

Algorithms for Approximate Subtropical Matrix Factorization

Sanjar Karaev · Pauli Miettinen

the date of receipt and acceptance should be inserted later

Abstract Matrix factorization methods are important tools in data mining and analysis. They can be used for many tasks, ranging from dimensionality reduction to visualization. In this paper we concentrate on the use of matrix factorizations for finding patterns from the data. Rather than using the standard algebra – and the summation of the rank-1 components to build the approximation of the original matrix – we use the subtropical algebra, which is an algebra over the nonnegative real values with the summation replaced by the maximum operator. Subtropical matrix factorizations allow “winner-takes-it-all” interpretations of the rank-1 components, revealing different structure than the normal (nonnegative) factorizations. We study the complexity and sparsity of the factorizations, and present a framework for finding low-rank subtropical factorizations. We present two specific algorithms, called Capricorn and Cancer, that are part of our framework. They can be used with data that has been corrupted with different types of noise, and with different error metrics, including the sum-of-absolute differences, Frobenius norm, and Jensen–Shannon divergence. Our experiments show that the algorithms perform well on data that has subtropical structure, and that they can find factorizations that are both sparse and easy to interpret.

Keywords Tropical algebra · Max-times algebra · Matrix factorizations · Data mining

1 Introduction

Finding simple patterns that can be used to describe the data is one of the main problems in data mining. The data mining literature knows many different techniques for this general task, but one of the most common pattern finding techniques rarely gets

Part of this work was done while P.M. was with Max-Planck-Institut für Informatik.

S. Karaev
Max-Planck-Institut für Informatik, Saarland Informatics Campus, Saarbrücken, Germany
E-mail: sanjar.karaev@mpi-inf.mpg.de

P. Miettinen
School of Computing, University of Eastern Finland, Kuopio, Finland
E-mail: pauli.miettinen@uef.fi

classified as such. Matrix factorizations (or decompositions, these two terms are used interchangeably in this paper) represent the given input matrix \mathbf{A} as a product of two (or more) factor matrices, $\mathbf{A} \approx \mathbf{BC}$. This standard formulation of matrix factorizations makes their pattern mining nature less obvious, but let us write the matrix product \mathbf{BC} as a sum of rank-1 matrices, $\mathbf{BC} = \mathbf{F}_1 + \mathbf{F}_2 + \dots + \mathbf{F}_k$, where \mathbf{F}_i is the outer product of the i th column of \mathbf{B} and the i th row of \mathbf{C} . Now it becomes clear that the rank-1 matrices \mathbf{F}_i are the “simple patterns”, and the matrix factorization is finding k such patterns whose sum is a good approximation of the original data matrix.

This so-called “component interpretation” (Skillicorn 2007) is more appealing with some factorizations than with others. For example, the classical singular value decomposition (SVD) does not easily admit such an interpretation, as the components are not easy to interpret without knowing the earlier components. On the other hand, the motivation for the nonnegative matrix factorization (NMF) often comes from the component interpretation, as can be seen, for example, in the famous “parts of faces” figures of Lee and Seung (1999). The “parts-of-whole” interpretation is in the hearth of NMF: every rank-1 component adds something to the overall decomposition, and never removes anything. This aids with the interpretation of the components, and is also often claimed to yield sparse factors, although this latter point is more contentious (see e.g. Hoyer 2004).

Perhaps the reason why matrix factorization methods are not often considered as pattern mining methods is that the rank-1 matrices are summed together to build the full data. Hence, it is rare for any rank-1 component to explain any part of the input matrix alone. But the use of summation as a way to aggregate the rank-1 components can be considered to be “merely” a consequence of the fact that we are using the standard algebra. If we change the algebra – in particular, if we change how we define the summation – we change the operator used for the aggregation. In this work, we propose to use the *maximum* operator to define the summation over the nonnegative matrices, giving us what is known as the *subtropical algebra*. As the aggregation of the rank-1 factors is now the element-wise maximum, we obtain what we call the “winner-takes-it-all” interpretation: the final value of each element in the approximation is defined only by the largest value in the corresponding element in the rank-1 matrices. This can be considered a staple of the subtropical structure – for each element in the data we can find a single rank-1 pattern, the “winner”, that determines its value exactly. This is in contrast to the NMF structure, where each pattern would only make a “contribution” to the final value.

Not only does the subtropical algebra give us the intriguing winner-takes-it-all interpretation, it also provides guarantees about the sparsity of the factors, as we will show in Section 3.2. Furthermore, a different algebra means that we are finding different factorizations compared to NMF (or SVD). The emphasis here is on the word *different*: the factorizations can be better or worse in terms of the reconstruction error but the patterns they find are usually different to those found by NMF. It is also worth mentioning that the same dataset often has both kinds of structures in it, in which case subtropical and NMF patterns are complementary to each other, and depending on an application, one or the other can be more useful. One practical advantage of the subtropical methods though is that they tend to find more concise representation of patterns in the data, while NMF often splits them into several smaller components, making it harder to see the big picture.

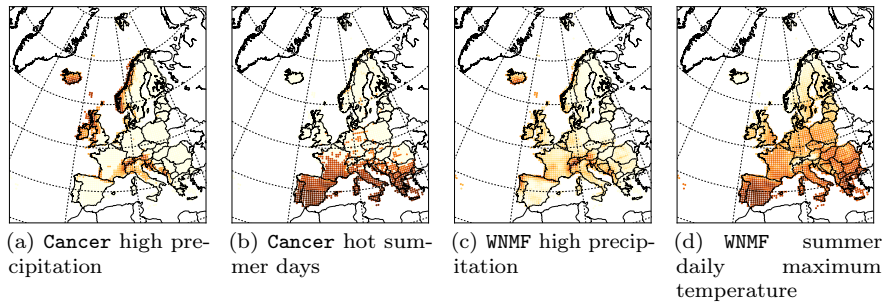


Fig. 1 Example results by **Cancer** and **WNMF** on the Worldclim dataset. For each method, two selected columns from the left-hand factor matrix is shown on a map. The values are normalized to the unit interval, and darker shades indicate higher values.

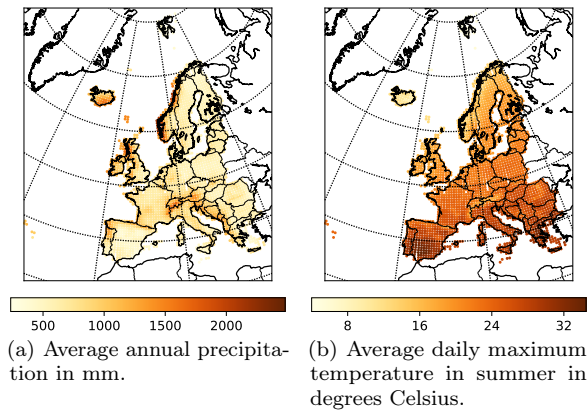


Fig. 2 Average climate data to be compared with the factors in Figure 1.

To illustrate this, and in general the kind of structure subtropical matrix factorization can reveal and how it is different from that of NMF, we show example results on the European climate data (Figure 1).

The data contains weather records for Europe between 1960 and 1990, and it was obtained from the global climate data repository.¹ The data has 2575 rows that correspond to 50-by-50 kilometer squares of land where measurements were made and 48 columns corresponding to observations. More precisely, the first 12 columns represent the average low temperature for each month, the next 12 columns the average high temperature, and the next 12 columns the daily mean. The last 12 columns represent the mean monthly precipitation for each month. We preprocessed every column of the data by first subtracting its mean, dividing by the standard deviation, and then subtracting its minimum value, so that the smallest value becomes 0. We compare the results of our subtropical matrix factorization algorithm, called **Cancer**, to those of an NMF algorithm, called **WNMF**, that obtained the best reconstruction error on this data (see Table 2 in Section 5). For both methods, we chose two factors: one that best

¹The raw data is available at <http://www.worldclim.org/>, accessed 18 July 2017.

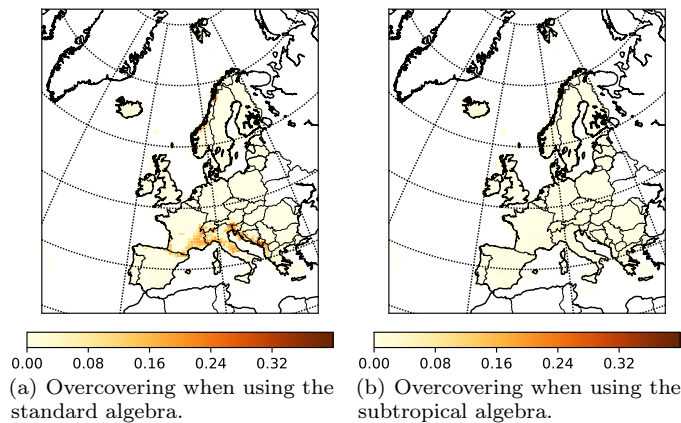


Fig. 3 Comparison of the overcovering when combining the **Cancer** factors from Figures 1(a)–(b) using the standard 3(a) and the subtropical 3(b) algebras. All values are divided by the average maximum value in the original matrix, negative values are ignored. Darker shades indicate higher values.

identifies the areas of high precipitation and another that reflects summer (i.e. June, July, and August) daily maximum temperatures. To be able to validate the results of the algorithms, we also include the average annual precipitation and average summer maximum temperature in Figures 2(a) and 2(b), respectively.

Both **Cancer** and **WMMF** identify the areas of high precipitation and their corresponding left-hand factors are shown in Figures 1(a) and 1(c), respectively. There are, however, two significant differences between their interpretations. First, **Cancer** emphasizes the wettest areas, while **WMMF** shows a much smoother transition similar to the original data. The second difference is that, unlike **Cancer**, **WMMF** does not identify either the UK or Norwegian coast as areas of (particularly) high precipitation. A potential explanation is that there are many overlaps with other factors (see Figure 15(b)), and hence having large values in any of them might lead to overcovering the original data. **Cancer**, as a subtropical method, does not suffer from this issue, as there the reconstructed precipitation is entirely based on the factor with the highest values.

In order to make the above argument more concrete, let us see what happens when we try to combine **Cancer**'s factors using the standard algebra instead of the subtropical one. Recall that if $\mathbf{A} = \mathbf{BC}$ is a rank- k matrix decomposition of \mathbf{A} , then we have $(\mathbf{A})_{ij} = \sum_{s=1}^k (\mathbf{F}_s)_{ij}$, where each pattern \mathbf{F}_s is an outer product of the s th column of \mathbf{B} and the s th row of \mathbf{C} . If for some l and t we have $(\mathbf{F}_l)_{ij} + (\mathbf{F}_t)_{ij} > \mathbf{A}_{ij}$, then also $(\mathbf{BC})_{ij} > \mathbf{A}_{ij}$ since all values are nonnegative. It is therefore generally undesirable for any subset of the patterns to overcover values in the original data, as there would be no way of decreasing these values by adding more patterns. As an example we will combine the patterns corresponding to **Cancer**'s factors from Figures 1(a)–(b). To obtain the actual rank-1 patterns we first need to compute the outer products of these factors with the corresponding rows of the right-hand side matrix. Now if we denote the obtained patterns by \mathbf{F}_1 and \mathbf{F}_2 , then the elements of the matrix $\max\{\mathbf{F}_1 + \mathbf{F}_2 - \mathbf{A}, 0\}$ show by how much the combination of \mathbf{F}_1 and \mathbf{F}_2 overcovers the original data \mathbf{A} . We now plot the average value of every row of the overcover matrix scaled by the average value in the original data (Figure 3(a)). Since each row corresponds to a location on the

map, it shows the average amount by which we would overcover the data, were we to use the standard algebra for combining the **Cancer**'s factors. It is evident that this method produces many values that are too high (mostly around Alps and other high precipitation areas). On the other hand, when we perform the same procedure using the subtropical algebra (Figure 3(b)), there is almost no overcovering.

A somewhat similar behaviour is seen with the maximal daily temperatures in summer. **WNMF** finds a factor that, with the exception of the Scandinavian peninsula, closely mimics the original data and maintains a smooth gradient of decreasing temperatures when moving towards north (Figure 1(d)). In contrast, **Cancer** identifies the areas where summer days are hot, while practically disregarding other parts (Figure 1(b)).

It is worth mentioning that, although UK and the coastal regions of Norway are not prominent in the **WNMF**'s factor shown above, they actually belong to some of its other factors (see Figure 15(b)). In other words, the high precipitation pattern is split into several parts and partially merged with other factors. This is likely a consequence of the pattern splitting nature of NMF mentioned earlier. On the other hand, using the subtropical structure, we were able to isolate the high precipitation pattern and present it in a single factor.

While the above discussion shows that the subtropical model can be a useful complement to NMF, it is generally difficult to claim that either of them is superior. For example **Cancer** generally provided a more concise representation of patterns in the climate data, outlining its most prominent properties, while **WNMF**'s strength was recovering the smooth transition between values.

Contributions and a roadmap. In this paper, we study the use of subtropical decompositions for data analysis.² We start by studying the theoretical aspects of the problem (Section 3), showing that the problem is NP-hard to even approximate, but also that sparse matrices have sparse dominated subtropical decompositions.

In Section 4, we develop a general framework, called **Equator**, for finding approximate, low-rank subtropical decompositions, and we will present two instances of this framework, tailored towards different types of data and noise, called **Capricorn** and **Cancer**. **Capricorn** assumes discrete data with noise that randomly flips the value to a random number, whereas **Cancer** assumes continuous-valued data with standard Gaussian noise.

Our experiments (Section 5) show that both **Capricorn** and **Cancer** work well on datasets that have the kind of noise they are designed for, and they outperform SVD and different NMF methods when data has subtropical structure. On real-world data, **Cancer** is usually the better of the two, although in terms of reconstruction error, neither of the methods can challenge SVD. On the other hand, we show that both **Cancer** and **Capricorn** return interpretable results that show different aspects of the data compared to factorizations made under the standard algebra.

2 Notation and Basic Definitions

Basic notation. Throughout this paper, we will denote a matrix by upper-case boldface letters (**A**), and vectors by lower-case boldface letters (**a**). The i th row of matrix **A** is

²This work is a combined and extended version of our preliminary papers that described these algorithms (Karaev and Miettinen 2016a,b).

denoted by \mathbf{A}_i and the j th column by \mathbf{A}^j . The matrix \mathbf{A} with the i th column removed is denoted by \mathbf{A}^{-i} , and \mathbf{A}_{-i} is the respective notation for \mathbf{A} with a removed row. Most matrices and vectors in this paper are restricted to the nonnegative real numbers $\mathbb{R}_+ = [0, \infty)$.

We use the shorthand $[n]$ to denote the set $\{1, 2, \dots, n\}$.

Algebras. In this paper we consider matrix factorization over so called *max-times* (or *subtropical*) algebra. It differs from the standard algebra of real numbers in that addition is replaced with the operation of taking the maximum. Also the domain is restricted to the set of nonnegative real numbers.

Definition 1 The *max-times* (or *subtropical*) algebra is a set \mathbb{R}_+ of nonnegative real numbers together with operations $a \boxplus b = \max\{a, b\}$ (addition) and $a \boxtimes b = ab$ (multiplication) defined for any $a, b \in \mathbb{R}_+$. The identity element for addition is 0 and for multiplication it is 1.

In the future we will use the notation $a \boxplus b$ and $\max\{a, b\}$ and the names *max-times* and *subtropical* interchangeably. It is straightforward to see that the max-times algebra is a *doid*, that is, a semiring with idempotent addition ($a \boxplus a = a$). It is important to note that subtropical algebra is anti-negative, that is, there is no subtraction operation.

A very closely related algebraic structure is the *max-plus* (*tropical*) algebra (see e.g. Akian et al 2007).

Definition 2 The *max-plus* (or *tropical*) algebra is defined over the set of extended real numbers $\mathbb{R} = \mathbb{R} \cup \{-\infty\}$ with operations $a \oplus b = \max\{a, b\}$ (addition) and $a \odot b = a + b$ (multiplication). The identity elements for addition and multiplication are $-\infty$ and 0, respectively.

The tropical and subtropical algebras are isomorphic (Blondel et al 2000), which can be seen by taking the logarithm of the subtropical algebra or the exponent of the tropical algebra (with the conventions that $\log 0 = -\infty$ and $\exp(-\infty) = 0$). Thus, most of the results we prove for subtropical algebra can be extended to their tropical analogues, although caution should be used when dealing with approximate matrix factorizations. The latter is because, as we will see in Theorem 4, the *reconstruction error* of an approximate matrix factorization under the two different algebras does not transfer directly.

Matrix products and ranks. The matrix product over the subtropical algebra is defined in the natural way:

Definition 3 The *max-times matrix product* of two matrices $\mathbf{B} \in \mathbb{R}_+^{n \times k}$ and $\mathbf{C} \in \mathbb{R}_+^{k \times m}$ is defined as

$$(\mathbf{B} \boxtimes \mathbf{C})_{ij} = \max_{s=1}^k \mathbf{B}_{is} \mathbf{C}_{sj}. \quad (1)$$

We will also need the matrix product over the *tropical* algebra.

Definition 4 For two matrices $\mathbf{B} \in \mathbb{R}^{n \times k}$ and $\mathbf{C} \in \mathbb{R}^{k \times m}$, their *tropical matrix product* is defined as

$$(\mathbf{B} \odot \mathbf{C})_{ij} = \max_{s=1}^k \{\mathbf{B}_{is} + \mathbf{C}_{sj}\}. \quad (2)$$

The *matrix rank* over the subtropical algebra can be defined in many ways, depending on which definition of the normal matrix rank is taken as the starting point. We will discuss different subtropical ranks in detail in Section 3.4. Here we give the main definition of the rank we are using throughout this paper, the so-called *Schein* (or *Barvinok*) *rank* of a matrix.

Definition 5 The *max-times* (*Schein* or *Barvinok*) *rank* of a matrix $\mathbf{A} \in \mathbb{R}_+^{n \times m}$ is the least integer k such that \mathbf{A} can be expressed as an element-wise maximum of k rank-1 matrices, $\mathbf{A} = \mathbf{F}_1 \boxplus \mathbf{F}_2 \boxplus \dots \boxplus \mathbf{F}_k$. Matrix $\mathbf{F} \in \mathbb{R}_+^{n \times m}$ has subtropical (Schein/Barvinok) rank of 1 if there exist column vectors $\mathbf{x} \in \mathbb{R}_+^n$ and $\mathbf{y} \in \mathbb{R}_+^m$ such that $\mathbf{F} = \mathbf{x}\mathbf{y}^T$. Matrices with subtropical Schein (or Barvinok) rank of 1 are called *blocks*.

When it is clear from the context, we will use the term *rank* (or *subtropical rank*) without other qualifiers to denote the subtropical Schein/Barvinok rank.

Special matrices. The final concepts we need in this paper are *pattern matrices* and *dominating matrices*.

Definition 6 A *pattern* of a matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$ is an n -by- m binary matrix \mathbf{P} such that $\mathbf{P}_{ij} = 0$ if and only if $\mathbf{A}_{ij} = 0$, and otherwise $\mathbf{P}_{ij} = 1$. We denote the pattern of \mathbf{A} by $p(\mathbf{A})$.

Definition 7 Let \mathbf{A} and \mathbf{X} be matrices of the same size, and let Γ be a subset of their indices. Then if for all indices $(i, j) \in \Gamma$, $\mathbf{X}_{ij} \geq \mathbf{A}_{ij}$, we say that \mathbf{X} *dominates* \mathbf{A} *within* Γ . If Γ spans the entire size of \mathbf{A} and \mathbf{X} , we simply say that \mathbf{X} *dominates* \mathbf{A} . Correspondingly, \mathbf{A} is said to be *dominated by* \mathbf{X} .

Main problem definition. Now that we have sufficient notation, we can formally introduce the main problem considered in the paper.

Problem 1 (Approximate subtropical rank- k matrix factorization) Given a matrix $\mathbf{A} \in \mathbb{R}_+^{n \times m}$ and an integer $k > 0$, find factor matrices $\mathbf{B} \in \mathbb{R}_+^{n \times k}$ and $\mathbf{C} \in \mathbb{R}_+^{k \times m}$ minimizing

$$E(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \|\mathbf{A} - \mathbf{B} \boxtimes \mathbf{C}\| . \quad (3)$$

Here we have deliberately not specified any particular norm. Depending on the circumstances, different matrix norms can be used, but in this paper we will consider the two most natural choices – the Frobenius and L_1 norms.

3 Theory

Our main contributions in this paper are the algorithms for the subtropical matrix factorization. But before we present them, it is important to understand the theoretical aspects of subtropical factorizations. We will start by studying the computational complexity of Problem 1, showing that it is NP-hard even to approximate. After that, we will show that the dominated subtropical factorizations of sparse matrices are sparse. Then we compare the subtropical factorizations to factorizations over other algebras, analyzing how the error of an approximate decomposition behaves when moving from tropical to subtropical algebra. Finally, we briefly summarize different ways to define the subtropical rank, and how these different ranks can be used to bound each other, and the Boolean rank of a binary matrix, as well.

3.1 Computational complexity

The computational complexity of different matrix factorization problems varies. For example, SVD can be computed in polynomial time (Golub and Van Loan 2012), while NMF is NP-hard (Vavasis 2009). Unfortunately, the subtropical factorization is also NP-hard.

Theorem 1 *Computing the max-times matrix rank is an NP-hard problem, even for binary matrices.*

The theorem is a direct consequence of the following theorem by Kim and Roush (2005):

Theorem 2 (Kim and Roush 2005) *Computing the max-plus (tropical) matrix rank is NP-hard, even for matrices that take values only from $\{-\infty, 0\}$.*

While computing the rank deals with exact decompositions, its hardness automatically makes any approximation algorithm with provable multiplicative guarantees unlikely to exist, as the following corollary shows.

Corollary 1 *It is NP-hard to approximate Problem 1 to within any polynomially computable factor.*

Proof Any algorithm that can approximate Problem 1 to within a factor α must find a decomposition of error $\alpha \cdot 0 = 0$ if the input matrix has exact max-times rank- k decomposition. As this implies solving the max-times rank, per Theorem 1 it is only possible if $P=NP$. \square

3.2 Sparsity of the factors

It is often desirable to obtain sparse factor matrices if the original data is sparse, as well, and the sparsity of its factors is frequently mentioned as one of the benefits of using NMF (see, e.g. Hoyer 2004). In general, however, the factors obtained by NMF might not be sparse, but if we restrict ourselves to *dominated* decompositions, Gillis and Glineur (2010) showed that the sparsity of the factors cannot be less than the sparsity of the original matrix.

The proof of Gillis and Glineur (2010) relies on the anti-negativity, and hence their proof is easy to adapt to max-times setting. Let the *sparsity* of an n -by- m matrix \mathbf{A} , $s(\mathbf{A})$, be defined as

$$s(\mathbf{A}) = \frac{nm - \eta(\mathbf{A})}{nm}, \quad (4)$$

where $\eta(\mathbf{A})$ is the number of nonzero elements in \mathbf{A} . Now we have

Theorem 3 *Let matrices $\mathbf{B} \in \mathbb{R}_+^{n \times k}$ and $\mathbf{C} \in \mathbb{R}_+^{k \times m}$ be such that their max-times product is dominated by an n -by- m matrix \mathbf{A} . Then the following estimate holds:*

$$s(\mathbf{B}) + s(\mathbf{C}) \geq s(\mathbf{A}). \quad (5)$$

Proof The proof follows that of Gillis and Glineur (2010). We first prove (5) for $k = 1$. Let $\mathbf{b} \in \mathbb{R}_+^n$ and $\mathbf{c} \in \mathbb{R}_+^m$ be such that $\mathbf{b}_i \mathbf{c}_j^T \leq \mathbf{A}_{ij}$ for all $i \in [n]$ and $j \in [m]$. Since $(\mathbf{bc}^T)_{ij} > 0$ if and only if $\mathbf{b}_i > 0$ and $\mathbf{c}_j > 0$, we have that

$$\eta(\mathbf{bc}^T) = \eta(\mathbf{b}) \eta(\mathbf{c}). \quad (6)$$

By (4) we have $\eta(\mathbf{bc}^T) = nm(1 - s(\mathbf{bc}^T))$, $\eta(\mathbf{b}) = n(1 - s(\mathbf{b}))$ and $\eta(\mathbf{c}) = m(1 - s(\mathbf{c}))$. Plugging these expressions into (6) we obtain $(1 - s(\mathbf{bc}^T)) = (1 - s(\mathbf{b}))(1 - s(\mathbf{c}))$. Hence, the sparsity in a rank-1 dominated approximation of \mathbf{A} is

$$s(\mathbf{b}) + s(\mathbf{c}) \geq s(\mathbf{bc}^T). \quad (7)$$

From (7) and the fact that the number of nonzero elements in \mathbf{bc}^T is no greater than in \mathbf{A} , it follows that

$$s(\mathbf{b}) + s(\mathbf{c}) \geq s(\mathbf{A}). \quad (8)$$

Now let $\mathbf{B} \in \mathbb{R}_+^{n \times k}$ and $\mathbf{C} \in \mathbb{R}_+^{k \times m}$ be such that $\mathbf{B} \boxtimes \mathbf{C}$ is dominated by \mathbf{A} . Then $\mathbf{B}_{il} \mathbf{C}_{lj} \leq \mathbf{A}_{ij}$ for all $i \in [n]$, $j \in [m]$, and $l \in [k]$, which means that for each $l \in [k]$, $\mathbf{B}^l \mathbf{C}_l$ is dominated by \mathbf{A} . To complete the proof observe that $s(\mathbf{B}) = k^{-1} \sum_{l=1}^k s(\mathbf{B}^l \mathbf{C}_l)$ and $s(\mathbf{C}) = k^{-1} \sum_{l=1}^k s(\mathbf{C}_l)$ and that for each l estimate (8) holds. \square

3.3 Relation to other algebras

Let us now study how the max-times algebra relates to other algebras, especially the standard, the Boolean, and the max-plus algebras. For the first two, we compare the ranks, and for the last, the reconstruction error.

Let us start by considering the Boolean rank of a binary matrix. The *Boolean* (Schein or Barvinok) rank is the following problem:

Problem 2 (Boolean rank) Given a matrix $\mathbf{A} \in \{0, 1\}^{n \times m}$, find the smallest integer k such that there exist matrices $\mathbf{B} \in \{0, 1\}^{n \times k}$ and $\mathbf{C} \in \{0, 1\}^{k \times m}$ that satisfy $\mathbf{A} = \mathbf{B} \circ \mathbf{C}$, where \circ is the *Boolean matrix product*,

$$(\mathbf{B} \circ \mathbf{C})_{ij} = \bigvee_{l=1}^k \mathbf{B}_{il} \mathbf{C}_{lj}.$$

Lemma 1 *If \mathbf{A} is a binary matrix, then its Boolean and subtropical ranks are the same.*

Proof We will prove the claim by first showing that the Boolean rank of a binary matrix is no less than the subtropical rank, and then showing that it is no larger, either. For the first direction, let the Boolean rank of \mathbf{A} be k , and let \mathbf{B} and \mathbf{C} be binary matrices such that \mathbf{B} has k columns and $\mathbf{A} = \mathbf{B} \circ \mathbf{C}$. It is easy to see that $\mathbf{B} \circ \mathbf{C} = \mathbf{B} \boxtimes \mathbf{C}$, and hence, the subtropical rank of \mathbf{A} is no more than k .

For the second direction, we will actually show a slightly stronger claim: Let $\mathbf{A} \in \mathbb{R}_+^{n \times m}$ and let $p(\mathbf{A})$ be its pattern. Then the Boolean rank of $p(\mathbf{A})$ is never more than the subtropical rank of \mathbf{A} . As $p(\mathbf{A}) = \mathbf{A}$ for a binary \mathbf{A} , the claim follows. To prove the claim, let $\mathbf{A} \in \mathbb{R}_+^{n \times m}$ have subtropical rank of k and let $\mathbf{B} \in \mathbb{R}_+^{n \times k}$ and

$\mathbf{C} \in \mathbb{R}_+^{k \times m}$ be such that $\mathbf{A} = \mathbf{B} \boxtimes \mathbf{C}$. Let (i, j) be such that $\mathbf{A}_{ij} = 0$. By definition, $\max_{l=1}^k \mathbf{B}_{il} \mathbf{C}_{lj} = 0$, and hence

$$\max_{l=1}^k p(\mathbf{B})_{il} p(\mathbf{C})_{lj} = \bigvee_{l=1}^k p(\mathbf{B})_{il} p(\mathbf{C})_{lj} = 0. \quad (9)$$

On the other hand, if (i, j) is such that $\mathbf{A}_{ij} > 0$, then there exists l such that $\mathbf{B}_{il}, \mathbf{C}_{lj} > 0$ and consequently,

$$\max_{l=1}^k p(\mathbf{B})_{il} p(\mathbf{C})_{lj} = \bigvee_{l=1}^k p(\mathbf{B})_{il} p(\mathbf{C})_{lj} = 1. \quad (10)$$

Combining (9) and (10) gives us

$$p(\mathbf{A}) = p(\mathbf{B}) \circ p(\mathbf{C}), \quad (11)$$

showing that the Boolean rank of $p(\mathbf{A})$ is at most k . \square

Notice that Lemma 1 also furnishes us with another proof of Theorem 1, as computing the Boolean rank is NP-hard (see, e.g. Miettinen 2009). Notice also that while the Boolean rank of the pattern is never more than the subtropical rank of the original matrix, it can be much less. This is easy to see by considering a matrix with no zeroes: it can have arbitrarily large subtropical rank, but its pattern has Boolean rank 1.

Unfortunately, the Boolean rank does not help us with effectively estimating the subtropical rank, as its computation is an NP-hard problem. The standard rank is (relatively) easy to compute, but the standard rank and the max-times rank are incommensurable, that is, there are matrices that have smaller max-times rank than standard rank and others that have higher max-times rank than standard rank. Let us consider an example of the first kind,

$$\begin{pmatrix} 1 & 2 & 0 \\ 2 & 4 & 1 \\ 0 & 4 & 2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 2 & 1 \\ 0 & 2 \end{pmatrix} \boxtimes \begin{pmatrix} 1 & 2 & 0 \\ 0 & 2 & 1 \end{pmatrix}.$$

As the decomposition shows, this matrix has max-times rank of 2, while its normal rank is easily verified to be 3. Indeed, it is easy to see that the complement of the n -by- n identity matrix $\bar{\mathbf{I}}_n$, that is, the matrix that has 0s at the diagonal and 1s everywhere else, has max-times rank of $O(\log n)$ while its standard rank is n (the result follows from similar results regarding the Boolean rank, see, e.g. Miettinen 2009).

As we have discussed earlier, max-plus and max-times algebras are isomorphic, and consequently for any matrix $\mathbf{A} \in \mathbb{R}_+^{n \times m}$ its max-times rank agrees with the max-plus rank of the matrix $\log(\mathbf{A})$. Yet, the errors obtained in approximate decompositions do not have to (and usually will not) agree. In what follows we characterize the relationship between max-plus and max-times errors. We denote by $\bar{\mathbb{R}}$ the extended real line $\mathbb{R} \cup \{-\infty\}$.

Theorem 4 *Let $\mathbf{A} \in \bar{\mathbb{R}}^{n \times m}$, $\mathbf{B} \in \bar{\mathbb{R}}^{n \times k}$ and $\mathbf{C} \in \bar{\mathbb{R}}^{k \times m}$. Let $M = \exp\{N\}$, where*

$$N = \max_{\substack{i \in [n] \\ j \in [m]}} \left\{ \max\{ \mathbf{A}_{ij}, \max_{1 \leq d \leq k} \{ \mathbf{B}_{id} + \mathbf{C}_{dj} \} \} \right\}.$$

If an error can be bounded in max-plus algebra as

$$\|\mathbf{A} - \mathbf{B} \diamond \mathbf{C}\|_F^2 \leq \lambda, \quad (12)$$

then the following estimate holds with respect to the max-times algebra:

$$\|\exp\{\mathbf{A}\} - \exp\{\mathbf{B}\} \boxtimes \exp\{\mathbf{C}\}\|_F^2 \leq M^2 \lambda. \quad (13)$$

Proof Let $\alpha_{ij} = \max_{d=1}^k \{\mathbf{B}_{id} + \mathbf{C}_{dj}\}$. From (12) it follows that there exists a set of numbers $\{\lambda_{ij} \geq 0 : i \in [n], j \in [m]\}$ such that for any i, j we have $(A_{ij} - \alpha_{ij})^2 \leq \lambda_{ij}$ and $\sum_{ij} \lambda_{ij} = \lambda$. By the mean-value theorem, for every i and j we obtain

$$|\exp\{\mathbf{A}_{ij}\} - \exp\{\alpha_{ij}\}| = |\mathbf{A}_{ij} - \alpha_{ij}| \exp\{\alpha_{ij}^*\} \leq \sqrt{\lambda_{ij}} \exp\{\alpha_{ij}^*\},$$

for some $\min\{\mathbf{A}_{ij}, \alpha_{ij}\} \leq \alpha_{ij}^* \leq \max\{\mathbf{A}_{ij}, \alpha_{ij}\}$. Hence,

$$(\exp\{\mathbf{A}_{ij}\} - \exp\{\alpha_{ij}\})^2 \leq \lambda_{ij} (\exp\{\max\{\mathbf{A}_{ij}, \alpha_{ij}\}\})^2.$$

The estimate for the max-times error now follows from the monotonicity of the exponent:

$$\begin{aligned} \|\exp\{\mathbf{A}\} - \exp\{\mathbf{B}\} \boxtimes \exp\{\mathbf{C}\}\|_F^2 &\leq \sum_{ij} (\exp\{\alpha_{ij}^*\})^2 \lambda_{ij} \\ &\leq \sum_{ij} (\exp\{\max\{\mathbf{A}_{ij}, \alpha_{ij}\}\})^2 \lambda_{ij} \leq M^2 \lambda, \end{aligned}$$

proving the claim. \square

3.4 Different subtropical matrix ranks

The definition of the subtropical rank we use in this work is the so-called Schein (or Barvinok) rank (see Definition 5). Like in the standard linear algebra, this is not the only possible way to define the (subtropical) rank. Here we will review few other forms of subtropical rank that can allow us to bound the Schein/Barvinok rank of a matrix. Unless otherwise mentioned, the definitions are by Guillon et al (2015); naturally results without citations are ours. Following Guillon et al, we will present the definitions in this section over the tropical algebra. Recall that due to isomorphism, these definitions transfer directly to the subtropical case.

We begin with the tropical equivalent of the subtropical Schein/Barvinok rank:

Definition 8 The *tropical Schein/Barvinok rank* of a matrix $\mathbf{A} \in \overline{\mathbb{R}}^{n \times m}$, denoted $\text{rank}_{\text{S/B}}(\mathbf{A})$, is defined to be the least integer k such that there exist matrices $\mathbf{B} \in \overline{\mathbb{R}}^{n \times k}$ and $\mathbf{C} \in \overline{\mathbb{R}}^{k \times m}$ for which $\mathbf{A} = \mathbf{B} \diamond \mathbf{C}$.

Analogous to the standard case, we can also define the rank as the number of linearly independent rows or columns. The following definition of linear independence of a family of vectors in a tropical space is due to Gondran and Minoux (1984b).

Definition 9 A set of vectors $\mathbf{x}_1, \dots, \mathbf{x}_k$ from $\overline{\mathbb{R}}^n$ is called *linearly dependent* if there exist disjoint sets $I, J \subset [k]$ and scalars $\{\lambda_i\}_{i \in I \cup J}$, such that $\lambda_i \neq -\infty$ for all i and

$$\max_{i \in I} \{\lambda_i + \mathbf{x}_i\} = \max_{j \in J} \{\lambda_j + \mathbf{x}_j\}. \quad (14)$$

Otherwise the vectors $\mathbf{x}_1, \dots, \mathbf{x}_k$ are called *linearly independent*.

This gives rise to the so-called *Gondran–Minoux ranks*:

Definition 10 The *Gondran–Minoux row (column) rank* of a matrix $\mathbf{A} \in \overline{\mathbb{R}}^{n \times m}$ is defined as the maximal k such that \mathbf{A} has k independent rows (columns). They are denoted by $\text{rank}_{G-M;rw}(\mathbf{A})$ and $\text{rank}_{G-M;cl}(\mathbf{A})$ respectively.

Another way to characterize the rank of the matrix is to consider the space its rows or columns can span.

Definition 11 A set $X \subset \overline{\mathbb{R}}^n$ is called *tropically convex* if for any vectors $\mathbf{x}, \mathbf{y} \in X$ and scalars $\lambda, \mu \in \overline{\mathbb{R}}$, we have $\max\{\lambda + \mathbf{x}, \mu + \mathbf{y}\} \in X$.

Definition 12 The *convex hull* $H(\mathbf{x}_1, \dots, \mathbf{x}_k)$ of a finite set of vectors $\{\mathbf{x}_i\}_{i=1}^k \in \overline{\mathbb{R}}^n$ is defined as follows

$$H(\mathbf{x}_1, \dots, \mathbf{x}_k) = \left\{ \max_{i=1}^k \{\lambda_i + \mathbf{x}_i\} : \lambda_i \in \overline{\mathbb{R}} \right\} .$$

Definition 13 The *weak dimension* of a finitely generated tropically convex subset of $\overline{\mathbb{R}}^n$ is the cardinality of its minimal generating set.

We can define the rank of the matrix by looking at the weak dimension of the (tropically) convex hull its rows or columns span.

Definition 14 The *row rank* and the *column rank* of a matrix $\mathbf{A} \in \overline{\mathbb{R}}^{n \times m}$ are defined as the weak dimensions of the convex hulls of the rows and the columns of \mathbf{A} respectively. They are denoted by $\text{rank}_{rw}(\mathbf{A})$ and $\text{rank}_{cl}(\mathbf{A})$.

None of the above definitions coincide (see Akian et al 2009), unlike in the standard algebra. We can, however, have a partial ordering of the ranks:

Theorem 5 (Guillon et al 2015; Akian et al 2009) *Let $\mathbf{A} \in \overline{\mathbb{R}}^{n \times m}$. Then the following relations are true for the above definitions of the rank of \mathbf{A} :*

$$\left. \begin{array}{l} \text{rank}_{G-M;rw}(\mathbf{A}) \\ \text{rank}_{G-M;cl}(\mathbf{A}) \end{array} \right\} \leq \text{rank}_{S/B}(\mathbf{A}) \leq \left\{ \begin{array}{l} \text{rank}_{rw}(\mathbf{A}) \\ \text{rank}_{cl}(\mathbf{A}) \end{array} \right\} . \quad (15)$$

The row and column ranks of an n -by- n tropical matrix can be computed in $O(n^3)$ time (Butkovič 2010), allowing us to bound the Schein/Barvinok rank from above. Unfortunately, no efficient algorithm for the Gondran–Minoux rank is known. On the other hand, Guillon et al (2015) presented what they called the *ultimate tropical rank* that lower-bounds the Gondran–Minoux rank and can be computed in time $O(n^3)$. We can also check if a matrix has full Schein/Barvinok rank in time $O(n^3)$ (see Butkovič and Hevry 1985), even if computing any other value is NP-hard.

These bounds, together with Lemma 1 yield the following corollary regarding the bounding of the *Boolean rank* of a square matrix:

Corollary 2 *Given an n -by- n binary matrix \mathbf{A} , its Boolean rank can be bound from below, using the ultimate rank, and from above, using the tropical column and row ranks, in time $O(n^3)$.*

4 Algorithms

There are some unique challenges in doing subtropical matrix factorization, that stem from the lack of linearity and smoothness of the max-times algebra. One of such issues is that dominated elements in a decomposition have no impact on the final result. Namely, if we consider the subtropical product of two matrices $\mathbf{B} \in \mathbb{R}_+^{n \times k}$ and $\mathbf{C} \in \mathbb{R}_+^{k \times m}$, we can see that each entry $(\mathbf{B} \boxtimes \mathbf{C})_{ij} = \max_{s \in [k]} \mathbf{B}_{is} \mathbf{C}_{sj}$ is completely determined by a single element with index $\arg \max_{s \in [k]} \mathbf{B}_{is} \mathbf{C}_{sj}$. This means that all entries t with $\mathbf{B}_{it} \mathbf{C}_{tj} < \max_{s \in [k]} \mathbf{B}_{is} \mathbf{C}_{sj}$ do not contribute at all to the final decomposition. To see why this is a problem, observe that many optimization methods used in matrix factorization algorithms rely on local information to choose the direction of the next step (e.g. various forms of gradient descent). In the case of the subtropical algebra, however, the local information is practically absent, and hence we need to look elsewhere for effective optimization techniques.

A common approach to matrix decomposition problems is to update factor matrices alternately, which utilizes the fact that the problem $\min_{\mathbf{B}, \mathbf{C}} \|\mathbf{A} - \mathbf{BC}\|_F$ is biconvex. Unfortunately, the subtropical matrix factorization problem does not have the biconvexity property, which makes alternating updates less useful.

Here we present a different approach that, instead of doing alternating factor updates, constructs the decomposition by adding one rank-1 matrix at a time, following the idea by Kolda and O'Leary (2000). The corresponding algorithm is called **Equator** (Algorithm 1).

First observe that the max-times product can be represented as an elementwise maximum of rank-1 matrices (blocks)

$$\mathbf{B} \boxtimes \mathbf{C} = \max_{s=1}^k \mathbf{B}^s \mathbf{C}_s. \quad (16)$$

Hence, Problem 1 can be split into k subproblems of the following form: given a rank- $(l-1)$ decomposition $\mathbf{B} \in \mathbb{R}_+^{n \times (l-1)}$, $\mathbf{C} \in \mathbb{R}_+^{(l-1) \times m}$ of a matrix $\mathbf{A} \in \mathbb{R}_+^{n \times m}$, find a column vector $\mathbf{b} \in \mathbb{R}_+^{n \times 1}$ and a row vector $\mathbf{c} \in \mathbb{R}_+^{1 \times m}$ such that the error

$$\|\mathbf{A} - \max\{\mathbf{B} \boxtimes \mathbf{C}, \mathbf{bc}\}\| \quad (17)$$

is minimized. We assume by definition that the rank-0 decomposition is an all zero matrix of the same size as \mathbf{A} . The problem of rank- k subtropical matrix factorization is then reduced to solving (17) k times. One should of course remember that this scheme is just a heuristic and finding optimal blocks on each iteration does not guarantee converging to a global minimum.

One prominent issue with the above approach is that an optimal rank- $(k-1)$ decomposition might not be very good when considered as a part of a rank- k decomposition. This is because for smaller ranks we generally have to cover the data more crudely, whereas when the rank increases we can afford to use smaller and more refined blocks. In order to deal with this problem, we find and then update the blocks repeatedly, in a cyclic fashion. That means that after discovering the last block, we go all the way back to block one. The input parameter M defines the number of full cycles we make.

On a high level **Equator** works as follows. First the factor matrices are initialized to all zeros (line 2). Since the algorithm makes iterative changes to the current solutions that might in some cases lead to worsening of the results, it also stores the best reconstruction error and the corresponding factors found so far. They are initialized

Algorithm 1 Equator

Input: $A \in \mathbb{R}_+^{n \times m}$, $k > 0$, $M > 0$
Output: $B^* \in \mathbb{R}_+^{n \times k}$, $C^* \in \mathbb{R}_+^{k \times m}$

```

1: function Equator( $A, k, M$ )
2:    $B \leftarrow 0^{n \times k}$ ,  $C \leftarrow 0^{k \times m}$ 
3:    $B^* \leftarrow B$ ,  $C^* \leftarrow C$ 
4:    $bestError \leftarrow E(A, B, C)$ 
5:   for  $count \leftarrow 1$  to  $k \times M$  do
6:      $l \leftarrow (count - 1) \pmod k + 1$  ▷ Index of the current block
7:      $[B^l, C_l] \leftarrow \text{UpdateBlock}(A, B, C, count)$ 
8:     if  $E(A, B, C) < bestError$  then
9:        $B^* \leftarrow B$ ,  $C^* \leftarrow C$ 
10:       $bestError \leftarrow E(A, B, C)$ 
11:  return  $B^*, C^*$ 

```

with the starting solution on lines 3–4. The main work is done in the loop on lines 5–10, where on each iteration we update a single rank-1 matrix in the current decomposition using the `UpdateBlock` routine (line 7), and then check if the update improves the best result (lines 8–10).

We will present two versions of the `UpdateBlock` function, one called `Capricorn` and the other one `Cancer`. `Capricorn` is designed to work with discrete (or flipping) noise, when some of the elements in the data are randomly changed to different values. In this setting the level of noise is the proportion of the flipped elements relative to the total number of nonzeros. `Cancer` on the other hand is robust with continuous noise, when many elements are affected (e.g. Gaussian noise). We will discuss both of them in detail in the following subsections. In the rest of the paper, especially when presenting the experiments, we will use names `Capricorn` and `Cancer` not only for a specific variation of the `UpdateBlock` function, but also for the `Equator` algorithm that uses it.

4.1 Capricorn

We first describe `Capricorn`, which is designed to solve the subtropical matrix factorization problem in the presence of discrete noise, and minimizes the L_1 norm of the error matrix. The main idea behind the algorithm is to spot potential blocks by considering ratios of matrix rows. Consider an arbitrary rank-1 block $X = bc$, where $b \in \mathbb{R}_+^{n \times 1}$ and $c \in \mathbb{R}_+^{1 \times m}$. For any indices i and j such that $b_i > 0$ and $b_j > 0$, we have $X_j = \frac{b_j}{b_i} X_i$. This is a characteristic property of rank-1 matrices – all rows are multiples of one another. Hence, if a block X dominates some region Γ of a matrix A , then rows of A should all be multiples of each other within Γ . These rows might have different lengths due to block overlap, in which case the rule only applies to their common part.

`UpdateBlock` starts by identifying the index of the block that has to be updated at the current iteration (line 2). In order to find the best new block we need to take into account that some parts of the data have already been covered, and we must ignore them. This is accomplished by replacing the original matrix with a residual R that represents what there is left to cover. The building of the residual (line 3) reflects the winner-takes-it-all property of the max-times algebra: if an element of A is approximated by a smaller value, it appears as such in the residual; if it is approximated by a value that is at least as large, then the corresponding residual element is NaN ,

Algorithm 2 UpdateBlock (Capricorn)

Input: $\mathbf{A} \in \mathbb{R}_+^{n \times m}$, $\mathbf{B} \in \mathbb{R}_+^{n \times k}$, $\mathbf{C} \in \mathbb{R}_+^{k \times m}$, $count > 0$
Output: $\mathbf{b} \in \mathbb{R}_+^{n \times 1}$, $\mathbf{c} \in \mathbb{R}_+^{1 \times m}$
Parameters: $bucketSize > 0$, $\delta > 0$, $\theta > 0$, $\tau \in [0, 1]$

- 1: **function** UpdateBlock($\mathbf{A}, \mathbf{B}, \mathbf{C}, count$)
- 2: $l \leftarrow (count - 1) \pmod k + 1$ ▷ Index of the current block
- 3: $\mathbf{R}_{ij} \leftarrow \begin{cases} \mathbf{A}_{ij} & (\mathbf{B}^{-l} \boxtimes \mathbf{C}_{-l})_{ij} < \mathbf{A}_{ij} \\ NaN & \text{otherwise} \end{cases}$ ▷ Residual matrix
- 4: $idx \leftarrow \arg \max_i \sum_j r_{ij}$
- 5: $\mathbf{H} \leftarrow \text{CorrelationsWithRow}(\mathbf{R}, idx, bucketSize, \delta, \tau)$
- 6: $r \leftarrow \arg \max_i \sum_j h_{ij}$
- 7: $c \leftarrow \arg \max_j \sum_i h_{ij}$
- 8: $b_idx \leftarrow \{i : \mathbf{H}_{ic} = 1\}$
- 9: $c_idx \leftarrow \{i : \mathbf{H}_{ri} = 1\}$
- 10: $[\bar{\mathbf{b}}, \bar{\mathbf{c}}] \leftarrow \text{RecoverBlock}(\mathbf{R}, b_idx, c_idx)$
- 11: $\mathbf{b} \leftarrow \text{AddRows}(\bar{\mathbf{b}}, \mathbf{c}, \mathbf{A}, \theta, bucketSize, \delta)$
- 12: $\mathbf{c} \leftarrow \text{AddRows}(\bar{\mathbf{c}}^T, \mathbf{b}^T, \mathbf{A}^T, \theta, bucketSize, \delta)^T$
- 13: **return** \mathbf{b}, \mathbf{c}

indicating that this value is already covered. We then select a seed row (line 4), with an intention of growing a block around it. We choose the row with the largest sum as this increases the chances of finding the most prominent block. In order to find the best block \mathbf{X} that the seed row passes through, we first find a binary matrix \mathbf{H} that represents the pattern of \mathbf{X} (line 5). Next, on lines 6–9 we choose an approximation of the block pattern with index sets b_idx and c_idx , which define what elements of \mathbf{b} and \mathbf{c} should be nonzero. The next step is to find the actual values of elements within the block with the function `RecoverBlock` (line 10). Finally, we inflate the found core block with `ExpandBlock` (line 11).

The function `CorrelationsWithRow` (Algorithm 3) finds the pattern of a new block. It does so by comparing a given seed row to other rows of the matrix and extracting sets where the ratio of the rows is almost constant. As was mentioned before, if two rows locally represent the same block, then one should be a multiple of the other, and the ratios of their corresponding elements should remain level. `CorrelationsWithRow` processes the input matrix row by row using the function `FindRowSet`, which for every row outputs the most likely set of indices, where it is correlated with the seed row (lines 4–6). Since the seed row is obviously the most correlated with itself, we compensate for this by replacing its pattern with that of the second most correlated row (lines 7–8). Finally, we drop some of the least correlated rows after comparing their correlation value ϕ to that of the second most correlated row (after the seed row). The correlation function ϕ is defined as follows

$$\phi(\mathbf{H}, idx, i) = \frac{\langle \mathbf{H}_i, \mathbf{H}_{idx} \rangle}{\langle \mathbf{H}_i, \mathbf{H}_i \rangle + 1}. \quad (18)$$

The parameter τ is a threshold determining whether a row should be discarded or retained. The auxiliary function `FindRowSet` (Algorithm 4) compares two vectors and finds the biggest set of indices where their ratio remains almost constant. It does so by sorting the log-ratio of the input vectors into buckets of a fixed size and then choosing the bucket with the most elements. The notation $\mathbf{u} ./ \mathbf{v}$ on line 2 means elementwise ratio of vectors \mathbf{u} and \mathbf{v} .

Algorithm 3 CorrelationsWithRow

Input: $\mathbf{R} \in \mathbb{R}_+^{n \times m}$, $idx \in [n]$, $bucketSize > 0$, $\delta > 0$, $\tau \in [0, 1]$
Output: $\mathbf{H} \in \{0, 1\}^{n \times m}$

- 1: **function** CorrelationsWithRow(\mathbf{R} , idx , $bucketSize$, δ , τ)
- 2: turn all *NaN* elements of \mathbf{R} to 0
- 3: $\mathbf{H} \leftarrow \mathbf{0}^{n \times m}$
- 4: **for** $i \leftarrow 1$ **to** n **do**
- 5: $V_i \leftarrow \text{FindRowSet}(\mathbf{R}_{idx}, \mathbf{R}_i, bucketSize, \delta)$
- 6: $\mathbf{H}(i, V_i) \leftarrow 1$
- 7: $s \leftarrow \arg \max_{i: i \neq idx} \sum_j h_{ij}$
- 8: $\mathbf{H}_{idx} \leftarrow \mathbf{H}_s$
- 9: **for** $i \leftarrow 1$ **to** n **do**
- 10: **if** $\phi(\mathbf{H}, idx, i) < \phi(\mathbf{H}, idx, s) - \tau$ **then**
- 11: $\mathbf{H}_i \leftarrow 0$
- 12: **return** \mathbf{H}

Algorithm 4 FindRowSet

Input: $\mathbf{u} \in \mathbb{R}_+^m$, $\mathbf{v} \in \mathbb{R}_+^m$, $bucketSize > 0$, $\delta > 0$
Output: $V \subset [m]$

- 1: **function** FindRowSet(\mathbf{u} , \mathbf{v} , $bucketSize$, δ)
- 2: $\mathbf{r} \leftarrow \log(\mathbf{u} ./ \mathbf{v})$
- 3: $nBuckets \leftarrow \lceil (\max\{\mathbf{r}\} - \min\{\mathbf{r}\}) / \delta \rceil$
- 4: **for** $i \leftarrow 0$ **to** $nBuckets$ **do**
- 5: $V_i \leftarrow \{idx \in [m] : \min\{\mathbf{r}\} + i\delta \leq r_{idx} < \min\{\mathbf{r}\} + (i + 1)\delta\}$
- 6: $V \leftarrow \arg \max\{|V_i| : i = 1, \dots, nBuckets\}$
- 7: **if** $|V| < bucketSize$ **then**
- 8: $V \leftarrow \emptyset$
- 9: **return** V

It accepts two additional parameters: *bucketSize* and δ . If the largest bucket has fewer than *bucketSize* elements, the function will return an empty set – this is done because very small patterns do not reveal much structure and are mostly accidental. The width of the buckets is determined by the parameter δ .

At this point we know the pattern of the new block, that is, the locations of its non-zeros. To fill in the actual values, we consider the submatrix defined by the pattern, and find the best rank-1 approximation of it. We do this using the **RecoverBlock** function (Algorithm 5). It begins by setting all elements outside of the pattern to 0 as they are irrelevant to the block (line 2). Then it chooses one row to represent the block (lines 3–4), which will be used to find a good rank-1 cover.

Finally, we find the optimal column vector for the block by computing the best weights to be used for covering different rows of the block with its representing row (line 5). Here we optimize with respect to the Frobenius norm, rather than L_1 matrix norm, since it allows to solve the optimization problem in closed form.

Since blocks often heavily overlap, we are susceptible to finding only fragments of patterns in the data – some parts of a block can be dominated by another block and subsequently not recognized. Hence, we need to expand found blocks to make them complete. This is done separately for rows and columns in the method called **AddRows** (Algorithm 6), which, given a starting block $\mathbf{X} = \mathbf{bc}$ and the original matrix \mathbf{A} , tries to add new nonzero elements to \mathbf{b} . It iterates through all rows of \mathbf{A} and adds those that would make a positive impact on the objective without unnecessarily overcovering the data. In order to decide whether a given row should be added, it first extracts a set V_i

Algorithm 5 RecoverBlock**Input:** $\mathbf{R} \in \mathbb{R}_+^{n \times m}$, $bIdx \subset [n]$, $cIdx \subset [m]$ **Output:** $\mathbf{b} \in \mathbb{R}_+^{n \times 1}$, $\mathbf{c} \in \mathbb{R}_+^{1 \times m}$

- 1: **function** RecoverBlock(\mathbf{R} , $bIdx$, $cIdx$)
- 2: turn \mathbf{R} to 0 except elements with indices ($bIdx$, $cIdx$)
- 3: $p \leftarrow \text{RowRepresentingBlock}(\mathbf{R}, bIdx)$
- 4: $\mathbf{c} \leftarrow \mathbf{R}_p$
- 5: $\mathbf{b} \leftarrow \arg \min_{\mathbf{t} \in \mathbb{R}_+^{n \times 1}} \|\mathbf{R} - \mathbf{t}\mathbf{c}\|_F$
- 6: **return** \mathbf{b} , \mathbf{c}

Algorithm 6 AddRows**Input:** $\mathbf{b} \in \mathbb{R}_+^{n \times 1}$, $\mathbf{c} \in \mathbb{R}_+^{1 \times m}$, $\mathbf{A} \in \mathbb{R}_+^{n \times m}$, $\theta > 0$, $bucketSize > 0$, $\delta > 0$ **Output:** $\mathbf{b} \in \mathbb{R}_+^{n \times 1}$

- 1: **function** AddRows(\mathbf{b} , \mathbf{c} , \mathbf{A} , θ , $bucketSize$, δ)
- 2: $b_idx \leftarrow \{t : \mathbf{b}_t > 0\}$
- 3: **for** $i \in [n] \setminus b_idx$ **do**
- 4: $V_i \leftarrow \text{FindRowSet}(\mathbf{c}, \mathbf{R}_i, bucketSize, \delta)$
- 5: **if** $V_i = \emptyset$ **then**
- 6: **continue**
- 7: $\alpha \leftarrow \text{mean}(\mathbf{R}_{iV_i} ./ \mathbf{c}_{V_i})$
- 8: $impact \leftarrow \frac{\sum_{s \in V_i} \max\{0, \alpha c_s - \mathbf{A}_{is}\}}{\sum_{s \in V_i} \mathbf{A}_{is} - |\mathbf{A}_{is} - \alpha c_s|}$
- 9: **if** $impact \leq \theta$ **then**
- 10: $\mathbf{b}_i \leftarrow \alpha$
- 11: **return** \mathbf{b}

of indices where this row is a multiple of the row vector \mathbf{c} of the block (if they are not sufficiently correlated, then the row does not belong to the block) (line 4). A row is added if the evaluation of the following function (line 8)

$$\psi(\alpha) = \frac{\sum_{s \in V_i} \max\{0, \alpha c_s - \mathbf{A}_{is}\}}{\sum_{s \in V_i} \mathbf{A}_{is} - |\mathbf{A}_{is} - \alpha c_s|} \quad (19)$$

is below the threshold θ . In (19) the numerator measures by how much the new row would overcover the original matrix, and the denominator reflects the improvement in the objective compared to a zero row.

Parameters. **Capricorn** has four parameters in addition to the common parameters in the Equator framework: $bucketSize > 0$, $\delta > 0$, $\theta > 0$, and $\tau \in [0, 1]$. The first one, $bucketSize$ determines the minimum number of elements in two rows that must have “approximately” the same ratio for them to be considered for building a block. The parameter δ defines the bucket width when computing row correlations. When expanding a block, θ is used to decide whether to add a row (or column) to it – the decision is positive whenever the expression (19) is at most θ . Finally τ is used during the discovery of correlated rows. The value of τ belongs to the closed unit interval, and the higher it is, the more rows will be added.

4.2 Cancer

We now present our second algorithm, **Cancer**, which is a counterpart of **Capricorn** specifically designed to work in the presence of high levels of continuous noise. The

reason why **Capricorn** cannot deal with continuous noise is that it expects the rows in a block to have an “almost” constant elementwise ratio, which is not the case when too many entries in the data are disturbed. For example, even low levels of Gaussian noise would make the ratios vary enough to hinder **Capricorn**’s ability to spot blocks. With **Cancer** we take a new approach which is based on polynomial approximation of the objective. We also replace the L_1 matrix norm, which was used as an objective for **Capricorn**, with the Frobenius norm. The reason for that is that when the noise is continuous, its level is defined as the total deviation of the noisy data from the original, rather than a count of the altered elements. This makes the Frobenius norm a good estimator for the amount of noise. **Cancer** conforms to the general framework of **Equator** (Algorithm 1), and differs from **Capricorn** only in how it finds the blocks and in the objective function.

Observe that in order to solve the problem (17) we need to find a column vector $\mathbf{b} \in \mathbb{R}_+^{n \times 1}$ and a row vector $\mathbf{c} \in \mathbb{R}_+^{1 \times m}$ such that they provide the best rank-1 approximation of the input matrix given the current factorization. The objective function is not convex in either \mathbf{b} or \mathbf{c} and is generally hard to optimize directly, so we have to simplify the problem, which we do in two steps. First, instead of doing full optimization of \mathbf{b} and \mathbf{c} simultaneously, we update only a single element of one of them at a time. This way the problem is reduced to single variable optimization. Even then the objective is hard to minimize, and we replace it with a polynomial approximation, which is easy to optimize directly.

The **Cancer** version of the **UpdateBlock** function is described in Algorithm 7. It alternately updates the vectors \mathbf{b} and \mathbf{c} using the **AdjustOneElement** routine. Both \mathbf{b} and \mathbf{c} will be updated $\lfloor f(n+m)/2 \rfloor$ times. **UpdateBlock** starts by finding the index of the block that has to be changed (line 2). Since the purpose of **UpdateBlock** is to find the best rank-1 matrix to replace the current block, we also need to compute the reconstructed matrix without it, which is done on line 3. We then find the number of times **AdjustOneElement** will be called (line 4) and change the degree of polynomials used for objective function approximation (line 5). This is needed because high degree polynomials are better at finalizing a solution that is already reasonably good, but tend to overfit the data and cause the algorithm to get stuck in local minima at the beginning. It is therefore beneficial to start with polynomials of lower degrees and then gradually increase it. The actual changes to \mathbf{b} and \mathbf{c} happen in the loop (lines 7–9), where we update them using **AdjustOneElement**.

The **AdjustOneElement** function (Algorithm 8) updates a single entry in either a column vector \mathbf{b} or a row vector \mathbf{c} . Let us consider the case when \mathbf{b} is fixed and \mathbf{c} varies. In order to decide which element of \mathbf{c} to change, we need to compare the best changes to all m entries and then choose the one that yields the most improvement to the objective. A single element \mathbf{c}_l only has an effect on the error along the column l . Assume that we are currently updating block with index q and let \mathbf{N} denote the reconstruction matrix without this block, that is $\mathbf{N} = \mathbf{B}^{-q} \boxtimes \mathbf{C}_{-q}$. Minimizing $E(\mathbf{A}, \mathbf{B}, \mathbf{C})$ with respect to \mathbf{c}_l is then equivalent to minimizing

$$\gamma(\mathbf{A}_l, \mathbf{N}_l, \mathbf{b}, \mathbf{c}_l) = \sum_{i=1}^n (\mathbf{A}_{il} - \max\{\mathbf{N}_{il}, \mathbf{b}_i \mathbf{c}_l\})^2. \quad (20)$$

Instead of minimizing (20) directly, we use polynomial approximation in the **PolyMin** routine (line 4). It returns the (approximate) error err and the value x achieving that. The polynomial approximation is obtained by evaluating the objective function at

Algorithm 7 UpdateBlock (Cancer)

Input: $A \in \mathbb{R}_+^{n \times m}$, $B \in \mathbb{R}_+^{n \times k}$, $C \in \mathbb{R}_+^{k \times m}$, $count > 0$
Output: $\mathbf{b} \in \mathbb{R}_+^{n \times 1}$, $\mathbf{c} \in \mathbb{R}_+^{1 \times m}$
Parameters: $t > 2$, $0 < f < 1$

- 1: **function** UpdateBlock($A, B, C, count$)
- 2: $l \leftarrow (count - 1) \pmod{k} + 1$ ▷ Index of the current block
- 3: $N \leftarrow B^{-l} \boxtimes C_{-l}$ ▷ Reconstructed matrix without the l -th block
- 4: $niters \leftarrow \lfloor f(n + m)/2 \rfloor$
- 5: $deg \leftarrow 2 + \lfloor (count - 1)/k \rfloor \pmod{t}$
- 6: $\mathbf{b} \leftarrow B^l$, $\mathbf{c} \leftarrow C_l$
- 7: **for** $iter \leftarrow 1$ **to** $niters$ **do**
- 8: $\mathbf{c} = \text{AdjustOneElement}(A, N, \mathbf{b}, \mathbf{c}, deg)$
- 9: $\mathbf{b} = \text{AdjustOneElement}(A^T, N^T, \mathbf{c}^T, \mathbf{b}^T, deg)^T$
- 10: **return** \mathbf{b}, \mathbf{c}

Algorithm 8 AdjustOneElement

Input: $A \in \mathbb{R}_+^{n \times m}$, $N \in \mathbb{R}_+^{n \times m}$, $\mathbf{b} \in \mathbb{R}_+^{n \times 1}$, $\mathbf{c} \in \mathbb{R}_+^{1 \times m}$, $deg \geq 2$
Output: $\mathbf{c} \in \mathbb{R}_+^{1 \times m}$

- 1: **function** AdjustOneElement($A, N, \mathbf{b}, \mathbf{c}, deg$)
- 2: **for** $j \leftarrow 1$ **to** m **do**
- 3: $baseError \leftarrow \sum_{i=1}^n (A_{ij} - \max\{N_{ij}, \mathbf{b}_i \mathbf{c}_j\})^2$
- 4: $[err, \mathbf{x}_i] \leftarrow \text{PolyMin}(A^j, N^j, \mathbf{b}, deg)$
- 5: $\mathbf{u}_i \leftarrow baseError - err$
- 6: $i \leftarrow$ the index i of largest value of \mathbf{u}
- 7: $\mathbf{c}_i \leftarrow \mathbf{x}_i$
- 8: **return** \mathbf{c}

$deg + 1$ points generated uniformly at random from the interval $[0, 5]$ and then fitting a polynomial to the obtained values. The upper bound of 5 does not have any special meaning, rather it was chosen by trial and error. `PolyMin` is a heuristic and does not necessarily find the global minimum of the objective function. Moreover, in rare cases it might even cause an increase in the objective value. In such cases it would, in theory, make sense to just keep the value prior to the update, as in that case the objective at least does not increase. However in practice this phenomenon helps to get out of local minima. Since we are only interested in the improvement of the objective achieved by updating a single entry of \mathbf{c} , we compute the improvement of the objective after the change (line 5). After trying every column of \mathbf{c} , we update only the column that yield the largest improvement.

The function γ that we need to minimize in order to find the best change to the vector \mathbf{c} in `AdjustOneElement` is hard to work with directly since it is not convex, and also not smooth because of the presence of the maximum operator. To alleviate this, we approximate the error function γ with a polynomial g of degree deg . Notice that when updating \mathbf{c}_l , other variables of γ are fixed and we only need to consider function $\gamma'(x) = \gamma(A_l, N_l, \mathbf{b}, x)$. To build g we sample $deg + 1$ points from $(0, 1)$ and fit g to the values of γ' at these points. We then find the $x \in \mathbb{R}_+$ that minimizes $g(x)$ and return $g(x)$ (the approximate error) and x (the optimal value).

Parameters. `Cancer` has two parameters, $t > 2$ and $0 < f < 1$, that control its execution. The first one, t , is the maximum allowed degree of polynomials used for approximation of the objective, which we set to 16 in all our experiments. The second parameter, f , determines the number of single element updates we make to the row

and column vectors of a block in `UpdateBlock`. To demonstrate that the chosen values of the parameters are reasonable, we performed a grid search for various parameter values (see Figure 4 in Section 5).

Generalized Cancer. The `Cancer` algorithm can be adapted to optimize other objective functions. Its general polynomial approximation framework allows for a wide variety of possible objectives, the only constraint being that they have to be additive (we call a function $E(\mathbf{A}, \mathbf{R})$ *additive* if there exists a mapping $\phi: \mathbb{R}_+ \times \mathbb{R}_+ \rightarrow \mathbb{R}_+$ such that for all $\mathbf{A} \in \mathbb{R}_+^{n \times m}$ and $\mathbf{R} \in \mathbb{R}_+^{n \times m}$ we have $E(\mathbf{A}, \mathbf{R}) = \sum_{ij} \phi(\mathbf{A}_{ij}, \mathbf{R}_{ij})$). Some examples of such functions are L_1 and Frobenius matrix norms, as well as Kullback–Leibler and Jensen–Shannon divergences. In order to use the generalized form of `Cancer` one simply has to replace the Frobenius norm with another cost function wherever the error is evaluated.

4.3 Time complexity

The main work in `Equator` is performed inside the `UpdateBlock` routine, which is called Mk times. Since M is a constant parameter, the complexity of `Equator` is k times the complexity of `UpdateBlock`. In the following we find the theoretical bounds on the execution time of `UpdateBlock` for both `Capricorn` and `Cancer`.

Capricorn. In the case of `Capricorn` there are three main contributors to `UpdateBlock` (Algorithm 2): `CorrelationsWithRow`, `RecoverBlock`, and `AddRows`.

`CorrelationsWithRow` compares every row to the seed row, each time calling `FindRowSet`, which in turn has to process all m elements of both rows. This results in the total complexity of `CorrelationsWithRow` being $O(nm)$. To find the complexity of `RecoverBlock`, first observe that any “pure” block \mathbf{X} can be represented as $\mathbf{X} = \mathbf{bc}$, where $\mathbf{b} \in \mathbb{R}_+^{n' \times 1}$ and $\mathbf{c} \in \mathbb{R}_+^{1 \times m'}$ with $n' \leq n$ and $m' \leq m$. `RecoverBlock` selects \mathbf{c} from the rows of \mathbf{X} and then finds the corresponding column vector \mathbf{b} that minimizes $\|\mathbf{X} - \mathbf{bc}\|_F$. In order to select the best row, we have to try each of the n' candidates, and since finding the corresponding \mathbf{b} for each of them takes time $O(n'm')$, this gives the runtime of `RecoverBlock` as $O(n')O(n'm') = O(n^2m)$. The most computationally expensive parts of `AddRows` are `FindRowSet` (line 4), finding the mean (line 7), and computing the impact (line 8), which all run in $O(m)$ time. All of these operations have to be repeated $O(n)$ times, and hence the runtime of `AddRows` is $O(nm)$. Thus, we can now estimate the complexity of `UpdateBlock` to be $O(nm) + O(n^2m) + O(nm) = O(n^2m)$, which leads to the total runtime of `Capricorn` to be $O(n^2mk)$.

Cancer. Here `UpdateBlock` (Algorithm 7) is a loop that calls `AdjustOneElement` $\lfloor f(n+m) \rfloor$ times. In `AdjustOneElement` the contributors to the complexity are computing the base error (line 3) and a call to `PolyMin` (line 4). Both of them are performed n or m times depending on whether we supplied the column vector \mathbf{b} or the row vector \mathbf{c} to `AdjustOneElement`. Finding the base error takes time $O(m)$ for \mathbf{b} and $O(n)$ for \mathbf{c} . The complexity of `PolyMin` boils down to that of evaluating the max-times objective at $deg + 1$ points and then minimizing a degree deg polynomial. Hence, `PolyMin` runs in time $O(m)$ or $O(n)$ depending on whether we are optimizing \mathbf{b} or \mathbf{c} , and the complexity of `AdjustOneElement` is $O(nm)$.

Since `AdjustOneElement` is called $\lfloor f(n+m)/2 \rfloor$ times and f is a fixed parameter, this gives the complexity $O((n+m)nm)$ for `UpdateBlock` and $O((n+m)nmk) = O(\max\{n, m\}nmk)$ for `Cancer`.

Empirical evaluation of the time complexity is reported in Section 5.3.

5 Experiments

We tested both **Capricorn** and **Cancer** on synthetic and real-world data. In addition we also compare against a variation of **Cancer** that optimizes the Jensen–Shannon divergence, which we call **CancerJS**. The purpose of the synthetic experiments is to evaluate the properties of the algorithm in controlled environments where we know the data has the max-times structure. They also demonstrate on what kind of data each algorithm excels and what their limitations are. The purpose of the real-world experiments is to confirm that these observations also hold true in real-world data, and to study what kinds of data sets actually have max-times structure. The source code of **Capricorn** and **Cancer** and the scripts that run the experiments in this paper are freely available for academic use.³

Parameters of Cancer. Both variations of **Cancer** use the same set of parameters. For the synthetic experiments we used $M = 14$, $t = 16$, and $f = 0.1$. For the real world experiments we set $t = 16$, $f = 0.1$, and $M = 40$ (except for **Eigenfaces**, where we used $M = 50$ and **Bas1LP**, where we set $M = 8$). Increasing M , which controls the number of cycles of execution of **Cancer**, almost invariably improves the results. At some point though, the gains become marginal, and the value of $M = 40$ is chosen so as to reach the point where increasing M further would not yield much improvement. Sometimes though, this moment can be reached faster – for example the smaller choice of M for **Bas1LP** is motivated by the fact that **Cancer** quickly reached a point where it could no longer make significant progress, despite **Bas1LP** being the largest dataset. The relationship of the other two parameters and the quality of decomposition is more complex. We see in Figure 4(a) that the dependence on f and t is not monotone, and it is hard to pinpoint the best combination exactly. Moreover, the optimal values can differ depending on the dataset; for example, Figure 4(b) features an almost monotone dependence on f that flattens out before f reaches 0.1. From our experience, however, the values of $t = 16$ and $f = 0.1$ seem to be a good choice.

Parameters of Capricorn. In both synthetic and real-world experiments we used the following default set of parameters: $M = 4$, $bucketSize = 3$, $\delta = 0.01$, $\theta = 0.5$, and $\tau = 0.5$. As with **Cancer**, there is a complex dependency of the results on the parameters, but the values chosen above seem to produce good results in most cases. We do not show a comparison table, as we did with **Cancer**, due to a bigger number of parameters.

5.1 Other methods

We compared our algorithms against **SVD** and five versions of **NMF**. For **SVD**, we used Matlab’s built-in implementation. The first form of **NMF** is a sparse **NMF** algorithm by Hoyer (2004),⁴ which we call **SNMF**. It defines the sparsity of a vector $\mathbf{x} \in \mathbb{R}_+^n$ as

$$\text{sparsity}_H(\mathbf{x}) = \frac{\sqrt{n} - (\sum_i |\mathbf{x}_i|) / \sqrt{\sum_i \mathbf{x}_i^2}}{\sqrt{n} - 1}, \quad (21)$$

³<http://cs.uef.fi/~pauli/tropical/>

⁴<https://github.com/aludnam/MATLAB/tree/master/nmfpack>, accessed 18 July 2017

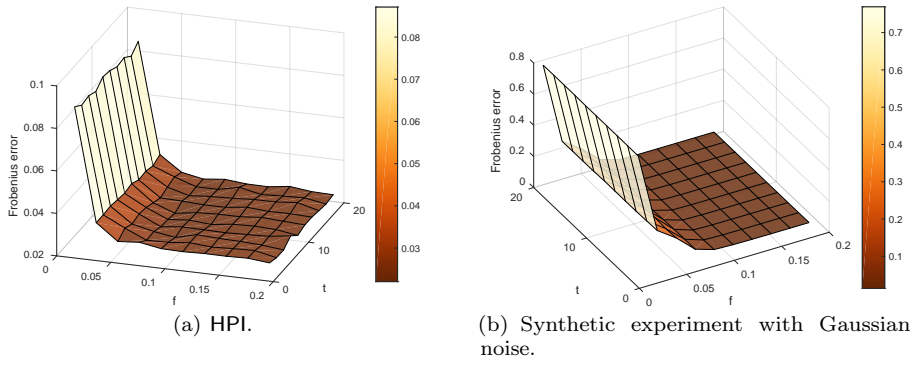


Fig. 4 Results of **Cancer** with different values of parameters t and f .

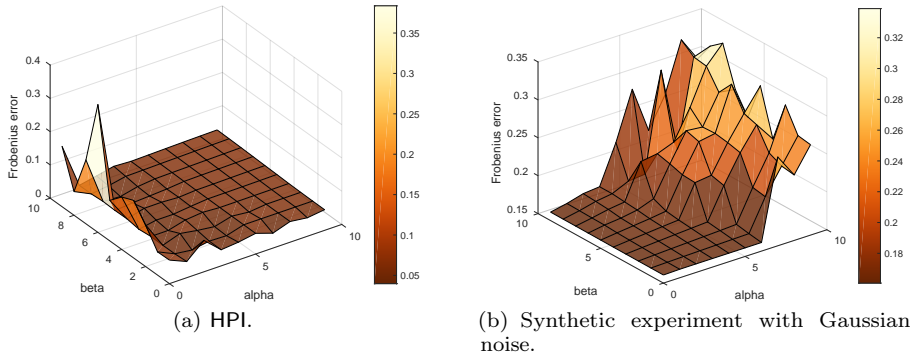


Fig. 5 Results of **ALSR** with different regularization parameters for the two factors matrices, which are denoted by α and β respectively.

and returns factorizations where the sparsity of the factor matrices is user-controllable. Note that the above definition of sparsity is different from the one we use elsewhere (see Equation (4)). In order to run **SMMF** we used the sparsity of **Cancer**'s factors (as defined by (21)) as its sparsity parameter. We also compare against a standard alternating least squares algorithm called **ALS** (Cichocki et al 2009). Next we have two versions of **NMF** that are essentially the same as **ALS**, but they use L_1 regularization for increased sparsity (Cichocki et al 2009), that is, they aim at minimizing

$$\|A - BC\|_F + \alpha \|B\|_1 + \beta \|C\|_1 .$$

The first method is called **ALSR** and uses regularizer coefficient $\alpha = \beta = 1$, and the other, called **ALSR5**, has regularizer coefficient $\alpha = \beta = 5$. It is natural to ask how **ALSR** would fare with different values of parameters. In Figure 5 we perform a grid search for the best parameter combination. While the experiment with **HPI** has a very uneven surface without much structure apart from a couple of spikes, the synthetic dataset demonstrates that high values of α and β can have serious adverse effects on the reconstruction error. It therefore seems safest to set $\alpha = \beta = 0$, which corresponds to the **ALS** method. It is worth mentioning that in many of our experiments larger values of α and β resulted in factors becoming close to zero, or some elements in the

factors getting enormous values due to numeric instability. This was the case for some other real-world experiments, such as **4News**, which is another indication to use the parameter values of $\alpha = \beta = 0$.

The last NMF algorithm, **WNMF** by Li and Ngom (2013), is designed to work with missing values in the data.

5.2 Synthetic experiments

The purpose of synthetic experiments is to prove the concept, that is that our algorithms are capable of identifying the max-times structure when it is there. In order to test this, we first generate the data with the pure max-times structure, then pollute it with some level of noise, and finally run the methods. The noise-free data is created by first generating random factors of some density with nonzero elements drawn from a uniform distribution on the $[0, 1]$ interval and then multiplying them using the max-times matrix product.

We distinguish two types of noise. The first one is the discrete (or tropical) noise, which is introduced in the following way. Assume that we are given an input matrix \mathbf{A} of size n -by- m . We first generate an n -by- m noise matrix \mathbf{N} with elements drawn from a uniform distribution on the $[0, 1]$ interval. Given a level of noise l , we then turn $\lfloor (1-l)nm \rfloor$ random elements of \mathbf{N} to 0, so that its resulting density is l . Finally, the noise is applied by taking elementwise maximum between the original data and the noise matrix $\mathbf{F} = \max\{\mathbf{A}, \mathbf{N}\}$. This is the kind of noise that **Capricorn** was designed to handle, so we expect it to be better than **Cancer** and other comparison algorithms.

We also test against continuous noise, as it is arguably more common in the real world. For that we chose Gaussian noise with 0 mean, where the noise level is defined to be its standard deviation. Since adding this noise to the data might result in negative entries, we truncate all values in a resulting matrix that are below zero.

Unless specified otherwise, all matrices in the synthetic experiments are of size 1000-by-800 with true max-times rank 10. All results presented in this section are averaged over 10 instances. For reconstruction error tests, we compared our algorithms **Capricorn**, **Cancer**, and **CancerJS** against **SVD**, **NMF**, **SNMF**, **ALS**, **ALSR**, and **ALSR5**. The error is measured as the relative Frobenius norm $\|\tilde{\mathbf{A}} - \mathbf{A}\|_F / \|\mathbf{A}\|$, where \mathbf{A} is the data and $\tilde{\mathbf{A}}$ its approximation, as that is the measure both **SVD** and **NMF** aim at minimizing. We also report the sparsity s of factor matrices obtained by algorithms, which is defined as a fraction of zero elements in the factor matrices, see (4). for an n -by- m matrix \mathbf{A} . For the experiments with tropical noise, the reconstruction errors are reported in Figure 6 and factor sparsity in Figure 7. For the Gaussian noise experiments, the reconstruction errors and factor sparsity are shown in Figure 8 and Figure 9 respectively.

Varying density with tropical noise. In our first experiment we studied the effects of varying the density of the factor matrices in presence of the tropical noise. We changed the density of the factors from 10% to 100% with an increment of 10%, while keeping the noise level at 10%. Figure 6(a) shows the reconstruction error and Figure 7(a) the sparsity of the obtained factors. **Capricorn** is consistently the best method, obtaining almost perfect reconstruction; only when the density approaches 100% does its reconstruction error deviate slightly from 0. This is expected since the data was generated with the tropical (flipping) noise that **Capricorn** is designed to optimize. Compared to **Capricorn** all other methods clearly underperform, with **Cancer**

being the second best. With the exception of **ALSR5**, all NMF methods obtain results similar to those of **SVD**, while having a somewhat higher reconstruction error than **Cancer**. That **SVD** and NMF methods (except **ALSR5**) start behaving better at higher levels of density indicates that these matrices can be explained relatively well using standard algebra. **Capricorn** and **Cancer** also have the highest sparsity of factors, with **Capricorn** exhibiting a decrease in sparsity as the density of the input increases. This behaviour is desirable since ideally we would prefer to find factors that are as close to the original ones as possible. For NMF methods there is a trade-off between the reconstruction error and the sparsity of the factors – the algorithms that were worse at reconstruction tend to have sparser factors.

Varying tropical noise. The amount of noise is always with respect to the number of nonzero elements in a matrix, that is, for a matrix \mathbf{A} with $\kappa(\mathbf{A})$ nonzero elements and noise level α , we flip $\alpha\kappa(\mathbf{A})$ elements to random values. There are two versions of this experiment – one with factor density 30% and the other with 60%. In both cases we varied the noise level from 0% to 110% with increments of 10%. Figure 6(b) and Figure 6(c) show the respective reconstruction errors and Figure 7(b) and Figure 7(c) the corresponding sparsities of the obtained factors. In the low-density case, **Capricorn** is consistently the best method with essentially perfect reconstruction for up to 80% of noise. In the high-density case, however, the noise has more severe effects, and in particular after 60% of noise, **Cancer**, **SVD**, and all versions of NMF are better than **Capricorn**. The severity of the noise is, at least partially, explained by the fact that in the denser data we flip more elements than in sparser data: for example when the data matrices are full, at 50% of noise, we have already replaced half of the values in the matrices with random values. Further, the quick increase of the reconstruction error for **Capricorn** hints strongly that the max-times structure of the data is mostly gone at these noise levels. **Capricorn** also produces clearly the sparsest factors for the low density case, and is mostly tied with **Cancer** and **ALSR5** when the density is high. It should be noted, however, that **ALSR5** generally has the highest reconstruction error among all the methods, which suggests that its sparse factors come at the cost of recovering little structure from the data.

Varying rank with tropical noise. Here we test the effects of the (max-times) rank, with the assumption that higher-rank matrices are harder to reconstruct. The true max-times rank of the data varied from 2 to 20 with increments of 2. There are three variations of this experiment: with 30% factor density and 10% noise (Figure 6(d)), with 30% factor density and 50% noise (Figure 6(e)), and with 60% factor density and 10% noise (Figure 6(f)). The corresponding sparsities are shown on Figures 7(d)–(f). **Capricorn** has a clear advantage for all settings, obtaining nearly perfect reconstruction. **Cancer** is generally second best, except for the high noise case, where it is mostly tied with a bunch of NMF methods. It also has a relatively high variance. To see why this happens, consider that **Cancer** always updates one element in factor matrices at a time. This update is completely dependent on values on a single row (or column) and is sensitive to the spikes that tropical noise introduces to some elements. Interestingly, on the last two plots the reconstruction error actually drops for **Cancer**, **SVD**, and NMF-based methods. This is a strong indication that at this point they no longer can extract meaningful structure in the data, and the improvement of the reconstruction error is largely due to uniformization of the data caused by high density and high noise levels.

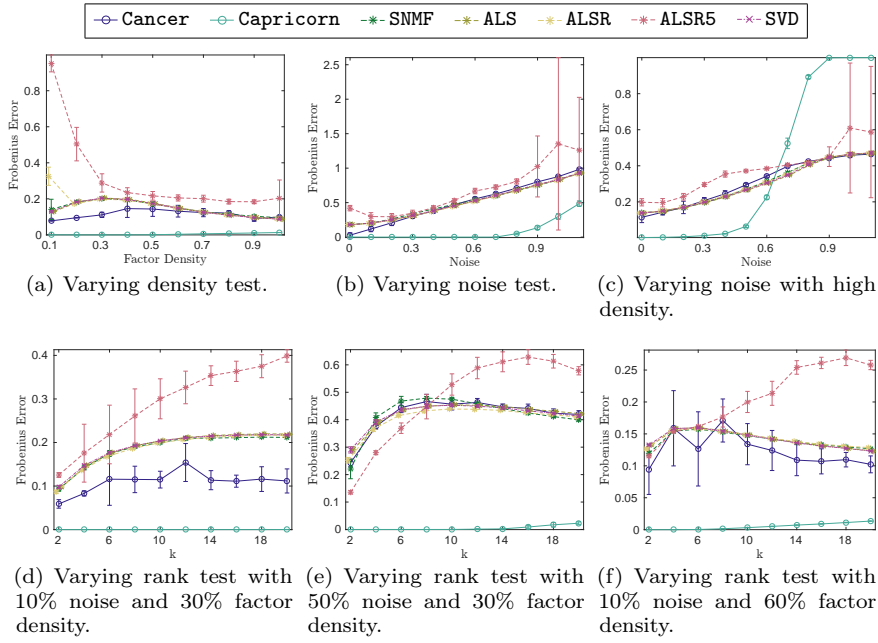


Fig. 6 Reconstruction errors on synthetic data with tropical noise. x -axis is the parameter varied and y -axis is the relative Frobenius norm. All results are averages over 10 random matrices and the width of the error bars is twice the standard deviation.

Varying Gaussian noise. Here we investigate how the algorithms respond to different levels of Gaussian noise, which was varied from 0 to 0.14 with increments of 0.01. A level of noise is a standard deviation of the Gaussian noise used to generate the noise matrix as described earlier. The factor density was kept at 50%. The results are given on Figure 8(a) (reconstruction error) and Figure 9(a) (sparsity of factors).

Here **Cancer** is generally the best method in reconstruction error, and second in sparsity only to **Capricorn**. The only time it loses to any method is when there is no noise, and **Capricorn** obtains a perfect decomposition. This is expected since **Capricorn** is by design better at spotting pure subtropical structure.

Varying density with Gaussian noise. In this experiment we studied what effects the density of factor matrices used in data generation has on the algorithms' performance. For this purpose we varied the density from 10% to 100% with increments of 10% while keeping the other parameters fixed. There are two versions of this experiment, one with low noise level of 0.01 (Figures 8(b) and 9(b)), and a more noisy case at 0.08 (Figures 8(c) and 9(c)).

Cancer provides the least reconstruction error in this experiment, being clearly the best until the density is 0.7, from which point on it is tied with **SVD** and the NMF-based methods (the only exception being the least-dense high-noise case, where **ALSR** obtains a slightly better reconstruction error). **Capricorn** is the worst by a wide margin, but this is not surprising, as the data does not follow its assumptions. On the other hand, **Capricorn** does produce generally the sparsest factorization, but these are of little use given its bad reconstruction error. **Cancer** produces the sparsest factors from the

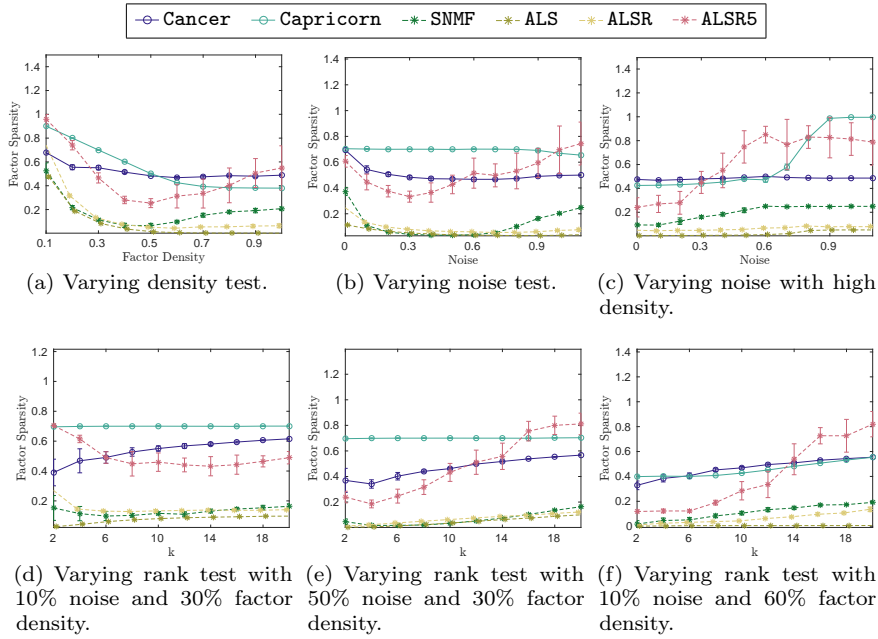


Fig. 7 Sparsity (fraction of zeroes) of the factor matrices for synthetic data with tropical noise. x -axis is the parameter varied and y -axis is the sparsity of the factors. The markers are averages of 10 random matrices and the width of the error bars is twice the standard deviation.

remaining methods, except in the first few cases where **ALSR5** is sparser (and worse in reconstruction error), meaning that **Cancer** produces factors that are both the most accurate and very sparse.

Varying rank with Gaussian noise. The purpose of this test is to study the performance of algorithms on data of different max-times ranks. We varied the true rank of the data from 2 to 20 with increments of 2. The factor density was fixed at 50% and Gaussian noise at 0.01. The results are shown on Figure 8(d) (reconstruction error) and Figure 9(d) (sparsity of factors). The results are similar to those considered above, with **Cancer** returning the most accurate and second sparsest factorizations.

Optimizing the Jensen–Shannon divergence. By default **Cancer** optimizes the Frobenius reconstruction error, but it can be replaced by an arbitrary additive cost function. We performed experiments with Jensen–Shannon divergence, which is given by the formula

$$J(\mathbf{A}, \mathbf{B}) = \sum_{ij} \mathbf{A}_{ij} \log \left(\frac{2\mathbf{A}_{ij}}{\mathbf{A}_{ij} + \mathbf{B}_{ij}} \right) + \mathbf{B}_{ij} \log \left(\frac{2\mathbf{B}_{ij}}{\mathbf{A}_{ij} + \mathbf{B}_{ij}} \right). \quad (22)$$

It is easy to see that (22) is an additive function, and hence can be plugged into **Cancer**. Figure 10 shows how this version of **Cancer** compares to other methods. The setup is the same as in the corresponding experiments on Figure 8. In all these experiments it is apparent that this version of **Cancer** is inferior to that optimizing the Frobenius error, but is generally on par with SVD and NMF-based methods. Also for the varying density

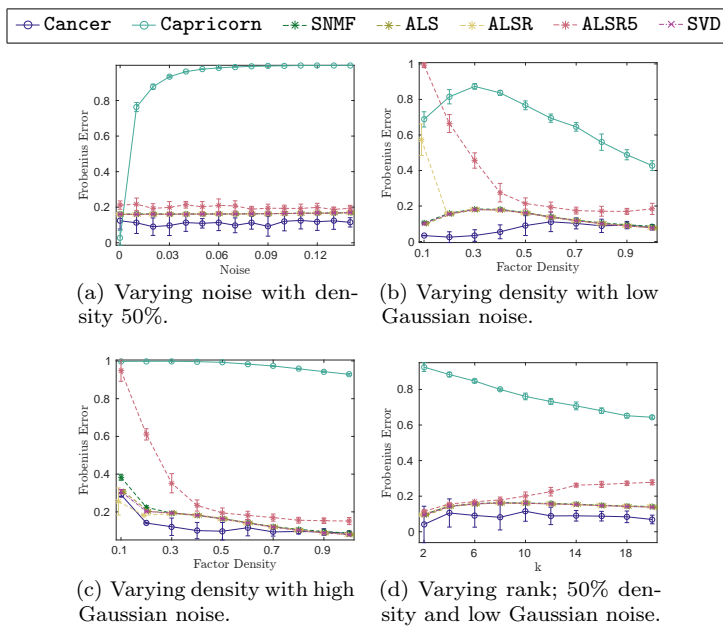


Fig. 8 Reconstruction error (Frobenius norm) for synthetic data with Gaussian noise. The markers are averages of 10 random matrices and the width of the error bars is twice the standard deviation.

test (Figure 10(b)) it produces better reconstruction errors than SVD and all the NMF methods, until the density reaches 50%, after which they become tied.

Prediction. In this experiment we choose a random holdout set and remove it from the data (elements of this set are marked as missing values). We then try to learn the structure of the data from its remaining part using the algorithms, and finally test how well they predict the values inside the holdout set. The factors are drawn uniformly at random from the set of integers in an interval $[0, a]$ with a predefined density of 30%, and then multiplied using the subtropical matrix product. We use two different values of a for each experiment, 10 and 3. With $a = 10$ input matrices have values in the range $[0, 100]$, and when $a = 3$, the range is $[0, 9]$. We then apply noise to the obtained matrices and feed them to the algorithms. Since all input matrices are integer-valued, and since the recovered data produced by the algorithms can be continuous-valued, we round it to the nearest integer. We report two measures of the prediction quality – prediction rate, which is defined as the fraction of correctly guessed values in the hold-out set, and root mean square error (RMSE). We tested this setup with both tropical noise (Figure 11) and Gaussian noise (Figure 12).

Capricorn gives by far the best prediction rate when using the higher $[0, 100]$ range of values in input matrices (Figures 11(a) and 12(a)). Especially interesting is that it also beats all other methods in the presence of Gaussian noise. In terms of RMSE it generally lands somewhere in the middle of the pack among various NMF methods. Such a large difference between these measures is caused by Capricorn not really being an approximation algorithm. It extracts subtropical patterns where they exist, while ignoring parts of the data where they cannot be found. This results in it either predicting

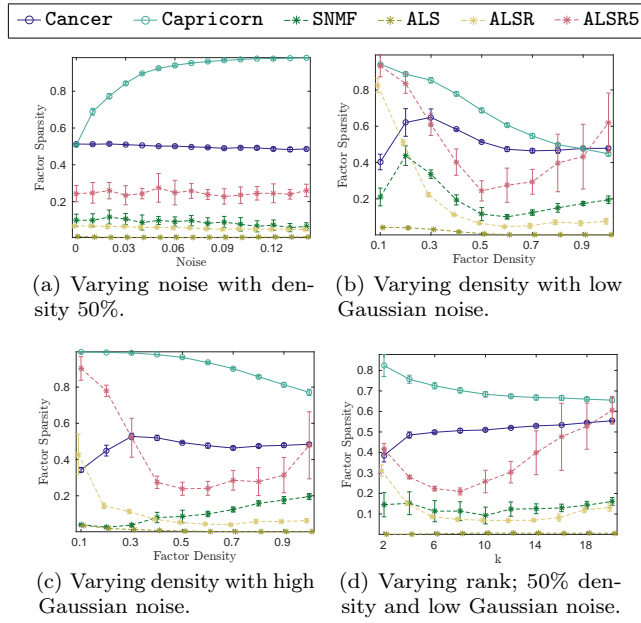


Fig. 9 Sparsity (fraction of zeroes) of the factor matrices for synthetic data with Gaussian noise. The markers are averages of 10 random matrices and the width of the error bars is twice the standard deviation.

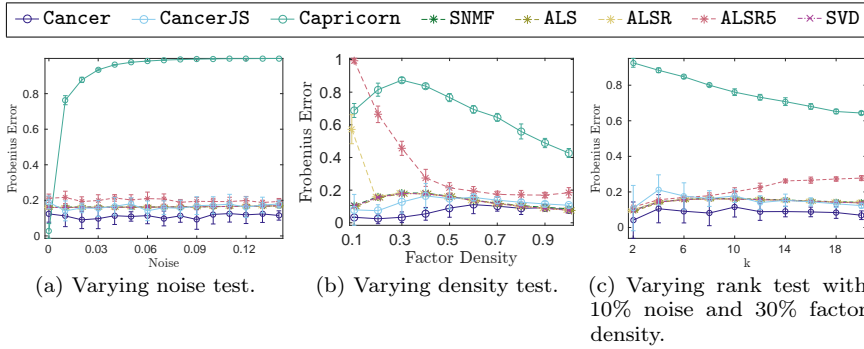


Fig. 10 Comparison of **Cancer** with Jensen–Shannon objective and other methods on synthetic data with Gaussian noise. x -axis is the parameter varied and y -axis is the relative Frobenius error. All results are averages over 10 random matrices and the width of the error bars is twice the standard deviation.

the integer values exactly or missing by a wide margin. With the $[0, 9]$ range of values the results of **Capricorn** become worse, which is especially evident with Gaussian noise. Although this behaviour might seem counterintuitive, it is simply a consequence of noise having a larger effect when values in the data are smaller. **Cancer** shows the opposite behaviour to **Capricorn** in that it benefits from smaller value range and Gaussian noise, where it consistently outperforms all other methods. Unlike **Capricorn**, **Cancer** approximates values in input data, which allows it to get a high number of hits with the

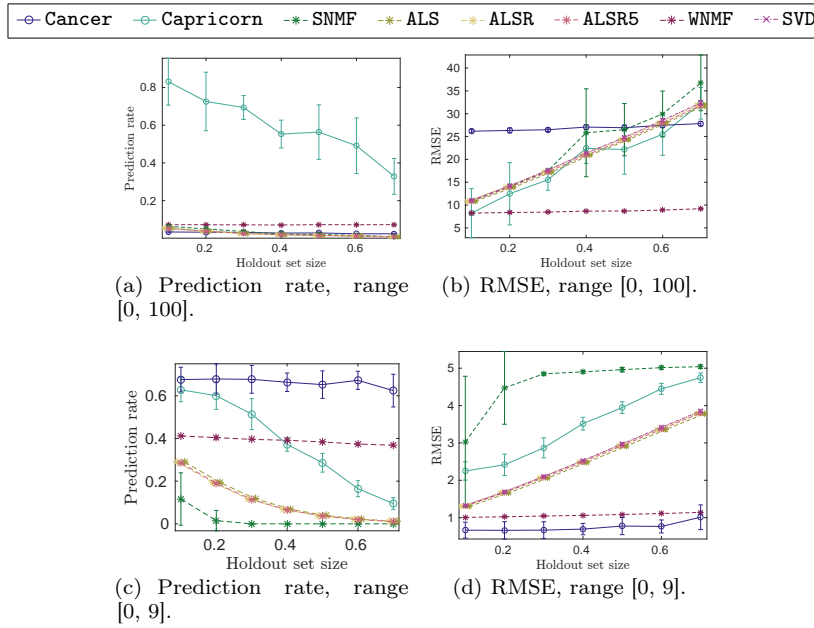


Fig. 11 Prediction rate on synthetic data with tropical noise. The x -axis represents the size of the holdout set. The y -axis is the correct prediction rate in Figures 11(a) and 11(c), and RMSE in Figures 11(b) and 11(d). The range is the interval that the values in input matrices are restricted to. All results are averages over 10 random matrices and the width of the error bars is twice the standard deviation.

$[0, 9]$ range after the rounding. On the $[0, 100]$ interval though, it is liable to guessing many values incorrectly since a much higher level of precision is required. For many prediction tasks, like predicting user ratings, **Cancer**'s approach seems more useful as input values are usually drawn from a relatively small range (for example, in **Movielens**, all ratings are from $[0, 5]$). Other competing methods generally do not perform well, with the exception of **SVD** winning the first place with RMSE measure for the high range experiments (Figures 11(b) and 12(b)). It illustrates once again that **SVD** is a good approximation method but does not help its prediction accuracy. In all other experiments the first place is held by either **Capricorn** or **Cancer**. As a general guideline, when choosing between **Capricorn** and **Cancer** for value prediction, one should consider that **Cancer** usually gives a superior performance, while **Capricorn** tends to be better for exact guessing of values having a wider range.

Discussion. The synthetic experiments confirm that both **Capricorn** and **Cancer** are able to recover matrices with max-times structure. The main practical difference between them is that **Capricorn** is designed to handle the tropical (flipping) noise, while **Cancer** is meant for the data that is perturbed with white (Gaussian) noise. While **Capricorn** is clearly the best method when the data has only the flipping noise – and is capable of tolerating very high noise levels – its results deteriorate when we apply Gaussian noise. Hence, when the exact type of noise is not known a priori, it is advisable to try both methods. It is also important to note that **Cancer** is actually a framework of algorithms as it can optimize various objectives. In order to demonstrate that, we performed

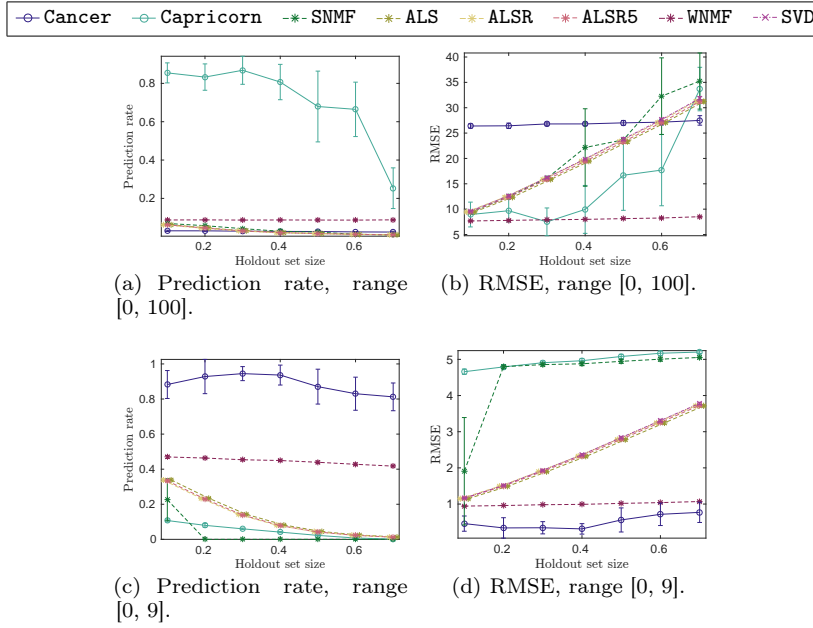


Fig. 12 Prediction rate on synthetic data with Gaussian noise. The x -axis represents the size of the holdout set. The y -axis is the correct prediction rate in Figures 12(a) and 12(c), and RMSE in Figures 12(b) and 12(d). The range is the interval that the values in input matrices are restricted to. All results are averages over 10 random matrices and the width of the error bars is twice the standard deviation.

experiments with Jensen–Shannon divergence as objective and obtained results that are, while inferior to **Cancer** that optimizes the Frobenius error, still slightly better than the rest of the algorithms. Overall we can conclude that **SVD** and the NMF-based methods generally cannot recover the structure from subtropical data, that is, we cannot use existing methods as a substitute to find the max-times structure neither for the reconstruction nor for the prediction tasks.

5.3 Real-world experiments

The main purpose of the real-world experiments is to study to which extent **Capricorn** and **Cancer** can find max-times structure from various real-world data sets. Having established with the synthetic experiments that both algorithms are capable of finding the structure when it is present, here we look at what kind of results they obtain in the real-world data.

It is probably unrealistic to expect real-world data sets to have “pure” max-times structure, as in the synthetic experiments. Rather, we expect **SVD** to be the best method (in reconstruction error’s sense), and our algorithms to obtain reconstruction error comparable to the NMF-based methods. We will also verify that the results from the real-world data sets are intuitive.

The datasets

Bas1LP represents a linear program.⁵ It is available from the University of Florida Sparse Matrix Collection⁶ (Davis and Hu 2011).

Trec12 is a brute force disjoint product matrix in tree algebra on n nodes.⁷ It can be obtained from the same repository as **Bas1LP**.

Worldclim contains weather records for various locations in Europe (full description can be found in Section 1).

NPAS is a nerdiness personality test that uses different attributes to determine the level of nerdiness of a person.⁸ It contains answers by 1418 respondents to a set of 36 questions that asked them to self-assess various statements about themselves on a scale of 1 to 7. We preprocessed the input matrix by dividing each column by its standard deviation and subtracting its mean. To make sure that the data is nonnegative, we subtracted the smallest value of the obtained normalized matrix from every its element.

Eigenfaces is a subset of the Extended Yale Face collection of face images (Georghiadis et al 2000). It consists of 32-by-32 pixel images under different lighting conditions. We used a preprocessed data by Xiaofei He et al.⁹ We selected a subset of pictures with lighting from the left and then preprocessed the input matrix by first subtracting from every column its smallest element and then dividing it by its standard deviation.

4News is a subset of the 20Newsgroups dataset,¹⁰ containing the usage of 800 words over 400 posts for 4 newsgroups.¹¹ Before running the algorithms we represented the dataset as a TF-IDF matrix, and then scaled it by dividing each entry by the greatest entry in the matrix.

HPI is a land registry house price index.¹² Rows represent months, columns are locations, and entries are residential property price indices. We preprocessed the data by first dividing each column by its standard deviation and then subtracting its minimum, so that each column has minimum 0.

Movielens is a collection of user ratings for a set of movies. The original dataset¹³ consists of 100 000 ratings from 1000 users on 1700 movies, with ratings ranging from 1 to 5. In order to be able to perform cross-validation on it, we had to preprocess **Movielens** by removing users that rated fewer than 10 movies and movies that were rated less than 5 times. After that we were left with 943 users, 1349 movies and 99 287 ratings.

The basic properties of these data sets are listed in Table 1.

⁵Submitted to the matrix repository by Csaba Meszaros.

⁶<http://www.cise.ufl.edu/research/sparse/matrices/>, accessed 18 July 2017

⁷Submitted by Nicolas Thiery.

⁸The dataset can be obtained on the online personality website http://personality-testing.info/_rawdata/NPAS-data.zip, accessed 18 July 2017.

⁹<http://www.cad.zju.edu.cn/home/dengcai/Data/FaceData.html>, accessed 18 July 2017

¹⁰<http://qwone.com/~jason/20Newsgroups/>, accessed 18 July 2017

¹¹The authors are grateful to Ata Kabán for pre-processing the data, see Miettinen (2009).

¹²Available at <https://data.gov.uk/dataset/land-registry-house-price-index-backgroud-tables/>, accessed 18 July 2017

¹³Available at <http://grouplens.org/datasets/movielens/100k/>, accessed 18 July 2017

Table 1 Real world datasets properties.

Dataset	Rows	Columns	Density
Bas1LP	9825	5411	1.1%
Trec12	2726	551	10.0%
Worldclim	2575	48	99.9%
NPAS	1418	36	99.6%
Eigenfaces	1024	222	97.0%
4News	400	800	3.5%
HPI	253	177	99.5%
Movielens	943	1349	7.8%

Table 2 Reconstruction error for various real-world datasets.

$k =$	Worldclim	NPAS	Eigenfaces	4News	HPI	Movielens	Trec12	Bas1LP
	10	10	40	20	15	10	25	25
Cancer	0.071	0.240	0.204	0.556	0.027	0.756	0.864	0.813
Capricorn	0.392	0.395	0.972	0.987	0.217	1.003	0.998	0.912
SNMF	0.046	0.225	0.178	0.546	0.023	0.745	0.841	0.749
ALS	0.087	0.227	0.313	0.538	0.074	0.749	0.828	0.733
ALSR	0.122	0.226	0.294	1.000	0.045	0.748	0.827	0.733
ALSR5	0.081	0.233	0.291	1.000	0.063	0.748	0.826	0.733
WNMF	0.034	0.221	0.169	0.545	0.021	0.741	0.824	0.733
SVD	0.025	0.209	0.140	0.533	0.015	0.728	0.802	0.722

Quantitative results: reconstruction error, sparsity, convergence, and runtime

The following experiments are meant to test **Cancer** and **Capricorn**, and how they compare to other methods, such as **SVD** and **NMF**. Table 2 provides the relative Frobenius reconstruction errors for various real-world data sets, as well as ranks used for factorization.¹⁴ Since there is no ground truth for these datasets, the ranks are chosen based mainly on the size of the data and our intuition on what the true rank should be. **SVD** is, as expected, consistently the best method, followed by **WNMF** and **SNMF**. **Cancer** generally lands in the middle of the pack of the **NMF** methods, which suggests that it is capable of finding max-times structure that is comparable to what **NMF**-based methods provide. Consequently, we can study the max-times structure found by **Cancer**, knowing that it is (relatively) accurate. On the other hand, **Capricorn** has a high reconstruction error. The discrepancy between **Cancer**’s and **Capricorn**’s results indicates that the datasets used cannot be represented using “pure” subtropical structure. Rather, they are either a mix of **NMF** and subtropical patterns or have relatively high levels of continuous noise.

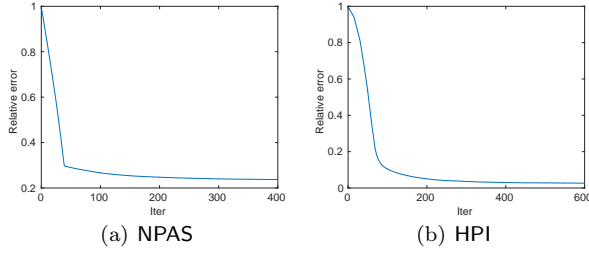
The sparsity of the factors for real-world data is presented in Table 3 (we do not include the sparsities for **SVD** and **WNMF** as they were all 0). Here, **Cancer** often returns the second-sparsest factors (behind only **Capricorn**), but with **4News** and **HPI**, **ALSR** and **ALSR5** obtains sparser decompositions.

We also studied the convergence behavior of **Cancer** using some of the real-world data sets. The results can be seen in Figure 13, where we plot the relative error with respect to the iterations over the main for-loop in **Cancer**. As we can see, in both cases

¹⁴The values are different than those presented by Karaev and Miettinen (2016b) because we used Frobenius error instead of L_1 and counted all elements towards the error, not just nonnegative ones.

Table 3 Factor sparsity for various real-world datasets.

	Worldclim	NPAS	Eigenfaces	4News	HPI	Movielens	Trec12	Bas1LP
$k =$	10	10	40	20	15	10	25	25
Cancer	0.645	0.528	0.571	0.812	0.422	0.666	0.838	0.951
Capricorn	0.795	0.733	0.949	0.991	0.685	0.957	0.988	0.978
SNMF	0.383	0.330	0.403	0.499	0.226	0.543	0.758	0.738
ALS	0.226	0.120	0.434	0.513	0.331	0.420	0.573	0.634
ALSR	0.275	0.117	0.480	1.000	0.729	0.438	0.681	0.748
ALSR5	0.549	0.189	0.648	1.000	0.622	0.481	0.743	0.811

**Fig. 13** Convergence rate of **Cancer** for two real-world datasets. Each iteration is a single run of **UpdateBlock**, that is if a factorization has rank k , then one full cycle would correspond to k iterations.**Table 4** The average runtime in seconds and standard deviation of the algorithms for various real-world datasets. The results were calculated based on 5 restarts of each method.

	Worldclim	NPAS	4News	HPI
Cancer	20116.000 \pm 15.14	6023.000 \pm 25.00	25520.000 \pm 60.44	924.000 \pm 8.00
Capricorn	205.870 \pm 1.39	87.000 \pm 1.30	165.960 \pm 7.12	41.000 \pm 0.72
SNMF	115.100 \pm 0.53	72.000 \pm 1.50	195.570 \pm 1.76	64.000 \pm 0.51
ALS	0.194 \pm 0.08	0.374 \pm 0.14	3.649 \pm 1.45	0.156 \pm 0.12
ALSR	0.187 \pm 0.02	0.280 \pm 0.04	4.684 \pm 0.74	0.309 \pm 0.13
WNMF	2.240 \pm 0.20	1.201 \pm 0.11	5.164 \pm 0.87	1.288 \pm 0.10
SVD	0.598 \pm 0.04	0.155 \pm 0.04	0.142 \pm 0.04	0.027 \pm 0.01

Cancer has obtained a good reconstruction error already after few full cycles, with the remaining runs only providing minor improvements. We can deduce that **Cancer** quickly reaches an acceptable solution.

To give some idea about the speed performance of the algorithms, we ran each of the competing methods on some of the real-world datasets. The runtime of each algorithm (in seconds) is shown in Table 4, where we report its mean value and the standard deviation averaged over 5 runs. All tests were performed on a Linux machine with Intel Xeon E5530 CPU with 16 2.40 GHz cores, although **Cancer** and **Capricorn** were only utilizing one core. As we can see, the simplest methods, such as **SVD** and **ALS**, are also the fastest, while more involved algorithms, such as **Cancer** or **SNMF**, take much longer to run. It is worth noting that **Cancer** and **Capricorn** are written in Matlab, and their performance can be potentially significantly improved by implementing time critical parts in C or another low-level programming language.

Prediction

Here we investigate how well both **Capricorn** and **Cancer** can predict missing values in the data. We used three real-world datasets, a user-movie rating matrix **Movielens**, a brute force disjoint product matrix in tree algebra **Trec12** and **Bas1LP**, that represents a linear program. All these matrices are integer valued, and hence we will also round the results of all methods to the nearest integer. We compare the results of our methods against **WMMF** and **SVD**. The choice of **WMMF** is motivated by its ability to ignore missing elements in the input data and its generally good performance on the previous tests. There is only one caveat: **WMMF** sometimes produces very high spikes for some elements in the matrix. They do not cause too much problem with prediction, but they seriously deteriorate the results of **WMMF** with respect to various distance measures. For this reason we always ignore such elements. While this comparison method is obviously not completely fair towards other methods, it can serve as a rough upper bound for what performance is possible with NMF-based algorithms. Comparing against other methods is obviously not fair as they are not designed to deal with missing values, but we will still present the results of **SVD** for completeness.

On **Movielens** we perform standard cross-validation tests, where a random selection of elements is chosen as a holdout set and removed from the data. The data has 943 users, each having rated from 19 to 648 movies. A holdout set is chosen by sampling uniformly at random 5 ratings from each user. We run the algorithms, while treating the elements from the holdout set as missing values, and then compare the reconstructed matrices to the original data. This procedure is repeated 10 times.

To get a more complete view on how good the predictions are, we report various measures of quality: Frobenius error, root mean square error (RMSE), reciprocal rank, Spearman’s ρ , mean absolute error (MAE), Jensen–Shannon divergence (JS), optimistic reciprocal rank, Kendall’s τ , and prediction accuracy. The prediction accuracy allows us to see if the methods are capable of recovering the missing user ratings. The remaining tests can be divided into two categories. The first one, which comprises Frobenius error, root mean square error, mean absolute error, and Jensen–Shannon divergence, aims to quantify the distance between the original data and the reconstructed matrix. The second group of tests finds the correlation between rankings of movies for each user. It includes Spearman’s ρ , Kendall’s τ , reciprocal rank, and optimistic reciprocal rank. All these measures are well known, with perhaps only the reciprocal rank requiring some explanation. Let us first denote by U the set of all users. In the following, for each user $u \in U$ we only consider the set of movies $M(u)$ that this user has rated that belong to the holdout set. The ratings by user u induce a natural ranking on $M(u)$. On the other hand, the algorithms produce approximations $r'(u, m)$ to the true ratings $r(u, m)$, which also induce a corresponding ranking of the movies. The reciprocal rank is a convenient way of comparing the rankings obtained by the algorithms to the original one. For any user $u \in U$, denote by $H(u)$ a set of movies that this user ranked the highest (that is $H(u) = \{m \in M(u) : r(u, m) = \max_{m' \in M(u)} r(u, m')\}$). The reciprocal rank for user u is now defined as

$$RR(u) = \frac{1}{\min_{m \in H} R(u, m)}, \quad (23)$$

where $R(u, m)$ is the rank of the movie m within $M(u)$ according to the rating approximations given by the algorithm in question. Now the mean reciprocal rank is defined as

the average of the reciprocal ranks for each individual user $MRR = \frac{1}{|U|} \sum_{u \in U} RR(u)$. When computing the ranks $R(u, m)$, all tied elements receive the same rank, which is computed by averaging. That means that if, say, movies m_1 and m_2 have tied ranks of 2 and 3, then they both receive the rank of 2.5. An alternative way is to always assign the smallest possible rank. In the above example both m_1 and m_2 will receive rank 2. When ranks $R(u, m)$ are computed like this, the equation (23) defines the optimistic reciprocal rank.

For each test, Table 5 shows the mean and the standard deviation of the results of each algorithm. In addition we report the p -value based on the Wilcoxon signed-rank test. It shows if an advantage of one method over another is statistically significant. We say that a method A is significantly better than method B if the p -value is < 0.05 . It is unreasonable to report the p -value for every method pair – instead we only show p -values involving the best method. For each method, the value given next to it is the p -value for this method and the best method.

Cancer is significantly better for the Frobenius error, root mean square error, mean absolute error, Jensen–Shannon divergence, and accuracy. For the remaining tests the results are less clear, with **Cancer** winning on the reciprocal rank, **Capricorn** taking the optimistic reciprocal rank, and **WNMF** being better on Spearman’s ρ and Kendall’s τ tests. It should be noted though, that the victories of **WNMF** on Spearman’s ρ and Kendall’s τ tests, as well as **Cancer**’s on the reciprocal rank, are not statistically significant as the p -values are quite high. In summary, our experiments show that **Cancer** is significantly better in tests that measure the direct distance between the original and the reconstructed matrices, as well as the prediction accuracy, whereas for the ranking experiments it is difficult to give any of the algorithms an edge.

For **Trec12** and **Bas1LP** we also perform cross-validation, where on each fold we take 10% of the nonzero values in the data as a holdout set and then try to predict them. In total there are 5 folds. Unlike **Movielens**, **Trec12** and **Bas1LP** datasets are not so readily interpretable as they were generated from abstract mathematical data structures. Ranking, in particular, makes no sense, and hence we do not perform any ranking associated experiments. The results for **Trec12** are shown in Table 6 and for **Bas1LP** in Table 7. It is apparent that on **Trec12** **WNMF** performs significantly better than any other method, being better in all metrics. As discussed earlier, however, that should be taken with a grain of salt as we ignore the elements where **WNMF** produced unreasonably large values. Without this preprocessing its results are much worse than those of **Cancer**. This presents evidence, although not conclusive, that the **Trec12** dataset has less subtropical structure than **Movielens**. The p -value is 0.004 in for all metrics, which is the result of a particular number of folds (5) that we used. The fact that we have this number everywhere in the table simply indicates that **WNMF** was better than any other method on every fold with respect to all measures. With **Bas1LP** the roles reverse, and this time **Cancer** is clearly the best method, winning according to all metrics and on all folds, just as **WNMF** did on **Trec12**.

Interpretability of the results

The crux of using max-times factorizations instead of standard (nonnegative) ones is that the factors (are supposed to) exhibit the “winner-takes-it-all” structure instead of the “parts-of-whole” structure. To demonstrate this, we analysed results in four different datasets: **Eigenfaces**, **NPAS**, **Worldclim**, and **Mammals**. The **Mammals** dataset is explained below.

Table 5 Comparison between the predictive power of different methods on the **Movielens** data. The arrow after the value indicates whether higher or lower values are preferable. The p -values are computed using the Wilcoxon signed-rank test.

	Frobenius		RMSE	
	value(↓)	p -value	value(↓)	p -value
Cancer	0.2876 ± 0.003		1.0802 ± 0.011	
Capricorn	0.6993 ± 0.024	0.0001	2.6267 ± 0.085	0.0001
WNMF	0.2989 ± 0.003	0.0001	1.1227 ± 0.012	0.0001
SVD	0.7336 ± 0.002	0.0001	2.7558 ± 0.014	0.0001
	Recip. rank		Spearman's ρ	
	value(↑)	p -value	value(↑)	p -value
Cancer	0.7451 ± 0.010		0.3071 ± 0.015	0.5749
Capricorn	0.5601 ± 0.017	0.0001	0.2354 ± 0.017	0.0001
WNMF	0.7395 ± 0.004	0.0521	0.3084 ± 0.012	
SVD	0.7217 ± 0.008	0.0004	0.2445 ± 0.013	0.0001
	MAE		JS	
	value(↓)	p -value	value(↓)	p -value
Cancer	0.8203 ± 0.008		0.0201 ± 0.000	
Capricorn	2.0518 ± 0.106	0.0001	0.2826 ± 0.026	0.0001
WNMF	0.8555 ± 0.008	0.0001	0.0209 ± 0.000	0.0057
SVD	2.4756 ± 0.014	0.0001	0.1153 ± 0.001	0.0001
	Recip. rank opt.		Kendall's τ	
	value(↑)	p -value	value(↑)	p -value
Cancer	0.7451 ± 0.010	0.0001	0.2659 ± 0.013	0.4251
Capricorn	0.8547 ± 0.010		0.2127 ± 0.016	0.0001
WNMF	0.7395 ± 0.004	0.0001	0.2679 ± 0.010	
SVD	0.7217 ± 0.008	0.0001	0.2111 ± 0.012	0.0001
	Accuracy			
	value(↑)	p -value		
Cancer	0.3968 ± 0.008			
Capricorn	0.2053 ± 0.019	0.0001		
WNMF	0.3828 ± 0.006	0.0011		
SVD	0.0588 ± 0.003	0.0001		

We plotted the left factor matrices for the **Eigenfaces** data for **Cancer** and **ALS** in Figure 14. At first, it might look like **ALS** provides more interpretable results, as most factors are easily identifiable as faces. This, however, is not a very interesting result: we already knew that the data has faces, and many factors in the **ALS**'s result are simply some kind of "prototypical" faces. The results of **Cancer** are harder to identify on the first sight. Upon closer inspection, though, one can see that they identify areas that are lighter in the different images, that is, have higher grayscale values. These factors tell us the variances in the lighting in the different photos, and can reveal information we did not know a priori. In addition almost every one of **Cancer**'s factors contains one or two main feature of the face (such as nose, left eye, right cheek, etc.). In other words, while **NMF**'s patterns are for the most part close to fully formed faces, **Cancer** finds independent fragments that indicate the direction of the lighting and (or) contain some of the main features of a face.

Table 6 Comparison between the predictive power of different methods on the **Trec12** data. The arrow after the value indicates whether higher or lower values are preferable. The p -values are computed using the Wilcoxon signed-rank test.

	Frobenius		RMSE	
	value(↓)	p -value	value(↓)	p -value
Cancer	0.4824 ± 0.016	0.0040	2.3124 ± 0.085	0.0040
Capricorn	0.7827 ± 0.023	0.0040	3.7521 ± 0.131	0.0040
WNMF	0.4374 ± 0.006		2.0925 ± 0.041	
SVD	0.6005 ± 0.003	0.0040	2.8784 ± 0.032	0.0040
	MAE		JS	
	value(↓)	p -value	value(↓)	p -value
Cancer	1.5852 ± 0.050	0.0040	0.0675 ± 0.005	0.0040
Capricorn	2.3871 ± 0.099	0.0040	0.2929 ± 0.028	0.0040
WNMF	1.2138 ± 0.011		0.0367 ± 0.000	
SVD	1.8413 ± 0.013	0.0040	0.0786 ± 0.001	0.0040
	Accuracy			
	value(↑)	p -value		
Cancer	0.2315 ± 0.010	0.0040		
Capricorn	0.1918 ± 0.019	0.0040		
WNMF	0.3996 ± 0.004			
SVD	0.2061 ± 0.002	0.0040		

Table 7 Comparison between the predictive power of different methods on the **Bas1LP** data. The arrow after the value indicates whether higher or lower values are preferable. The p -values are computed using the Wilcoxon signed-rank test.

	Frobenius		RMSE	
	value(↓)	p -value	value(↓)	p -value
Cancer	0.3690 ± 0.018		1.1462 ± 0.065	
Capricorn	0.5741 ± 0.054	0.0040	1.7822 ± 0.161	0.0040
WNMF	0.4113 ± 0.014	0.0040	1.2748 ± 0.038	0.0040
SVD	0.5003 ± 0.002	0.0040	1.5534 ± 0.019	0.0040
	MAE		JS	
	value(↓)	p -value	value(↓)	p -value
Cancer	0.3286 ± 0.014		0.0228 ± 0.001	
Capricorn	0.6712 ± 0.094	0.0040	0.1208 ± 0.037	0.0040
WNMF	0.3932 ± 0.006	0.0040	0.0268 ± 0.000	0.0040
SVD	0.9391 ± 0.006	0.0040	0.0919 ± 0.000	0.0040
	Accuracy			
	value(↑)	p -value		
Cancer	0.8841 ± 0.002			
Capricorn	0.7111 ± 0.050	0.0040		
WNMF	0.8562 ± 0.001	0.0040		
SVD	0.2837 ± 0.002	0.0040		



(a) Cancer



(b) ALS

Fig. 14 Cancer finds the dominant patterns from the Eigenfaces data. Pictured are the left factor matrices for the Eigenfaces data.

Further, as seen in Table 2, **Cancer** obtains a better reconstruction error than **ALS** with this data, confirming that these factors are indeed useful to recreate the data.

In order to interpret NPAS we first observe that each column represents a single personality attribute. Denote by \mathbf{A} the obtained approximation of the original matrix. For each rank-1 factor \mathbf{X} and each column \mathbf{A}_i we define the score $\sigma(i)$ as the number of elements in \mathbf{A}_i that are determined by \mathbf{X} . By sorting attributes in descending order of $\sigma(i)$ we obtain relative rankings of the attributes for a given factor. The results are shown in Table 8. The first factor clearly shows introverted tendencies, while the second one can be summarized as having interests in fiction and games.

Figure 15 shows all of the factors for the *Worldclim* data, as obtained by **Cancer** and **WNMF** (the best NMF-method in Table 2). Figures 15(a) and 15(b) show left-hand sides of factors found by **Cancer** and **WNMF**, respectively, plotted on the map. Darker colours indicate higher values, that can be interpreted as “more important”. The right-hand side factors are presented in Figures 15(c) and 15(d), respectively. Here, each row corresponds to a factor, and each column to a single observation column from the original data (that is columns 1–12 represent average low temperatures for each month,

Table 8 Top three attributes for the first two factors of NPAS.

Factor 1	Factor 2
I am more comfortable with my hobbies than I am with other people	I have played a lot of video games
I gravitate towards introspection	I collect books
I sometimes prefer fictional people to real ones	I care about super heroes

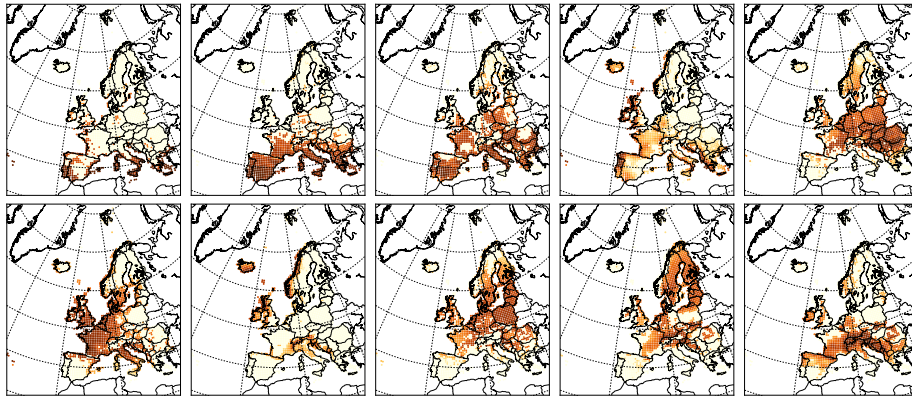
columns 13–24 average high temperatures, columns 25–36 daily means, and columns 37–48 average monthly precipitation). Again, higher values can be seen as having more importance. Recall that a pattern is formed by taking an outer product of a single left-hand factor and the corresponding right-hand factor. It is easy to see that largest (and thus the most important) values in a pattern are those that are products of high values in both right-hand side and left-hand side factors.

The **WNMF** factors have less high values (dark colours – all factors are normalized to the unit interval). For **Cancer**, there are more large values in each factor. This highlights the difference between the subtropical and the normal algebra: in normal algebra, if you sum two large values, the result is even larger, whereas in subtropical algebra, the result is no larger than the largest of the summands. In decompositions, this means that **WNMF** cannot have overlapping high values in its factors; instead it has to split its factors to mostly non-overlapping parts. **Cancer**, on the other hand, can have overlap, and hence its factors can share some phenomena. For instance, the seventh factor of **Cancer** clearly indicates areas of high precipitation (cf. Figure 1). The same phenomenon is split into many factors by **WNMF** (at least third, sixth, and seventh factor), mostly explaining areas with higher precipitation at different parts of the year. While many elements in the right-hand side factors of **Cancer** are nonzero, that does not mean that all of them are of equal importance. Because some of them are dominated by larger features, they do not influence the final outcome. Generally, since larger values are more likely to make a contribution than smaller ones, they should be considered more important when interpreting the data.

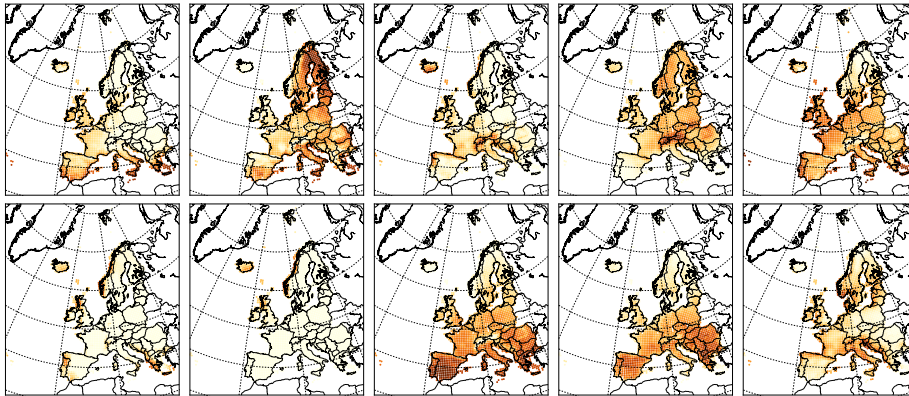
The possibility of the factors to overlap is not always desired, but in some applications it can be seen to be almost necessary. Consider, for example, mammal species’ co-location data. This dataset, called **Mammals**, is a matrix whose rows and columns correspond to locations in Europe, and for every column-row pair, the corresponding entry represents the degree to which the sets of mammals inhabiting them overlap. This dataset¹⁵ was obtained from the original binary location-species matrix (see Mitchell-Jones et al 1999) by multiplying it with its transpose and then normalizing by dividing each column by its maximal element. The obtained matrix has 2670 rows and columns and density 91%. Due to its special nature, we use it only in this experiment to provide intuition about the subtropical factorizations.

The factors obtained by **Cancer** with the **Mammals** data are depicted in Figure 16, where we can see that many of these factors cover the central parts of the European Plain, extending a bit south to cover most of Germany. There are, naturally, many mammal species that inhabit the whole European Plain, and the east–west change is gradual. This gradual change is easier to model in subtropical algebra, as we do not have to worry about the sums of the factors getting too large. Factors 1–6 model various

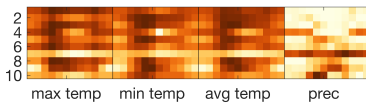
¹⁵Available for research purposes from the Societas Europaea Mammalogica at <http://www.european-mammals.org>



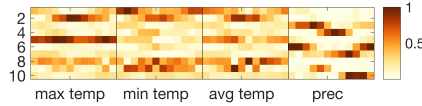
(a) Cancer left-hand factors.



(b) WNMF left-hand factors.



(c) Cancer right-hand factors.



(d) WNMF right-hand factors.

Fig. 15 Cancer factors in the Worldclim data. The factor vectors are normalized to take values from the unit interval and darker shades indicate higher values.

aspects of the east–west change, emphasizing either the south–west, central, or eastern parts of the plain. Similarly, the ninth factor explains mammal species found in the UK and southern Scandinavia, while the tenth factor covers species found in Scotland, Scandinavia, and Baltic countries, indicating that these areas have roughly the same biome. If we compare these results to those of WNMF (Figure 17), then it becomes evident that the latter tries to find relatively disjoint factors and avoids the factor overlap whenever possible. This is because in NMF any feature that is nonzero at a given data point is always “active” in a sense that it contributes to the final value. That being said, WNMF does find some interesting patterns, such as rather distinct factors representing France and Scandinavian peninsula.

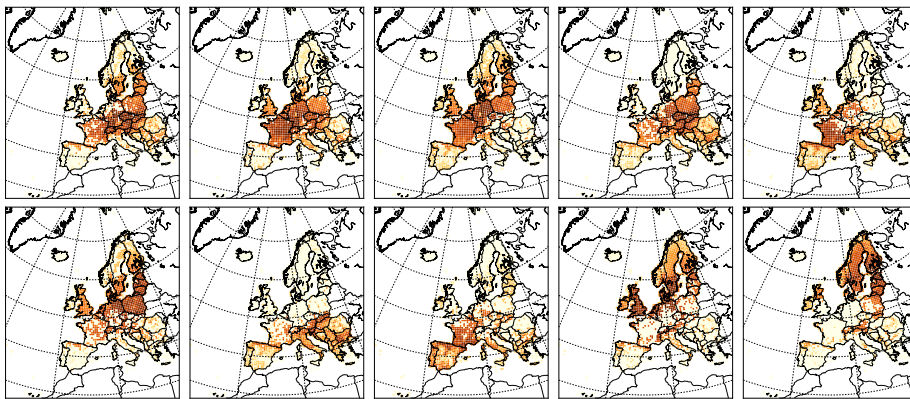


Fig. 16 Values in the factors by Cancer in the Mammals data plotted on a map. Every factor is normalized to take values from the unit interval and darker shades indicate higher values.

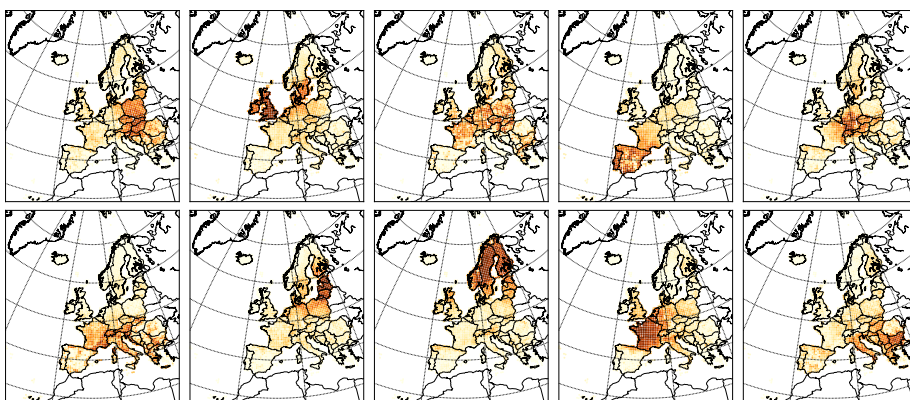


Fig. 17 Values in the factors by WNMF in the Mammals data plotted on a map. Every factor is normalized to take values from the unit interval and darker shades indicate higher values.

6 Related Work

Here we present earlier research that is related to the subtropical matrix factorization. We start by discussing classic methods, such as SVD and NMF, that have long been used for various data analysis tasks, and then continue with approaches that use idempotent structures. Since the tropical algebra is very closely related to the subtropical algebra, and since there has been a lot of research on it, we dedicate the last subsection to discuss it in more detail.

6.1 Matrix factorization in data analysis

Matrix factorization methods play a crucial role in data analysis as they help to find low-dimensional representations of the data and uncover the underlying latent structure. A classic example of a real-valued matrix factorization is the singular value decomposition (SVD) (see e.g. Golub and Van Loan 2012), which is very well known

and finds extensive applications in various disciplines, such as signal processing and natural language processing. The SVD of a real n -by- m matrix \mathbf{A} is a factorization of the form $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$, where $\mathbf{U} \in \mathbb{R}^{n \times n}$ and $\mathbf{V} \in \mathbb{R}^{m \times m}$ are orthogonal matrices, and $\mathbf{\Sigma} \in \mathbb{R}^{n \times m}$ is a rectangular diagonal matrix with nonnegative entries. An important property of SVD is that it provides the best low-rank approximation of a given matrix with respect to the Frobenius norm (Golub and Van Loan 2012), giving rise to the so called truncated SVD. This property is frequently used to separate important parts of data from the noise. For example, it was used by Jha and Yadava (2011) to remove the noise from sensor data in electronic nose systems. Another prominent usage of the truncated SVD is in dimensionality reduction (see for example Sarwar et al 2000; Deerwester et al 1990).

Despite SVD being so ubiquitous, there are some restrictions to its usage in data mining due to possible presence of negative elements in the factors. In many applications negative values are hard to interpret, and thus other methods have to be used. Nonnegative matrix factorization (NMF) is a way to tackle this problem. For a given nonnegative real matrix \mathbf{A} , the NMF problem is to find a decomposition of \mathbf{A} into two matrices $\mathbf{A} \approx \mathbf{B}\mathbf{C}$ such that \mathbf{B} and \mathbf{C} are also nonnegative. Its applications are extensive and include text mining (Pauca et al 2004), document clustering (Xu et al 2003), pattern discovery (Brunet et al 2004), and many other. This area drew considerable attention after a publication by Lee and Seung (1999), where they provided an efficient algorithm for solving the NMF problem. It is worth mentioning that even though the paper by Lee and Seung is perhaps the most famous in NMF literature, it was not the first one to consider this problem. Earlier works include Paatero and Tapper (1994) (see also Paatero 1997), Paatero (1999), and Cohen and Rothblum (1993). Berry et al (2007) provide an overview of NMF algorithms and their applications. There exist various flavours of NMF that impose different constraints on the factors; for example Hoyer (2004) used sparsity constraints. Though both NMF and SVD perform approximations of a fixed rank, there are also other ways to enforce compact representation of data. For example, in maximum-margin matrix factorization constraints are imposed on the norms of factors. This approach was exploited by Srebro et al (2004), who showed it to be a good method for predicting unobserved values in a matrix. The authors also indicate that posing constraints on the factor norms, rather than on the rank, yields a convex optimization problem, which is easier to solve.

6.2 Idempotent semirings

The concept of the subtropical algebra is relatively new, and as far as we know, its applications in data mining are not yet well studied. Indeed, its only usage for data analysis that we are aware of was by Weston et al (2013), where it was used as a part of a model for collaborative filtering. The authors modeled users as a set of vectors, where each vector represents a single aspect about the user (e.g. a particular area of interest). The ratings are then reconstructed by selecting the highest scoring prediction using the max operator. Since their model uses max as well as the standard plus operation, it stands on the border between the standard and the subtropical worlds.

Boolean algebra, despite being limited to the binary set $\{0, 1\}$, is related to the subtropical algebra by virtue of having the same operations, and is thus a restriction of the latter to $\{0, 1\}$. By the same token, when both factor matrices are binary, their subtropical product coincides with the Boolean product, and hence the Boolean matrix

factorization can be seen as a degenerate case of the subtropical matrix factorization problem. The dioid properties of the Boolean algebra can be checked trivially. The motivation for the Boolean matrix factorization comes from the fact that in many applications data is naturally represented as a binary matrix (e.g. transaction databases), which makes it reasonable to seek decompositions that preserve the binary character of the data. The conceptual and algorithmic analysis of the problem was done by Miettinen (2009), with the focus mainly on the data mining perspective of the problem. For a linear algebra perspective see Kim (1982), where the emphasis is put on the existence of exact decompositions. A number of algorithms have been proposed for solving the BMF problem (Miettinen et al 2008; Lu et al 2008; Lucchese et al 2014; Karaev et al 2015).

6.3 Tropical algebra

Another close cousin of the max-times algebra is the max-plus, or so called tropical algebra, which uses plus in place of multiplication. It is also a dioid due to the idempotent nature of the max operation. As was mentioned earlier, the two algebras are isomorphic, and hence many of the properties are identical (see Sections 2 and 3 for more details).

Despite the theory of the tropical algebra being relatively young, it has been thoroughly studied in recent years. The reason for this is that it finds extensive applications in various areas of mathematics and other disciplines. An example of such a field is the discrete event systems (DES) (Cassandras and Lafortune 2008), where the tropical algebra is ubiquitously used for modeling (see e.g. Baccelli et al 1992; Cohen et al 1999). Other mathematical disciplines where the tropical algebra plays a crucial role are optimal control (Gaubert 1997), asymptotic analysis (Dembo and Zeitouni 2010; Maslov 1992; Akian 1999), and decidability (Simon 1978, 1994).

Research on tropical matrix factorization is of interest to us because of the above mentioned isomorphism between the two algebras. However, as was explained in Section 3, the approximate matrix factorizations are not directly transferable as the errors can differ dramatically. It should be mentioned that in the general case the problem of the tropical matrix factorization is NP-hard (see e.g. Shitov 2014). De Schutter and De Moor (2002) demonstrated that if the max-plus algebra is extended in such a way that there is an additive inverse for each element, then it is possible to solve many of the standard matrix decomposition problems. Among other results the authors obtained max-plus analogues of QR and SVD. They also claimed that the techniques they propose can be readily extended to other types of classic factorizations (e.g. Hessenberg and LU decomposition).

The problem of solving tropical linear systems of equations arises naturally in numerous applications, and is also closely related to matrix factorization. In order to illustrate this connection, assume that we are given a tropical matrix $\mathbf{A} \in \overline{\mathbb{R}}^{n \times m}$ and one of the factors $\mathbf{B} \in \overline{\mathbb{R}}^{n \times k}$. Then the other factor $\mathbf{C} \in \overline{\mathbb{R}}^{k \times m}$ can be found by solving the following set of problems

$$\mathbf{C}_j = \arg \min_{\mathbf{c} \in \overline{\mathbb{R}}^k} \|\mathbf{B} \diamond \mathbf{c} - \mathbf{A}_j\|_F, \quad j = 1, \dots, m. \quad (24)$$

Each problem in (24) requires “approximately” solving a system of tropical linear equations. The minus operation in (24) does not belong to the tropical semiring, so the

approximation here should be understood in terms of minimizing the classical distance. The general form of tropical linear equations

$$\mathbf{Ax} \oplus \mathbf{b} = \mathbf{Cx} \oplus \mathbf{d} \quad (25)$$

is not always solvable (see e.g. Gaubert 1997); however various techniques exist for checking the existence of the solution for particular cases of (25).

For equations of the form $\mathbf{Ax} = \mathbf{b}$ the feasibility can be established for example through the so called *matrix residuation*. There is a general result that for an n -by- m matrix \mathbf{A} over a complete idempotent semiring, the existence of the solution can be checked in $O(nm)$ time (see Gaubert 1997). Although the tropical algebra is not complete, there is an efficient way of finding if the solution exists (Cuninghame-Green 1979; Zimmermann 2011). It was shown by Butkovič (2003) that this type of tropical equations is equivalent to the set cover problem, which is known to be NP-hard. This directly affects the max-times algebra through the above-mentioned isomorphism and makes the problem of precisely solving max-times linear systems of the form $\mathbf{Ax} = \mathbf{b}$ infeasible for high dimensions.

Homogeneous equations $\mathbf{Ax} = \mathbf{Bx}$ can be solved using the *elimination* method, which is based on the fact that the set of solutions of a homogeneous system is a finitely generated semimodule (Butkovič and Hegedüs 1984) (independently rediscovered by Gaubert 1992). If only a single solution is required, then according to Gaubert (1997), a method by Walkup and Borriello (1998) is usually the fastest in practice.

Now let \mathbf{A} be a tropical square matrix of size $n \times n$. For complete idempotent semirings a solution to the equation $\mathbf{x} = \mathbf{Ax} \oplus \mathbf{b}$ is given by $\mathbf{x} = \mathbf{A}^* \mathbf{b}$ (see e.g. Salomaa and Soittola 2012), where the operator \mathbf{A}^* is defined as

$$\mathbf{A}^* = \bigoplus_{k=1}^{\infty} \mathbf{A}^k.$$

Since the tropical semiring is not complete (it is missing the ∞ element), \mathbf{A}^* can not always be computed. However, when there are no positive weight circuits in the graph defined by \mathbf{A} , then we have $\mathbf{A}^* = \mathbf{A}^0 \oplus \dots \oplus \mathbf{A}^{n-1}$, and all entries of \mathbf{A}^* belong to the tropical semiring (Baccelli et al 1992). Computing the operator \mathbf{A}^* takes time $O(n^3)$ (see e.g. Gondran and Minoux 1984a; Gaubert 1997).

Another important direction of research is the eigenvalue problem $\mathbf{Ax} = \lambda \mathbf{x}$. Tropical analogues of the Perron–Frobenius theorem (see e.g. Vorobyev 1967; Maslov 1992), and Collatz–Wielandt formula (Bapat et al 1995; Gaubert 1992) were developed. For a general overview of the results in the $(\max, +)$ spectral theory, see for example Gaubert (1997).

Tropical algebra and tropical geometry were used by Gärtner and Jaggi (2008) to construct a tropical analogue of an SVM. Unlike in the classical case, tropical SVMs are localized, in the sense that the kernel at any given point is not influenced by all the support vectors. Their work also utilizes the fact that tropical hyperplanes are somewhat more complex than their counterparts in the classical geometry, which makes it possible to do multiple category classification with a single hyperplane.

7 Conclusions

Subtropical low-rank factorizations are a novel approach for finding latent structure from nonnegative data. The factorizations can be interpreted using the winner-takes-it-all

interpretation: the value of the element in the final reconstruction depends only on the largest of values in the corresponding elements of the rank-1 components (compare that to NMF, where the value in the reconstruction is the *sum* of the corresponding elements). That the factorizations are different does not necessarily mean that they are better in terms of reconstruction error, although they can yield lower reconstruction error than even SVD. It does mean, however, that they find different structure from the data. This is an important advantage, as it allows the data analyst to use both the classical factorizations and the subtropical factorizations to get a broader understanding of the kinds of patterns that are present in the data.

Working in the subtropical algebra is harder than in the normal algebra, though. The various definitions for the rank, for example, do not agree, and computing many of them – including the subtropical Schein rank, which is arguably the most useful one for data analysis – is computationally hard. That said, our proposed algorithms, **Capricorn** and **Cancer**, can find the subtropical structure when it is present in the data. Not every data has subtropical structure, though, and due to the complexity of finding the optimal subtropical factorization we cannot distinguish between the cases where our algorithms fail to find the latent subtropical structure, and where it does not exist. Based on our experiments with synthetic data, our hypothesis is that the failure of finding a good factorization more probably indicates the lack of the subtropical structure rather than the algorithms' failure. Naturally, more experiments using data with known subtropical structure should improve our confidence of the correctness of the hypothesis.

The presented algorithms are heuristics. Developing algorithms that achieve better reconstruction error is naturally an important direction of future work. In our **Equator** framework, this hinges on the task of finding the rank-1 components. In addition, the scalability of the algorithms could be improved. A potential direction could be to take into account the sparsity of the factor matrices in dominated decompositions. This could allow one to concentrate only on the non-zero entries in the factor matrices.

The connection between Boolean and (sub-)tropical factorizations raises potential directions for future work. The continuous framework could allow for easier optimization in the Boolean algebra. Also, the connection allows us to model combinatorial structures (e.g. cliques in a graph) using subtropical matrices. This could allow for novel approaches on finding such structures using continuous subtropical factorizations.

References

- Akian M (1999) Densities of idempotent measures and large deviations. *Trans Amer Math Soc* 351(11):4515–4543, DOI 10.1090/S0002-9947-99-02153-4
- Akian M, Bapat R, Gaubert S (2007) Max-plus algebra. In: Hogben L (ed) *Handbook of Linear Algebra*, Chapman & Hall/CRC, Boca Raton
- Akian M, Gaubert S, Guterman A (2009) Linear independence over tropical semirings and beyond. *Contemp Math* 495:1–38
- Baccelli F, Cohen G, Olsder GJ, Quadrat JP (1992) *Synchronization and linearity: An algebra for discrete event systems*. John Wiley & Sons, Hoboken, New Jersey, DOI 10.2307/2583959
- Bapat R, Stanford DP, Van den Driessche P (1995) Pattern properties and spectral inequalities in max algebra. *SIAM J Matrix Anal Appl* 16(3):964–976, DOI 10.1137/S0895479893251782
- Berry MW, Browne M, Langville AN, Pauca VP, Plemmons RJ (2007) Algorithms and applications for approximate nonnegative matrix factorization. *Comput Stat Data Anal* 52(1):155–173, DOI 10.1016/j.csda.2006.11.006
- Blondel VD, Gaubert S, Tsitsiklis JN (2000) Approximating the spectral radius of sets of matrices in the max-algebra is NP-hard. *IEEE Trans Autom Control* 45(9):1762–1765, DOI 10.1109/9.880644

- Brunet JP, Tamayo P, Golub TR, Mesirov JP (2004) Metagenes and molecular pattern discovery using matrix factorization. *Proc Natl Acad Sci USA* 101(12):4164–4169, DOI 10.1073/pnas.0308531101
- Butkovič P (2003) Max-algebra: The linear algebra of combinatorics? *Linear Algebra Appl* 367:313–335, DOI 10.1016/S0024-3795(02)00655-9
- Butkovič P (2010) Max-linear systems: Theory and algorithms. Springer Science & Business Media, New York, DOI 10.1007/978-1-84996-299-5
- Butkovič P, Hegedüs G (1984) An elimination method for finding all solutions of the system of linear equations over an extremal algebra. *Ekon-Mat Obzor* 20(2):203–215
- Butkovič P, Hevery F (1985) A condition for the strong regularity of matrices in the minimax algebra. *Discrete Appl Math* 11(3):209–222, DOI 10.1016/0166-218X(85)90073-3
- Cassandras CG, Lafortune S (2008) Introduction to discrete event systems, 2nd edn. Springer, Berlin, DOI 10.1007/978-0-387-68612-7
- Cichocki A, Zdunek R, Phan AH, Amari S (2009) Nonnegative matrix and tensor factorizations: Applications to exploratory multi-way data analysis and blind source separation. John Wiley & Sons, Chichester, DOI 10.1002/9780470747278
- Cohen G, Gaubert S, Quadrat JP (1999) Max-plus algebra and system theory: Where we are and where to go now. *Annu Rev Control* 23:207–219, DOI 10.1016/S1367-5788(99)90091-3
- Cohen JE, Rothblum UG (1993) Nonnegative ranks, decompositions, and factorizations of non-negative matrices. *Linear Algebra Appl* 190:149–168, DOI 10.1016/0024-3795(93)90224-C
- Cuninghame-Green RA (1979) Minimax algebra. Springer, Berlin, DOI 10.1007/978-3-642-48708-8
- Davis TA, Hu Y (2011) The University of Florida sparse matrix collection. *ACM Trans Math Soft* 38(1):1–25, DOI 10.1145/2049662.2049663
- De Schutter B, De Moor B (2002) The QR decomposition and the singular value decomposition in the symmetrized max-plus algebra revisited. *SIAM Rev* 44(3):417–454, DOI 10.1137/S00361445024039
- Deerwester S, Dumais ST, Furnas GW, Landauer TK, Harshman R (1990) Indexing by latent semantic analysis. *J Am Soc Inf Sci* 41:391–407, DOI 10.1002/(SICI)1097-4571(199009)41:6%3C391::AID-ASI1%3E3.0.CO;2-9
- Dembo A, Zeitouni O (2010) Large deviations techniques and applications, 2nd edn. Springer, Berlin, DOI 10.1007/978-3-642-03311-7
- Gärtner B, Jaggi M (2008) Tropical support vector machines. Tech. Rep. ACS-TR-362502-01
- Gaubert S (1992) Théorie des systèmes linéaires dans les dioïdes. PhD thesis, Ecole nationale supérieure des mines de Paris
- Gaubert S (1997) Methods and applications of $(\max,+)$ linear algebra. In: 14th Annual Symposium on Theoretical Aspects of Computer Science (STACS), Springer, pp 261–282, DOI 10.1007/BFb0023465
- Georghiadis AS, Belhumeur PN, Kriegman DJ (2000) From few to many: Generative models for recognition under variable pose and illumination. In: 4th IEEE International Conference on Automatic Face and Gesture Recognition (FG), pp 277–284, DOI 10.1109/AFGR.2000.840647
- Gillis N, Glineur F (2010) Using underapproximations for sparse nonnegative matrix factorization. *Pattern Recogn* 43(4):1676–1687, DOI 10.1016/j.patcog.2009.11.013
- Golub GH, Van Loan CF (2012) Matrix computations, 4th edn. Johns Hopkins University Press, Baltimore
- Gondran M, Minoux M (1984a) Graphs and algorithms. John Wiley & Sons, New York
- Gondran M, Minoux M (1984b) Linear algebra in dioids: A survey of recent results. *North-Holland Math Stud* 95:147–163, DOI 10.1016/S0304-0208(08)72960-8
- Guillon P, Izhakian Z, Mairesse J, Merlet G (2015) The ultimate rank of tropical matrices. *J Algebra* 437:222–248, DOI 10.1016/j.jalgebra.2015.02.026
- Hoyer PO (2004) Non-negative matrix factorization with sparseness constraints. *J Mach Learn Res* 5:1457–1469
- Jha SK, Yadava R (2011) Denoising by singular value decomposition and its application to electronic nose data processing. *IEEE Sens J* 11(1):35–44, DOI 10.1109/JSEN.2010.2049351
- Karaev S, Miettinen P (2016a) Cancer: Another algorithm for subtropical matrix factorization. In: European Conference on Machine Learning and Principles of Knowledge Discovery in Databases (ECML PKDD), pp 576–592, DOI 10.1007/978-3-319-46227-1_36
- Karaev S, Miettinen P (2016b) Capricorn: An algorithm for subtropical matrix factorization. In: 16th SIAM International Conference on Data Mining (SDM), pp 702–710, DOI

- 10.1137/1.9781611974348.79
- Karaev S, Miettinen P, Vreeken J (2015) Getting to know the unknown unknowns: Destructive-noise resistant boolean matrix factorization. In: 15th SIAM International Conference on Data Mining (SDM), pp 325–333, DOI 10.1137/1.9781611974010.37
- Kim KH (1982) Boolean matrix theory and applications. Marcel Dekker, New York
- Kim KH, Roush FW (2005) Factorization of polynomials in one variable over the tropical semiring. Tech. Rep. math/0501167, arXiv
- Kolda T, O’Leary D (2000) Algorithm 805: Computation and uses of the semidiscrete matrix decomposition. *ACM Trans Math Softw* 26(3):415–435, DOI 10.1145/358407.358424
- Lee DD, Seung HS (1999) Learning the parts of objects by non-negative matrix factorization. *Nature* 401(6755):788–791, DOI 10.1038/44565
- Li Y, Ngom A (2013) The non-negative matrix factorization toolbox for biological data mining. *Source Code Biol Med* 8(1):1–15, DOI 10.1186/1751-0473-8-10
- Lu H, Vaidya J, Atluri V (2008) Optimal boolean matrix decomposition: Application to role engineering. In: 24th IEEE International Conference on Data Engineering (ICDE), pp 297–306, DOI 10.1109/ICDE.2008.4497438
- Lucchese C, Orlando S, Perego R (2014) A unifying framework for mining approximate top- k binary patterns. *IEEE Trans Knowl Data Eng* 26(12):2900–2913, DOI 10.1109/TKDE.2013.181
- Maslov V (1992) Idempotent analysis. American Mathematical Society, Providence
- Miettinen P (2009) Matrix decomposition methods for data mining: Computational complexity and algorithms. PhD thesis, University of Helsinki
- Miettinen P, Mielikäinen T, Gionis A, Das G, Mannila H (2008) The discrete basis problem. *IEEE Trans Knowl Data Eng* 20(10):1348–1362
- Mitchell-Jones A, Amori G, Bogdanowicz W, Krystufek B, Reijnders PH, Spitzenberger F, Stubbe M, Thissen J, Vohralik V, Zima J (1999) The atlas of European mammals. Academic Press, London
- Paatero P (1997) Least squares formulation of robust non-negative factor analysis. *Chemometr Intell Lab* 37(1):23–35, DOI 10.1016/S0169-7439(96)00044-5
- Paatero P (1999) The multilinear engine—table-driven, least squares program for solving multilinear problems, including the n -way parallel factor analysis model. *J Comp Graph Stat* 8(4):854–888, DOI 10.1080/10618600.1999.10474853
- Paatero P, Tapper U (1994) Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics* 5(2):111–126, DOI 10.1080/10618600.1999.10474853
- Pauca VP, Shahnaz F, Berry MW, Plemmons RJ (2004) Text mining using nonnegative matrix factorizations. In: 4th SIAM International Conference on Data Mining (SDM), pp 22–24, DOI 10.1137/1.9781611972740.45
- Salomaa A, Soittola M (2012) Automata-theoretic aspects of formal power series. Springer Science & Business Media, New York
- Sarwar B, Karypis G, Konstan J, Riedl J (2000) Application of dimensionality reduction in recommender system – a case study. Tech. rep., GroupLens Research Group
- Shitov Y (2014) The complexity of tropical matrix factorization. *Adv Math* 254:138–156, DOI 10.1016/j.aim.2013.12.013
- Simon I (1978) Limited subsets of a free monoid. In: 19th IEEE Annual Symposium on Foundations of Computer Science (FOCS), pp 143–150, DOI 10.1109/SFCS.1978.21
- Simon I (1994) On semigroups of matrices over the tropical semiring. *Inform Theor Appl* 28(3-4):277–294, DOI 10.1051/ita/1994283-402771
- Skillcorn D (2007) Understanding complex datasets: Data mining with matrix decompositions. Data Mining and Knowledge Discovery, Chapman & Hall/CRC, Boca Raton, DOI 10.1007/s00362-008-0147-y
- Srebro N, Rennie J, Jaakkola TS (2004) Maximum-margin matrix factorization. In: 17th Advances in Neural Information Processing Systems (NIPS), pp 1329–1336
- Vavasis SA (2009) On the complexity of nonnegative matrix factorization. *SIAM J Optim* 20(3):1364–1377, DOI 10.1137/070709967
- Vorobyev N (1967) Extremal algebra of positive matrices. *Elektron Informationsverarbeitung und Kybernetik* 3:39–71
- Walkup EA, Borriello G (1998) A general linear max-plus solution technique. In: Gunawardena J (ed) Idempotency, Cambridge University Press, Cambridge, pp 406–415, DOI 10.1017/CBO9780511662508.024

- Weston J, Weiss RJ, Yee H (2013) Nonlinear latent factorization by embedding multiple user interests. In: 7th ACM Conference on Recommender Systems (RecSys), pp 65–68, DOI 10.1145/2507157.2507209
- Xu W, Liu X, Gong Y (2003) Document clustering based on non-negative matrix factorization. In: 26th Annual International ACM SIGIR Conference (SIGIR), pp 267–273, DOI 10.1145/860435.860485
- Zimmermann U (2011) Linear and combinatorial optimization in ordered algebraic structures. Elsevier, Amsterdam