

# Scalable Boolean Tensor Factorizations using Random Walks

Dóra Erdős\*    Pauli Miettinen†

October 21, 2013

## Abstract

Tensors are becoming increasingly common in data mining, and consequently, tensor factorizations are becoming more and more important tools for data miners. When the data is binary, it is natural to ask if we can factorize it into binary factors while simultaneously making sure that the reconstructed tensor is still binary. Such factorizations, called Boolean tensor factorizations, can provide improved interpretability and find Boolean structure that is hard to express using normal factorizations. Unfortunately the algorithms for computing Boolean tensor factorizations do not usually scale well. In this paper we present a novel algorithm for finding Boolean CP and Tucker decompositions of large and sparse binary tensors. In our experimental evaluation we show that our algorithm can handle large tensors and accurately reconstructs the latent Boolean structure.

## 1 Introduction

Tensors, and their factorizations, are getting increasingly popular in data mining. Many real-world data sets can be interpreted as ternary (or higher arity) relations (e.g. sender, receiver, and date in correspondence or object, relation, and subject in RDF data bases or natural language processing). Such relations have a natural representations as 3-way (or higher order) tensors. A data miner who is interested in finding some structure from such a tensor would normally use tensor decomposition methods, commonly either CANDECOP/PARAFAC (CP) or Tucker decomposition (or variants thereof). In both of these methods, the goal is to (approximately) reconstruct the input tensor as a sum of simpler elements (e.g. rank-1 tensors) with the hope that these simpler elements would reveal the latent structure of the data.

The type of these simpler elements plays a crucial role on determining what kind of structure the decomposition will reveal. For example, if the elements

---

\*Boston University, Boston, MA, USA

†Max-Planck-Institut für Informatik, Saarbrücken, Germany

contain arbitrary real numbers, we are finding general linear relations; if the numbers are non-negative, we are finding parts-of-whole representations. In this paper, we study yet another type of structure: that of *Boolean tensor factorizations* (BTF). In BTF, we require the data tensor to be binary, and we also require any matrices and tensors that are part of the decomposition to be binary. Further, instead of normal addition, we use Boolean *or*, that is, we define  $1 + 1 = 1$ . The type of structure found under BTF is different to the type of structure found under normal algebra (non-negative or otherwise). Intuitively, if there are multiple “reasons” for a 1 in the data, under normal algebra and non-negative values, for example, we explain this 1 using a sum of smaller values, but under Boolean algebra, any of these reasons alone is sufficient, and there is no penalty for having multiple reasons. For a concrete example, consider a data that contains noun phrase–verbal phrase–noun phrase patterns extracted from textual data. Underlying this data are the true facts: which entities are connected to which entities by which relations. We see the noun phrase–verbal phrase–noun phrase  $(n_1, v, n_2)$  triple if 1) there is a fact  $(e_1, r, e_2)$ , that is, entity  $e_1$  is connected to entity  $e_2$  via relation  $r$ ; and 2)  $n_1$  is one of the phrases for  $e_1$ ,  $n_2$  is a phrase for  $e_2$ , and  $v$  is a phrase for  $r$ . It does not matter if there is a different “core triple”  $(e'_1, r', e'_2)$  that could also generate the same observed triple as long as there is at least one of them. This kind of model is exactly the Boolean Tucker decomposition (see Section 3), and we will show in the experiments (Section 5) how it performs in this type of data.

We want to emphasize that we do not consider BTF as a replacement of other tensor factorization methods even if the data is binary. Rather, we consider it as an addition to the data miner’s toolbox, letting her to explore another type of structure.

But how do we find a Boolean factorization of a given tensor? There exists algorithms for BTF (e.g. [1, 11, 13]), but they do not scale well. Our main contribution in this paper is to present a scalable algorithm for finding Boolean CP and Tucker decompositions. Further, we apply the minimum description length principle to automatically select the size of the decomposition.

Our algorithm can be divided into many phases. The main work is done by the WALK’N’MERGE algorithm (Section 3), but to obtain proper Boolean CP or Tucker decomposition, we need to apply some post-processing to the output of WALK’N’MERGE (explained in Section 4). We present our experiments in Section 5 and discuss related work in Section 6. Before all of this, however, we present some important definitions.

## 2 Definitions

Before we can present our algorithm, we will explain our notation and formally define the tensor factorization problems we are working with. At the end of this section, we introduce two important concepts, blocks and convex hulls, that will be used extensively in the algorithm

## 2.1 Notation

Throughout this paper vectors are indicated as bold-face lower-case letters ( $\mathbf{v}$ ), matrices as bold-face upper-case letters ( $\mathbf{M}$ ), and tensors as bold-face upper-case calligraphic letters ( $\mathcal{T}$ ). We present the notation for 3-way tensors, but it can be extended to  $N$ -way tensors in a straight forward way. Element  $(i, j, k)$  of a 3-way tensor  $\mathcal{X}$  is denoted either as  $x_{ijk}$  or as  $(\mathcal{X})_{ijk}$ . A colon in a subscript denotes taking that mode entirely; for example, if  $\mathbf{X}$  is a matrix,  $\mathbf{x}_i$  denotes the  $i$ th row of  $\mathbf{X}$  (for a shorthand, we use  $\mathbf{x}_j$  to denote the  $j$ th column of  $\mathbf{X}$ ). For a 3-way tensor  $\mathcal{X}$ ,  $\mathbf{x}_{:jk}$  is the  $(j, k)$  mode-1 (column) fiber,  $\mathbf{x}_{i:k}$  the  $(i, k)$  mode-2 (row) fiber, and  $\mathbf{x}_{ij:}$  the  $(i, j)$  mode-3 (tube) fiber. Furthermore,  $\mathbf{X}_{::k}$  is the  $k$ th frontal slice of  $\mathcal{X}$ . We use  $\mathbf{X}_k$  as a shorthand for the  $k$ th frontal slice.

For a tensor  $\mathcal{X}$ , the number of non-zero elements in it is denoted by  $|\mathcal{X}|$ . The Frobenius norm of a 3-way tensor  $\mathcal{X}$ ,  $\|\mathcal{X}\|$ , is defined as  $\sqrt{\sum_{i,j,k} x_{ijk}^2}$ . If  $\mathcal{X}$  is binary, i.e. takes values only from  $\{0, 1\}$ ,  $|\mathcal{X}| = \|\mathcal{X}\|^2$ .

The *tensor sum* of two  $n$ -by- $m$ -by- $l$  tensors  $\mathcal{X}$  and  $\mathcal{Y}$  is the element-wise sum,  $(\mathcal{X} + \mathcal{Y})_{ijk} = x_{ijk} + y_{ijk}$ . The *Boolean tensor sum* of binary tensors  $\mathcal{X}$  and  $\mathcal{Y}$  is defined as  $(\mathcal{X} \vee \mathcal{Y})_{ijk} = x_{ijk} \vee y_{ijk}$ .

The outer product of vectors in  $N$  modes is denoted by  $\boxtimes$ . That is, if  $\mathbf{a}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$  are vectors of length  $n$ ,  $m$ , and  $l$ , respectively,  $\mathcal{X} = \mathbf{a} \boxtimes \mathbf{b} \boxtimes \mathbf{c}$  is an  $n$ -by- $m$ -by- $l$  tensor with  $x_{ijk} = a_i b_j c_k$ . A tensor that is an outer product of three vectors has *tensor rank* 1.

Finally, if  $\mathcal{X}$  and  $\mathcal{Y}$  are binary  $n$ -by- $m$ -by- $l$  tensors, we say that  $\mathcal{Y}$  *contains*  $\mathcal{X}$  if  $x_{ijk} = 1$  implies  $y_{ijk} = 1$  for all  $i, j$ , and  $k$ . This relation defines a partial order of  $n$ -by- $m$ -by- $l$  binary tensors, and it is therefore understood that when we say that  $\mathcal{X}$  is the *smallest*  $n$ -by- $m$ -by- $l$  binary tensor for which some property  $P$  holds, we mean that there exists no other  $n$ -by- $m$ -by- $l$  binary tensors for which  $P$  holds and that are contained in  $\mathcal{X}$ .

## 2.2 Ranks and Factorizations

With the basic notation explained, we first define the CP decomposition and tensor rank under the normal algebra, after which we explain how the Boolean concepts differ. After that we define the Boolean Tucker decomposition.

**Tensor Rank and CP Decomposition.** The so-called CP factorization,<sup>1</sup> we are studying in this paper is defined as follows

**Problem 1** (CP decomposition). Given tensor  $\mathcal{X}$  of size  $n$ -by- $m$ -by- $l$  and an integer  $r$ , find matrices  $\mathbf{A}$  ( $n$ -by- $r$ ),  $\mathbf{B}$  ( $m$ -by- $r$ ), and  $\mathbf{C}$  ( $l$ -by- $r$ ) such that they minimize

$$\left\| \mathcal{X} - \sum_{i=1}^r \mathbf{a}_i \boxtimes \mathbf{b}_i \boxtimes \mathbf{c}_i \right\|^2. \quad (1)$$

<sup>1</sup>The name is short for two names given to the same decomposition: CANDECOMP [2] and PARAFAC [7].

Notice that the  $i$ th columns of the *factor matrices*  $\mathbf{A}$ ,  $\mathbf{B}$ , and  $\mathbf{C}$  define a rank-1 tensor  $\mathbf{a}_i \boxtimes \mathbf{b}_i \boxtimes \mathbf{c}_i$ . In other words, the CP decomposition expresses the given tensor as a sum of  $r$  rank-1 tensors.

Using the CP decomposition, we can define the tensor rank analogous to the matrix (Schein) rank as the smallest  $r$  such that the tensor can be exactly decomposed into a sum of  $r$  rank-1 tensors. Note that, unlike the matrix rank, computing the tensor rank is NP-hard [8].

**The Boolean Tensor Rank and Decompositions.** The Boolean versions of tensor rank and CP decomposition are rather straight forward to define given their normal counterparts. One only needs to change the summation to  $1 + 1 = 1$ . Notice that this does not change the definition of a rank-1 tensor (or vector outer product). Thence, a 3-way Boolean rank-1 tensor is a tensor that is an outer product of three binary vectors.

**Definition 1** (Boolean tensor rank). The *Boolean rank* of a 3-way binary tensor  $\mathcal{X}$ ,  $\text{rank}_B(\mathcal{X})$ , is the least integer  $r$  such that there exist  $r$  rank-1 binary tensors with

$$\mathcal{X} = \bigvee_{i=1}^r \mathbf{a}_i \boxtimes \mathbf{b}_i \boxtimes \mathbf{c}_i . \quad (2)$$

The Boolean CP decomposition follows analogously. Instead of subtraction, we take the element-wise exclusive or (denoted by  $\oplus$ ), and instead of sum of squared values, we simply count the number of non-zero elements in the residual. Note, however, that with all-binary data, our error function is equivalent to the squared Frobenius error.

**Problem 2** (Boolean CP decomposition). Given an  $n$ -by- $m$ -by- $l$  binary tensor  $\mathcal{X}$  and an integer  $r$ , find binary matrices  $\mathbf{A}$  ( $n$ -by- $r$ ),  $\mathbf{B}$  ( $m$ -by- $r$ ), and  $\mathbf{C}$  ( $l$ -by- $r$ ) such that they minimize

$$\left| \mathcal{X} \oplus \left( \bigvee_{i=1}^r \mathbf{a}_i \boxtimes \mathbf{b}_i \boxtimes \mathbf{c}_i \right) \right| . \quad (3)$$

Analogous to the normal CP decomposition, the Boolean CP decomposition can be seen as a (Boolean) sum of  $r$  binary rank-1 tensors. Unsurprisingly, both finding the Boolean rank of a tensor and finding its minimum-error rank- $r$  Boolean CP decomposition are NP-hard [13].

**Boolean Tucker decompositions.** Given a (binary) tensor, its Tucker decomposition contains a *core tensor* and three factor matrices. The number of rows in the factor matrices are defined by the dimensions of the original tensor while the number of columns in them are defined by the dimensions of the core tensor. In case of the Boolean Tucker decomposition, all involved tensors and matrices are required to be binary, and the arithmetic is again done over the Boolean semi-ring. The Boolean Tucker decomposition is defined formally as follows.

**Problem 3** (Boolean Tucker decomposition). Given an  $n$ -by- $m$ -by- $l$  binary tensor  $\mathcal{X} = (x_{ijk})$  and three integers  $p$ ,  $q$ , and  $r$ , find the minimum-error  $(p, q, r)$

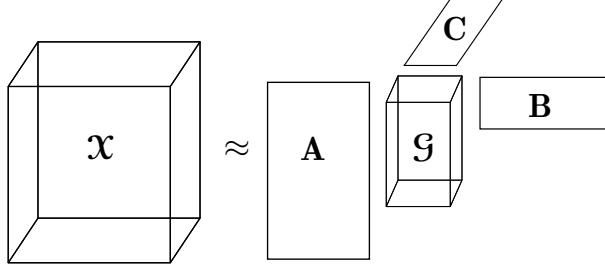


Figure 1: Tucker tensor decomposition.

*Boolean Tucker decomposition* of  $\mathcal{X}$ , that is, tuple  $(\mathcal{G}, \mathbf{A}, \mathbf{B}, \mathbf{C})$ , where  $\mathcal{G}$  is a  $p$ -by- $q$ -by- $r$  binary *core tensor* and  $\mathbf{A}$  ( $n$ -by- $p$ ),  $\mathbf{B}$  ( $m$ -by- $q$ ), and  $\mathbf{C}$  ( $l$ -by- $r$ ) are binary *factor matrices*, such that  $(\mathcal{G}, \mathbf{A}, \mathbf{B}, \mathbf{C})$  minimizes

$$\sum_{i,j,k} \left( x_{ijk} \oplus \left( \bigvee_{\alpha=1}^p \bigvee_{\beta=1}^q \bigvee_{\gamma=1}^r g_{\alpha\beta\gamma} a_{i\alpha} b_{j\beta} c_{k\gamma} \right) \right). \quad (4)$$

For a schematic view of Tucker decomposition, see Figure 1.

### 2.3 Blocks, Convex Hulls, and Factorizations

Let  $\mathcal{X}$  be a binary  $n$ -by- $m$ -by- $l$  tensor and let  $X \subseteq [n]$ ,  $Y \subseteq [m]$ , and  $Z \subseteq [l]$ , where  $[x] = \{1, 2, \dots, x\}$ . A *block* of  $\mathcal{X}$  is a  $|X|$ -by- $|Y|$ -by- $|Z|$  sub-tensor  $\mathcal{B}$  that is formed by taking the rows of  $\mathcal{X}$  defined by  $X$ , columns defined by  $Y$ , and tubes defined by  $Z$ . Block  $\mathcal{B}$  is *monochromatic* if all of its values are 1. We will often (implicitly) embed  $\mathcal{B}$  to  $n$ -by- $m$ -by- $l$  tensor by filling the missing values with 0s. If  $\mathcal{B}$  is monochromatic it is (embedded or not) a rank-1 tensor. If  $\mathcal{B}$  is not monochromatic, we say it is *dense*.

Now let the sets  $I$ ,  $J$ , and  $K$  be such that they contain the indices of all the non-zero slices of  $\mathcal{X}$ . That is,  $I = \{i : x_{ijk} = 1 \text{ for some } j, k\}$ ,  $J = \{j : x_{ijk} = 1 \text{ for some } i, k\}$ , and  $K = \{k : x_{ijk} = 1 \text{ for some } i, j\}$ . The *convex hull* of  $\mathcal{X}$  is a binary  $n$ -by- $m$ -by- $l$  tensor  $\mathcal{Y}$  that has 1 in every position defined by the Cartesian product of  $I$ ,  $J$ , and  $K$ ,  $I \times J \times K = \{(i, j, k) : i \in I, j \in J, k \in K\}$ .

The following lemma will explain the connection between monochromatic blocks (rank-1 tensors) and convex hulls. We utilize this lemma throughout our algorithms by searching for convex blocks rather than explicitly rank-1 tensors in the data.

**Lemma 1.** *Let  $\mathcal{X}$  be a binary  $n$ -by- $m$ -by- $l$  tensor. Then the convex hull of  $\mathcal{X}$  is the smallest  $n$ -by- $m$ -by- $l$  rank-1 binary tensor that contains  $\mathcal{X}$ .*

*Proof.* Let us start by showing that the convex hull of  $\mathcal{X}$  is indeed a rank-1 tensor. To that end, let  $I$ ,  $J$ , and  $K$  be the sets of indices of slices of  $\mathcal{X}$  that have 1s in them (i.e.  $I = \{i : x_{ijk} = 1 \text{ for some } j, k\}$  and similarly for  $J$  and  $K$ ). If  $\mathcal{Y}$  is the convex hull of  $\mathcal{X}$ , by definition  $y_{ijk} = 1$  if and only if  $(i, j, k) \in I \times J \times K$ . Let us now define three binary vectors,  $\mathbf{a}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$  (of dimensions  $n$ ,  $m$ , and  $l$ , respectively). Let  $a_i = 1$  if and only if  $i \in I$ ,  $b_j = 1$  if and only if  $j \in J$ , and  $c_k = 1$  if and only if  $k \in K$ . Then the outer product  $\mathbf{a} \boxtimes \mathbf{b} \boxtimes \mathbf{c}$  has 1 at position  $(i, j, k)$  if and only if  $(i, j, k) \in I \times J \times K$ , that is  $\mathcal{Y} = \mathbf{a} \boxtimes \mathbf{b} \boxtimes \mathbf{c}$ .

That  $\mathcal{Y}$  contains  $\mathcal{X}$  is straight forward to see. This means we only have to prove that there exists no other tensor that is rank-1, contains  $\mathcal{X}$ , and is contained in  $\mathcal{Y}$ . Assume, for a contradiction, that  $\mathcal{Z} \neq \mathcal{Y}$  is such. Then, it has to be that there is a location  $(i, j, k)$  for which  $x_{ijk} = z_{ijk} = 0$  but  $y_{ijk} = 1$ . As  $\mathcal{Z}$  is rank-1, we can represent it as  $\mathcal{Z} = \mathbf{a} \boxtimes \mathbf{b} \boxtimes \mathbf{c}$  for some  $\mathbf{a}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$ . As  $z_{ijk} = 0$ , it must be that  $a_i b_j c_k = 0$ , that is, one of the three elements is 0. Let  $c_k = 0$  (other cases are analogous). This means that the slice  $\mathcal{Z}_{::k}$  is empty. As  $\mathcal{Z}$  contains  $\mathcal{X}$ , also  $\mathcal{X}_{::k}$  must be empty. But this is a contradiction, as  $y_{ijk} = 1$  only if  $k \in K$ , and therefore  $\mathcal{X}_{::k}$  cannot be empty.  $\square$

As a corollary to Lemma 1 we get that  $\mathcal{X}$  is rank-1 if and only if it is its own convex hull.

**Blocks and factorizations.** The key observation underlying our algorithms is the fact that both the CP as the Tucker decomposition can be thought of as a decomposition of the data tensor  $\mathcal{X}$  into some combination of rank-1 sub-tensors. While this is obvious in case of the CP decomposition, it is also easy to see for the Tucker: every triplet of factors  $a_{\alpha}$ ,  $b_{\beta}$ , and  $c_{\gamma}$  where  $g_{\alpha\beta\gamma}$  is non-zero defines such a rank-1 tensor. The main idea of our algorithm, which we will explain next, is to find dense blocks from the input data, construct their convex hulls, and build the Boolean CP or Tucker factorization from the resulting rank-1 tensors.

### 3 The Walk'n'Merge Algorithm

In this section we present the main part of our algorithm, WALK'N'MERGE, that aims to find the dense blocks from which we build the factorizations (how that is done is explained in the next section). The WALK'N'MERGE algorithm contains two phases. The first phase, RANDOMWALK, aims at finding and removing the most prominent blocks quickly from the tensor. The second phase, BLOCKMERGE, uses these blocks together with smaller, easier-to-find monochromatic blocks and tries to merge them into bigger blocks.

#### 3.0.1 Random walk algorithm

In this phase we represent the tensor  $\mathcal{X}$  with a graph  $G(V, E)$  that is defined as follows. For every  $x_{ijk} = 1$  we have a node  $v_{ijk} \in V$ . Two nodes  $v_{ijk}$  and  $v_{pqr}$  are connected by an edge  $(v_{ijk}, v_{pqr})$  if  $(i, j, k)$  and  $(p, q, r)$  differ in exactly one coordinate.

---

**Algorithm 1** Random walk algorithm to find blocks.

---

**Input:**  $\mathcal{X}$ ,  $d$ , walk\_length, num\_walks, freq

**Output:**  $\mathcal{B}_1, \mathcal{B}_2 \dots \mathcal{B}_k$

- 1: create graph  $G(V, E)$  from  $\mathcal{X}$
- 2: **while**  $V$  is not empty **do**
- 3:      $v \leftarrow$  random node from  $V$
- 4:     visitedNodes  $\leftarrow (v, count_v = 1)$
- 5:     **for** num\_walks number of times **do**
- 6:          $v_{vis} \leftarrow$  random node from visitedNodes
- 7:         **for** walk\_length number of times **do**
- 8:              $v' \leftarrow$  random neighbor of  $v_{vis}$
- 9:             visitedNodes  $\leftarrow (v', count_{v'} ++)$
- 10:      $\mathcal{B} \leftarrow$  empty block
- 11:     **for**  $v \in$  visitedNodes **do**
- 12:         **if**  $count_v > \text{freq}$  **then**
- 13:              $\mathcal{B} \leftarrow v$
- 14:      $V \setminus \text{convex\_hull}(\mathcal{B})$
- 15:     block  $\mathcal{B}$  is the convex hull of nodes in  $\mathcal{B}$
- 16:     **if** density of  $\mathcal{B} > d$  **then**
- 17:         add  $\mathcal{B}$  to blocks
- 18: **return** blocks

---

Observe that a node  $v_{ijk}$  is connected to all nodes in  $V$  that are in the same fiber as  $v_{ijk}$  in any mode of  $\mathcal{X}$ . Moreover, a monochromatic block in  $\mathcal{X}$  corresponds to a subgraph of  $G$  with radius at most 3. In case of noisy data, blocks are not perfectly monochromatic and some of the nodes in  $V$  might be missing. Still, if the blocks are fairly dense, the radius of the corresponding subgraph is not too big. More precisely, if  $v_{ijk}$  is a node that participates in a block of density  $d$ , the probability of a random neighbor of  $v_{ijk}$  also participating in that block is  $\frac{d}{d+d'}$ , where  $d'$  is the density of the full tensor. This observation implies that if the blocks are significantly denser than the noisy part of a tensor, then a random neighbor of a node inside block  $\mathcal{B}$  is with high probability also in  $\mathcal{B}$ . Our first algorithm exploits this property by performing short random walks in  $G$ . The intuition is that if such a walk hits a node in a block, then with high probability the consecutive hops in this walk are also hitting the block.

The pseudo code for our RANDOMWALK algorithm is given in Algorithm 1. RANDOMWALK takes as an input the data tensor  $\mathcal{X}$  and parameters controlling the length and number of the random walks, and the minimum density of the resulting blocks. After creating the graph  $G(V, E)$  it finds a block  $\mathcal{B}$  in every iteration of the algorithm by means of executing random walks. Nodes that have been assigned to  $\mathcal{B}$  are removed from  $V$ , resulting in a smaller graph  $G'(V - V_{\mathcal{B}}, E')$  on which the subsequent random walks are executed.

The block  $\mathcal{B}$  is found by way of executing a number of random walks on  $G$ . The first walk is initiated from a random node in  $V$ . For every node we maintain

a counter for the number of times any of the walks has visited that node. For any consecutive walk, we pick a random starting point among those nodes that already have a positive counter. This ensures that once we hit a block  $\mathcal{B}$  with a walk, the consecutive walks start with higher and higher probability from within that block. The length of the walks is given as an input to the algorithm. In order to traverse as big part of  $\mathcal{B}$  as possible, we make many short walks. Since we know that nodes corresponding to a dense block  $\mathcal{B}$  have with high probability a higher visit count than nodes corresponding to noise, we abandon all nodes with visit counts less than the average.

In order to make sure that the block we find is a rank-1 tensor (and to include those nodes we might have missed in the random walk), we take  $\mathcal{B}$  to be the convex hull of the discovered frequent nodes. Finally, we accept  $\mathcal{B}$  only if it has density above a user-specified threshold  $d$ . Before proceeding with the next iteration of RANDOMWALK we remove all nodes corresponding to  $\mathcal{B}$ , regardless of whether  $\mathcal{B}$  was accepted.

**Running time of RandomWalk.** The crux of this algorithm is that the running time of every iteration of the algorithm is fixed and depends only on the number and length of the walks. How often we have to re-start the walks depends on how quickly we remove the nodes from the graph, but the worst-case running time is bound by  $O(|V|) = O(|\mathcal{X}|)$ . However, if  $\mathcal{X}$  contains several dense blocks, then the running time is significantly less, since all nodes corresponding to cells in the block are removed at the same time.

**Paralellization.** RANDOMWALK is easily paralellizable, as we can start the random walk iterations (in Line 3 in Algorithm 1) from several (non-neighboring) nodes at the same time. In this case it may happen that some indices are chosen in multiple blocks. We don't mind that (as  $\mathcal{X}$  may contain partially overlapping blocks) and simply return all resulting blocks.

### 3.0.2 BlockMerge Algorithm

The RANDOMWALK algorithm is a fast method, but it is only able to reliably find the most prominent blocks. If a block is too small, the random walks might visit it as a part of a bigger sparse (and hence rejected) block. It can also happen that while most part of a block is found by RANDOMWALK, due to the randomness in the algorithm, some of its slices are not discovered.

Therefore we present the second part of our algorithm, BLOCKMERGE, that executes two tasks. First it finds smaller monochromatic blocks that for some reason are undiscovered. After finding the smaller blocks, the algorithm has a merging phase, where it tries to merge some of the newly found blocks and the dense blocks found by the RANDOMWALK algorithm. The output of BLOCKMERGE is a set of dense blocks.

The BLOCKMERGE algorithm is akin to normal bottom-up frequent itemset mining algorithms in that it starts with elementary blocks and advances by merging these elementary blocks into bigger blocks, although without the benefit of anti-monotonicity.



The input for BLOCKMERGE is the same data tensor  $\mathcal{X}$  given to the RANDOMWALK algorithm, the blocks already found, and the minimum density  $d$ . As its first step, the algorithm will find all *non-trivial* monochromatic blocks of  $\mathcal{X}$  that are not yet included in any of the blocks found earlier. A monochromatic block is non-trivial if its volume and dimensions are above some user-defined thresholds (e.g. all modes have at least 2 dimensions). We find these non-trivial blocks in a greedy fashion. We start with *singletons*: elements  $x_{ijk} = 1$  that do not belong into any block. We pick one of them,  $x_{ijk}$ , and find all singletons that share at least one coordinate with it. Among these singletons we do an exhaustive search to find all monochromatic non-trivial blocks containing  $x_{ijk}$ . As a result, for every cell that is included in a non-trivial block in  $\mathcal{X}$ , we find at least one monochromatic block it is included in, but we may not find all of them. In our implementation we maintain some practical data indices based on the coordinates defining the cells  $x_{ijk}$  so that looking up neighbors of a cell takes at most  $O(n + m + l)$  time. Since the singleton blocks that remain after the initialization step could not be incorporated in any of the non-trivial blocks, we regard them as noise, and will not consider them for merging to any other block.

The second part of the BLOCKMERGE algorithm is to try and merge the remaining blocks so that we get larger (usually not monochromatic, but still dense) blocks. Each block  $\mathcal{B}$  is defined by three sets of indices,  $I$ ,  $J$ , and  $K$ , giving the row, column, and tube indices of this block. When we merge two blocks,  $\mathcal{B}$  and  $\mathcal{C}$ , with indices given by  $(I_{\mathcal{B}}, J_{\mathcal{B}}, K_{\mathcal{B}})$  and  $(I_{\mathcal{C}}, J_{\mathcal{C}}, K_{\mathcal{C}})$ , respectively, the resulting block  $\mathcal{B} \boxplus \mathcal{C}$  has its indices given by  $(I_{\mathcal{B}} \cup I_{\mathcal{C}}, J_{\mathcal{B}} \cup J_{\mathcal{C}}, K_{\mathcal{B}} \cup K_{\mathcal{C}})$ . (This is equivalent on taking the convex hull of  $\mathcal{B} \vee \mathcal{C}$ , ensuring again that the block is rank-1.)

The way we merge two blocks means that the resulting block can, and typically will, include elements that were not in either of the merged blocks. Therefore, when deciding whether to merge two blocks, we must look how well we do in those areas that are not in either of the blocks. To that end, we will again employ the user-defined density parameter  $d$ . We will only merge two blocks if the joint density of 1s and elements already included in the other blocks in the area not in either of merged blocks is higher than  $d$ .

To present the above consideration more formally, let  $\mathcal{A}$  and  $\mathcal{B}$  be the two blocks we are currently considering to merge.  $\overline{\mathcal{A} \vee \mathcal{B}} = (\mathcal{A} \boxplus \mathcal{B}) \setminus (\mathcal{A} \cup \mathcal{B})$  is the area (monochromatic sub-tensor) in  $\mathcal{A} \boxplus \mathcal{B}$  that is not in either  $\mathcal{A}$  or  $\mathcal{B}$ , and  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_R$  are the rest of the non-trivial blocks we have build so far, then what we compute is the density of 1s in  $\bigvee_{r=1}^R \mathcal{D}_r \vee \mathcal{X}$  in those locations that are 1s in  $\overline{\mathcal{A} \vee \mathcal{B}}$ , that is,

$$\frac{\sum_{i,j,k} ((\overline{\mathcal{A} \vee \mathcal{B}})_{ijk} (\bigvee_{r=1}^R \mathcal{D}_r \vee \mathcal{X})_{ijk})}{\sum_{i,j,k} (\overline{\mathcal{A} \vee \mathcal{B}})_{ijk}}. \quad (5)$$

The reason for including the other blocks  $\mathcal{D}_r$  in the equation is that we do not want to pay multiple times for the same error. Recall that our representation of  $\mathcal{X}$  after  $\mathcal{X}$  and  $\mathcal{Y}$  are merged will be  $\bigvee_{r=1}^R \mathcal{D}_r \vee (\mathcal{X} \boxplus \mathcal{Y})$ , and hence, if we already have expressed some 0 of  $\mathcal{X}$  by 1 in one of the  $\mathcal{D}_r$ 's, this error is already

done, and cannot be revoked. Similarly, whatever error we will do in  $\mathcal{X}$  or  $\mathcal{Y}$ , we will still do in  $\mathcal{X} \boxplus \mathcal{Y}$ , and therefore we only consider the area not included in either of the merged tensor.

Now the only remaining question is how to select which blocks to merge. A simple answer would be to try all possible pairs and select the best. That, however, would require us to compute quadratic number of possible merges, which in practice is too much. Instead we restrict our attention to pairs of blocks that share coordinates in at least one mode, that is, if  $(I_{\mathcal{B}}, J_{\mathcal{B}}, K_{\mathcal{B}})$  and  $(I_{\mathcal{C}}, J_{\mathcal{C}}, K_{\mathcal{C}})$  are as above, we would consider merging  $\mathcal{B}$  and  $\mathcal{C}$  only if at least one of the sets  $I_{\mathcal{B}} \cap I_{\mathcal{C}}$ ,  $J_{\mathcal{B}} \cap J_{\mathcal{C}}$ , or  $K_{\mathcal{B}} \cap K_{\mathcal{C}}$  is non-empty. We call a pair of blocks for which  $I_{\mathcal{B}} \cap I_{\mathcal{C}} = J_{\mathcal{B}} \cap J_{\mathcal{C}} = K_{\mathcal{B}} \cap K_{\mathcal{C}} = \emptyset$  *independent*.

It is worth asking will this restriction mean we will not find all the meaningful blocks. We argue that it does not. The intuition is the following. Let  $\mathcal{B}$  and  $\mathcal{C}$  be the two blocks we should merge but that are independent and let their index sets be as above. If we would merge them, the majority of the volume of the new block would be outside of  $\mathcal{B}$  or  $\mathcal{C}$  ( $(|I_{\mathcal{B}}| + |I_{\mathcal{C}}|)(|J_{\mathcal{B}}| + |J_{\mathcal{C}}|)(|K_{\mathcal{B}}| + |K_{\mathcal{C}}|) - |I_{\mathcal{B}}||J_{\mathcal{B}}||K_{\mathcal{B}}| - |I_{\mathcal{C}}||J_{\mathcal{C}}||K_{\mathcal{C}}|$ , to be exact). If this area is very sparse, then so will be the whole block, and we should not have merged the two original block, after all. But if parts of that area are dense, we should find there another block that shares co-ordinates with both  $\mathcal{B}$  and  $\mathcal{C}$ . If that block is large and dense enough, we will merge it with either  $\mathcal{B}$  or  $\mathcal{C}$ , at which point these two blocks do share co-ordinates, and we will consider them for merging.

This, then, is how we proceed: for every block, count how good a merge it would be with every other block with shared coordinates, select the best merge and execute it, put the merged block back at the bottom of the list of blocks to consider and pick up the next block from the list until no new merges are possible. This means that we execute as many merges as possible in a single sweep of the list of the blocks, as opposed to making a merge and starting again from the begin of the list, as we consider this the faster way to perform the merges. An overview of the whole merging algorithm is presented in Algorithm 2.

This part of the algorithm can be implemented parallel as well; first pick a block  $\mathcal{B}$  and choose all blocks that share a cell with  $\mathcal{B}$ . These blocks are the candidates to merge with  $\mathcal{B}$  in this iteration. Now, among the remaining blocks that are not candidates we choose another  $\mathcal{B}'$  and find the candidates for  $\mathcal{B}'$ . We repeat this until there are no unchosen blocks left. The processing of the candidate lists to find potential merges can then be executed in parallel. Observe that the set of candidates for different blocks may overlap. We don't regard this as a problem and if this happens, then (provided density constraints are met) the block is simply merged to multiple blocks.

**Parallelization.** The merging phase of the BLOCKMERGE algorithm can easily be parallelized as well. Since we only merge blocks that are not independent, it is an obvious choice to parallelize the merging procedure of independent blocks. In every iteration we find a maximal set of independent blocks in a greedy fashion; we pick a block  $\mathcal{B}_1$ , then we pick a block  $\mathcal{B}_2$  from those that are independent of  $\mathcal{B}_1$ , etc. We then consider possible merges for  $\mathcal{B}_1, \mathcal{B}_2 \dots$  with the remaining blocks in parallel. Note that any block can be considered for merge in more

---

**Algorithm 2** BLOCKMERGE algorithm for merging blocks.

---

**Input:** Data  $\mathcal{X}$ , threshold  $d$ , blocks  $B = \{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_r\}$  from random walk

**Output:** Final blocks  $\mathcal{B}_1, \mathcal{B}_2 \dots \mathcal{B}_k$

- 1: find all non-trivial monochromatic blocks  $\mathcal{B}$  of size at least 2-by-2-by-2 not included in blocks in  $B$
- 2: **for**  $\mathcal{B}$  is a non-trivial monochromatic block **do**
- 3:     add  $\mathcal{B}$  to  $B$
- 4: let  $Q$  be a queue of all the blocks in  $B$
- 5: **while**  $Q$  is not empty **do**
- 6:      $\mathcal{B} \leftarrow Q.\text{pop}$
- 7:     **for all**  $\mathcal{C}$  that shares co-ordinates with  $\mathcal{B}$  in at least one mode **do**
- 8:         compute the density of  $\mathcal{B} \boxplus \mathcal{C}$
- 9:         **if** density  $> d$  **then**
- 10:              $Q.\text{push}(\mathcal{B} \boxplus \mathcal{C})$
- 11:             replace  $\mathcal{B}$  and  $\mathcal{C}$  in  $B$  with  $\mathcal{B} \boxplus \mathcal{C}$
- 12:             **break**
- 13: **return**  $B$

---

than one of the threads and as a result may end up being merged with several different blocks.

**Running time of the BlockMerge algorithm.** Let the densest fiber in  $\mathcal{X}$  have  $b = \max\{n, m, l\} \times d$  ones. Observe that any nontrivial monochromatic block is defined exactly by 2 of its cells. Thus for a cell  $x_{ijk}$  we can compute all nontrivial monochromatic blocks containing it in  $b^2$  time by checking all blocks defined by pairs of ones in fibers  $i, j$  and  $k$ . This checking takes constant time. Hence, the first part of the algorithm takes  $O(Bb^2)$  time if there are  $B$  trivial blocks in the data. In worst case  $B = |\mathcal{X}|$ . The second part of the algorithm is the actual merging of blocks. If there are  $D$  blocks at the begin of this phase, we will try at most  $\binom{D}{2}$  merges. The time it takes to check whether to merge depends on the size of the two blocks involved. Executing the merge  $\mathcal{A} = \mathcal{B} \boxplus \mathcal{C}$  takes at most  $|\mathcal{A}|$  time. In worst case  $|\mathcal{A}| = |\mathcal{X}|$ . As a result, a very crude upper bound on the running time can be given as  $O(|\mathcal{X}|(b^2 + D^3))$ .

## 4 From Blocks to Factorizations

The WALK'N'MERGE algorithm only returns us a set of rank-1 tensors, corresponding to dense blocks in the original tensor. To obtain the final decompositions, we will have to do some additional post-processing.

### 4.1 Ordering and Selecting the Final Blocks for the CP-decomposition

We can use all the blocks returned by WALK'N'MERGE to obtain a Boolean CP factorization. The rank of this factorization, however, cannot be controlled, as

it is the number of blocks WALK'N'MERGE returned. Furthermore, it can be that it is better to not use all these blocks but only a subset of them. Ideally, therefore, we would like to be able to select a subset of the blocks such that together they give the CP-decomposition that minimizes the error. It turns out, however, that even this simple selection task is computationally very hard.

*Proposition 4.1.* Given a binary  $n$ -by- $m$ -by- $l$  tensor  $\mathcal{X}$ , and a set  $B$  of  $r$  binary rank-1 tensors of the same size (blocks),  $B = \{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_r\}$ , it is NP-hard to select  $B^* \subset B$  such that

$$\left| \mathcal{X} \oplus \bigvee_{\mathcal{B} \in B^*} \mathcal{B} \right| \quad (6)$$

is minimized. Furthermore, for any  $\varepsilon > 0$ , it is quasi-NP-hard to approximate (6) to within  $\Omega\left(2^{(4 \log r)^{1-\varepsilon}}\right)$  and NP-hard to approximate it to within  $\Omega\left(2^{\log^{1-\varepsilon} |\mathcal{X}|}\right)$ .

*Proof.* For the proof, we need the following result: Consider the Basis Usage (BU) problem [14], where we are given a binary  $n$ -dimensional vector  $\mathbf{a}$  and a binary  $n$ -by- $r$  matrix  $\mathbf{B}$ , and the task is to find a binary  $r$ -dimensional vector  $\mathbf{x}$  such that we minimize the Hamming distance between  $\mathbf{a}$  and  $\mathbf{B} \circ \mathbf{x}$  (where  $\circ$  is the matrix product with Boolean addition). This problem is NP-hard to approximate within  $\Omega\left(2^{\log^{1-\varepsilon} |\mathcal{X}|}\right)$  and quasi-NP-hard to approximate within  $\Omega\left(2^{(4 \log r)^{1-\varepsilon}}\right)$  [12].

The BU problem is equivalent to the problem of selecting the blocks: Take the tensor  $\mathcal{X}$  and write it as a long ( $nml$ -dimensional) binary vector. This will be the vector  $\mathbf{a}$  of the BU problem. Vectorize the blocks in  $B$  in the same way; these will be the columns of  $\mathbf{B}$  in the BU problem. Now, the Boolean product  $\mathbf{B} \circ \mathbf{x}$  is equivalent to taking those columns of  $\mathbf{B}$  for which the corresponding row of  $\mathbf{x}$  is 1 and taking their Boolean sum. But this is the same as selecting some of the blocks in  $B$  and taking their Boolean sum. Furthermore, the error metrics are the same (number of element-wise disagreements). This shows that we can reduce the block selection problem to the BU problem. For the other direction it suffices to note that a vector is a tensor, and therefore, the BU problem is merely a special case of the block selection problem.  $\square$

Given Proposition 4.1, we cannot hope for always finding the optimal solution. But luckily the same proposition also tells us how to solve the block selection problem given that we know how to solve the BU problem. Therefore we will use the greedy algorithm proposed in [14]: We will always select the block that has the highest gain given the already-selected blocks. The gain of a block is defined as the number of not-yet-covered 1s of  $\mathcal{X}$  minus the number of not-yet-covered 0s of  $\mathcal{X}$  covered by this block, and an element  $x_{ijk}$  is covered if  $b_{ijk} = 1$  for some already-selected block.

The greedy algorithm has the benefit that it gives us an ordering of the blocks, so that if the user wants a rank- $k$  decomposition, we can simply return the first  $k$  blocks, instead of having to re-compute the ordering.

## 4.2 The MDL Principle and Encoding the Data for the CP decomposition

The greedy algorithm in the previous section returns an ordering of the columns of matrices  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{C}$  of the CP-decomposition. However, this still does not tell us the optimal rank of the decomposition. In order to choose the best rank  $r$  for the decomposition we apply the *Minimum Description Length* (MDL) principle [17] to the encoding of the obtained decomposition. In this section we explain how this is done.

**Minimum Description Length Principle.** The intuition behind the MDL principle is that the best model is the one that allows us to compress the data best. For our application that means that we should choose the rank  $r$  of the CP decomposition in such a way that the size of the resulting compression is minimal.

To compute the encoding length of the data, we use the two-part (or crude) MDL: if  $\mathcal{D}$  is our data (the data tensor) and  $\mathcal{M}$  is a model of it (often called *hypothesis* in the MDL literature), we aim to minimize  $L(\mathcal{M}) + L(\mathcal{D} | \mathcal{M})$ , where  $L(\mathcal{M})$  is the number of bits we need to encode  $\mathcal{M}$  and  $L(\mathcal{D} | \mathcal{M})$  is the number of bits we need to encode the data *given* the model  $\mathcal{M}$ .

In our application, the model  $\mathcal{M}$  is the Boolean CP decomposition of the data tensor. As MDL requires us to explain the data exactly, we also need to encode the differences between the data and its (approximate) decomposition; this is the  $\mathcal{D} | \mathcal{M}$  part.

The intuition of using the MDL principle lies in the following simple observation: When we move from the rank- $r$  to the rank- $(r + 1)$  decomposition defined by  $\mathbf{A} \boxtimes \mathbf{B} \boxtimes \mathbf{C}$  two things happen. First, the size of the factor matrices increases (and so does  $L(\mathcal{M})$ ). Second, (hopefully) the reconstruction error decreases (and so does  $L(\mathcal{D} | \mathcal{M})$ ). Hence our goal is to find the rank  $r$  where the trade off between the encoding of  $\mathcal{M}$  and  $L(\mathcal{D} | \mathcal{M})$  is optimal.

We will now explain how we compute the encoding length. For this, we modify the *Typed XOR Data-to-Model encoding* for encoding Boolean matrix factorizations [15]. But first, let us emphasize two details. First, we are not interested on the actual encoding lengths; rather, we are interested on the *change* on the encoding lengths between two models. We can therefore omit all the parts that will not change between two models. Second, we are not interested on creating actual encodings, only computing the encoding lengths. We are therefore perfectly happy with fractional bits and will omit the rounding to full bits for the sake of simplicity. Also, the base of the logarithm does not matter (as long as we use the same base for all logarithms); the reader can consider all the logarithms in this chapter taken on base 2.

We will first explain how to encode the model  $\mathcal{M}$ , that is, the tuple  $(\mathbf{A}, \mathbf{B}, \mathbf{C})$  that defines a Boolean CP decomposition of a 3-way binary tensor. The first thing we need to encode is the size of the original tensor,  $n$ ,  $m$ , and  $l$  and the rank  $r$  of the decomposition. For this, we can use any universal code for nonnegative integers, such as the Elias Delta code [5], taking  $\Theta(\log x + 2 \log \log x)$  bits per integer  $x$ . In practice we can omit the numbers  $n$ ,  $m$ , and  $l$  and only compute

the length of  $r$ , as the former do not change between two decompositions of the same data tensor.

We encode the factor matrix  $\mathbf{A}$  (other factor matrices follow analogously). We note first that the size of  $\mathbf{A}$  has already been encoded in the size of the data tensor and  $r$ . Let us assume  $\mathbf{A}$  has  $r$  factors (i.e. columns) and  $n$  rows. We encode each factor  $\mathbf{a}_i$  (which is just a binary vector) separately by enumerating over all  $n$ -dimensional binary vectors with  $|\mathbf{a}_i|$  1s in some fixed order, and storing just the index of the vector we want to encode in this enumeration. As there are  $\binom{n}{|\mathbf{a}_i|}$  such binary vectors, storing this index takes  $\log \binom{n}{|\mathbf{a}_i|}$  bits. (Note that we do not need to do the actual enumeration, as we only need to know the number of bits storing the number would take.) To be able to reverse this computation, we need to encode the number  $|\mathbf{a}_i|$ ; this takes  $\log n$  bits, and so in total a single factor takes  $\log \binom{n}{|\mathbf{a}_i|} + \log n$  bits and the whole matrix  $p \log n + \sum_{i=1}^p \binom{n}{|\mathbf{a}_i|}$ .

With the length of encoding the model computed, we still need to compute  $L(\mathcal{D} | \mathcal{M})$ , that is, the difference between the approximation induced by the decomposition and the actual data. Following [15], we split this difference into two groups: false positives (elements that are 1 in the approximation but 0 in the data) and false negatives (elements that are 0 in the approximation but 1 in the data). We can represent the false positives using a binary  $n$ -by- $m$ -by- $k$  tensor  $\mathcal{F}_+$  that has 1 in each of the positions that are false positives in the approximation and 0 elsewhere. We can encode this tensor by unfolding it into a long binary vector and using the same approach we used to encode the factors. The size of the tensor has already been encoded (it is the same size as the data). The naïve upper bound to the number of 1s in  $\mathcal{F}_+$  is  $nmk$ , but in fact we know that we can only make a false positive if the approximation is 1. Therefore, if the number of 1s in the approximation is  $|\tilde{\mathcal{D}}|$ , we can encode the number of 1s in  $\mathcal{F}_+$  using  $\log |\tilde{\mathcal{D}}|$  bits. Using the same numbering scheme as above, we still need  $\log \binom{nmk}{|\mathcal{F}_+|}$  bits to encode the contents of the tensor.

We can encode the false negative tensor  $\mathcal{F}_-$  analogously, except that the upper bound for 1s is  $nmk - |\tilde{\mathcal{D}}|$ . In summary we have that  $L(\mathcal{D} | \mathcal{M})$  is

$$\log |\tilde{\mathcal{D}}| + \log \binom{nmk}{|\mathcal{F}_+|} + \log(nmk - |\tilde{\mathcal{D}}|) + \log \binom{nmk}{|\mathcal{F}_-|}.$$

Having the encoding in place, we can simply compute the change of description length for every rank  $1 \leq r \leq B$  and return  $r$  where this value is minimized. The corresponding (truncated) matrices  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{C}$  are the factors of the final CP decomposition that our algorithm returns.

### 4.3 Encoding the Data for the Tucker decomposition

Similar to obtaining a CP decomposition from the blocks returned by WALK'N'-MERGE these blocks also define a trivial Tucker decomposition of the same tensor. The factor matrices  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{C}$  are defined the same way as for the CP. The core  $\mathcal{G}$  of the Tucker decomposition is a  $B$ -by- $B$ -by- $B$  size tensor with ones in its hyperdiagonal.

Our goal is to obtain a more compact decomposition starting from this trivial one by merging some of the factors and adjusting the dimensions and content of the core accordingly. We want to allow the merge of two factors even if it would increase the error slightly. But how to define when error is increasing too much and merge should not be made? To solve that problem, we again use the MDL principle.

**Encoding the Boolean Tucker decomposition.** The model  $\mathcal{M}$  we want to encode is the Boolean Tucker decomposition of the data tensor, that is, a tuple  $(\mathcal{G}, \mathbf{A}, \mathbf{B}, \mathbf{C})$ . Encoding the size of the data tensor as well as the content of the factor matrices is done in the same way as for the CP decomposition. As the size of the core tensor determines the size of the factor matrices, we do not need to encode it separately. To encode the core tensor, we need to encode its dimensions  $p$ ,  $r$ , and  $q$ . For this, we again use the Elias delta coding. The actual core we encode similarly to how we encoded the error tensors with the CP factorization, that is, we unfold the core into a long binary vector and encode that vector using its index in the enumeration. This takes  $\log pqr + \log \binom{pqr}{|\mathcal{G}|}$  bits. Again, remember that we do not need to compute the actual index, only how many bits storing it would take.

Finally the positive and negative error tensors are identical to the ones in the CP decomposition and hence are encoded in the same way.

**Applying the MDL principle.** Given the encoding scheme we can use a straight forward heuristic to obtain the final Tucker decomposition starting from the trivial one determined by the output of WALK’N’MERGE. In every mode and for every pair of factors we compute the description length of the resulting decompositions if we were to merge these two factors. Ideally we would compute all possible merging sequences and pick the one with the highest overall gain in encoding length. This is of course infeasible, hence we follow a greedy heuristic and apply every merge that yields an improvement. An overview of this procedure is given in Algorithm 3. We use the notation  $\text{MDL}(\mathcal{G}, \mathbf{A}, \mathbf{B}, \mathbf{C})$  to indicate the encoding length of a Tucker decomposition.  $\text{MDL}(\mathcal{G}, \mathbf{A}, \mathbf{B}, \mathbf{C}, f_1, f_2)$  indicates the encoding length if factors  $f_1$  and  $f_2$  would be merged.

One question is what the merged factor should be. Let us assume we are considering merging factors  $f_1$  and  $f_2$ . Trivial solutions would be to either take the union ( $f_1 \cup f_2$ ) or the intersection ( $f_1 \cap f_2$ ) of the indices in the two factors. We found that both approaches perform poorly. Instead we apply a greedy heuristic that makes a decision for every index in the union. The basis of the merged factor is  $f_1 \cap f_2$ . If the intersection of the factors is empty, we move on and don’t merge them. If it is not, then for every element in the symmetric difference we make a greedy decision whether to include it in the merged factor or not. For this we compute the change in the encoding length of the whole decomposition with or without that element. Bear in mind that in order to compute this, we have to check every block (thus combination of factors as indicated by the current core tensor) that this factor participates in. If we are able to find a merged factor that decreases the overall encoding length, then we always execute this merge. The algorithm finishes when there is no merges

---

**Algorithm 3** Reducing the size of the Boolean Tucker decomposition with help of the MDL principle.

---

**Input:** Data  $\mathcal{X}$ , threshold  $d$ , blocks  $B = \{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_r\}$  from random walk

**Output:**  $\mathcal{G}, \mathbf{A}, \mathbf{B}, \mathbf{C}$  of the Tucker decomposition

- 1: create trivial Tucker decomposition  $\mathcal{G}, \mathbf{A}, \mathbf{B}, \mathbf{C}$
- 2:  $Len \leftarrow \text{MDL}(\mathcal{G}, \mathbf{A}, \mathbf{B}, \mathbf{C})$
- 3: **repeat**
- 4:   **for all**  $\mathbf{a}_i, \mathbf{a}_j \in \mathcal{A}$  **do**
- 5:      $newLen \leftarrow \text{MDL}(\mathcal{G}, \mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{a}_i, \mathbf{a}_j)$
- 6:     **if**  $newLen < Len$  **then**
- 7:        $Len \leftarrow newLen$
- 8:       merge( $\mathbf{a}_i, \mathbf{a}_j$ )
- 9:   **for all**  $\mathbf{b}_i, \mathbf{b}_j \in \mathcal{B}$  **do**
- 10:      $newLen \leftarrow \text{MDL}(\mathcal{G}, \mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{b}_i, \mathbf{b}_j)$
- 11:     **if**  $newLen < Len$  **then**
- 12:        $Len \leftarrow newLen$
- 13:       merge( $\mathbf{b}_i, \mathbf{b}_j$ )
- 14:   **for all**  $\mathbf{c}_i, \mathbf{c}_j \in \mathcal{C}$  **do**
- 15:      $newLen \leftarrow \text{MDL}(\mathcal{G}, \mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{c}_i, \mathbf{c}_j)$
- 16:     **if**  $newLen < Len$  **then**
- 17:        $Len \leftarrow newLen$
- 18:       merge( $\mathbf{c}_i, \mathbf{c}_j$ )
- 19: **until** no more merges are performed

---

executed anymore.

## 5 Experimental Evaluation

We evaluated our algorithms with both synthetic and real-world data.

### 5.1 Other methods and Evaluation Criteria

To the best of our knowledge, this paper is the first to present a scalable Boolean CP decomposition algorithm. Therefore, we cannot compare our algorithm against other *Boolean* CP decomposition algorithms with the kind of data sets we are interested about. We did try the BCP\_ALS algorithm [13] (implementation from the author), but it ran out of memory in all but single dataset. Therefore we cannot report results with it.

Instead, we used two real-valued scalable CP decomposition methods: namely CP-APR [4] (implementation from the Matlab Tensor Toolbox v2.5<sup>2</sup>) and PAR-CUBE [16]<sup>3</sup>. CP-APR is an alternating Poisson regression algorithm that is

<sup>2</sup><http://www.sandia.gov/~tgkolda/TensorToolbox/>

<sup>3</sup><http://www.cs.cmu.edu/~epapalex/>



specifically developed for sparse (counting) data (which can be expected to follow the Poisson distribution) with the goal of returning sparse factors. The aim for sparsity and, to some extent, considering the data as a counting data, make this method suitable for comparison; on the other hand, it aims to minimize the (generalized) K–L divergence, not squared error, and binary data is not Poisson distributed.<sup>4</sup>

The other method we compare against, PARCUBE, uses clever sampling to find smaller sub-tensors. It then solves the CP decomposition in this sub-tensor, and merges the solutions back into one. We used a non-negative variant of PARCUBE that expects non-negative data, and returns non-negative factor matrices. PARCUBE aims to minimize the squared error.

To compute the error, we used the Boolean error function (3) for WALK’N’-MERGE and the squared error function (1) for the comparison methods. This presents yet another apples-versus-oranges comparison: on one hand, the squared error can help the real-valued methods, as it scales all errors less than 1 down; on the other hand, small errors cumulate unlike with fully binary data. To alleviate this problem, we also rounded the reconstructed tensors from CP\_APR and PARCUBE to binary tensors. Instead of simply rounding from 0.5, we tried different rounding thresholds between 0 and 1 and selected the one that gave the lowest (Boolean) reconstruction error. With some of the real-world data, we were unable to perform the rounding for the full representation due to time and memory limitations. For these data sets, we estimated the rounded error using stratified sampling, where we sampled 10 000 1s and 10 000 0s from the data, computed the error on these, and scaled the results.

## 5.2 Synthetic Data

We start by evaluating our algorithms with synthetic data. Our algorithm is aimed to reconstruct the latent structure from large and sparse binary tensors and therefore we tested the algorithms with such data. We generated sparse 1000-by-1500-by-2000 synthetic binary tensor as follows: We first fixed parameters for the Boolean rank of the tensor and the noise to apply. We generated three (sparse) factor matrices to obtain the noise-free tensor. As we assume that the rank-1 tensors in the real-world data are relatively small (e.g. synonyms of an entity), the rank-1 tensors we use were approximately of size 16-by-16-by-16, with each of them overlapping with another block. We then added noise to this tensor. We separate the noise in two types: additive noise flips elements that are 0 to 1 while destructive noise flips elements that are 1 in the noise-free tensor to 0. The amount of noise depends on the number of 1s in the noise-free data, that is 10% of destructive noise means that we delete 10% of the 1s, and 20% of additive noise means that we add 20% more 1s.

We varied three parameters – rank, additive noise, destructive noise, and overlap of the latent blocks – and created five random copies for each set

---

<sup>4</sup>Sampling Poisson distribution can give a binary matrix, but it cannot be forced to give one.

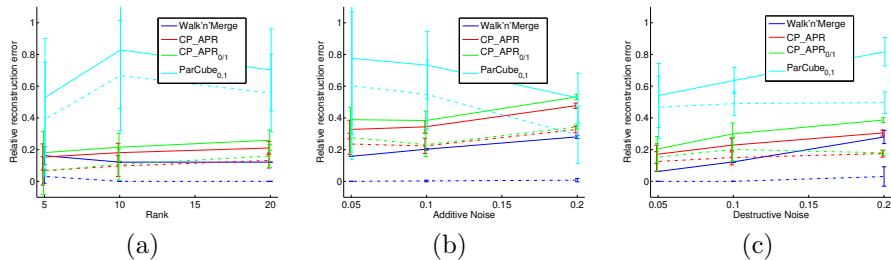


Figure 2: Results on synthetic data sets using CP-type decompositions. (a) Varying rank. (b) Varying additive noise. (c) Varying destructive noise. Solid lines present the relative reconstruction error w.r.t. input tensor; dashed lines present it w.r.t. the original noise-free tensor. All points are mean values over five random datasets and the width of the error bars is twice the standard deviation.

parameters. We measured the quality of the factorizations using the sum of squared differences (1) for continuous-valued methods and the number of disagreements (3) for binary methods. We normalized the errors by the number of non-zeros in the data (e.g. the sum of squared values, as the data is binary). We compared the reconstruction error against both the input data (with noise) and the original noise-free data. Our goal, after all, is to recover the latent structure, not the noise. The rank of the decomposition was set to the true rank of the data for all methods. For WALK’N’MERGE we set the merging threshold to  $1 - (n_d + 0.05)$ , where  $n_d$  was the amount of destructive noise, the length of the random walks was set to 5, and we only considered blocks of size 4-by-4-by-4 or larger. The results for varying rank and different types of noise are presented in Figure 2. Varying the amount of overlap did not have any effect on the results of WALK’N’MERGE, and we omit the results. Results for PARCUBE were consistently worse than anything else and they are omitted from the plots.

**Rank.** For the first experiment (Figure 2(a)) we varied the rank while keeping the additive and destructive noise at 10%. With rank-5 decomposition, WALK’N’MERGE fits to the input data slightly worse than CP\_APR (unrounded) but clearly better than CP\_APR<sub>0/1</sub> (rounded) and PARCUBE<sub>0/1</sub>, the latter being clearly the worse with all ranks. For larger ranks, WALK’N’MERGE is clearly better than variations of CP\_APR. Note that here rank is both the rank of the data and the rank of the decomposition. When comparing the fit to the original data (dashed lines), WALK’N’MERGE is consistently better than the variants of CP\_APR or PARCUBE<sub>0/1</sub>, to the extent that it achieves perfect results for ranks larger than 5.

**Additive noise.** In this experiment, rank was set to 10, destructive noise to 10%, and additive noise was varied. Results are presented in Figure 2(b). In all results, WALK’N’MERGE is consistently better than any other method, and always recovers the original tensor perfectly.

Table 1: Data set properties

Data set	Rows	Columns	Tubes	Density
Enron	146	146	38	0.0023
TracePort	501	10266	8622	$2.51 \times 10^{-7}$
Facebook	63891	63890	228	$9.42 \times 10^{-7}$

**Destructive noise.** For this experiment, rank was again set to 10 and additive noise to 10% while the amount of destructive noise was varied (Figure 2(c)). The results are similar to those in Figure 2(b), although it is obvious that the destructive noise has the most significant effect on the quality of the results.

**Discussion.** In summary, the synthetic experiments show that when the Boolean structure is present in the data, WALK’N’MERGE is able to find it – in many cases even exactly. That CP-APR is not able to do that should not come as a surprise as it does not try to find such structure. That PARCUBE<sub>0/1</sub> is almost consistently the worse is slightly surprising (and the results from the unrounded PARCUBE were even worse). From Figure 2(b) we can see that the results of PARCUBE<sub>0/1</sub> start improving when the amount of additive noise increases. This hints that PARCUBE’s problems are due to its sampling approach not performing well on these extremely sparse tensors.

## 5.3 Real-World Data

### 5.3.1 Datasets

To assess the quality of our algorithm, we tested it with three real-world data sets, namely Enron, TracePort, and Facebook. The Enron data<sup>5</sup> contains information about who sent e-mail to whom (rows and columns) per months (tubes). The TracePort data set<sup>6</sup> contains anonymized passive traffic traces (source and destination IP and port numbers) from 2009. The Facebook data set<sup>7</sup> [20] contains information about who posted a message on whose wall (rows and columns) per weeks (tubes). Basic properties of the data sets are given in Table 1.

### 5.3.2 CP Factorization

We start by reporting the reconstruction errors with CP decompositions using the same algorithms we used with the synthetic data. The results can be seen in Table 2. For Enron, we used single rank ( $r = 12$ ) and for the other two, we used two ranks:  $r = 15$  and whichever gave the smallest reconstruction error by WALK’N’MERGE (after ordering the blocks). In case of the Facebook data, WALK’N’MERGE obtained minimum error of 611 561, but no other method was

<sup>5</sup><http://www.cs.cmu.edu/~enron/>

<sup>6</sup>[http://www.caida.org/data/passive/passive\\_2009\\_dataset.xml](http://www.caida.org/data/passive/passive_2009_dataset.xml)

<sup>7</sup>The data is publicly available from the authors of [20], see <http://socialnetworks.mpi-sws.org>

Table 2: Reconstruction errors rounded to the nearest integer. Numbers prefixed with \* are obtained using sampling.

Algorithm	Enron	TracePort		Facebook
	$r = 12$	$r = 15$	$r = 1370$	$r = 15$
WALK’N’MERGE	1 753	10 968	7 613	612 314
PARCUBE	2 089	33 741	$4 \cdot 10^{55}$	$8 \cdot 10^{140}$
PARCUBE <sub>0/1</sub>	1 724	11 189	$* 2 \cdot 10^7$	* 1 788 874
CP-APR	1 619	11 069	5 230	626 349
CP-APR <sub>0/1</sub>	1 833	11 121	* 1 886	* 626 945

able to finish within 48 hours with the higher rank ( $r = 3233$ ) and we omit the results from the table and only report the errors for  $r = 15$ .

The smallest of the data sets, **Enron**, reverses the trend we saw with the synthetic data: now WALK’N’MERGE is no more the best, as both CP-APR and PARCUBE<sub>0/1</sub> obtain slightly better reconstruction errors. This probably indicates that the data does not have strong Boolean CP type structure. In case of **TracePort** and  $k = 15$  however, WALK’N’MERGE is again the best, if only slightly. With  $r = 1370$ , WALK’N’MERGE improves, but CP-APR and especially CP-APR<sub>0/1</sub> improve even more, obtaining significantly lower reconstruction errors. The very high rank probably lets CP-APR to better utilize the higher expressive power of continuous factorizations, thus explaining the significantly improved results. For **Facebook**, we only report the  $r = 15$  results as the other methods were not able to handle the rank-3300 factorization that gave WALK’N’MERGE its best results. For this small rank, the situation is akin to **TracePort** with  $r = 15$  in that WALK’N’MERGE is the best followed directly with CP-APR. PARCUBE’s errors were off the charts with both **TracePort** ( $r = 1370$ ) and **Facebook**; we suspect that the extreme sparsity (and high rank) fooled its sampling algorithm.

Observing the results of WALK’N’MERGE, we noticed that the resulting blocks were typically very small (e.g. 3-by-3-by-2). This is understandable given the extreme sparsity of the data. For example, the **TracePort** data does not contain any 2-by-2-by-2 monochromatic submatrix. On the other hand, the small factors fit to our intuition of the data. Consider, for example, the **Facebook** data: a monochromatic block corresponds to a set of people who all write to everybody’s walls in the other group of people in certain days. Even when we relax the constrain to dense blocks, it is improbable that these groups would be very big.

**Running time.** Final important question is the running time of the algorithm. The running time of WALK’N’MERGE depends on one hand on the structure of the input tensor (number, but also location, of non-zeros) and on the other hand, on the parameters used (number of random walks, their length, minimum density threshold, and how big a block has to be to be non-trivial). It is therefore hard to provide any systematic study of the running times. But to give some idea, we report the running times for the **Facebook** data, as that

is the biggest data set we used. The fastest algorithm for  $k = 15$  was PARCUBE, finishing in a matter of minutes (but note that it gave very bad results). Second-fastest was WALK’N’MERGE. We tried different density thresholds  $d$ , effecting the running time. The fastest was  $d = 0.2$ , when WALK’N’MERGE took 85 minutes, the slowest was  $d = 0.70$ , taking 277 minutes, and the average was 140 minutes. CP-APR was in between these extremes, taking 128 minutes for one run. Note, however, that WALK’N’MERGE didn’t return just the  $r = 15$  decomposition, but in fact all decompositions up to  $r = 3300$ . Neither PARCUBE or CP-APR was able to handle so large ranks with the Facebook data.

### 5.3.3 Tucker Decomposition

We did some further experiments with the Boolean Tucker decomposition. For the Enron dataset we obtained a decomposition with a core of size 9-by-11-by-9 from the MDL step. While this might feel small, the reconstruction error was 1775, i.e. almost as good as the best BCP decomposition. (Recall that MDL does not try to optimize the reconstruction error, but the encoding length.)

With the Tucker decomposition, we also used a fourth semi-synthetic data set, YPSS.<sup>8</sup> This data set contains noun phrase–context pattern–noun phrase triples that are observed (surface) forms of subject entity–relation–object entity triples. With this data our goal is to find a Boolean Tucker decomposition such that the core  $\mathcal{G}$  corresponds to the latent subject–relation–object triples and the factor matrices tell us which surface forms are used for which entity and relation. A detailed analysis of the fact-recovering power of the Tucker decomposition applied to the YPSS dataset can be found in [6]. The size of the data is 39 500-by-8 000-by-21 000 and it contains 804 000 surface term triplets.

The running time of WALK’N’MERGE on YPSS was 52 minutes, and computing the Tucker decomposition took another 3 hours.

An example of a factor of the subjects would be {`claudio de lorimier, de lorimier, louis, jean-baptiste`}, corresponding to Claude-Nicolas-Guillaume de Lorimier, a Canadian politician and officer from the 18th Century (and his son, Jean-Baptiste). An example of an object-side factor is {`borough of lachine, villa st. pierre, lachine quebec`}, corresponding to the borough of Lachine in Quebec, Canada (town of St. Pierre was merged to Lachine in 1999). Finally, an example of a factor in the relations is {`was born was , [[det]] born in`}, with an obvious meaning. In the Boolean core  $\mathcal{G}$  the element corresponding to these three factors is 1, meaning that according to our algorithm, de Lorimier was born in Lachine, Quebec – as he was.

### 5.3.4 Discussion

Unlike with synthetic data, with real-world data we cannot guarantee that the data has Boolean structure. And if the data does not have the Boolean structure, there does not exist any good BTF. Yet, with most of our experiments, WALK’N’MERGE performs very well, both in quantitative and qualitative analysis.

<sup>8</sup>The data set is available at <http://www.mpi-inf.mpg.de/~pmiETTIN/btf/>.

Considering running times, WALK’N’MERGE is comparative to CP-APR with most datasets.

## 6 Related Work

Normal tensor factorizations are well-studied, dating back to the late Twenties. The two popular decomposition methods, Tucker and CP, were proposed in Sixties [19] and Seventies [2, 7], respectively. The topic has nevertheless attained growing interest in recent years, both in numerical linear algebra and computer science communities. For a comprehensive study of recent work, see [10], and the recent work on scalable factorizations [16].

One field of computer science that has adopted tensor decompositions is computer vision and machine learning. The interest to non-negative tensor factorizations stems from these fields [9, 18].

The theory of Boolean tensor factorizations was studied in [13], although the first algorithm for Boolean CP factorization was presented in [11]. A related line of data mining research has also studied a specific type of Boolean CP decomposition, where no 0s can be presented as 1s (e.g. [3]). For more on these methods and their relation to Boolean CP factorization, see [13].

## 7 Conclusions

We have presented WALK’N’MERGE, an algorithm for computing the Boolean tensor factorization of large and sparse binary tensors. Analysing the results of our experiments sheds some light on the strengths and weaknesses of our algorithm. First, it is obvious that it does what it was designed to do, that is, finds Boolean tensor factorizations of large and sparse tensors. But it has its caveats, as well. The random walk algorithm, for example, introduces an element of randomness, and it seems that it benefits from larger tensors. The algorithm, and its running time, is also somewhat sensible to the parameters, possibly requiring some amount of tuning.

## References

- [1] Radim Bělohávek, Cynthia Glodeanu, and Vilém Vychodil. Optimal Factorization of Three-Way Binary Data Using Triadic Concepts. *Order*, March 2012.
- [2] J. Douglas Carroll and Jih-Jie Chang. Analysis of individual differences in multidimensional scaling via an N-way generalization of ‘Eckart-Young’ decomposition. *Psychometrika*, 35(3):283–319, 1970.
- [3] Loïc Cerf, Jérémy Besson, Céline Robardet, and Jean-François Boulicaut. Closed patterns meet n-ary relations. *ACM Trans. Knowl. Discov. Data*, 3(1), 2009.

- [4] Eric C Chi and Tamara G Kolda. On Tensors, Sparsity, and Nonnegative Factorizations. *SIAM J. Matrix Anal. Appl.*, 33(4):1272–1299, December 2012.
- [5] Peter Elias. Universal codeword sets and representations of the integers. *IEEE Trans. Inform. Theory*, 21(2):194–203, March 1975.
- [6] Dóra Erdős and Pauli Miettinen. Discovering Facts with Boolean Tensor Tucker Decomposition. In *CIKM '13*, 2013.
- [7] Richard A. Harshman. Foundations of the PARAFAC procedure: Models and conditions for an ‘explanatory’ multimodal factor analysis. Technical report, UCLA Working Papers in Phonetics, 1970.
- [8] Johan Håstad. Tensor rank is NP-complete. *J. Algorithms*, 11(4):644–654, December 1990.
- [9] Yong-Deok Kim and Seungjin Choi. Nonnegative Tucker Decomposition. In *CVPR '07*, pages 1–8, 2007.
- [10] Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM Review*, 51(3):455–500, 2009.
- [11] Iwin Leenen, Iven Van Mechelen, Paul De Boeck, and Seymour Rosenberg. INDCLAS: A three-way hierarchical classes model. *Psychometrika*, 64(1):9–24, March 1999.
- [12] Pauli Miettinen. On the positive-negative partial set cover problem. *Inform. Process. Lett.*, 108(4):219–221, 2008.
- [13] Pauli Miettinen. Boolean Tensor Factorizations. In *ICDM '11*, pages 447–456, 2011.
- [14] Pauli Miettinen, Taneli Mielikäinen, Aristides Gionis, Gautam Das, and Heikki Mannila. The Discrete Basis Problem. *IEEE Trans. Knowl. Data Eng.*, 20(10):1348–1362, October 2008.
- [15] Pauli Miettinen and Jilles Vreeken. MDL4BMF: Minimum description length for boolean matrix factorization. Technical Report MPI-I–2012–5–001, Max-Planck-Institut für Informatik, June 2012.
- [16] Evangelos E Papalexakis, Christos Faloutsos, and Nicholas D Sidiropoulos. ParCube: Sparse Parallelizable Tensor Decompositions. In *ECML PKDD '12*, pages 521–536, 2012.
- [17] J Rissanen. Modeling by shortest data description. *Automatica*, 14(5):465–471, September 1978.
- [18] Amnon Shashua and Tamir Hazan. Non-negative tensor factorization with applications to statistics and computer vision. In *ICML '05*, 2005.

- [19] Ledyard R Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31(3):279–311, 1966.
- [20] Bimal Viswanath, Alan Mislove, Meeyoung Cha, and Krishna P Gummadi. On the Evolution of User Interaction in Facebook. In *WOSN '09*, pages 37–42, 2009.