problem is NP-hard, we cannot construct any polynomial-time algorithm with a multiplicative guarantee. Hence, we need to

resort to heuristics. The most natural heuristic is a greedy heuristic, where given a  $(k-1)$ -sized factorization we look for a  $k$ -sized factorization by adding one row and one column to  $X$  and  $Y$ . Note that these rows need to be selected carefully such that  $X$  and  $Y$  remain unimodal, and we also need to maintain the permutation(s).

Unfortunately, Theorem 1 states that we cannot even find the best solution for  $k = 1$  in polynomial-time. Fortunately, we can solve quickly a subproblem, where we have fixed one side.

**Problem 3** (Ordered BMF step, `OBMFSTEP`). Given a binary matrix  $D$  of size  $n$ -by- $m$  and two unimodal matrices,  $X'$  of size  $k$ -by- $n$  and  $Y'$  of size  $(k-1)$ -by- $m$ , find the decomposition  $X^T \circ Y$  solving `OBMF` such that  $X = X'$  and  $Y$  is obtained by adding one new row to  $Y'$ .

We can use `OBMFSTEP` as follows. Assume that we have already found  $(k-1)$ -by- $m$  matrices  $X$  and  $Y$ . We first extend  $X$  with a new row using a given seed, and find the optimal new row for  $Y$  (strategy for such selection is given later using `OBMFSTEP`). We fix the discovered row, and use `OBMFSTEP` to find the corresponding row for  $X$ . Since we solve each step optimally, the error will never increase. We stop when the error stops decreasing. Note that we will need to provide a seed for the initial row in  $X$ . Here, we test several possible seeds  $S$ , and select the best. We experiment with several options in experiments, but the default is that  $S$  is equal to all singleton columns. The pseudo-code for the algorithm is given in Algorithm 1.

The remainder of this section is about solving `OBMFSTEP` in linear time. Almost the same approach will also work for the cyclic version, `COBMFSTEP`; we will point the minute difference.

### 3.2 Expressing permutations with pq-trees

The complicated aspect of `OBMFSTEP` is that we need to make sure that the new matrix is unimodal. Luckily, we can use pq-trees, a classic structure that allows us to express every permutation for which a set of binary vertices remain unimodal. In this section we will give a brief review of pq-trees and the two main properties that are relevant to us.

Assume that we are given a universe  $U$ ; in our case this will be either rows or columns of the input matrix. A pq-tree is a tree with each leaf corresponding to  $u \in U$ .

---

**Algorithm 1:** Greedy iterative algorithm for estimating `OBMF`. The algorithm takes as input the dataset, the desired dimension  $k$ , and the seed set  $S$  used for selecting the first candidate for a column.

---

```

1  $X \leftarrow$  matrix of size 0-by- $n$ ;
2  $Y \leftarrow$  matrix of size 0-by- $m$ ;
3 foreach  $i = 1 \dots k$  do
4   foreach  $s \in S$  do
5      $c \leftarrow s$ ;
6     while error decreases do
7        $r \leftarrow$  best row for fixed columns  $[Y; c]$ ;
8        $c \leftarrow$  best column for fixed rows  $[X; r]$ ;
9      $X \leftarrow [X; r]$ ;
10     $Y \leftarrow [Y; c]$ ;

```

---

There are two types of non-leaf nodes, these types will dictate what permutations we can perform on the children. We can permute children of p-node in any order whereas the order of the children of q-node is fixed but we can flip the direction. The leaves of the permuted tree will then indicate an order. We will denote such orders by  $order(T)$ , where  $T$  is the pq-tree.

Two seminal results are important to us. The first result states that there is a pq-tree  $T$  such that  $order(T)$  are exactly the orders under which a set of binary vertices remain unimodal.

**Theorem 3** (Booth and Lueker [3]). *Given a universe  $U$  and  $k$  sets  $S_i \subseteq U$ , there is a pq-tree  $T$  such that  $order(T)$  are exactly the permutations of  $U$  under which each  $S_i$  is contiguous.*

The second result states that we can efficiently update the pq-tree.

**Theorem 4** (Booth and Lueker [3]). *Assume that we have a pq-tree  $T$  over a universe  $U$  and a set  $S \subseteq T$ . Let  $P$  be the set of all permutations of  $U$  where  $S$  is contiguous. If  $order(T) \cap P \neq \emptyset$ , then there is an  $O(|U|)$ -time algorithm that constructs a tree  $T'$  such that  $order(T') = order(T) \cap P$ . If  $order(T) \cap P = \emptyset$ , then the same algorithm detects a failure.*

The detailed description of the algorithm for updating the pq-tree can be found in [3].

### 3.3 Finding the optimal row

In this section we describe the algorithm that solves OBMFSTEP. Assume that we have a pq-tree  $T$  representing the permutations of columns in  $\mathbf{D}$  allowed by the previously discovered rows in  $\mathbf{Y}'$ . When dealing with pq-trees it is notationally easier to deal with sets rather than with vectors. Naturally every binary vector  $\mathbf{y}$  can be represented as a set  $S = \{i : y_i = 1\}$ .

Let us define  $U$  to be the column indices of  $\mathbf{D}$ ; these are exactly the leaves of  $T$ . We say that a set  $S \subseteq U$  is *compatible* with a pq-tree  $T$ , if there is an order in  $order(T)$  where  $S$  is contiguous. Obviously, compatible sets  $S$  correspond exactly to suitable new rows in  $\mathbf{Y}$ .

We can express OBMFSTEP as an instance of the following problem.

**Problem 4 (OPTSET).** Given a universe  $U$ , weights  $w(u)$  for each  $u \in U$ , and a pq-tree  $T$  over the universe  $U$ , find a set  $S$  that is compatible with  $T$  and maximizes the total weight  $\sum_{u \in S} w(u)$ .

Recall that  $u \in U$  corresponds to a column index of  $\mathbf{D}$ . Define  $w(u)$  to be the gain in the error-function if we were to use  $u$  in our new row for  $\mathbf{Y}$ . More formally, let  $\mathbf{x}$  be the fixed counterpart in  $\mathbf{X}$  for the new row in  $\mathbf{Y}$ . Let  $p$  be the number of ones in  $\mathbf{D}$  at rows  $\mathbf{x}$  and column  $u$  that are not yet covered by the previous factors. Let  $n$  be the number of zeros in  $\mathbf{D}$  at rows  $\mathbf{x}$  and column  $u$  that are not yet covered by the previous factors. We define  $w(u) = p - n$ . Solving OPTSET with these weights solves OBMFSTEP.

In order to solve COBMFSTEP, we solve OPTSET using  $w(u) = p - n$ , as above, yielding a set, say  $S_1$ . In addition, we also solve OPTSET using  $w(u) = n - p$ , yielding a set, say  $S_2$ . Then, we use either  $S_1$  or  $U \setminus S_2$ , whichever yields a better gain.

In order to solve OPTSET, we need an additional definition: Let  $S$  be a compatible set of a pq-tree  $T$ . If there is a permutation in  $order(T)$  with the first or the last element in  $S$ , we call  $S$  a *border-compatible* set.

Let  $T$  be a pq-tree. To solve OPTSET we will compute 3 counters for a node  $v$  in  $T$ , namely,  $inner(v)$ ,  $border(v)$ , and  $total(v)$ . The counter  $total(v)$  corresponds to the total weight of leaves under  $v$ , while the counter  $inner(v)$  corresponds to the best  $S$  that is compatible with the subtree starting at  $v$ . Finally,  $border(v)$  corresponds to the best  $S$  that is border-compatible with the subtree starting at  $v$ .

We should stress that, strictly by definition,  $inner(v)$  can represent an empty set, whereas  $total(v)$  and  $border(v)$  should be never empty, even if they produce a negative value. Thus,  $inner(v) \geq 0$  but  $border(v)$  and  $total(v)$  can have negative values. Moreover, it is possible that  $border(v)$  represents every leaf of  $v$ , in which case,  $border(v) = total(v)$ .

Naturally, we want to compute  $inner(r)$ , where  $r$  is the root of  $T$ . To obtain this value we compute each value iteratively, children first. We also maintain the lists of the children that were responsible for producing the optimal value. These lists are clear from the proofs of the following lemmata. This allows us to extract the optimal  $S$ .

First, note that computing  $total(v)$  is trivial since  $total(v) = \sum_{c \in ch(v)} total(c)$ . If  $v$  is a leaf-node, then  $border(v) = total(v)$  and  $inner(v) = \max(0, total(v))$ .

The next two lemmata establish how to compute the counters for q-nodes.

**Lemma 5.** Let  $v$  be a q-node and let  $c_1, \dots, c_\ell$  be its children. Then

$$border(v) = \max(x, y), \quad \text{where}$$

$$x = \max_i border(c_i) + \sum_{j=1}^{i-1} total(c_j),$$

$$y = \max_i border(c_i) + \sum_{j=i+1}^{\ell} total(c_j) \quad .$$

**Lemma 6.** Let  $v$  be a q-node and let  $c_1, \dots, c_\ell$  be its children. Then

$$inner(v) = \max(x, y), \quad \text{where}$$

$$x = \max_i inner(c_i),$$

$$y = \max_{i < j} border(c_i) + border(c_j) + \sum_{\ell=i+1}^{j-1} total(c_\ell) \quad .$$

Our next step is to compute the counters for p-nodes. For that we need to define the following helper function: given a node  $v$  we define  $g(v) = border(v) - \max(total(v), 0)$ . We will use  $g(v)$  in the next two lemmata describing on how to compute the counters for p-node.

**Lemma 7.** Let  $v$  be a p-node and let  $c_1, \dots, c_\ell$  be its children. Define  $b = \max g(c_i)$ . Then

$$border(v) = b + \sum_i \max(total(c_i), 0) \quad .$$

Note that since we require the set responsible for  $border(v)$  be non-empty, it is possible that  $border(v) < 0$ . This can happen only if  $b < 0$  and every child  $w$  of  $v$  has  $total(w) < 0$ .

**Lemma 8.** *Let  $v$  be a  $p$ -node and let  $c_1, \dots, c_\ell$  be its children. Define  $b_1$  and  $b_2$  be the top-2 values of  $g(c_i)$ . Then*

$$\begin{aligned} inner(v) &= \max(x, y), \quad \text{where} \\ x &= \max_i inner(c_i), \\ y &= \max(b_1, 0) + \max(b_2, 0) + \sum_i \max(total(c_i), 0). \end{aligned}$$

Note that using these lemmas every counter can be trivially solved in linear time, except for  $inner(v)$ , where  $v$  is  $q$ -node. To compute  $inner(v)$  in linear time, it is enough if we can solve

$$border(c_j) + \max_{i < j} border(c_i) + \sum_{\ell=i+1}^{j-1} total(c_\ell)$$

in constant time for a fixed  $j$ . Luckily, we can rewrite this function as

$$border(c_j) + \left( \sum_{\ell=1}^{j-1} total(c_\ell) \right) + \max_{i < j} t(i, j),$$

where

$$t(i, j) = \max_{i < j} border(c_i) - \sum_{\ell=1}^i total(c_\ell) \quad .$$

Let  $i(j)$  to be the optimal  $i$  for a fixed  $j$ . Since

$$\max_{i < j} t(i, j) = \max \left( t(j-1, j) \max_{i < j-1} t(i, j) \right),$$

we have either  $i(j) = i(j-1)$  or  $i(j) = j-1$ . If we were to test each  $j$  consecutively, then this allows us to compute  $i(j)$  in constant time: we simply compare the solution  $i = j-1$  to the best previous solution  $i(j-1)$ .

In summary, each counter of  $v$  can be computed in  $O(ch(v))$ . Thus we need  $O(\ell)$ , where  $\ell$  is the number of nodes in  $T$ . Since  $\ell \in O(|U|)$ , we can compute the counters in  $O(|U|)$  time, where  $|U|$  is the number of columns in  $D$ .

When computing the counters we also store which children were responsible for this value. Once we have computed  $inner(r)$ , where  $r$  is the root of the tree, we can backtrack to obtain the optimal  $S$ . This can be also done in linear time.

Computing the weights  $w$  in  $OPTSET$  can be done in  $O(p)$  time, where  $p$  is the number of 1s in the dataset  $D$  of size  $n$ -by- $m$ . Consequently,  $OBMFSTEP$  can be done in  $O(p + n + m)$  time.

## 4 Symmetric decomposition

We now propose an extension for symmetric matrices.

### 4.1 Definition

If  $D$  is symmetric (e.g. an adjacency matrix of an undirected graph), we have the following problem:

**Problem 5** (Symmetric  $OBMF$ ,  $OBMF_{sym}$ ). Given a binary matrix  $D$  and an integer  $k \in \mathbb{N}$ , find two binary matrices  $X$  and  $Y$  such that  $[X; Y]$  is unimodal, that minimize the number of disagreements

$$(4) \quad \left\| D - ((X^T \circ Y) \vee (Y^T \circ X)) \right\|_F^2 .$$

We define similarly  $COBMF_{sym}$ , a cyclic and symmetric variant of  $OBMF$ .

The unimodality condition in  $OBMF_{sym}$  states that we should be able to permute  $X$  and  $Y$  with the same permutation so that the rows are in form of  $[0, \dots, 0, 1, \dots, 1, 0, \dots, 0]$ .

Notice that we do not use the more common symmetric decomposition  $D \approx X^T \circ X$  as this would lead to necessarily having the blocks around the diagonal.

### 4.2 Algorithm

The discovery algorithm for symmetric  $OBMF$  is similar. Like with the regular  $OBMF$ , we use a greedy algorithm as an iterative step for discovering new rows.

The first difference is that we maintain only one pq-tree, corresponding to the rows in both  $X$  and  $Y$ .

The second difference is that – as  $X^T \circ Y$  and  $Y^T \circ X$  can have overlapping 1s – maximizing  $OPTSET$  does not necessarily produce the optimal row. Instead, we can show

that solving `OPTSET`, with the weights as described in the previous section, minimizes  $\|D - X^T \circ Y\|_F^2 + \|D - Y^T \circ X\|_F^2$ . It follows easily that minimizing this function yields a 2-approximation for finding optimal counterpart row.

## 5 Experimental evaluation

In this section we study how well the algorithms from Sections 3 and 4.2 work with synthetic and real-world data. We denote the algorithms with the same names as the problems they are solving, and differentiate the algorithms from the problems via the font. That is, `obmf` is the algorithm for `OBMF`, and so on. The algorithms are implemented in C++, and we make the source code and synthetic experiments freely available.<sup>1</sup>

### 5.1 Resilience to Noise

We start by evaluating the algorithms’ resilience to noise. To that end, we synthesized random matrices of size  $95 \times 95$  with block structure (6 blocks of size  $20 \times 20$  along the diagonal, with 5 overlapping rows and columns) and corrupted those matrices with flipping a varying amounts of entries. The amount of flipped entries varied from 0% to 50% (of total elements) and we compared the quality of the results to both the noise-free matrix and noisy matrix. The results are shown in Figure 1.

With lower levels of noise (35% for `obmf` and `cobmf` and 25% for the symmetric variants), the reconstruction of the original data is more accurate. With higher levels of noise, the noise has destroyed so much of the structure that the algorithms start fitting to the noise only, with a clear reduction of the quality versus the original data.

It is also worth noticing that `obmf` obtains exact decompositions when the data has no noise; the other methods introduce a slight error even in these cases emphasizing their more complex setting.

### 5.2 Scalability

In this section we test how well `obmf` scales to larger data sets and how well it benefits from multiple cores. These experiments were executed on a server with 40 cores of

<sup>1</sup>[https://cs.uef.fi/~pauli/bmf/ordered\\_bmf/](https://cs.uef.fi/~pauli/bmf/ordered_bmf/)

Intel Xeon E7-4870 processors running at 2.4 GHz. The algorithm was compiled using GCC 8.1.0 and the parallel code uses the OpenMP library.

To test the scalability, we generated  $n$ -by- $n$  square matrices with  $n = 2^i$  for  $i = 9, \dots, 13$ . All matrices have a density of approximately 24%. The results are presented in Figure 2a.

The algorithm shows very good scalability over the full range, although it does get slower when the data size increases from  $2^{12}$  to  $2^{13}$ . It should be noted, though, that as the density is constant, the number of non-zeros in the matrices increases as the square of the matrix size. Hence, `obmf` exhibits linear growth with respect to the number of non-zero elements.

Algorithm 1 is almost embarrassingly parallel over the different seeds vectors. Hence, we parallelized the test of different seeds, and tested how the algorithm behaves with increased number of cores. The results are in Figure 2b, where we can see that the speed-up is essentially linear up to 4 cores, slightly slower until 16 cores, and only marginal gains are available when increasing the number of cores to 32, indicating that at the algorithm has become memory bus constrained.

Overall, the experiments show that the algorithm scales very well, and is able to benefit from modern multi-core computers. We study further speed-up options later in Section 5.3.2.

### 5.3 Experiments with Real-World Data

We now turn to real-world data sets. We used six different real-world data sets, selected to offer a wide variety of different types of data. The data sets we used are as follows. *Les Misérables* is a standard benchmark data<sup>2</sup> of the characters of Victor Hugo’s novel *Les Misérables*. *Paleo* is a palaeontological data<sup>3</sup> in the form of a locations-by-genera matrix, giving information where different fossils have been found. *Newsgroups* is a subset of the famous 20Newsgroups data<sup>4</sup> consisting four newsgroups and 100 terms. *Terms* the terms-by-terms co-occurrence matrix based on *Newsgroups*. *Locations* is locations-by-locations matrix indicating mammal species co-location in the northern hemisphere: the data has a 1 in element  $(i, j)$  if loca-

<sup>2</sup><http://moreno.ss.uci.edu/data.html>

<sup>3</sup>NOW 030717, <http://www.helsinki.fi/science/now/>

<sup>4</sup><http://qwone.com/~jason/20Newsgroups/>

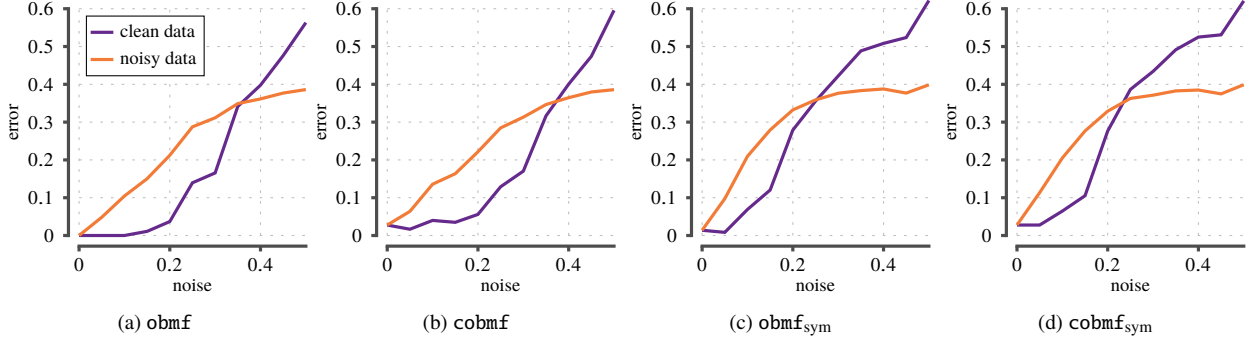


Figure 1: Error as a function of noise. Here the error is the proportion of disagreements between the reconstructed matrix and either the noise-free or the noisy matrix. The decomposition was done using the noisy matrix.

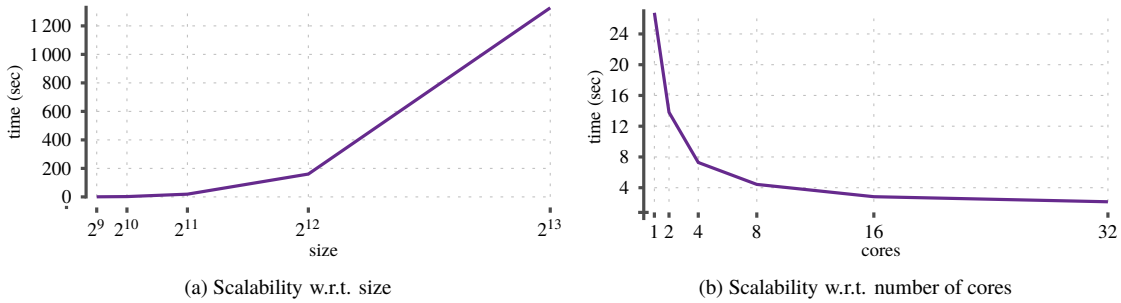


Figure 2: Scalability with respect to the size and number of cores.

tions  $i$  and  $j$  have at least five mammals in common. The data is based on the IUC Red List data<sup>5</sup>. The final data set, *Mammals*, contains a species-by-species co-inhabitation matrix.<sup>6</sup> The data set properties are summarized in Table 1.

To the best of our knowledge, this is the first work to address the ordered Boolean matrix factorization problem. To understand what kind of an effect the ordering constraint has to the reconstruction error, we compare our results with those of *asso* [15]. The *asso* algorithm is a well-known method for computing the standard Boolean matrix factorization. We used an implementa-

Table 1: Properties of real-world data sets. Rank indicates the rank used in the decomposition.

data	rows	cols	% of 1s	sym.	rank
<i>Les Misérables</i>	77	77	8.57	Yes	10
<i>Paleo</i>	124	139	11.48	No	10
<i>Newsgroups</i>	100	348	6.30	No	10
<i>Terms</i>	100	100	48.54	Yes	10
<i>Locations</i>	3203	3203	8.42	Yes	50
<i>Mammals</i>	194	194	58.04	Yes	10

tion available from the author<sup>7</sup> and set the rank for *asso* the same as for our algorithms, and used threshold values

<sup>5</sup><http://www.iucnredlist.org/technical-documents/spatial-data>

<sup>6</sup>Available for research purposes from the Societas Europaea Mammalogica at <http://www.european-mammals.org>

<sup>7</sup><https://cs.uef.fi/~pauli/basso/basso-0.5.tar.gz>

$\tau = \{0.2, 0.4, 0.6, 0.8\}$ .

For symmetric data sets, we also computed the symmetric Boolean factorization. This was done by first computing the standard  $X^T \circ Y$  factorization, and then testing whether  $X^T \circ X$  or  $Y^T \circ Y$  gives smaller reconstruction error and using that one. This version of `asso` is denoted `assosym`.

### 5.3.1 Reconstruction errors

We first compute the reconstruction errors for the various data sets. To facilitate the comparisons, we report the *relative* reconstruction error

$$\frac{\|D - X^T \circ Y\|_F^2}{\|D\|_F^2}.$$

The results of all datasets are given in Table 2.

In case of asymmetric decompositions, `asso` is – as expected, as its factor matrices are not restricted to unimodal or cyclic – almost always slightly better than either `obmf` or `cobmf`. This difference is, however, very small in many data sets (only 8% in *Les Misérables* and 0.50% in *Paleo*). A remarkable exception is the *Mammals* data, where `asso` is in fact worse than either `obmf` or `cobmf`. As the data set is the densest of the ones we tested, it is possible that `asso` was unable to obtain good candidates from it with the rounding thresholds we tried.

There is almost no difference between `obmf` and `cobmf` in the terms of reconstruction error in these data sets. Usually, `obmf` is on par or slightly better than `cobmf`, except again in *Mammals*, where `cobmf` is slightly better. The asymmetric data sets, *Paleo* and *Newsgroups*, cause the highest reconstruction errors at over 70%. It should be noted, though, that also `asso` has similarly high errors with these data sets, indicating that they might not have strong Boolean low-rank structure.

In symmetric decompositions, the relationship between the ordered BMF algorithms and `asso` is reversed, with `assosym` being often the worse method (with the exception of *Terms*). This is not very surprising, given that `asso` is not designed for symmetric decompositions. The errors are slightly worse than with the asymmetric algorithms, highlighting the complexity of finding the symmetric decompositions.

### 5.3.2 Changing the seeds

In the above experiments, we used the columns as the seeds  $S$  for the algorithm (cf. Algorithm 1). This slows the algorithm down, as it has to attempt all of the potential seeds. In this section we study if we can improve the running time without hurting the reconstruction error by sampling only some of the columns for the seed set  $S$ .

In particular, we sampled 10% of the columns uniformly at random to create the seed set. As the algorithm scales linearly with the number of seeds, this provides an order of magnitude speed-up. To test the quality, we repeated the sampling ten times and report the average relative reconstruction errors and standard deviations in Table 3.

The first thing to notice in Table 3 are the low standard deviations; less than 3% in almost all data sets. The reconstruction errors are also only slightly higher than those in Table 2; for instance, `obmf` with *Paleo* has only 6% higher error on average when using random sampling. In most cases the speed-up obtained by the sampling is significant compared to the loss in accuracy.

## 5.4 Visualizing the Graphs

One of the motivations for the ordered BMF is that it allows the convenient visualization of the graphs using edge bundles (or ribbons) between nodes that are placed in a circle. In this section we explore some of these visualizations and explain what we can learn from the respective data sets using them. In the following plots, the edge bundles and the ordering are obtained from the factorization. Further visualizations can be found in Appendix B.

**The *Les Misérables* data:** The visualization of the *Les Misérables* data is presented in Figure 3. Most edge bundles form a circular segment indicating that all of the nodes under the segment are connected to each other (the characters appear in the same parts of the book). Some of the bundles are contained in other bundles, indicating important subset of characters. Multiple bundles intersect on a node at south-east of the circle called *Valjean* – the protagonist of the book.

**The *Mammals* data:** The second data set is the *Mammals* data, in Figure 4. For a clearer visualization, we only consider 134 species that do not appear too frequently in the data, as such species are neighbours of every other



Table 2: Relative errors with asymmetric (left) and symmetric (right) algorithms on real-world data.

	<i>Les Mis</i>	<i>Paleo</i>	<i>News</i>	<i>Terms</i>	<i>Locations</i>	<i>Mammals</i>		<i>Les Mis</i>	<i>Terms</i>	<i>Locations</i>	<i>Mammals</i>
obmf	0.36	0.71	0.74	0.32	0.40	0.26	obmf <sub>sym</sub>	0.40	0.35	0.48	0.27
cobmf	0.36	0.72	0.74	0.32	0.40	0.26	cobmf <sub>sym</sub>	0.41	0.36	0.45	0.27
asso	0.33	0.71	0.72	0.29	0.34	0.26	asso <sub>sym</sub>	0.66	0.33	1.02	0.33

Table 3: Average relative errors and standard deviation with random columns as seeds for asymmetric algorithms on real-world data. Ten random samples.

	<i>Les Misérables</i>	<i>Paleo</i>	<i>Newsgroups</i>	<i>Terms</i>	<i>Locations</i>	<i>Mammals</i>
obmf	0.51 ± 0.08	0.76 ± 0.02	0.83 ± 0.02	0.32 ± 0.00	0.53 ± 0.00	0.26 ± 0.00
cobmf	0.46 ± 0.05	0.76 ± 0.02	0.83 ± 0.02	0.32 ± 0.01	0.51 ± 0.00	0.26 ± 0.00

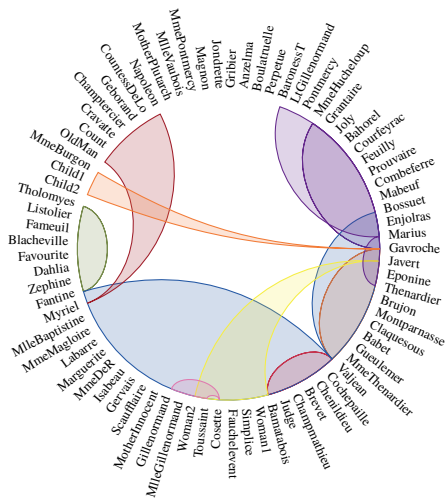


Figure 3: Visualization of the *Les Misérables* data with the ribbons and ordering from cobmf.

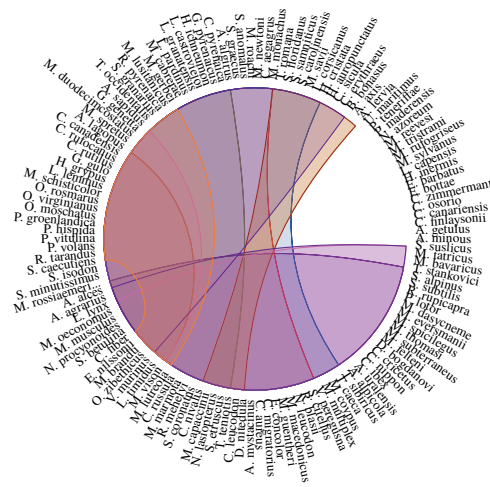


Figure 4: Visualization of the *Mammals* data with the ribbons and ordering from obmf.

species in graph. The edge bundles in Figure 4 are essentially rotating around the middle. This probably corresponds to the change of fauna when moving from north to south. The change is gradual, hence two consecutive edge bundles have a significant overlap, but over longer distance, the change in the fauna becomes more obvious and the edge bundles are more disjoint. This gives a good intuition about the structure of the data.

## 6 Related Work

*Boolean matrix factorization* (BMF) has received increasing interest in the data analysis community [2, 9–17], proving to be a versatile tool for analyzing Boolean matrices. Many different algorithms have been proposed, including algorithms based on candidate creation and selection [12, 15], proximal alternations [10], and message passing [16], to name but a few. It has also found applications in diverse fields, such as bioinformatics [5], information

extraction [4], and lifted inference [18]. To the best of our knowledge, however, the ordering constraint is not studied in earlier work related to Boolean matrix factorization.

*Tiling* databases [6] can be seen as a restricted version of BMF, where the factorization cannot express any 0s as 1. *Geometric tiling* [8] is a variation thereof, where the tiles have to be consecutive. The main difference to our work is a different optimization function, [8] uses log-likelihood, and that it assumes that the order is already given, for example, by spectral ordering, whereas we discover the order on the fly.

A binary matrix has the *consecutive ones property* (C1P) if its columns can be permuted so that all rows have all 1s consecutively. The pq-trees can be used to check for the C1P [3] and Atkins et al. [1] propose spectral ordering algorithm. The spectral ordering approach is used in [8] to permute the data for finding the geometric tiles.

## 7 Conclusions

Ordered Boolean matrix factorization (OBMF) and its variations (COBMF, OBMF<sub>sym</sub>) are restricted versions of Boolean matrix factorization, requiring the factors to have the consecutive ones property (or be cyclic, in case of COBMF). This restriction facilitates the interpretation of the factorization, in particular in the case of the edge bundle visualizations of graphs, as we saw in Section 5.4. On the other hand, the restriction yields higher reconstruction errors, though our experiments show that the difference to state-of-the-art Boolean matrix factorization algorithm is usually very small.

In this paper we laid the theoretical foundations of the OBMF problem and its variations, and proposed algorithms based on the pq-trees. An important part of the proposed algorithm is the choice of the seed vectors. In this paper, we mostly used all columns of the data as the seed, though the experiments in Section 5.3.2 show that sampling the columns could work equally well. An interesting question for the future is whether other methods for selecting the seeds would yield better reconstruction errors.

In the problem setting of this paper, the user provides the rank of the decomposition and the goal is to minimize the reconstruction error over the rank- $k$  OBMF decompositions. A common variant in the Boolean matrix factorization world is to make the rank a free variable and replace

the target function with measure that penalizes for higher ranks (see, e.g. [10, 12, 14]). The *Minimum Description Length* principle is a common approach. The ordered nature of our factor matrices could help with finding more efficient MDL decompositions, as the factor matrices are easier to compress using run-length encoding or similar approaches.

## References

- [1] J. E. Atkins, E. G. Boman, and B. Hendrickson. A Spectral Algorithm for Seriation and the Consecutive Ones Problem. *SIAM J. Comput.*, 28(1):297–310, 1998.
- [2] R. Bělohávek and V. Vychodil. Discovery of optimal factors in binary data via a novel method of matrix decomposition. *J. Comput. Syst. Sci.*, 76(1):3–20, 2010.
- [3] K. S. Booth and G. S. Lueker. Testing for the consecutive ones property, interval graphs, and graph planarity using pq-tree algorithms. *J. Comput. Syst. Sci.*, 13(3):335–379, 1976.
- [4] E. Cergani and P. Miettinen. Discovering relations using matrix factorization methods. In *CIKM '13*, pages 1549–1552, 2013.
- [5] G. Corrado, T. Tebaldi, G. Bertamini, F. Costa, A. Quattrone, G. Viero, and A. Passerini. PTRcombiner: mining combinatorial regulation of gene expression from post-transcriptional interaction maps. *BMC Genomics*, 15(1), Apr. 2014.
- [6] F. Geerts, B. Goethals, and T. Mielikäinen. Tiling databases. In *DS '04*, pages 278–289, 2004.
- [7] N. Gillis and S. A. Vavasis. On the Complexity of Robust PCA and  $\ell_1$ -norm Low-Rank Matrix Approximation. *arXiv*, 2015.
- [8] A. Gionis, H. Mannila, and J. K. Seppänen. Geometric and Combinatorial Tiles in 0–1 Data. In *PKDD '04*, pages 173–184, 2004.
- [9] S. Hess and K. Morik. C-SALT: Mining Class-Specific ALTERations in Boolean Matrix Factorization. In *ECMLPKDD '17*, pages 547–563, 2017.
- [10] S. Hess, K. Morik, and N. Piatkowski. The PRIMPING routine—Tiling through proximal alternating linearized minimization. *Data Min. Knowl. Discov.*, 31(4):1090–1131, May 2017.

- [11] S. Karavev, P. Miettinen, and J. Vreeken. Getting to Know the Unknown Unknowns: Destructive-Noise Resistant Boolean Matrix Factorization. In *SDM '15*, pages 325–333, 2015.
- [12] C. Lucchese, S. Orlando, and R. Perego. A Unifying Framework for Mining Approximate Top-k Binary Patterns. *IEEE Trans. Knowl. Data Eng.*, 26(12): 2900–2913, Dec. 2013.
- [13] S. Maurus and C. Plant. Ternary Matrix Factorization: problem definitions and algorithms. *Knowl. Inf. Syst.*, 46(1):1–31, Jan. 2016.
- [14] P. Miettinen and J. Vreeken. MDL4BMF: Minimum Description Length for Boolean Matrix Factorization. *ACM Trans. Knowl. Discov. Data*, 8(4):–31, Oct. 2014.
- [15] P. Miettinen, T. Mielikäinen, A. Gionis, G. Das, and H. Mannila. The Discrete Basis Problem. *IEEE Trans. Knowl. Data Eng.*, 20(10):1348–1362, Oct. 2008.
- [16] S. Ravanbakhsh, B. Póczos, and R. Greiner. Boolean Matrix Factorization and Noisy Completion via Message Passing. In *ICML '16*, 2016.
- [17] T. Rukat, C. C. Holmes, M. K. Titsias, and C. Yau. Bayesian Boolean Matrix Factorisation. In *ICML '17*, pages 2969–2978, July 2017.
- [18] G. van den Broeck and A. Darwiche. On the Complexity and Approximation of Binary Evidence in Lifted Inference. In *NIPS '13*, pages 2868–2876, 2013.

## A Proofs

*Proof of Theorem 1.* In this case, we are looking for a decomposition of format  $\mathbf{D} \approx \mathbf{x}^T \mathbf{y}$ , where  $\mathbf{D} \in \{0, 1\}^{n \times m}$ ,  $\mathbf{x} \in \{0, 1\}^n$ , and  $\mathbf{y} \in \{0, 1\}^m$ . Notice that (i) whether we use normal or Boolean algebra does not matter in this case; and (ii) we can always find the ordering after we have found the decomposition, as we only need to order the vectors  $\mathbf{x}$  and  $\mathbf{y}$ . But this problem, the *rank-1 binary matrix factorization* problem, is known to be NP-hard [7], finalizing the proof.  $\square$

*Proof of Theorem 2.* The decision problem is obviously in NP.

We prove the hardness by reduction from HAMILTON PATH, where we are given a graph  $G = (V, E)$  and asked whether there is a hamiltonian path, that is, a path visiting every vertex exactly once.

Assume that we are given a graph  $G = (V, E)$  with  $n$  vertices and  $m$  edges. Assume that we have some arbitrary order on the vertices  $V = v_1, \dots, v_n$ , and on the edges  $E = e_1, \dots, e_m$ .

Let us define  $\mathbf{D}$  first. The dataset will be of size  $(n + m + 1)$ -by- $(3m + 1)$ . To define the matrix, we split the rows in two parts  $R = r_1, \dots, r_n$  and  $S = s_0, \dots, s_m$ , containing respectively  $n$  and  $m$  rows. Similarly, we split the columns in 3 parts,  $X = x_1, \dots, x_m$ ,  $Y = y_1, \dots, y_m$ ,  $Z = z_0, \dots, z_m$ .

The 1s in  $\mathbf{D}$  are as follows. for each edge  $e_\ell = (v_i, v_j)$ , we set the cells  $(r_i, x_\ell)$   $(r_j, x_\ell)$   $(r_i, y_\ell)$   $(r_j, y_\ell)$  to be 1. For two adjacent edges  $e_\ell$  and  $e_{\ell+1}$ , we set the cells  $(s_\ell, y_\ell)$   $(s_\ell, z_\ell)$   $(s_\ell, x_{\ell+1})$ . Finally, we set  $(s_0, x_1)$ ,  $(s_0, z_0)$ , and  $(s_m, y_m)$ ,  $(s_m, z_m)$  to be 1. The remaining values are 0.

We argue that there is a zero-error solution for OBMF using  $k = 3m - n + 2$  if and only there is a hamiltonian path.

Let us prove the easy direction: assume that there is a hamiltonian path. To that end, let us permute the rows and columns  $\mathbf{D}$  such that the factor matrices do not have gap zeros. Permute  $\mathbf{D}$  as follows: Set the column order as  $z_0, x_1, y_1, z_1, x_2, y_2, \dots$ . Order the rows in  $R$  according to the hamiltonian path, followed by the rows in  $S$ . We denote the resulting matrix by  $\mathbf{D}'$ . There is a zero-solution if the ones in  $\mathbf{D}'$  are a union of  $k$  contiguous blocks. The  $k$  blocks are as follows:  $m + 1$  blocks covering individual rows in  $S$ ,  $n - 1$  blocks covering edges along the hamiltonian path (this can be done since the corresponding rows in  $R$  and the corresponding columns in  $X$  and  $Y$  are adjacent), and  $2(m - n + 1)$  blocks to cover the remaining edges, 2 blocks per edge. This covers all 1s using  $m+1+n-1+2(m-n+1) = k$  blocks.

Let us prove the other direction. Assume that there is zero-error solution, and let  $\mathbf{D}'$  be the permuted version of  $\mathbf{D}$  with no gap zeros. Then the ones in  $\mathbf{D}'$  must be a union of  $k$  contiguous blocks. For a column index  $i$ , we define  $f_i$  to be the number of blocks started at the  $i$ th column. Let us also define  $g_i$  to be the number of blocks ended at  $i$ th columns. Trivially,  $\sum_i f_i + g_i = 2k$ .

We say that an edge  $(v_i, v_j) \in E$  is *active* if  $i$  and  $j$  are adjacent in  $\mathbf{D}'$ . Let  $h$  be the total number of active edges.

Note that we have  $h \leq n - 1$ . Assume for a moment that  $h = n - 1$  and let  $w_1, \dots, w_n$  be the vertices ordered according to the order of  $R$  in  $\mathbf{D}'$ . Since  $h = n - 1$ , we are forced to have  $(w_i, w_{i+1}) \in E$ . This implies that  $w_1, \dots, w_n$  is a hamiltonian path.

We will now argue that  $h \geq n - 1$ .

Consider two adjacent columns at  $i$  and  $i + 1$ . If none of the columns are in  $Z$ , then both columns contain 1 that is not in the other column. This forces  $g_i + f_{i+1} \geq 2$ . The same argument holds if both columns are in  $Z$ .

Assume that the  $j$ th column is in  $X$  and  $(j + 1)$ th column is in  $Z$ . Assume that  $g_i + f_{i+1} = 1$ . Let  $a$  and  $b$  be the rows in  $R$  that are active in the  $j$ th columns. Since  $Z$  does not have active rows, the block(s) covering  $a$  and  $b$  must terminate, and since  $g_i \geq 1$ , we have only block, implying that  $a$  and  $b$  are adjacent. The same result holds if we replace  $X$  with  $Y$  or permute the order of the two columns. To summarize, if  $g_i + f_{i+1} = 1$ , then either  $i$ th or the  $(i + 1)$ th column corresponds to an active edge.

In addition, we must have  $f_1 \geq 1$  and  $g_{3m+1} \geq 1$  as these columns have 1s. This leads to

$$\begin{aligned} 2(3m - n + 2) &= 2k = \sum_{i=1}^{3m+1} f_i + g_i \\ &= f_1 + g_{3m+1} + \sum_{i=1}^{3m} f_{i+1} + g_i \\ &\geq 2(3m + 1) - 2h, \end{aligned}$$

proving the result.  $\square$

*Proof of Lemma 5.* Let  $S$  be the optimal border-compatible set. Then there is  $i$  such that  $S$  is a union of the best border-compatible set of  $c_i$  and either the union of all leaves in  $c_1, \dots, c_{i-1}$  or  $c_{i+1}, \dots, c_\ell$ .  $\square$

*Proof of Lemma 6.* Let  $S$  be the optimal compatible set. Then  $S$  is either included completely within one child, or there are indices  $i < j$  such that  $S$  is a union of the best border-compatible sets of  $c_i, c_j$ , and the union of all leaves in  $c_{i+1}, \dots, c_{j-1}$ .  $\square$

*Proof of Lemma 7.* Let  $S$  be the optimal border-compatible set. Then there is  $i$  such that  $S$  is a union of the best border-compatible set of  $c_i$  and the union of all leaves of some children.  $\square$

Let  $w$  be a child of  $v$ , if  $total(w) \geq 0$ , then having the leaves of  $w$  in  $S$  has positive gain. Let  $P$  be these children. The total gain corresponds of having these children is  $\sum_i \max(total(v), 0)$ .

We need to transform one of the children to a partial. Let  $w$  be a child of  $v$ . If  $total(w) < 0$ , then  $v \notin P$  and adding  $w$  will have a gain of  $border(w)$ . If  $total(w) \geq 0$ , then  $v \in P$ , and transforming  $w$  from a fully-covered node to a partial node will have a gain of  $border(w) - total(w)$ . In summary, the gain is equal to  $g(w)$ . Thus, selecting the vertex with the maximal  $g(w)$  should be the partial child in  $S$ .  $\square$

*Proof of Lemma 8.* Let  $S$  be the optimal compatible set. Then  $S$  is either included completely within one child, or  $S$  is a union of some children and possibly up to two of the best border-compatible sets for some  $c_i$  and  $c_j$ .

Let  $w$  be a child of  $v$ , if  $total(w) \geq 0$ , then having the leaves of  $w$  in  $S$  has positive gain. Let  $P$  be these children. The total gain corresponds of having these children is  $\sum_i \max(total(v), 0)$ .

As shown in the proof of Lemma 7,  $b_1$  and  $b_2$  correspond the top-2 border-compatible sets. It may happen that  $b_1$  or  $b_2$  are negative, in which case we simply do not add them to  $S$ . Thus the total gain of border-compatible sets is  $\max(b_1, 0) + \max(b_2, 0)$ .  $\square$

## B Further Visualizations

$\square$  Here we present for the *Terms* and *Locations* data sets.

**The Terms data** The visualization of the *Terms* data, in Figure 5, is markedly different from Figure 3. Here, most bundles overlap each other. This indicates that many of these terms are used together in different posts. Yet, we can also identify specialized groups of terms. At the left of Figure 5, we have a blue bundle, from *mission* to *nasa*, that contains terms used when discussing space programs. This overlaps with a larger orange bundle, from *chip* to *tap*, containing terms related to cryptography.

**The Locations data** For the *Locations* data, in Figure 6, we cannot print any labels, as the data consists of 3203 geographical locations. For these results, we did a rank-10 decomposition. Most of the edge bundles again form

