

**Nonlinear Mixture-wise Expansion Approach to  
Underdetermined Blind Separation of Nonnegative  
Dependent Sources**

Journal:	<i>Journal of Chemometrics</i>
Manuscript ID:	CEM-13-0069.R1
Wiley - Manuscript type:	Research Article
Date Submitted by the Author:	n/a
Complete List of Authors:	Kopriva, Ivica; Ruđer Bošković Institute, Laser and Atomic Research and Development Jerić, Ivanka; Ruđer Bošković Institute, Organic Chemistry and Biochemistry Brkljačić, Lidija; Ruđer Bošković Institute, Organic Chemistry and Biochemistry
Keyword:	Underdetermined blind source separation, Dependent sources, Reproducible kernel Hilbert spaces, Empirical kernel maps, Nonnegative matrix factorization

SCHOLARONE™  
Manuscripts

View

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

# Nonlinear Mixture-wise Expansion Approach to Underdetermined Blind Separation of Nonnegative Dependent Sources

Ivica Kopriva<sup>1\*</sup>, Ivanka Jerić<sup>2</sup>, and Lidija Brkljačić<sup>2</sup>

Ruder Bošković Institute, Bijenička cesta 54, HR-10000, Zagreb, Croatia

<sup>1</sup>Division of Laser and Atomic Research and Development

phone: +385-1-4571-286, fax:+385-1-4680-104

e-mail: ikopriva@irb.hr

<sup>2</sup>Division of Organic Chemistry and Biochemistry

e-mail: ijeric@irb.hr, Lidija.Brkljacic@irb.hr

## Abstract

Underdetermined blind separation of nonnegative dependent sources consists in decomposing set of observed mixed signals into greater number of original nonnegative and dependent component (source) signals. That is an important problem for which very few algorithms exist. It is also practically relevant for contemporary metabolic profiling of biological samples, such as biomarker identification studies, where sources (a.k.a. pure components or analytes) are aimed to be extracted from mass spectra of complex

1  
2  
3  
4 multicomponent mixtures. This paper presents method for underdetermined blind  
5  
6 separation of nonnegative dependent sources. The method performs nonlinear mixture-  
7  
8 wise mapping of observed data in high-dimensional reproducible kernel Hilbert space  
9  
10 (RKHS) of functions and sparseness constrained nonnegative matrix factorization  
11  
12 (NMF) therein. Thus, original problem is converted into new one with increased  
13  
14 number of mixtures, increased number of dependent sources and higher-order (error)  
15  
16 terms generated by nonlinear mapping. Provided that amplitudes of original components  
17  
18 are sparsely distributed, that is the case for mass spectra of analytes, sparseness  
19  
20 constrained NMF in RKHS yields, with significant probability, improved accuracy  
21  
22 relative to the case when the same NMF algorithm is performed on original problem.  
23  
24 The method is exemplified on numerical and experimental examples related  
25  
26 respectively to extraction of ten dependent components from five mixtures and to  
27  
28 extraction of ten dependent analytes from mass spectra of two to five mixtures.  
29  
30 Thereby, analytes mimic complexity of components expected to be found in biological  
31  
32 samples.  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

43 *Key words:* Underdetermined blind source separation, Dependent sources, Reproducible  
44  
45 kernel Hilbert spaces, Empirical kernel maps, Nonnegative matrix factorization.  
46  
47  
48  
49  
50

## 51 **1. INTRODUCTION**

52  
53  
54  
55 Blind source separation (BSS) refers to extraction of unknown source signals from  
56  
57 observed mixture signals only [1-4]. Within BSS framework a nonnegative BSS  
58  
59 (NBSS), where both mixing and source matrix are nonnegative, has drawn significant  
60

1  
2  
3  
4 attention recently yielding algorithms such as nonnegative independent component  
5 analysis (NICA) [5], nonnegative matrix factorization (NMF) [3, 6-8], convex  
6 analysis/geometry [9-11], nonnegative least correlated component analysis (nLCA)  
7 [12], determinant based sparseness measure approach to NBSS [13], and sparse  
8 component analysis (SCA) that combines data clustering and  $\ell_1$ -minimization [14, 15].

9  
10  
11  
12  
13  
14  
15  
16 A challenge for NBSS algorithms set by real world problems is characterized by more  
17 sources than mixtures available, i.e. NBSS problem is underdetermined (uNBSS),  
18 whereas sources are dependent. Such problems, associated with research related to  
19 health, food and environment, set motivation for development of the uNBSS algorithm  
20 to be presented herein. For example, 326 analytes were quantified in extracts of  
21 *Arabidopsis thaliana* leaf tissue [16], while the independent gas chromatography-mass  
22 spectrometry (GC-MS) study of *Arabidopsis thaliana* leaves detected 497 unique  
23 chemical components [17]. Metabolic profiling, that is seen as one of the most  
24 challenging tasks in chemical biology [18], aims to identify and quantify small-  
25 molecule analytes (a.k.a. pure components or sources) present in biological samples,  
26 typically urine, serum or tissue extract. Thereby, number of analytes can be large. For  
27 example, analysis of human adult urinary metabolome by liquid chromatography-mass  
28 spectrometry (LC-MS) revealed presence of 1484 components, while 384 of them were  
29 characterized by matching their spectra with references stored in libraries [19]. Great  
30 majority of algorithms developed for separation of dependent sources are incapable to  
31 deal with uNBSS problem, [5, 6, 9-14]. As opposed to them, few algorithms capable to  
32 handle uNBSS problem with dependent sources include [7, 8, 15, 20]. Hence, we  
33 propose new method for uNBSS problem with nonnegative dependent sources. It is a  
34 preprocessing method that performs nonlinear mixture-wise mapping of observed data  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4 and sparseness constrained NMF in high-dimensional mapped space. As properly  
5  
6 pointed out in [13], performances of many algorithms depends on optimal usage of  
7  
8 parameters required to be known *a priori*, such as balance parameter that regulates  
9  
10 influence of sparseness constraint [15, 20], or number of overlapping components that  
11  
12 exist in mixtures [8]. These parameters are difficult to select optimally in practice. To  
13  
14 the best of our knowledge the nonnegative matrix underapproximation (NMU)  
15  
16 algorithm is the only one that can handle uNBSS problem with dependent sources and  
17  
18 does not require *a priori* information from the user. Therefore, we propose herein to  
19  
20 combine nonlinear preprocessing transform (NPT) with the NMU algorithm in mapped  
21  
22 high-dimensional space. Hence, the NPT-NMU algorithm. The NPT-NMU is  
23  
24 exemplified on numerical and experimental problems. Nevertheless, proposed  
25  
26 preprocessing method can be used in combination with other sparseness constrained  
27  
28 NMF algorithms such as NMF algorithm with  $\ell_0$ -constraints (NMF\_L0) [8].  
29  
30  
31  
32  
33  
34  
35

36 The rest of the paper is organized as follows. Section 2 introduces instantaneous  
37  
38 (memoryless) linear mixture model, commonly used in chemometrics, defines uNBSS  
39  
40 problem and presents theory upon which proposed NPT approach is based. Section 3  
41  
42 describes experiments performed on synthetic and MS mixtures. Results of comparative  
43  
44 performance analysis between NMU, NMF-L0, NPT-NMU and NPT-NMF-L0  
45  
46 algorithms are discussed in Section 4. The NMF-L0 algorithm has been used as a  
47  
48 reference since it is known that  $\ell_0$ -constraints yield best results in the case of dependent  
49  
50 (overlapping) sources [21, 22]. In numerical, and even experimental, examples it was  
51  
52 possible to set optimally parameter related to number of overlapping sources.  
53  
54  
55  
56  
57  
58 Concluding remarks are given in Section 5.  
59  
60

## 2. THEORY AND ALGORITHM

Aimed application of proposed method is in extraction of analytes from multicomponent mixtures of mass spectra. MS is chosen due to its increasing importance in clinical chemistry, safety and quality control as well as biomarker discovery and validation. Identification of analytes is often achieved by matching experimental spectra to the ones stored in the library [23]. For an example the NIST and Wiley-Interscience universal spectral library [24], contains more than 800 000 mass spectra (corresponding to more than 680 000 compounds). Thus, we also assume that library of reference mass spectra is available to evaluate quality of components extracted by the proposed method.<sup>1</sup> Although various analytical methods are available for the separation of individual compounds from mixtures, ideal separation cannot be always accomplished, especially when dealing with complex samples [19]. There are also analytes that are prone to chemical decomposition and thus cannot be isolated [25]. Furthermore, when two or more analytes elute from chromatography column close to each other in time their peaks overlap partially or completely [26]. Thus, instead of analytes, their mixture will be compared with the reference pure components in the library. This sets motivation for development of algorithm for uNBSS problem with dependent sources.

### 2.1. Underdetermined nonnegative blind source separation with dependent sources

---

<sup>1</sup> Please note that any BSS algorithm when applied to experimental data requires some kind of expert knowledge to evaluate the separation results. Herein the library of pure components is such an "expert". The same concept is also in use in hyperspectral image analysis.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Linear mixture model (LMM) is commonly used in chemometrics [27-30] in general, and in MS in particular [29, 30]. It is the model upon which linear instantaneous BSS methods are based [1-4]. In the absence of additive noise the model reads as:

$$\mathbf{X} = \mathbf{A}\mathbf{S} \quad (1)$$

where  $\mathbf{X} \in \mathbb{R}_{0+}^{N \times T}$  represents matrix of acquired nonnegative mass spectra such that each row  $\{\mathbf{x}_n\}_{n=1}^N$  of  $\mathbf{X}$  contains one recorded multicomponent mixture mass spectra comprised of intensity values at  $T$   $m/z$  channels.  $\mathbf{A} \in \mathbb{R}_{0+}^{N \times M}$  represents mixture matrix, whereas each column vector  $\{\mathbf{a}_m\}_{m=1}^M$  represents concentration profile of the corresponding analyte across the  $N$  mixture spectra.  $\mathbf{S} \in \mathbb{R}_{0+}^{M \times T}$  is a matrix with the rows  $\{\mathbf{s}_m\}_{m=1}^M$  representing mass spectra of the unknown number of  $M$  analytes present in the mixture spectra  $\mathbf{X}$ . Thereby, the number of analytes  $M$  can be less than, equal to and greater than the number of recorded mixtures spectra  $N$ . This, respectively, leads to over-, even- and under-determined BSS problems in which case it is assumed that information about concentration of analytes (stored in the mixing matrix  $\mathbf{A}$ ) is not known to the BSS algorithm. That is, it is expected from BSS method to estimate matrix of analytes  $\mathbf{S}$  by having at disposal matrix with recorded mixtures spectra  $\mathbf{X}$  only. Due to high complexity of spectra of biological samples and, quite often, small number of recorded spectra available, it is certain that in such analyses corresponding BSS problem will be (highly) underdetermined:  $M > N$ . Thus, in this paper the following assumptions are made on LMM (1):

$$A1) 0 \leq s_{mt} < 1 \quad \forall m=1, \dots, M \quad t=1, \dots, T, ^2$$

$$A2) a_{nm} \geq 0 \quad \forall n=1, \dots, N \quad m=1, \dots, M \quad \text{and} \quad \|\mathbf{a}_{\cdot m}\|_2 = 1 \quad \forall m=1, \dots, M, ^3$$

$$A3) M > N$$

$$A4) M \ll T, ^4$$

Due to A1) and A2) it is clear that  $\mathbf{X} \geq \mathbf{0}$  as well. Furthermore, components spectra will overlap implying that at some  $m/z$  coordinates multiple components will be present.

This implies for column vectors of  $\mathbf{S}$ :  $\{\|\mathbf{s}_{\cdot t}\|_0 \leq K\}_{t=1}^T$ , where  $\|\mathbf{s}_{\cdot t}\|_0$  denotes  $\ell_0$  quasi-norm that counts number of nonzero entries of  $\mathbf{s}_{\cdot t}$ . Thus,  $K$  stands for maximal number of analytes that can be present at the particular  $m/z$  coordinate. Hence, sources  $\{\mathbf{s}_{m\cdot}\}_{m=1}^M$  will be statistically dependent. The uNBSS problem (1) is ill-posed due to the fact that matrix factorization suffers from indeterminacies:  $\mathbf{X} = \mathbf{A}\mathbf{S} = \mathbf{A}\mathbf{B}^{-1}\mathbf{B}\mathbf{S}$  for some invertible  $M \times M$  square matrix  $\mathbf{B}$ . Hence, it has an infinite number of solutions. Meaningful solutions are characterized by the permutation and scaling indeterminacies in which case  $\mathbf{B} = \mathbf{P}\mathbf{\Lambda}$ , where  $\mathbf{P}$  represents permutation and  $\mathbf{\Lambda}$  represents diagonal scaling matrix.

<sup>2</sup> Provided that A1) is not satisfied it can be satisfied by scaling  $\mathbf{X}$  with a constant  $c$ :  $\mathbf{X} \rightarrow \mathbf{X}/c$ . The

conservative scaling strategy that always guarantees A1) is given with:  $c = \arg \max_t \{\|\mathbf{x}_{\cdot t}\|_1\}_{t=1}^T$ . However,

scaling by  $c = \arg \max_{n,t} \{\mathbf{X}_{nt}\}_{n,t=1}^{N,T}$  will satisfy A1) in great majority of occasions.

<sup>3</sup> Due to the scaling indeterminacy that is inherent to the BSS problem magnitude of the mixing vectors cannot be guaranteed. Therefore, A2) constraint is assumed commonly in BSS.

<sup>4</sup> This technical assumption is necessary to ensure that resolution of the spectrometer is high enough to enable discrimination between components the number of which is expected to be large.



1  
2  
3  
4 However, constraints are necessary to be imposed on  $\mathbf{A}$  and/or  $\mathbf{S}$  to obtain solution of  
5  
6 uNBSS problem (1) that is unique up to permutation and scaling indeterminacies. The  
7  
8 necessary constraint is sparseness of analytes spectra  $\{\mathbf{s}_m\}_{m=1}^M$ . Sparseness constraint  
9  
10 implies that in relation to  $N$  and  $M$  the maximal number of analytes  $K$  present at the  
11  
12 particular  $m/z$  coordinate is small enough. However,  $K$  is application dependent. When  
13  
14 number of sources present in the mixture is large,  $K$  will grow. Compressive sensing  
15  
16 theory has established condition between  $N$ ,  $M$  and  $K$  necessary to obtain unique  
17  
18 solution for underdetermined system of linear equations:  $\{\mathbf{x}_t = \mathbf{A}\mathbf{s}_t\}_{t=1}^T$  assuming that  $\mathbf{A}$   
19  
20 is known and random with the entries distributed according to Gaussian or Bernoulli  
21  
22 distributions. For  $\ell_1$ -constrained solutions  $\{\mathbf{s}_t\}_{t=1}^T$  number of measurements  $N$   
23  
24 necessary to obtain unique solution with probability one is given with:  $N \approx K \log(M/K)$   
25  
26 [31]. When  $\ell_p$ -constraint,  $0 \leq p \leq 1$ , is used instead, condition on number of  
27  
28 measurements  $N$  is given with:  $N \geq C_1(p)K + pC_2(p)K \log(M/K)$  [21], where  $C_1$  and  $C_2$   
29  
30 are constants that depend on choice of the norm  $p$ . Hence,  $\lim_{p \rightarrow 0} N \geq C_1(0)K$ , i.e. when  $p=0$   
31  
32 number of measurements  $N$  does not depend on  $M$ . That explains good results of  $\ell_0$ -  
33  
34 constrained algorithms for solving (1) [8, 23], when compared against  $\ell_1$ -constrained  
35  
36 algorithms when  $K$  is increasing. However, when (1) is associated with uNBSS problem  
37  
38 in chemometrics  $\mathbf{A}$  is not random but deterministic, i.e. it is a concentration matrix. To  
39  
40 the best of our knowledge there is only one result related to condition necessary for  
41  
42 unique solution of the underdetermined system of equations:  $\{\mathbf{x}_t = \mathbf{A}\mathbf{s}_t\}_{t=1}^T$  when  $\mathbf{A}$  is  
43  
44 deterministic. It is shown in [32] that for cyclic polynomial matrix  $\mathbf{A}$  it applies:  
45  
46  $N = O(K^2)$ . That is significantly worse than  $N \approx K \log(M/K)$  [31], for random  $\mathbf{A}$ . When  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

analysis of complex mixtures takes place, where number of sources can be large and consequently  $K$  will grow, it can be necessary to ensure large number of mixtures spectra  $N$  in order to obtain solution of the uNBSS problem (1) that is, possibly, unique up to scaling and permutation indeterminacies. However, when  $N$  is associated with biological samples it can virtually be impossible to satisfy this requirement. Therefore, we propose nonlinear transformation of LMM (1) into quasi-linear model with increased number of measurements.

## 2.2. Nonlinear transform of linear mixture model

We propose mixture-wise nonlinear transform of LMM (1):  $\{\mathbf{x}_t \mapsto \phi(\mathbf{x}_t) \in \mathbb{R}_{0+}^{\bar{N}}\}_{t=1}^T$ ,

such that  $\bar{N} \gg N$ . We would like  $\bar{N}$  to be (very) large and possibly even infinite. The mapping has the following structure:

$$\phi(\mathbf{x}_t) = \left[ \left\{ c_{q_1 \dots q_N} x_{1t}^{q_1} \dots x_{Nt}^{q_N} \right\}_{q_1, \dots, q_N=0}^{\bar{N}} \right]^T \text{ such that } \sum_{n=1}^N q_n \leq \bar{N}, \quad \forall t = 1, \dots, T. \quad (2)$$

In (2)  $\{c_{q_1 \dots q_N}\}$  are real constants that are mapping dependent. By taking into account

that  $x_{nt} = \sum_{m=1}^M a_{nm} s_{mt}$ , (2) can be written as:

$$\phi(\mathbf{x}_t) = c_0 \mathbf{e}_1 + \mathbf{B} \begin{bmatrix} 0 \\ \mathbf{s}_t \end{bmatrix} + \mathbf{B}_{HOT} \begin{bmatrix} 0 \\ \mathbf{0}_{M \times 1} \\ \mathbf{s}_{tHOT} \end{bmatrix} \quad \forall t = 1, \dots, T \quad (3)$$

where  $HOT$  stands for higher order (nonlinear) terms introduced by mapping  $\phi(\mathbf{x}_t)$ ,  $\mathbf{e}_1$

is a unit vector in  $\mathbb{R}^{\bar{N}}$ ,  $\mathbf{0}_{M \times 1}$  is column vector with zero entries and  $\mathbf{s}_{tHOT}$  is  $\bar{N} - M - 1$

column vector comprised of  $\left\{ s_{1t}^{q_1} \times \dots \times s_{Mt}^{q_M} \right\}_{q_1, \dots, q_M=2}^{\bar{N}}$  and  $\sum_{m=1}^M q_m \leq \bar{N}$ . Provided that

LMM (1) is related to MS data analysis,  $\{s_{mt}\}_{m=1}^M$  represent analytes in mixtures mass spectra at the particular  $m/z$  coordinate, i.e.  $t$  corresponds to  $m/z$ . Then, all the cross-terms  $s_{1t}^{q_1} \times \dots \times s_{Mt}^{q_M}$  will be zero if only one analyte is not present at this coordinate.

Thus,  $\mathbf{s}_{t,HOT}$  in (3) will simplify to  $\mathbf{s}_{t,HOT} \approx [s_{1t}^2 \dots s_{Mt}^2 \dots s_{1t}^{\bar{N}} \dots s_{Mt}^{\bar{N}}]^T$ . Due to assumption A1),  $\{0 \leq s_{mt} < 1\}_{m=1}^M$ , many higher order terms in  $\mathbf{s}_{t,HOT}$  will go to zero as power term increases. Speed of decay depends on distribution of amplitudes. For sparse distributions, such as those encountered in MS, it is reasonable to expect that only several HOTS of each source will be significantly greater than zero. For an example, for amplitude  $s_{mt}=0.5$ , the 10th order power is  $9.7 \times 10^{-4}$ . Nevertheless, powers of  $\{s_{m\cdot}\}_{m=1}^M$  will represent new sources that are statistically dependent with the original ones. Thus we can write (3) as:

$$\phi(\mathbf{x}_{\cdot t}) \approx c_0 \mathbf{e}_1 + \bar{\mathbf{B}} \begin{bmatrix} 0 \\ \bar{\mathbf{s}}_{\cdot t} \end{bmatrix} \quad \forall t = 1, \dots, T \quad (4)$$

where  $\bar{\mathbf{s}}_{\cdot t} = [\mathbf{s}_{\cdot t} \ \mathbf{s}_{\cdot t,HOT}]^T$  and  $\bar{\mathbf{B}}$  combines on appropriate way  $\mathbf{B}$  and  $\mathbf{B}_{HOT}$ . Model (4) can be written in matrix format yielding:

$$\phi(\mathbf{X}) \approx \underbrace{\begin{bmatrix} c_0 \mathbf{e}_1 & \dots & c_0 \mathbf{e}_1 \end{bmatrix}}_{\times T \text{ times}} + \bar{\mathbf{B}} \begin{bmatrix} 0 \\ \bar{\mathbf{S}} \end{bmatrix} \quad (5)$$

where  $\phi(\mathbf{X}) \in \mathbb{R}_{0+}^{\bar{N} \times T}$ ,  $\bar{\mathbf{B}} \in \mathbb{R}_{0+}^{\bar{N} \times P+1}$  and  $\bar{\mathbf{S}} \in \mathbb{R}_{0+}^{P \times T}$ . Hence, the uNBSS problem (1) characterized by triplet  $(N, M, K)$  is converted into new problem (5) characterized by triplet  $(\bar{N}, P, Q)$  where  $P > M$  stands for number of dependent sources in (5) and  $Q > K$

stands for number of overlapping sources in (5). Provided that amplitudes of the sources are sparsely distributed it is justified to expect that:

$$(\bar{N}/N) \gg (P/M) \text{ as well as } (\bar{N}/N) \gg (Q/K). \quad (6)$$

In the light of the uniqueness condition related analysis presented in [32], sparseness constrained factorization of (5) will with significant probability yield, depending on fulfillment of (6), increased accuracy when compared against the same factorization method used for the uNBSS problem (1). The difficulty with factorization of problem (5) is that  $\bar{N}$  can be large or even infinite, in which case factorization becomes computationally intractable. To alleviate this difficulty a special type of nonlinear mapping  $\phi$  is selected such that space induced by it is reproducing kernel Hilbert space (RKHS) of functions. To this end we introduce the following definitions and theorems.

**Definition 2.2.1.** A real function  $\kappa: \mathbb{R}^N \times \mathbb{R}^N \rightarrow \mathbb{R}$  is positive semi-definite if it is symmetric and satisfies for any finite set of points  $\{\mathbf{x}_i\}_{i=1}^T$  in  $\mathbb{R}^N$  and real numbers

$$\{\alpha_i\}_{i=1}^T : \sum_{i,j=1}^T \alpha_i \alpha_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \geq 0.$$

**Theorem 2.2.1.** The Moore-Aronszajn theorem, [33]. Given any nonnegative definite function  $\kappa(\mathbf{x}, \mathbf{y})$  there exists a uniquely determined RKHS  $H_\kappa$  consisting of real valued

functions on set  $\mathbf{X} \subset \mathbb{R}^N$   $f: \mathbf{X} \rightarrow \mathbb{R}$  such that: (i)  $\forall \mathbf{x} \in \mathbf{X}, \kappa(\circ, \mathbf{x}) \in H_\kappa$ ; (ii)

$\forall \mathbf{x} \in \mathbf{X}, \forall f \in H_\kappa, f(\mathbf{x}) = \langle f, \kappa(\circ, \mathbf{x}) \rangle_{H_\kappa}$ . Here,  $\langle \circ, \circ \rangle$  denotes the inner product

associated with  $H_\kappa$ .

**Definition 2.2.2.** Replacing  $f(\mathbf{x})$  in (ii) in Theorem 2.2.1. by  $\kappa(\circ, \mathbf{x})$ , it follows

$\kappa(\mathbf{x}_t, \mathbf{x}) = \langle \kappa(\circ, \mathbf{x}_t), \kappa(\circ, \mathbf{x}) \rangle_{H_\kappa}$ . By selecting the nonlinear map as  $\phi(\mathbf{x}) = \kappa(\circ, \mathbf{x})$  it follows  $\kappa(\mathbf{x}_t, \mathbf{x}) = \langle \phi(\mathbf{x}_t), \phi(\mathbf{x}) \rangle_{H_\kappa}$ . That is known as *kernel trick*. The nonlinear mapping  $\phi(\mathbf{x}_t)$  is called explicit feature map (EFM).

Practical importance of the *kernel trick* is enormous since it substitutes evaluation of inner product of possibly infinite dimensional mappings  $\langle \phi(\mathbf{x}_t), \phi(\mathbf{x}) \rangle_{H_\kappa}$  in  $H_\kappa$  by evaluation of kernel function  $\kappa(\mathbf{x}_t, \mathbf{x})$  in the space spanned by empirical set of patterns  $\mathbf{X}$ . To substitute EFM-based nonlinear mappings in (5) by implicit kernel-based mappings we need to define empirical kernel map (EKM). To this end we use the following definition, see also definition 2.15 in [34].

**Definition 2.2.3.** Empirical kernel map. For a given set of patterns  $\{\mathbf{v}_d \in \mathbb{R}^N\}_{d=1}^D \subset \mathbf{X}$ ,

$D \in \mathbb{N}$ , we call  $\psi: \mathbb{R}^N \rightarrow \mathbb{R}^D$ , where

$\mathbf{x}_t \mapsto \kappa(\circ, \mathbf{x}_t) \Big|_{\{\mathbf{v}_d\}_{d=1}^D} = [\kappa(\mathbf{v}_1, \mathbf{x}_t), \dots, \kappa(\mathbf{v}_D, \mathbf{x}_t)]^T \quad \forall t = 1, \dots, T$  the empirical kernel

map with respect to  $\{\mathbf{v}_d\}_{d=1}^D$ .

Hence, EKM  $\psi(\mathbf{x}_t)$  is obtained by projecting EFM  $\phi(\mathbf{x}_t)$  associated with kernel

$\kappa(\circ, \mathbf{x}_t)$  on a  $D$ -dimensional subspace in RKHS spanned by  $\{\phi(\mathbf{v}_d) \in \mathbb{R}^{\bar{N}}\}_{d=1}^D$ :

$$\psi(\mathbf{x}_t) = [\phi(\mathbf{v}_{.1}) \dots \phi(\mathbf{v}_{.D})]^T \phi(\mathbf{x}_t) = \begin{bmatrix} \kappa(\mathbf{x}_t, \mathbf{v}_{.1}) \\ \dots \\ \kappa(\mathbf{x}_t, \mathbf{v}_{.D}) \end{bmatrix} \quad \forall t=1, \dots, T. \quad (7)$$

If (4) is substituted in (7) we obtain:

$$\begin{aligned} \psi(\mathbf{x}_t) &= [\phi(\mathbf{v}_{.1}) \dots \phi(\mathbf{v}_{.D})]^T \phi(\mathbf{x}_t) \approx [\phi(\mathbf{v}_{.1}) \dots \phi(\mathbf{v}_{.D})]^T \left[ c_0 \mathbf{e}_1 + \bar{\mathbf{B}} \begin{bmatrix} 0 \\ \bar{\mathbf{s}}_t \end{bmatrix} \right] \\ &\approx c_0 \underbrace{\begin{bmatrix} \kappa(\mathbf{e}_1, \mathbf{v}_{.1}) \\ \dots \\ \kappa(\mathbf{e}_1, \mathbf{v}_{.D}) \end{bmatrix}}_{\mathbf{c}} + \underbrace{[\phi(\mathbf{v}_{.1}) \dots \phi(\mathbf{v}_{.D})]^T}_{\bar{\mathbf{B}}} \bar{\mathbf{B}} \begin{bmatrix} 0 \\ \bar{\mathbf{s}}_t \end{bmatrix} \approx \mathbf{c} + \bar{\mathbf{B}} \begin{bmatrix} 0 \\ \bar{\mathbf{s}}_t \end{bmatrix} \quad \forall t=1, \dots, T \end{aligned}$$

Hence, we can write (7) in the matrix version as:

$$\psi(\mathbf{X}) \approx \mathbf{C} + \bar{\mathbf{B}} \begin{bmatrix} \mathbf{0}_{1 \times T} \\ \bar{\mathbf{S}} \end{bmatrix} \quad (8)$$

where  $\psi(\mathbf{X}) \in \mathbb{R}_{0+}^{D \times T}$ ,  $\mathbf{C} \in \mathbb{R}_{0+}^{D \times T}$ ,  $\bar{\mathbf{B}} \in \mathbb{R}_{0+}^{D \times P+1}$  and  $\bar{\mathbf{S}} \in \mathbb{R}_{0+}^{P \times T}$ .  $\mathbf{C}$  in (8) represents bias term

and does not play a role in parts based decomposition of  $\Psi(\mathbf{X})$  that is enforced by sparseness constrained NMF. Hence, the uNBSS problem (1) characterized by triplet  $(N, M, K)$  is converted into new problem (8) characterized by triplet  $(D, P, Q)$  where  $P > M$  stands for number of dependent sources in (8) and  $Q > K$  stands for number of overlapping sources in (8). Analogously to (6), provided that amplitudes of the sources are sparsely distributed it is justified to expect that:

$$(D/N) \gg (P/M) \text{ as well as } (D/N) \gg (Q/K). \quad (9)$$

In the light of the uniqueness condition related analysis presented in [32], sparseness constrained factorization of (8) will with significant probability yield, depending on

fulfillment of (9), increased accuracy when compared against the same factorization method used for the uNBSS problem (1). Thus nonnegativity and sparseness constrained factorization of (8) should extract original sources  $\{\mathbf{s}_{m\cdot}\}_{m=1}^M$  as well as their powers that actually are new sources that are dependent with the original ones. While in (5)  $\bar{N}$  is large or even infinite,  $D$  in (8) is finite. To perform projection implied by (7) a basis in the original empirical data set  $\mathbf{X}$  has to be constructed  $\mathbf{V} = \{\mathbf{v}_{\cdot d} \in \mathbb{R}_{0+}^N\}_{d=1}^D$  such that

$$\text{span}\{\phi(\mathbf{v}_{\cdot d})\}_{d=1}^D \approx \text{span}\{\phi(\mathbf{x}_{\cdot t})\}_{t=1}^T \quad (10)$$

where *span* denotes a vector space spanned by particular set of vectors, i.e. it is expected that basis vectors span the same vector space that is spanned by empirical set of patterns. The basis  $\mathbf{V}$  can be constructed on several ways, for example using data clustering whereas cluster centers represent basis vectors. Hence, basis construction can be computationally challenging problem for itself. This, however, can be avoided if each pattern vector  $\{\mathbf{x}_{\cdot t}\}_{t=1}^T$  is chosen as a basis vector, i.e.  $\mathbf{V}=\mathbf{X}$ . Then condition (10) is satisfied perfectly. In this case, however, dimension of the projected sub-space  $D$  equals the number of the  $m/z$  channels  $T$ . Hence, the matrix  $\psi(\mathbf{X})$  implied by (8) will have dimensions  $T \times T$ . For low-resolution mass spectrometry  $T$  is of the order of several thousands and matrix factorization problems implied by (8) are computationally tractable even on today's personal computers. When it comes to the kernel function  $\kappa(\circ, \mathbf{x}_{\cdot t})$  necessary to compute  $\psi(\mathbf{x}_{\cdot t})$  in (7), respectively (8) for matrix formulation, it is important that induced RKHS is high-dimensional. Although there are many kernel functions that satisfy this requirement we restrict ourselves herein to the one, arguably,

most often used kernel, [34]: the Gaussian kernel

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / \sigma^2\right) \text{ where } \sigma^2 \text{ denotes kernel bandwidth or variance. In}$$

principle, optimal value of  $\sigma^2$  depends on dimension  $N$  and has to be estimated through cross-validation. When A1) and A2) apply we have found empirically that variance is in the interval  $\sigma^2 \in [0.1, 1]$ . Higher order terms (HOT) present in  $\bar{\mathbf{S}}$  in (8) can also be seen as a noise that is generated by nonlinear transformation. Thus, error introduced by HOT can be reduced by applying entry-wise either soft thresholding nonlinearity on

$$\{\psi(\mathbf{X})_{dt}\}_{d,t=1}^{D,T} : \{\psi(\mathbf{X})_{dt} \rightarrow \eta_\tau(\psi(\mathbf{X})_{dt}) = \max(0, \psi(\mathbf{X})_{dt} - \tau)\}_{d,t=1}^{D,T}, [35], \text{ or hard}$$

$$\text{thresholding nonlinearity: } \{\psi(\mathbf{X})_{dt} \rightarrow \nu_\tau(\psi(\mathbf{X})_{dt}) = \psi(\mathbf{X})_{dt} 1_{\psi(\mathbf{X})_{dt} > \tau}\}_{d,t=1}^{D,T} \text{ where } 1_{\psi(\mathbf{X})_{dt} > \tau}$$

represents indicator function. Through numerical experiments under assumptions A1) and A2) we have empirically found that, if "de-noising" operator is applied, thresholding parameter ought to be set to  $\tau \approx 10^{-7}$ . Hence, from the user perspective the NPT-NMU algorithm is virtually almost parameter free.

### 2.3. Sparseness constrained factorization

The NMU algorithm [7], with a MATAB code available at [36], has been used to evaluate effectiveness of proposed nonlinear mixture expansion method. The NMU method performs factorization of (8) in a recursive manner extracting one component at a time. After identifying optimal rank-one solution  $(\bar{\mathbf{b}}_1, \bar{\mathbf{s}}_1)$  the rank-one factorization is performed on the residue matrix  $\psi(\mathbf{X}) \leftarrow \psi(\mathbf{X}) - \bar{\mathbf{b}}_1 \bar{\mathbf{s}}_1$ . To preserve non-negativity of  $\psi(\mathbf{X})$  an underapproximation constraint is imposed on  $\bar{\mathbf{B}}$  and  $\bar{\mathbf{S}}$ :  $\bar{\mathbf{B}}\bar{\mathbf{S}} \leq \psi(\mathbf{X})$ . This constraint yields localized parts-based decomposition where different basis elements



1  
2  
3  
4 describe disjoint parts of the input data  $\psi(\mathbf{X})$ . It has been proven in theorem 1 in [7] that  
5  
6  
7 sparseness (number of non-zero entries) of  $\bar{\mathbf{B}}$  and  $\bar{\mathbf{S}}$  is less than sparseness of  $\psi(\mathbf{X})$ . A  
8  
9  
10 main reason for preferring the NMU algorithm over other sparseness constrained NMF  
11  
12 algorithms is that there are no regularization constants that require a tuning procedure.

13  
14 When performing NMU-based factorization of matrix  $\psi(\mathbf{X})$  in (8), the unknown number  
15  
16  
17 of analytes  $P$  needs to be given to the algorithm as an input.  $P$  represents nonnegative  
18  
19  
20 rank of  $\psi(\mathbf{X})$  in (8) and its optimal selection is related the model order selection.

21  
22 However, the well-known difficulty with model selection methods is that very often,  
23  
24  
25 due to different theoretical assumptions, they yield (significantly) different result when  
26  
27  
28 applied to the same data. Due to this reason and also due to the reason of not to lose  
29  
30  
31 some of the components our strategy was to set an assumed value of  $P$  as:  $\hat{P} = D = T$ .

32  
33 Then, due to (8),  $P < \hat{P}$ . Then, the NMU algorithm will extract all the  $\hat{P}$  components  
34  
35  
36 contained in  $\psi(\mathbf{X})$  at an increased computational complexity, i.e. the  $T$ - $P$  rank-one  
37  
38  
39 factors will be computed unnecessarily. Nevertheless, that is the price worth being paid  
40  
41  
42 in order not to lose some components that potentially can lead to biomarker discovery.

43  
44 When all  $T$  components are extracted they are compared with the reference components  
45  
46  
47 from the library. We identify analytes candidates as those that are maximally correlated  
48  
49  
50 with components in the library. As a reference solution in the benchmark problem we  
51  
52  
53 have used solution obtained by applying the NMF\_L0 algorithm [8], to the original  
54  
55  
56 problem (1). The MATLAB code of the NMF\_L0 algorithm is available at [37].

57  
58 NMF\_L0 is based on natural sparseness measure, the  $\ell_0$ -pseudo-norm of the source or  
59  
60  
61 component matrix  $\mathbf{S}$ . The NMF\_L0 when applied to (1) requires *a priori* information  
62  
63  
64 on the number of components  $M$  and number of overlapping components  $K$ . In

numerical scenario both  $M$  and  $K$  are known while in experimental scenario selecting optimal (true) value of  $K$  can be difficult. Nevertheless, the NMF\_L0 can provide a good reference in validating worst-case and average performance of the NMU algorithm when applied on uNBSS problem (1) and as well as the NPT-NMU algorithm, that is the NMU algorithm applied on uNBSS problem (8). We summarize the NPT-NMU algorithm in the Algorithm 1.

### 3. EXPERIMENTS

Studies on numerical and experimental data reported below were executed on personal computer running under Windows 64-bit operating system with 64GB of RAM using Intel Core i7-3930K processor and operating with a clock speed of 3.2 GHz. Matlab 2012b environment has been used for programming.

#### 3.1. Numerical simulation

In numerical study we simulate LMM (1) with  $N=5$ ,  $M=10$ ,  $T=1000$  and  $K \in \{2, 3, 4\}$ . Each source is generated according to A1) and distributed according to probability density function of mixed state random variable, [38, 39]:

$$p(s_{mt}) = \rho \delta(s_{mt}) + (1 - \rho) \delta^*(s_{mt}) f(s_{mt}) \quad \forall m = 1, \dots, M \quad \forall t = 1, \dots, T \quad (11)$$

where  $\delta(s_{mt})$  is an indicator function and  $\delta^*(s_{mt}) = 1 - \delta(s_{mt})$  is its complementary function,

$\rho = \{P(s_{mt} = 0)\}_{t=1}^T$ . Hence,  $\{P(s_{mt} > 0)\}_{t=1}^T = 1 - \rho$ . We have generated sources with

probability of being zero  $\rho=0.9$ . The nonzero state of  $s_{mt}$  is distributed according to

$f(s_{mt})$ . We have chosen exponential distribution  $f(s_{mt}) = (1/\mu) \exp(-s_{mt}/\mu)$  to model

sparse distribution of the nonzero states such that the most probable outcomes were equal to  $\mu=0.1$  and  $\mu=0.01$ .<sup>5</sup> For these outcomes figures 1 and 2 respectively show values of normalized correlation coefficients

$$c_{mm} = \langle \mathbf{s}_{m\cdot}, \hat{\mathbf{s}}_{m\cdot} \rangle / \|\mathbf{s}_{m\cdot}\| \|\hat{\mathbf{s}}_{m\cdot}\|, \quad \forall m = 1, \dots, M \quad (12)$$

between true and separated sources versus Monte Carlo run index. Sources were separated by the NMU algorithm [7], the NMF\_L0 algorithm [8] and proposed NPT-NMU algorithm. Since for the NPT-NMU algorithm one run took roughly two hours, only 10 Monte Carlo runs were executed for each simulation scenario. Left column shows minimal value of the correlation coefficient attained by each of the algorithms while right column shows average value of the correlation coefficients for ten sources. The true values for  $M$  and  $K$  were reported to the NMF\_L0 algorithm and true value for  $M$  was reported to the NMU algorithm. The NPT-NMU algorithm assumed that  $M=T$ . The NMF\_L0 algorithm was run with the following parameter setup: reverse sparse nonnegative least square (rsNNLS) sparse coder and alternating nonnegative least square (ANLS) for dictionary update stage. Careful inspection of results presented in figures 1 and 2 suggests that NPT-NMU algorithm yields better accuracy than NMU algorithm in 30% of the runs, while NMU is better in 60% of the runs. While average performance of the NMF\_L0 algorithm is the best it yields the worse value of the minimal correlation coefficient.

### 3.2. MS measurements

<sup>5</sup> Even though the exponential distribution has support on the  $[0, \infty)$  interval setting  $\mu=0.1$  implies that with probability 0.9999546 realizations will be contained in  $[0, 1]$  interval. For  $\mu=0.01$  realizations will be contained in  $[0, 1]$  interval with a probability that is close to 1 with an error of  $3.72 \times 10^{-44}$ . Thus, this justifies a choice of exponential distribution to model sparse distribution of amplitudes  $s_{mi}$  on interval  $[0, 1]$ .

### 3.2.1. Chemicals

A library composed of mass spectra of ten amino acids, namely Ala, Asn, Asp, Gln, Glu, Leu, Lys, Phe, Pro and Val ( $C_1$ - $C_{10}$ ), was constructed. All amino acids and solvents were commercially available. Stock solutions of these amino acids (1 mg/mL) were prepared in 10% methanol. Working solutions (0.16 mg/mL) were prepared by diluting the stock solutions with 10% methanol. Five mixtures ( $X_1$ - $X_5$ ) were prepared by mixing different volumes of amino acid stock solutions according to Table S-1 given in Supplemental Information. Mass spectra of analytes were recorded by injection of 5  $\mu$ l of amino acid working solutions and mass spectra of five mixtures were obtained by injection of 15  $\mu$ l of mixture solutions prepared as described above, to the ion source. Mass spectra of analytes  $C_1$ - $C_{10}$  and mixtures  $X_1$ - $X_5$  are given in Supplemental Information (Figures S-1 and S-2).

### 3.2.2. Mass spectroscopy measurements

Electrospray ionization-mass spectrometry (ESI-MS) measurements operating in a positive ion mode were performed on a HPLC-MS triple quadrupole instrument equipped with an autosampler (Agilent Technologies, Palo Alto, CA, USA). The desolvation gas temperature was 300 °C with flow rate of 6.0 L/min. The fragmentor voltage was 135 V and capillary voltage was 4.0 kV. Mobile phase was 0.1% formic acid in 50% methanol and a flow rate of mobile phase was 0.2 ml/min. Mass spectra as total ion current spectra were recorded in  $m/z$  segment of 10-300. All data acquisition and processing was performed using Agilent MassHunter software.

### 3.2.3. Setting up an experiment

1  
2  
3  
4 Naturally occurring L-amino acids Ala, Asn, Asp, Gln, Glu, Leu, Lys, Phe, Pro and Val,  
5  
6  
7 Figure 3, were chosen for the construction of mass spectra library and preparation of  
8  
9 five mixtures used in validation of the proposed method. Using amino acids as testing  
10  
11 compounds is rationalized as follows: (a) they are metabolites often followed in  
12  
13 metabolomic studies, [40, 41], (b) their mass profile falls into relatively narrow  $m/z$   
14  
15 window thus mimicking complexity (overlapping) expected to be found in spectra of  
16  
17 real biological samples and (c) owing to the fragmentation often taking place in the MS  
18  
19 ion source, mass spectra of amino acids are enriched with numerous fragment ions  
20  
21 making separation problem even more challenging. Mass spectra of components C<sub>1</sub>-C<sub>10</sub>  
22  
23 are given in Figure S-1 in Supporting Information together with the assignment of the  
24  
25 most abundant fragment ions. Inspection of mass spectra clearly shows that some  
26  
27 fragment ions are present in spectra of different components. For example, the fragment  
28  
29 ion at  $m/z$  84 share components C<sub>4</sub>, C<sub>5</sub> and C<sub>7</sub>, while that at  $m/z$  116 components C<sub>2</sub>, C<sub>3</sub>  
30  
31 and C<sub>9</sub>. Moreover, difference in mass spectra of components C<sub>4</sub> (Gln) and C<sub>7</sub> (Lys) is  
32  
33 only in the intensity of fragment ions. Normalized cross-correlation coefficients of  
34  
35 components are shown in Table S-2 given in Supporting Information. As seen from  
36  
37 mass spectra, many of them (C<sub>4</sub> and C<sub>5</sub>, C<sub>5</sub> and C<sub>7</sub>, C<sub>3</sub> and C<sub>9</sub>, C<sub>2</sub> and C<sub>3</sub> as well as C<sub>2</sub>  
38  
39 and C<sub>9</sub>) are significantly correlated, while correlation between C<sub>4</sub> and C<sub>7</sub> is 0.9539.  
40  
41 Thus, blind extraction of these analytes from small number of given mixtures is a (very)  
42  
43 hard problem. It is also important to emphasize that mass spectra of mixtures were  
44  
45 obtained by direct injection of sample to the ion source (one minute run), *without*  
46  
47 chromatographic separation prior to MS analysis.  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

#### 4. RESULTS AND DISCUSSION

Sparseness constrained matrix factorization methods such as those used in our previous work [30] and the NMU method [7] failed to extract analytes from mixtures. That is explained by significant correlation (overlapping) between the analytes. The NMF\_L0 method [8] when applied to LMM (1) yielded decent results in extracting components from five mixtures. Correlation structure discussed in subsection 3.2.3 suggested that optimal value for  $K$  could be 4, 5 or 6. Thus, NMF\_L0 was cross-validated for values of  $K \in \{3,4,5,6\}$ . Note, however, that in truly experimental scenario correlation structure of analytes is unknown and selection of optimal  $K$  would require extensive cross-validation. Table S-3 shows the best results, in terms of maximal normalized correlation coefficients between extracted components and components in the library, obtained by applying the NMU and NMF\_L0 algorithms on mixture spectra according to the LMM (1). Arguably, the best result by NMF\_L0 is obtained for  $K=5$ . Since recorded mass spectra were composed of 2901  $m/z$  points, extraction of analytes according to model (8) has been reduced to NMU-, respectively NMF\_L0, based factorization of the  $2901 \times 2901$  matrix in mapped space. Before mapping the mixing matrix  $\mathbf{X}$  was scaled by  $\arg \max_t \{\|\mathbf{x}_{\cdot,t}\|_1\}_{t=1}^T$  as well as by  $\arg \max_{n,t} \{\mathbf{X}_{nt}\}_{n,t=1}^{N,T}$ . Gaussian kernel with  $\sigma^2 \in \{0.1, 1, 10\}$  has been used. Proposed NPT-NMU method managed to extract ten analytes with a reasonable accuracy from even two mixtures only. The accuracy improves by increasing the number of mixtures. In Table 1 selected results are shown in terms of maximal normalized correlation coefficient (12) obtained by the proposed NPT-NMU method as well as by NMF\_L0, while more comprehensive results are reported in Table S-4 in the Supporting Information. It took roughly two hours in described software environment

1  
2  
3  
4 to perform decomposition of each particular combination of mixture spectra. It was  
5  
6 especially demanding to cross-validate number of overlapping components  $Q$  in model  
7  
8 (8), and that is required by the NMF\_L0 algorithm, since value of  $Q$  depends on how  
9  
10 fast power terms of the original sources decay toward zero. The "best" result was  
11  
12 obtained for  $Q \approx 50$ , but that is significantly worse than obtained by the NMF\_L0 based  
13  
14 factorization of (1) for  $K=5$ . On the other side the NPT-NMU method yielded much  
15  
16 better result than NMF\_L0 method. Although quality of components extracted by the  
17  
18 NPT-NMU and NMF\_L0 methods was not perfect they were assigned uniquely to the  
19  
20 true ones in the library. This aspect is of practical importance in different areas such as  
21  
22 disease diagnosis, food quality control, environmental related studies that depend on  
23  
24 library matching. Mass spectra of analytes extracted by the proposed NPT-NMU  
25  
26 method from five mixtures are shown in Figure S-3 in Supporting Information. To take  
27  
28 into account scaling indeterminacy extracted analytes were scaled to 0-100 range  
29  
30 (dividing each extracted analyte by its maximal value and multiplying by 100).  
31  
32  
33  
34  
35  
36  
37

## 38 **5. CONCLUSION**

39  
40  
41 Problems such as metabolic profiling of biological samples aim to extract many  
42  
43 dependent (overlapping) analytes from small number of multicomponent mixtures mass  
44  
45 spectra. That results in underdetermined nonnegative blind source separation problem  
46  
47 (uNBSS) with dependent sources for which an algorithm is proposed. It performs  
48  
49 nonlinear mixture-wise mapping of observed data into reproducible kernel Hilbert space  
50  
51 (RKHS) of functions and sparseness constrained nonnegative matrix factorization  
52  
53 (NMF) in RKHS. For sparse signals such as those encountered in mass spectroscopy the  
54  
55 method yields, with significant probability, improved accuracy relative to the case when  
56  
57  
58  
59  
60

1  
2  
3  
4 the same NMF algorithm is performed on the original uNBSS problem. On demanding  
5 numerical and experimental problems the algorithm demonstrated capability to extract  
6  
7 ten dependent analytes from two to five mixtures. Thereby, extracted components were  
8  
9 assigned uniquely to the true ones in the library. That is practically important for  
10  
11 biomarker identification studies.  
12  
13  
14

## 15 16 17 **ACKNOWLEDGMENT**

18  
19  
20 This work has been supported through grant 9.01/232 "Nonlinear component analysis  
21  
22 with applications in chemometrics and pathology" funded by the Croatian Science  
23  
24 Foundation.  
25  
26  
27  
28  
29  
30

## 31 32 **REFERENCES**

- 33  
34 1. Hyvärinen A, Karhunen J, Oja E. *Independent Component Analysis*. John Wiley &  
35  
36 Sons, Inc.: New York, US, 2001.  
37  
38  
39 2. Cichocki A, Amari S. *Adaptive Blind Signal and Image Processing*. John Wiley: New  
40  
41 York, 2002.  
42  
43  
44 3. Cichocki A, Zdunek R, Phan A H, Amari, S I. *Nonnegative Matrix and Tensor*  
45  
46 *Factorizations*. John Wiley: Chichester, UK, 2009.  
47  
48  
49  
50 4. Comon P, Jutten C (eds). *Handbook of Blind Source Separation*. Academic Press:  
51  
52 Oxford, UK, 2010.  
53  
54  
55 5. Plumbley M, Oja E. A 'nonnegative PCA' algorithm for independent component  
56  
57 analysis. *IEEE Trans. Neural Netw.* 2004; **15** (1): 66-76.  
58  
59  
60



- 1  
2  
3  
4 6. Lee DD, Seunng HS. Learning the parts of objects by nonnegative matrix  
5  
6 factorization. *Nature* 1999; **401**: 788-791.  
7  
8  
9  
10 7. Gillis N, Glineur F. Using underapproximations for sparse nonnegative matrix  
11  
12 factorization. *Pattern Rec.* 2010; **43**: 1676-1687.  
13  
14  
15 8. Peharz R, Pernkopf, F. Sparse nonnegative matrix factorization with  $\ell^0$ -constraints.  
16  
17 *Neurocomputing.* 2012; **80**: 38-46.  
18  
19  
20 9. Chan T H, Ma W K, Chi C Y, Wang Y. Convex Analysis Framework for blind  
21  
22 separation of nonnegative sources. *IEEE Trans. Sig. Proc.* 2008; **56**(10): 5120-5134.  
23  
24  
25 10. Ambikapathi, A M, Chan T H, Ma W K, Chi C Y. Chance-Constrained Robust  
26  
27 Minimum-Volume Enclosing Simplex Algorithm for Hyperspectral Unmixing. *IEEE*  
28  
29 *Trans. Geosc. Remote Sens.* 2011; **49** (11): 4194-4209.  
30  
31  
32 11. Naanaa W, Nuzillard J M. A geometric approach to blind separation of  
33  
34 nonnegative and dependent sources. *Sig. Proc.* 2012; **92**: 2775-2784.  
35  
36  
37 12. Wand F Y, Chi C Y, Chan T H, Wang Y. Nonnegative least correlated component  
38  
39 analysis for separation of dependent sources by volume maximization. *IEEE Trans.*  
40  
41 *Pattern Anal. Mach. Intell.* 2010; **32** (5): 875-888.  
42  
43  
44 13. Yang Z, Xiang Y, Xie S, Ding S, Rong Y. Nonnegative Blind Source Separation by  
45  
46 Sparse Component Analysis Based on Determinant Measure. *IEEE Trans. Neural Netw.*  
47  
48 *and Learn. Sys.* 2012; **23** (10): 1601-1610.  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3  
4  
5 14. Sun Y, Xin J. Nonnegative Sparse Blind Source Separation for NMR Spectroscopy  
6  
7 by Data Clustering, Model Reduction, and  $\ell_1$  Minimization. *SIAM J. Imaging Sci.*  
8  
9 2012; **5** (3): 886-911.  
10  
11  
12 15. Kopriva I, Cichocki A. Blind decomposition of low-dimensional multi-spectral  
13  
14 image by sparse component analysis. *J. Chemometrics* 2009; **23** (11): 590-597.  
15  
16  
17  
18 16. Fiehn O, Kopka J, Dörmann P, Altmann T, Trethewey R N, Willmitzer L.  
19  
20 Metabolite profiling for plant functional genomics. *Nature Biotechnology* 2000; **18**:  
21  
22 1157-1161.  
23  
24  
25  
26 17. Jonsson P, Johansson A I, Gullberg J, Trygg J, Jiye A, Grung B, Marklund S,  
27  
28 Sjöström M, Antti H, Moritz T. High-throughput data analysis for detecting and  
29  
30 identifying differences between samples in GC/MS-based metabolomic analyses,"  
31  
32 *Analytical Chem.* 2005; **77**: 5635-5642.  
33  
34  
35  
36 18. Nicholson J K, Lindon J C. Systems biology: Metabonomics. *Nature* 2008; **455**  
37  
38 (7216): 1054-1056.  
39  
40  
41  
42 19. Roux A, Xu Y, Heilier J-F, Olivier M-F, Ezan E, Tabet J-C, Junot C. Annotation of  
43  
44 the human adult urinary metabolome and metabolite identification using ultra high  
45  
46 performance liquid chromatography coupled to a linear quadrupole ion trap-orbitrap  
47  
48 mass spectrometer. *Anal. Chem.* 2012; **84**: 6429-6437.  
49  
50  
51  
52 20. Cichocki A, Zdunek R, Amari S I. Hierarchical ALS Algorithms for Nonnegative  
53  
54 Matrix Factorization and 3D Tensor Factorization. *LNCS* 2007; **4666**: 169-176.  
55  
56  
57  
58 21. Chartran R, Staneva V. Restricted isometry properties and nonconvex compressive  
59  
60 sensing. *Inverse Problems* 2008; **24**: 035020 (14 pages).

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60
22. Mohimani H, Babie-Zadeh B, Jutten C. A Fast Approach for Overcomplete Sparse Decomposition Based on Smoothed  $\ell_0$  Norm. *IEEE Trans. Sig. Proc.* 2009; **57**(1): 289-301.
23. Abu-Farha M, Elisma F, Zhou H, Tian R, Asmer M S, Figeys D. Proteomics: from technology developments to biological applications. *Anal. Chem.* 2009; **81**: 4585-4599.
24. McLafferty F W, Stauffer D A, Loh S Y, Wesdemiotis C. Unknown Identification Using Reference Mass Spectra. Quality Evaluation of Databases. *J. Am. Soc. Mass. Spectrom.* 1999; **10**: 1229-1240.
25. Forseth R R, Schroeder F C. NMR-spectroscopic analysis of mixtures: from structure to function. *Curr. Opin. Chem. Biol.* 2011; **15** (1): 38-47.
26. Ni Y, Qiu Y, Jiang W, Suttlemyre K, Su M, Zhang W, Jia W, Du X. ADAP-GC 2.0 Deconvolution of Coeluting Metabolites from GC/TOF-MS Data for Metabolomic Studies. *Anal. Chem.* 2012; **84**: 6619-6629.
27. Nuzillard D, Bourg S, Nuzillard J M. Model-Free Analysis of Mixtures by NMR Using Blind Source Separation *J. Magn. Reson.* 1998; **133**: 358-363.
28. Visser E, Lee T W. An information-theoretic methodology for the resolution of pure component spectra without prior information using spectroscopic measurements. *Chemom. Int. Lab. Syst.* 2004; **70**: 147-155.
29. Kopriva I, Jerić I. Multi-component Analysis: Blind Extraction of Pure Components Mass Spectra using Sparse Component Analysis. *J. Mass Spectrom.* 2009; **44** (9): 1378-1388.

- 1  
2  
3  
4  
5 30. Kopriva I, Jerić I. Blind separation of analytes in nuclear magnetic resonance  
6 spectroscopy and mass spectrometry: sparseness-based robust multicomponent analysis.  
7  
8  
9 *Anal. Chem.* 2010; **82**: 1911-1920.  
10  
11  
12 31. Candès E, Tao T. Near optimal signal recovery from random projections: universal  
13 encoding strategy?. *IEEE Trans. Information Theory* 2006; **52**: 5406-5425.  
14  
15  
16  
17  
18 32. DeVore R A. Deterministic constructions of compressed sensing matrices. *Journal*  
19 *of Complexity* 2007; **23**: 918-925.  
20  
21  
22  
23 33. Aronszajn N. The theory of reproducing kernels. *Trans. of the Amer. Math. Soc.*  
24 1950; **68**: 337-404.  
25  
26  
27  
28  
29 34. Schölkopf B, Smola A. *Learning with kernels*. MIT Press, 2002.  
30  
31  
32 35. Donoho D L. De-Noising by Soft-Thresholding. *IEEE Trans. Inf. Theory* 1995; **41**  
33 (3): 613-627.  
34  
35  
36  
37  
38 36. The Nicolas Gillis Website. <https://sites.google.com/site/nicolasgillis/code> [7 March  
39 2013].  
40  
41  
42  
43 37. The Robert Peharz Website. <http://www3.spsc.tugraz.at/people/robert-peharz> [7  
44 March 2013].  
45  
46  
47  
48  
49 38. Bouthemy P, Piriou C H G, Yao J. Mixed-state auto-models and motion texture  
50 modeling. *J. Math Imaging Vision* 2006; **25**: 387-402.  
51  
52  
53  
54  
55 39. Caifa C, Cichocki A. Estimation of Sparse Nonnegative Sources from Noisy  
56 Overcomplete Mixtures Using MAP. *Neural Comput.* 2009; **21**: 3487-3518.  
57  
58  
59  
60

1  
2  
3  
4 40. Hisamatsu T, Okamoto S, Hashimoto M, Muramatsu T, Andou A, Uo M, Kitazume  
5  
6 M T, Matsuoka K, Yajima T, Inoue N, Kanai T, Ogata H, Iwao Y, Yamakado M, Sakai  
7  
8 R, Ono N, Ando T, Suzuki M, Hibi T. Novel, objective, multivariate biomarkers  
9  
10 composed of plasma amino acid profiles for the diagnosis and assessment of  
11  
12 inflammatory bowel disease. *PLoS ONE* 2012; **7**: 1-10.  
13  
14

15  
16  
17 41. Xu Y, Yang L, Yang F, Xiong Y, Wang Z, Hu Z. Metabolic profiling of fifteen  
18  
19 amino acids in serum of chemical-induced liver injured rats by hydrophilic interaction  
20  
21 liquid chromatography coupled with tandem mass spectrometry. *Metabolomics* 2012;  
22  
23 **8**(3): 475-483.  
24  
25  
26  
27  
28  
29  
30

### 31 **Figure Captions**

32  
33  
34 Figure 1. Normalized correlation coefficient vs. Monte Carlo run index between true  
35  
36 and extracted sources by algorithms: NMF\_L0 (squares), NMU (stars) and NPT-NMU  
37  
38 (circles). Left: minimal (worst) value; (right) mean value for ten sources. From top to  
39  
40 bottom - number of overlapping sources  $K$ : 2, 3, and 4. Most probable value of the  
41  
42 nonzero state, generated according to exponential distribution, equal to 0.1.  
43  
44  
45  
46

47  
48 Figure 2. Normalized correlation coefficient vs. Monte Carlo run index between true  
49  
50 and extracted sources by algorithms: NMF\_L0 (squares), NMU (stars) and NPT-NMU  
51  
52 (circles). Left: minimal (worst) value; (right) mean value for ten sources. From top to  
53  
54 bottom - number of overlapping sources  $K$ : 2, 3, and 4. Most probable value of the  
55  
56 nonzero state, generated according to exponential distribution, equal to 0.01.  
57  
58  
59

60 Figure 3. Structures of components  $C_1$ - $C_{10}$ .

1  
2  
3  
4  
5 **Table Captions**  
6  
7

8 Algorithm 1. The NPT-NMF (preferably NMU) algorithm.  
9

10  
11 **Required:**  
12

13  $\mathbf{X} \in \mathbb{R}_{0+}^{N \times T}$ . If A1) is not satisfied perform scaling

14  
15  
16  $\mathbf{X} \rightarrow \mathbf{X} / \arg \max_t \{\|\mathbf{x}_t\|_1\}^T$  or  $\mathbf{X} \rightarrow \mathbf{X} / \arg \max_{nt} \{\mathbf{X}_{nt}\}_{n,t=1}^{N,T}$ .  
17

- 18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60
1. Perform mapping  $\psi(\mathbf{X})$  in (7)/(8).
  2. Optionally, apply "de-noising" by soft- or hard thresholding operator entry-wise on  $\psi(\mathbf{X})$  in (8) with  $\tau=10^{-7}$ .
  3. Apply selected NMF algorithm to  $\mathbf{X}$  in (1) and  $\psi(\mathbf{X})$  in (8) to estimate, respectively,  $\mathbf{S}$  and  $\bar{\mathbf{S}}$ .
  4. Compare estimated  $\mathbf{S}$  and  $\bar{\mathbf{S}}$  with the reference components in the library to obtain the final estimate of  $\mathbf{S}$ .

Table 1. Maximal normalized correlation coefficients between analytes  $C_1$  to  $C_{10}$  and components extracted by the proposed NPT-NMU method. Columns from left to right: correlation coefficients; combinations of mixture spectra. The star '\*' denotes analytes in the library associated with the same extracted component. As expected highly correlated analytes  $C_4$  and  $C_7$  were associated with the same extracted component.

	$\mathbf{X}_1$ to $\mathbf{X}_5$	$\mathbf{X}_1$ to $\mathbf{X}_5$	$\mathbf{X}_{(1,3,4,5)}$	$\mathbf{X}_{(1,2,3)}$	$\mathbf{X}_{(3,4)}$
	NMF_L0	NPT-NMU	NPT-NMU	NPT-NMU	NPT-NMU
$c_{1,1}$	0.7269	0.8792	0.8486	0.7194	0.6232
$c_{2,2}$	0.9567	0.9370	0.8484	0.8855	0.8479
$c_{3,3}$	0.7448	0.9160	0.9142	0.6495	0.6889
$c_{4,4}$	0.8595	0.9816*	0.9008*	0.7474*	0.7308*
$c_{5,5}$	0.5616	0.6994	0.6107	0.5863	0.6461
$c_{6,6}$	0.9922	0.9844	0.9160	0.7958	0.9386
$c_{7,7}$	0.7117	0.9684*	0.8993*	0.7830*	0.7621*
$c_{8,8}$	0.6401	0.9869	0.9826	0.9671	0.9318
$c_{9,9}$	0.9924	0.9194	0.8746	0.9413	0.7998
$c_{10,10}$	0.9880	0.9398	0.9359	0.9826	0.8085
Kernel variance	Does not apply	1.0	1.0	1.0	0.1
Scaled by	$\max_{nt} \{\mathbf{X}_{nt}\}_{n,t}^{N,T}$	$\max_{nt} \{\mathbf{X}_{nt}\}_{n,t=1}^{N,T}$	$\max_{nt} \{\mathbf{X}_{nt}\}_{n,t=1}^{N,T}$	$\max_{nt} \{\mathbf{X}_{nt}\}_{n,t=1}^{N,T}$	$\max_t \{\ \mathbf{x}_{\cdot,t}\ \}_{t=1}^T$
"Denoising"	Does not apply	Hard thresholding $\tau=10^{-7}$	Hard thresholding $\tau=10^{-7}$	Hard thresholding $\tau=10^{-7}$	Hard thresholding $\tau=10^{-7}$

1  
2  
3  
4  
5  
6  
7 **Nonlinear Mixture-wise Expansion Approach to Underdetermined Blind**  
8  
9 **Separation of Nonnegative Dependent Sources**  
10  
11

12  
13  
14  
15 Ivica Kopriva<sup>1\*</sup>, Ivanka Jerić<sup>2</sup>, and Lidija Brkljačić<sup>2</sup>  
16

17  
18 Ruđer Bošković Institute, Bijenička cesta 54, HR-10000, Zagreb, Croatia  
19

20  
21 <sup>1</sup>Division of Laser and Atomic Research and Development  
22

23  
24  
25 phone: +385-1-4571-286, fax:+385-1-4680-104  
26

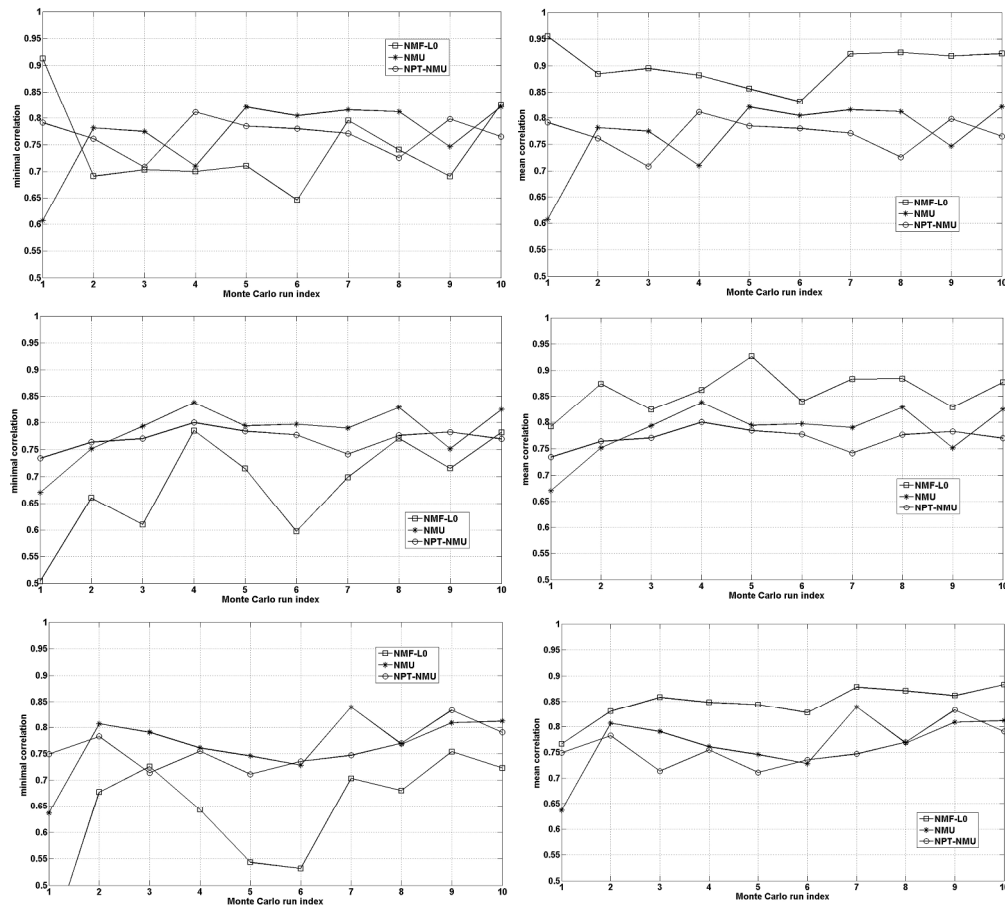
27  
28 e-mail: ikopriva@irb.hr  
29

30  
31 <sup>2</sup>Division of Organic Chemistry and Biochemistry  
32

33  
34 e-mail: ijeric@irb.hr, Lidija.Brkljacic@irb.hr  
35  
36  
37  
38  
39  
40  
41

42 **Summary abstract.** A method for underdetermined blind separation of nonnegative dependent  
43 sources is proposed. The method performs nonlinear mixture-wise mapping of observed data  
44 and sparseness constrained nonnegative matrix factorization (NMF) in high-dimensional  
45 mapped space. Proposed method can be applied with existing NMF algorithms to extract  
46 analytes from mass spectra of multicomponent mixtures in biomarker related studies of  
47 biological samples.  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

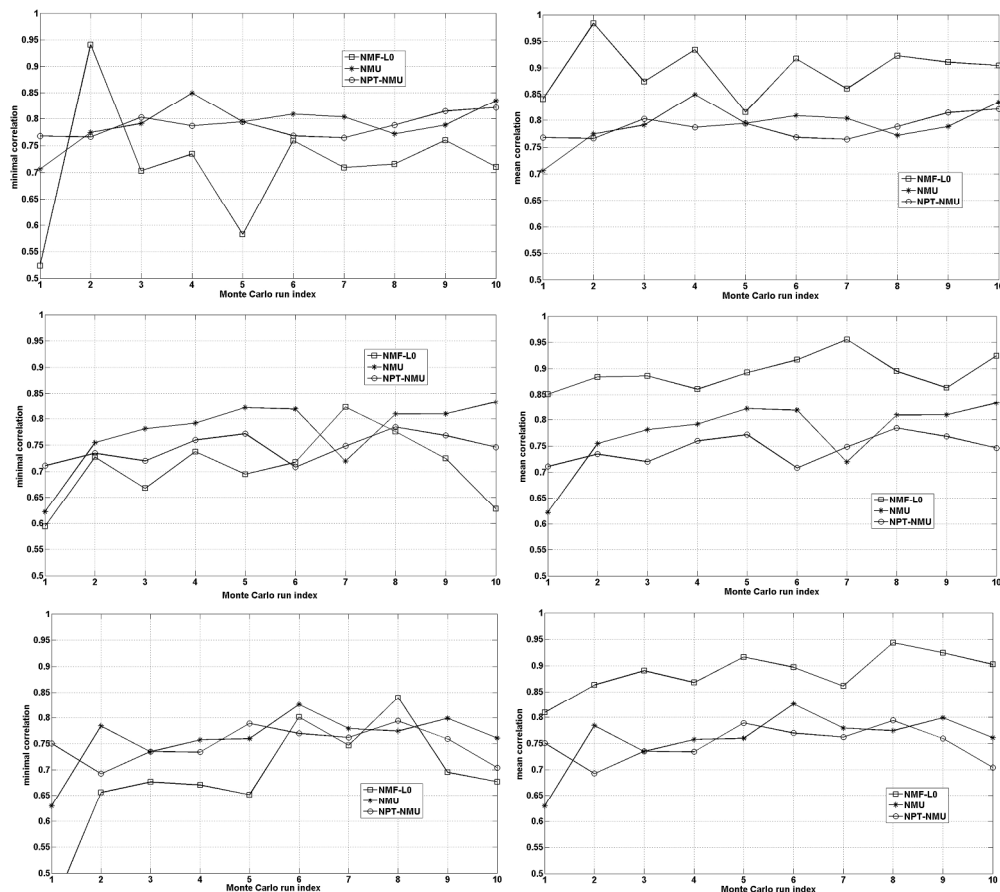




Normalized correlation coefficient vs. Monte Carlo run index between true and extracted sources by algorithms: NMF\_L0 (squares), NMU (stars) and NPT-NMU (circles). Left: minimal (worst) value; (right) mean value for ten sources. From top to bottom - number of overlapping sources  $K$ : 2, 3, and 4. Most probable value of the nonzero state, generated according to exponential distribution, equal to 0.1.

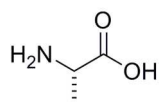
169x152mm (600 x 600 DPI)



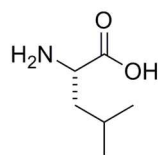


Normalized correlation coefficient vs. Monte Carlo run index between true and extracted sources by algorithms: NMF\_LO (squares), NMU (stars) and NPT-NMU (circles). Left: minimal (worst) value; (right) mean value for ten sources. From top to bottom - number of overlapping sources K: 2, 3, and 4. Most probable value of the nonzero state, generated according to exponential distribution, equal to 0.01.  
173x154mm (600 x 600 DPI)

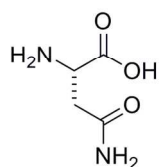




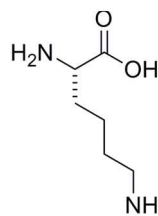
Ala (C1)



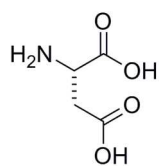
Leu (C6)



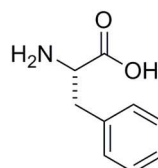
Asn (C2)



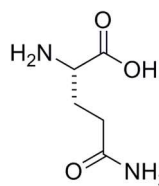
Lys (C7)



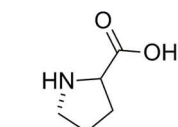
Asp (C3)



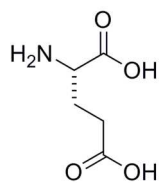
Phe (C8)



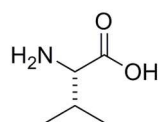
Gln (C4)



Pro (C9)



Glu (C5)



Val (C10)

Structures of components C1 to C10.  
58x179mm (300 x 300 DPI)