Contents

4 Non-negative matrix factorization 1 D. Brie, N. Gillis and S. Moussaoui 4.1 Introduction 1 4.1.1 Brief historical overview 2 Geometrical interpretation of NMF and the non-negative rank 2 4.2 Statistical formulation of NMF 9 4.3 4.3.1 Case of a Gaussian noise 9 4.3.2 Case of Poissonian noise 10 4.4 Uniqueness and admissible solutions of NMF 10 4.4.1 Uniqueness conditions 11 4.4.2 Finding the admissible solutions 12 4.5 Non-negative matrix factorization algorithms 16

- 4.5.1 Iterative factorization methods *16*
- 4.5.2 Constrained and penalized factorization methods 22
- 4.5.3 Geometrical approaches and separability 24
- 4.6 Applications of NMF chemical sensing. Two examples of reducing admissible solutions 25

| i

- 4.6.1 Polarized Raman spectroscopy: a data augmentation approach 25
- 4.6.2 Unmixing blurred Raman spectroscopy images 29
- 4.7 Conclusions 34

Index 41

ii |

4 Non-negative matrix factorization

D. Brie, N. Gillis and S. Moussaoui

Solving a source separation problem when the data at hand can be interpreted as linear instantaneous mixing of non-negative sources with non-negative mixing weights reduces to performing a non-negative factorization of the data matrix, which is referred to as non-negative matrix factorization (NMF). NMF has a long story originating from linear algebra and analytical chemistry and extensive developments has been recently achieved in the signal and image processing fields. The popularity of NMF is guided either by the mathematically challenging question of factorizing a matrix under non-negativity constraint and also by the need to explain observations as purely additive combination of non-negative factors or physically meaning quantities. This chapter addresses the concept of NMF, presents its foundations in terms of model setting and indeterminacy in addition to the main guidelines of existing factorization algorithms. The application of NMF to real situations of chemical data processing is illustrated with two examples of Raman spectroscopy measurements.

4.1 Introduction

The linear instantaneous mixing model assumes that P observations gathered in a vector $\boldsymbol{x}(t)$ can be modeled as the linear combination of R unknown sources

$$\boldsymbol{x}(t) = \boldsymbol{A}\,\boldsymbol{s}(t), \quad \forall t = 1, 2, \dots, T.$$
(4.1)

where t index refers to the observation variability parameter depending on the considered application. It can correspond for instance to time, frequency, wavelength, pixel index, etc. We will often use the matrix notation, merging the observations $\boldsymbol{x}(t)$ into the observation matrix $\boldsymbol{X} \in \mathbb{R}^{P \times T}$ and the source signals $\boldsymbol{s}(t)$ into the source matrix $\boldsymbol{S} \in \mathbb{R}^{R \times T}$. Moreover, we also consider that measured data are subject to measurement noise and errors. Therefore, given \boldsymbol{X} , the source separation aims at recovering the source matrix $\boldsymbol{S} \in \mathbb{R}^{R \times T}$ and the matrix of mixing coefficients $\boldsymbol{A} \in \mathbb{R}^{P \times R}$ such that

$$X = AS + E, \tag{4.2}$$

where the matrix E refers to the measurement noise. Moreover, we will only consider in the sequel the case where the matrices containing the source signals and the mixing coefficients are component-wise non-negative, that is, $A \succeq 0$ and $S \succeq 0$.

Finding non-negative matrices S and A allowing to reproduce the observation matrix X according to (4.2) is known as *non-negative matrix factorization* (NMF). Note that the smallest R such that such a decomposition of X exists is called the of X and is denoted rank₊{X}. Clearly, we have

 $\operatorname{rank}\{\boldsymbol{X}\} \leq \operatorname{rank}_{+}\{\boldsymbol{X}\} \leq \min(P,T).$

4.1.1

Brief historical overview

It is difficult to trace back the first time the model (4.1) with non-negativity constraints was introduced, as it is a rather natural model in many situations; see section 4.6 for the description of several applications. For example, Imbrie and Van Andel [1] used this model in 1960's for the analysis of mineral data. In the field of linear algebra, the first publications related to the mathematical formulation of the NMF problem concentrated their effort on the conditions for the existence and the uniqueness of such factorization [2, 3, 4, 5, 6]. The problem was named as non-negative rank factorization and defined as the factorization of a non-negative matrix into the product of two non-negative matrices. However, NMF in its current form (4.2) was introduced by Paatero and Tapper [7] and referred to as (PMF). In 1999, Lee and Seung [8] popularized NMF with a paper in Nature 'Learning the parts of objects by non-negative matrix factorization' where they applied it to the extraction of facial features in a set of facial images and to identify topics in a set of documents. Regarding the decomposition algorithms, pioneering contributions of Tauler, Kowalsi and Fleming [9], Paatero and Tapper [10, 11] proposed original algorithms to find an approximate factorization of a matrix in the case of spectroscopic data and noisy observations, by alternating non-negative least squares estimation in the former and penalized least squares estimation in the latter. They proposed a simple alternating optimization scheme (optimizing over A and S alternatively; see section 4.5) and applied it to air emission control. More recently, Lee and Seung [12] presented two algorithms based on multiplicative updates dedicated to non-negative matrix factorization (NMF); one for the Frobenius norm and one for the Kullback-Leibler divergence (see section 4.5). NMF and source separation with non-negativity constraint have since remained an active research topic using several factorization approaches and applications [13, 14, 15, 16].

4.2

Geometrical interpretation of NMF and the non-negative rank

Let us assume that there is no noise and that the linear mixture model is exact, that is, that each observation can be written as a non-negative linear combination of the R sources:

$$x_p(t) = \sum_{r=1}^R a_{pr} s_r(t) \quad (\forall t = 1, 2, \dots, T), (\forall p = 1, 2, \dots, P).$$
(4.3)

Given X, finding $A \succeq 0$ and $S \succeq 0$ such that X = AS is referred to as exact NMF. Note that, even in noisy settings, the geometrical interpretation that we will describe in this section for (4.3) is useful because the noisy observations are approximated by points satisfying this exact linear mixing model.

Let us get rid of the index t by denoting x_p the p-th row of X and s_r the r-th row of S. We have

$$x_p = \sum_{r=1}^{R} a_{pr} s_r \quad (\forall p = 1, \dots, P).$$
 (4.4)

Since the coefficients a_{pr} are non-negative, the rows of X belong to the convex cone generated by the rows of S, that is, by the set $\{s_1, s_2, \dots s_R\}$; see Figure 4.1 (left) for an illustration in the case R = 3.



Figure 4.1 Geometric illustration of exact NMF for T = R = 3 and P = 25. Both figures represent the same data set. The figure on the right is the normalization to unit ℓ_1 norm of the rows of X and S from the figure on the left.

Equivalently, we can scale the observations and the sources as follows:

$$\frac{\boldsymbol{x}_{p}}{||\boldsymbol{x}_{p}||_{1}} = \sum_{r=1}^{R} \underbrace{\left(a_{pr} \frac{||\boldsymbol{s}_{r}||_{1}}{||\boldsymbol{x}_{p}||_{1}}\right)}_{a'_{pr}} \underbrace{\frac{\boldsymbol{s}_{r}}{||\boldsymbol{s}_{r}||_{1}}}_{\boldsymbol{s}'_{r}}, \tag{4.5}$$

where $||\boldsymbol{y}||_1 = \sum_i |y_i|$. In this way, \boldsymbol{x}'_p and \boldsymbol{s}'_r belong to the unit simplex \mathcal{S} , defined as

$$\mathcal{S} = \left\{ \boldsymbol{y} \in \mathbb{R}^T \mid y_t \ge 0 \text{ for } 1 \le t \le T, \text{ and } \sum_{t=1}^T y_t = 1 \right\}.$$

Moreover, we must have $\sum_{r=1}^{R} a'_{pr} = 1$ for all p since $\mathbf{x}'_p = \sum_{r=1}^{R} a'_{pr} \mathbf{s}'_r \in S$ and $\mathbf{s}'_r \in S$. In other words, the vectors \mathbf{x}'_p $(1 \le p \le P)$ belong to the convex hull of the set $\{\mathbf{s}'_1, \mathbf{s}'_2, \dots, \mathbf{s}'_R\}$, defined as

$$\operatorname{conv}\{\boldsymbol{s}_1', \boldsymbol{s}_2', \dots \boldsymbol{s}_r'\} = \left\{ \sum_{r=1}^R a_{pr}' \boldsymbol{s}_r' \mid \boldsymbol{a} \in \mathcal{S} \right\} = \operatorname{conv}\{\boldsymbol{S}'\} = \{ \boldsymbol{a}'^T \boldsymbol{S}' \mid \boldsymbol{a} \in \mathcal{S} \},$$

see Figure 4.1 (right) for an illustration.

For simplicity, we assume in the remainder of this section that x_p and s_r are normalized to have unit ℓ_1 norm for all $1 \le p \le P$ and $1 \le r \le R$, that is, that the entries of x_p and s_r sum to one (note that we also assume without loss of generality that the observations x_p and the sources s_r are different from zero, otherwise we discard them). Hence, given a normalized $X \succeq 0$, finding $A \succeq 0$ and $S \succeq 0$ such that X = AS is equivalent to finding a set of sources $\{s_1, s_2, \ldots, s_r\}$ such that

 $\operatorname{conv} \{ oldsymbol{x}_1, oldsymbol{x}_2, \dots oldsymbol{x}_p \} \subseteq \operatorname{conv} \{ oldsymbol{s}_1, oldsymbol{s}_2, \dots oldsymbol{s}_r \} \subseteq \mathcal{S}.$

This is an instance of the so called *nested polytope problem* (NPP) in computational geometry;

see e.g. [17] and the references therein. NPP is defined as follows: given two nested polytopes, $\mathcal{P} \subset \mathcal{Q}$, the goal is to find, if possible, a set of R points $\{s_1, s_2, \ldots s_r\}$ such that its convex hull is nested between the two given polytopes \mathcal{P} and \mathcal{Q} , that is, such that

$$\mathcal{P} \subset \operatorname{conv}\{s_1, s_2, \dots s_r\} \subset \mathcal{Q},$$

see Figure 4.2 for an illustration in two dimensions where \mathcal{P} and \mathcal{Q} are squares, and $\operatorname{conv}\{s_1, s_2, s_3\}$ is a nested triangle. For our exact NMF problem above, we have $\mathcal{P} = \operatorname{conv}\{X\}$ and $\mathcal{Q} = \mathcal{S}$. Note that, the outer polytope \mathcal{Q} can have a higher dimension than the inner polytope \mathcal{P} . NPP is a very difficult geometric problem, being NP-hard already in dimension 3 [17], although a polynomial-time algorithm exists when the inner and outer polytopes have dimension two, that is, when they are polygons [18].

The one-to-one equivalence between the exact NMF and NPP was established in [19, 20, 21, 22]. In fact, it can also be shown that any NPP instance can be written as an exact NMF problem: given the outer polytope described with its T facets $\{y|a_i^{\mathsf{T}}y \leq b_i\} \ 1 \leq t \leq T$ and the P vertices $v_p \ 1 \leq p \leq P$ of the inner polytope, solving NPP is equivalent to solving exact NMF for the input matrix

$$X_{p,t} = b_i - \boldsymbol{a}_i^{\mathsf{T}} \boldsymbol{v}_p \ge 0 \quad \text{for all } 1 \le p \le P, 1 \le t \le T.$$

It is often assumed in practice that the rank of the input matrix X is equal to the number of sources R, that is, $R = \operatorname{rank}\{X\}$. In that case, the NPP problem corresponding to exact NMF can be simplified because the row space of X and Smust coincide, since X = AS and S has R rows. Therefore, the rows of S must be contained in the polytope $S \cap row(X)$, where row(X) denotes the row space of X. Hence we can restrict the outer polytope S to be $S \cap row(X)$, so that the dimensions of the inner and outer polytope coincide and are equal to $rank{X}-1$. Therefore, we are looking for a matrix S such that

$$\operatorname{conv}\{X\} \subseteq \operatorname{conv}\{S\} \subseteq S \cap \operatorname{row}(X). \tag{4.6}$$

We will refer to this variant of NPP as the restricted NPP, where the inner and outer polytopes have the same dimension (hence the nested polytope also has the same dimension). The one-to-one correspondence between exact NMF with $R = \operatorname{rank}\{X\}$ and restricted NPP was established in [19, 20]:

THEOREM 4.1 Given a normalized non-negative matrix $X \in \mathbb{R}^{P \times T}_+$, deciding whether rank $\{X\}$ = rank $_+\{X\}$ and computing a corresponding factorization X = AS with $A \in \mathbb{R}^{P \times R}_+$ and $S \in \mathbb{R}^{R \times T}_+$ is equivalent to deciding whether there exists a polytope with R vertices nested between conv $\{X\}$ and $S \cap row(X)$.

The theorem above can be generalized to check whether the $\operatorname{rank}_{+}\{X\} = \operatorname{rank}\{X\} + 1$. In fact, when $R = \operatorname{rank}\{X\} + 1$, it can be shown that the nested polytope can be assumed, without loss of generality, to have the same dimension as the inner polytope. In other words, it can be assumed without loss of generality that $\operatorname{rank}\{S\} = \operatorname{rank}\{X\}$. Therefore, checking whether $\operatorname{rank}_{+}\{X\} = \operatorname{rank}\{X\} + 1$ is equivalent to checking whether there exists a polytope with R + 1 vertices nested between $\operatorname{conv}\{X\}$ and $S \cap \operatorname{row}(X)$ [21, Corollary 2]. However, if $\operatorname{rank}_{+}\{X\} > \operatorname{rank}\{X\} + 1$, then in general $\operatorname{rank}\{S\} > \operatorname{rank}\{X\}$ and it can no longer be assumed that the nested and inner polytopes have the same dimension, that is, we will have that $\operatorname{conv}\{S\} \not\subset S \cap \operatorname{row}(X)$ [21] (see example 2 below).

We can use the equivalence between NPP and the computation of the non-negative rank described above to derive possible values of the non-negative rank of a non-negative matrix $X \in \mathbb{R}^{P \times T}$ in simple cases:

- For rank{X} = 1, it is clear that rank₊{X} = 1 since all rows of X are multiple of the same vector. Geometrically, the inner polytope corresponding to the NPP instance is a single point and the NPP instance is trivial.
- For rank{X} = 2, the inner polytope is a line segment. Therefore, the problem can easily be solved by identifying the two extreme points of that line segment hence rank₊{X} = 2 [3].
- For min(P,T) = rank{X}, we have rank₊{X} = rank{X} using the trivial decompositions X = XI = IX where I is the identity matrix of appropriate dimensions. Geometrically, this means that we either take the nested polytope conv{S} as the inner polytope conv{X} (for X = XI where S = X) or as the outer polytope S (for X = IX where S = I).

From the three results above, we have the following theorem from [?].

THEOREM 4.2 Let $X \in \mathbb{R}^{P \times T}_+$. If rank $\{X\} \leq 2 \text{ or } \min(P, T) = \operatorname{rank}\{X\}$, then rank $\{X\} = \operatorname{rank}\{X\}$.

When $\min(P,T) \ge 4$ and $\operatorname{rank}\{X\} \ge 3$, determining the non-negative rank becomes more difficult.

Let us analyze higher dimensional cases.

For rank{X} = 3 and min(P,T) ≥ 4, the restricted NPP instance is the problem to finding a polygon nested between two given polygons. As mentioned above, this problem can be solved in polynomial time, more precisely in O((P+T) log(min(P,T))) operations [18]. Therefore, when rank{X} = 3, one can decide in polynomial time whether rank₊{X} = 3. Moreover, because it does not help to try to find a higher dimensional nested polytope with only four vertices (this is the case R = rank{X} + 1; see above), this algorithm can also be used to decide in polynomial time whether rank₊{X} = 4.

Deciding whether rank₊{X} = 5 becomes more difficult (unless min(P, T) = 5) because the nested polytope might live in a higher dimension than conv{X}; see the second example below. In fact, it is important to realize that, when rank{X} \geq 3, the non-negative rank of X can be arbitrarily large. For example, for the matrix corresponding to the restricted NPP instances of the regular P-gon nested with itself (for which T = P and conv{X} = $S \cap row(X)$), we have [23]

 $\operatorname{rank}_{+}\{\boldsymbol{X}\} \geq \left\lceil \log_2(2P+2) \right\rceil,$

while rank $\{X\} = 3$; see also [24] the references therein for more details.

• For $R = \operatorname{rank}\{X\} = 4$ and $\min(P, T) \ge 5$, it is difficult to compute the non-negative rank [20] since three-dimensional restricted NPP instances are NP-hard [17] (note that this result does not fix a priori the number of vertices of the nested polytope).

However, checking whether rank $\{X\} = R$ for fixed R (that is, R is not part of the input) can be done in polynomial time in P and T requiring $\mathcal{O}((PT)^{R^2})$ operations [25]. This approach, although theoretically appealing, is not very useful in practice as it requires to solving systems of equations via quantifier elimination theory, and we were not able to identify a software able to solve problems already for P = T = 4 and R = 3, the first non-trivial case (see above).

We now illustrate the results of this section on two interesting examples.

Example 1: Nested squares [19]. We consider the smallest possible case where $\operatorname{rank}_{+}{X} > \operatorname{rank}{X}$: this requires that P and T are at least 4 (see above). Consider the matrix

$$\mathbf{X} = \frac{1}{4} \begin{pmatrix} 1+a & 1-a & 1+a & 1-a \\ 1-a & 1+a & 1+a & 1-a \\ 1-a & 1+a & 1-a & 1+a \\ 1+a & 1-a & 1-a & 1+a \end{pmatrix},$$
(4.7)

where $0 < a \le 1$ and rank $\{X\} = 3$. The restricted NPP instance corresponding to the exact NMF problem for X (Theorem 4.3) are two nested squares; see Figure 4.2.

The following can be shown [26]:



Figure 4.2 Restricted NPP instance of two nested squares corresponding to the exact NMF problem for X from (4.7) for $a = \sqrt{2} - 1$ and a = 1/4 [26]. The two triangles correspond to two exact NMF's for $0 < a \le \sqrt{2} - 1$.

For 0 < a < √2 − 1, the inner square is small enough so that there exists infinitely many triangles in between the two nested squares. This implies that rank₊{X} = 3. This also implies that the exact NMF of X is (highly) non-unique; see section 4.4 for more details.

(Note that for a = 0, the inner 'square' is a single point and rank $\{X\} = 1$.)

- For $a = \sqrt{2} 1$, there exists 8 different triangles nested between the two squares hence rank₊{**X**} = 3 (leading to 8 different exact factorizations, up to permutation and scaling, not infinitely many as above); see Figure 4.2 where two such triangles are represented (the other six solutions are rotations of these two).
- For a > √2 − 1, there does not exist any triangle between the two nested squares (note that, for a = 1, the two squares coincide) hence rank₊{X} = 4.

Example 2: Regular hexagon. A popular example is the non-negative matrix corresponding to the restricted NPP instance of the regular hexagon nested with itself. It is the smallest nontrivial case for which rank $\{S\} > \operatorname{rank}\{X\}$ is necessary to find the exact NMF of minimum rank when rank $\{X\} = 3$. In fact, for rank $\{S\} > \operatorname{rank}\{X\}$ to be necessary when rank $\{X\} = 3$, we need that rank $\{X\}$ is at least five (see above) hence we need $\min(P,T) \ge 6$ to have a nontrivial factorization.

Let us consider the following non-negative matrix

. ...

~

$$\boldsymbol{X} = \begin{pmatrix} 0 & 1 & 4 & 9 & 16 & 25 \\ 1 & 0 & 1 & 4 & 9 & 16 \\ 4 & 1 & 0 & 1 & 4 & 9 \\ 9 & 4 & 1 & 0 & 1 & 4 \\ 16 & 9 & 4 & 1 & 0 & 1 \\ 25 & 16 & 9 & 4 & 1 & 0 \end{pmatrix}$$

$$= \begin{pmatrix} 0 & 0 & 4 & 5 & 1 \\ 1 & 0 & 1 & 3 & 0 \\ 4 & 0 & 0 & 1 & 1 \\ 4 & 1 & 0 & 0 & 1 \\ 1 & 3 & 1 & 0 & 0 \\ 0 & 5 & 4 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 1 \\ 5 & 3 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 3 & 5 \\ 0 & 1 & 0 & 0 & 1 & 0 \end{pmatrix} = \boldsymbol{AS}, \quad (4.8)$$

~~ .

for which

$$\operatorname{rank}{X} = 3 < \operatorname{rank}{S} = 4 < \operatorname{rank}{X} = 5 < \min(P, T) = 6.$$

The restricted NPP instance (4.6) corresponding to X are two hexagons that coincide, with conv $\{X\} = S \cap row(X)$. Clearly, it is not possible to find a polygon in between the hexagon and itself with less than 6 vertices hence (i) rank₊ $\{X\} > 4$ (see above) and (ii) it is not possible to find an exact NMF with R = 5 and rank $\{S\} = 3$. However, there exists a three dimensional polytope conv $\{S\}$ in S with five vertices that contains conv $\{X\}$ hence rank₊ $\{X\} = 5$ and, in any decomposition of rank 5, it is required that rank $\{S\} > rank\{X\}$; see Figure 4.3 for an illustration.



Figure 4.3 Illustration of the polytope $\operatorname{conv}\{S\}$ with 5 vertices containing the hexagon $\operatorname{conv}\{X\} = S \cap \operatorname{row}(X)$ in the unit simplex; see Equation (4.8). Note that the outer polytope S is not shown here since it has dimension 5.

4.3 Statistical formulation of NMF

The mixing model (4.1) assumes that P observations $\{x_p(t), t = 1, ..., T\}_{p=1}^{P}$ are linear instantaneous combinations of R unknown source signals $\{s_r(t), t = 1, ..., T\}_{r=1}^{R}$. In the noisy case, this mixing model is expressed as

$$x_p(t) = \sum_{r=1}^R a_{pr} s_r(t) + e_p(t),$$
(4.9)

where the additive noise terms $\{e_p(t), t = 1, ..., T\}_{p=1}^{P}$ represent the measurement errors and the model errors.

4.3.1 Case of a Gaussian noise

The statistical distribution of each noise term $e_p(t)$ is assumed to be Gaussian with a zero mean and variance σ_p^2 . More generally, the distribution of the noise vector e(t) is represented by a multivariate Gaussian distribution with zero mean vector and a covariance matrix Σ_t . This matrix will be diagonal in the case of mutually independent noise components and its diagonal terms are equal to $[\sigma_1^2, \ldots, \sigma_P^2]$. Alternatively, one can consider a sample-dependent noise variance but the same value for all the observations

$$f(\boldsymbol{e}_t | \sigma_t) = \mathcal{N}(\boldsymbol{0}, \sigma_t^2 \boldsymbol{I}_P),$$

where I_P denotes the identity matrix and $\mathcal{N}(\mu, \Sigma)$ stands for the Gaussian distribution with mean μ and covariance Σ . In addition, by assuming that samples of each noise component are independent and identically distributed, one can write

$$f(\boldsymbol{E}|\boldsymbol{\Sigma}) = \prod_{t=1}^{T} \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Sigma}).$$

Consequently the likelihood can be constructed as

$$f(\boldsymbol{X}|\boldsymbol{S}, \boldsymbol{A}, \boldsymbol{\Sigma}_t) = \prod_{t=1}^T \mathcal{N}(\boldsymbol{x}(t) - \boldsymbol{A}\boldsymbol{s}(t), \boldsymbol{\Sigma}_t).$$
(4.10)

By taking its negative logarithm, this likelihood leads to a data fitting term

$$\mathcal{Q}(\boldsymbol{S}, \boldsymbol{A}) = \sum_{t=1}^{T} \left(\boldsymbol{x}(t) - \boldsymbol{A}\boldsymbol{s}(t) \right)^{t} \boldsymbol{\Sigma}_{t}^{-1} \left(\boldsymbol{x}(t) - \boldsymbol{A}\boldsymbol{s}(t) \right).$$
(4.11)

This objective function corresponds to a weighted least squares criterion. However, when $\Sigma_t = \sigma^2 I$, this criterion simplifies to the quadratic data fitting objective function used in most NMF algorithms [10, 8]; see Section 4.5.

4.3.2 Case of Poissonian noise

Poissonian noise model is more adequate in the case where measurements in $x_p(t)$ are obtained through a counting process. In this model, each observed data sample $x_p(t) = [\mathbf{X}]_{p,t}$ is assumed to be a realization of a Poisson process of mean $[\mathbf{AS}]_{p,t}$. The likelihood can therefore be expressed as:

$$f(x_p(t)|\mathbf{A}, \mathbf{S}) = \frac{([\mathbf{A}\mathbf{S}]_{p,t})^{x_p(t)}}{x_p(t)!} \exp(-[\mathbf{A}\mathbf{S}]_{p,t}).$$
(4.12)

Under the assumption of mutually independent observations and identically distributed samples, the likelihood f(X|S, A) can be deduced:

$$f(\boldsymbol{X}|\boldsymbol{S},\boldsymbol{A}) = \prod_{p=1}^{P} \prod_{t=1}^{T} f(x_p(t)|\boldsymbol{S},\boldsymbol{A}).$$

It can be noted that the data fitting term, $Q(S, A) = -\log f(X|S, A)$, resulting from this likelihood is

$$\mathcal{Q}(\boldsymbol{S}, \boldsymbol{A}) = \sum_{p=1}^{P} \sum_{t=1}^{T} \left([\boldsymbol{A}\boldsymbol{S}]_{p,t} - x_p(t) \log[\boldsymbol{A}\boldsymbol{S}]_{p,t} \right).$$
(4.13)

This criterion is an instance of those used in non-negative matrix factorization algorithms based on divergence measures such as Kulback-Leibler [8, 27]. Actually, the source separation approach using the maximum likelihood approach, whose principle is the maximization of f(X|S, A) or equivalently the minimization of $-\log f(X|S, A)$ allows to give a statistical formulation of NMF algorithms. More generally, adding statistical priors on matrices X and S can also formalized using Bayesian source separation methods.

4.4 Uniqueness and admissible solutions of NMF

Before trying to effectively solve any NMF problem, a key point is to answer some questions related to the model indeterminacy and to the uniqueness of the solution. Let us assume the existence of a non-negative factorization of the data matrix X into matrices A and S. Let us start with any pair (A, S) that fulfill the mixing model (4.2) and then introduce a non-singular $(p \times p)$ matrix R. A new pair (\tilde{A}, \tilde{S}) can be defined by

$$\tilde{\boldsymbol{A}} = \boldsymbol{A}\boldsymbol{R}^{-1} \text{ and } \tilde{\boldsymbol{S}} = \boldsymbol{R}\boldsymbol{S}, \tag{4.14}$$

with no modification of the recovered data matrix, i.e. $X = \tilde{A} \tilde{S}$. In the unconstrained case, this well known result shows the existence of an infinite number of exact factorizations of the matrix X. Matrix R is sometimes called *rotational ambiguity* matrix which includes classical scaling and ordering indeterminacies (See Chapter ?? of this book). In the case of NMF, a possible linear transformation should lead to transformed matrices \tilde{A} and \tilde{S} satisfying the non-negativity constraints

 $\tilde{A} \ge 0 \text{ and } \tilde{S} \ge 0.$ (4.15)

In that respect, three questions arise:

- what are the conditions on the actual source signals and mixing coefficients (factors S and A) ensuring the uniqueness of the factorization of X according to (4.2), and satisfying the non-negativity constraints (4.15)?
- 2) if the decomposition is not unique, what are the admissible (feasible) solutions?
- 3) among all the admissible solutions, can we define a more plausible one?

These questions highlight three aspects of the NMF problem that will be detailed in the sequel, with a special focus on a practical procedure for obtaining a set of admissible NMF solutions.

4.4.1 Uniqueness conditions

Many necessary and/or sufficient conditions for establishing the have been formulated in the literature. In what follows, we will suppose that :

$$\operatorname{rank}\{\boldsymbol{X}\} = \operatorname{rank}_{+}\{\boldsymbol{X}\} = R.$$

The first theorem on uniqueness was proposed by Chen [5]. It gives a necessary and sufficient condition to have a unique NMF. But, this condition does not give any numerical mean to check if a given non-negative matrix admits a unique nonnegative factorization. In addition, Park et *al.* [28] and Smilde et *al.*, [29] developed some sufficient uniqueness conditions well adapted to some specific applications but they seem not to be applicable for a general purpose.

First, we give a necessary condition [30], to have the NMF uniqueness. It states that both S and A should have a minimum number of zero entries.

THEOREM 4.3 If the NMF of X = AS is unique, then the following condition are fulfilled: (A1) $\forall (r \neq r'), \exists k \text{ such as } : s_r(k) = 0 \text{ and } s_{r'}(k) > 0.$

(A2) $\forall (r \neq r'), \exists \ell \text{ such as }: a_{\ell r} = 0 \text{ and } a_{\ell r'} > 0.$

Chen's uniqueness results is the starting point of Donoho et *al*. [31] which gives a sufficient uniqueness condition :

THEOREM 4.4 The NMF of X = AS is unique if the following conditions are satisfied:

• Separability: $\forall r, \exists k \text{ such that } : s_r(k) \neq 0 \text{ and } s_r(\ell) = 0, \forall \ell \neq k$

- Generative model: the set $\{1, \dots, P\}$ is partitioned into L groups $\mathcal{P}_1, \dots, \mathcal{P}_L$, each containing exactly R elements. $\forall p, \forall \ell$, there exists an element a_{pr} such that: $a_{pr} \neq 0, r \in \mathcal{P}_\ell$ and $a_{p\ell} = 0, \forall \ell \in \mathcal{P}_\ell, \ell \neq r$
- Complete Factorial Sampling: $\forall r_1 \in \mathcal{P}_1, \cdots, r_L \in \mathcal{P}_L, \exists k \text{ such that: } a_{kr_1} \neq 0, \cdots, a_{kr_\ell} \neq 0$

The work of Laurberg et *al.* [32] is starting from a different point of view which was initially proposed by [3] and proves the following sufficient uniqueness condition

THEOREM 4.5 The NMF of X = AS is unique if the following conditions are satisfied:

- Sufficiently spread: $\forall r, \exists k \text{ such that } : s_r(k) \neq 0 \text{ and } s_r(\ell) = 0, \forall \ell \neq k$
- Strongly Boundary Close: the matrix A satisfies the following conditions

 ∀r, ∃ℓ such that: a_{ℓr} = 0 and a_{ℓr'} ≠ 0, ∀r' ≠ r
 There exists a permutation matrix P such that ∀r, there exists a set k₁, ..., k_{R-k} satisfying [AP]_{r,kj} = 0, ∀j ≤ R k; and the matrix [AP]_{r+1:R,k1:kR-r} is invertible.

A recent and complete survey on the analysis of can be found in [33] where uniqueness conditions are also formulated for the case of symmetric NMF; see also Chapter ??. In fact, although sparsity of the input matrix is neither a necessary nor a sufficient condition for uniqueness, there is a link between uniqueness of NMF and sparsity of the latent variable; see for example Theorem 4.3. It is observed in practice that if the true latent factors are sparse, NMF usually tends to recover the correct solution, the geometric interpretation of NMF shows that sparser matrices lead to more well-posed NMF problems. However, in many applications including multivariate curve resolution and spectral data unmixing, there is at least one factor that is nonsparse. This motivates the development of approaches allowing to assess the extent of the possible results termed as solutions.

4.4.2

Finding the admissible solutions

The retained approach for finding the admissible solutions is based on the same idea as the one previously proposed in [34, 30]. It consists in finding a set of parametric transformation matrices $T(\theta)$ minimizing a criterion $C_{nneg}(\theta)$ based on a non-negativity measure. This criterion is defined as

$$\mathcal{C}_{\text{nneg}}(\boldsymbol{\theta}) = \|f(\boldsymbol{T}(\boldsymbol{\theta})\boldsymbol{S})\|_{F}^{2} + \|f(\boldsymbol{A}\,\boldsymbol{T}^{-1}(\boldsymbol{\theta}))\|_{F}^{2}.$$
(4.16)

where $f(x) = \min(x, 0)$. Parameters θ are defined so as to implicitly handle the scaling ambiguity. Such objective function is generally used in constrained optimization as an exterior penalty function since it assigns a high cost to solutions that do not fulfill the non-negativity constraint. The minimization of this criterion with respect to the parameter vector θ is performed using an unconstrained optimization

method such as the Neldar-Mead simplex algorithm with several different random initializations of the transformation matrix parameters θ . The retained solutions are those that cancel C_{nneg} , i.e. that exactly solve the constrained factorization problem. The experiment is repeated with different starting points to get several admissible values of the transformation matrix parameters.

Illustration in the case of two sources. For a decomposition rank R = 2, it was shown in [30] that the transformation matrix $T(\theta)$ reduces to the following parametric form:

$$\boldsymbol{T}(\boldsymbol{\theta}) = \begin{bmatrix} 1 - \theta_1 & \theta_1 \\ \theta_2 & 1 - \theta_2 \end{bmatrix},$$
(4.17)

with parameter vector such that $(\theta_1 + \theta_2) < 1$, to get rid of the ordering indeterminacy and ensuring invertibility of T. Analytical calculations detailed in [30] can be performed in the case of two sources and lead to bounds on the values of the parameters

$$\begin{cases} \theta_{1} \in \left[-\min_{t \in \mathbb{T}_{1}} \left\{\frac{s_{1}(t)}{s_{2}(t) - s_{1}(t)}\right\}, & \min_{\ell} \left\{\frac{a_{\ell 2}}{a_{\ell 1} + a_{\ell 2}}\right\}\right], \\ \theta_{2} \in \left[-\min_{t \in \mathbb{T}_{2}} \left\{\frac{s_{2}(t)}{s_{1}(t) - s_{2}(t)}\right\}, & \min_{\ell} \left\{\frac{a_{\ell 1}}{a_{\ell 1} + a_{\ell 2}}\right\}\right] \end{cases}$$
(4.18)

with $\mathbb{T}_1 = \{t \in \{1, ..., T\}; s_2(t) > s_1(t)\}$ and $\mathbb{T}_2 = \{t \in \{1, ..., T\}; s_1(t) > s_2(t)\}$. Figure 4.4 gives an illustration of a set of admissible solutions in the case of two spectral sources. This example is inspired from measurement data that can be obtained from the monitoring of kinetic chemical reactions using spectroscopy [35]. One can see here that acceptable NMF solutions deviate from the original sources and presents additional peaks, which in practice may lead to data interpretation errors when the separation is performed under the non-negativity constraint alone.

Illustration in the case of more than two sources. In the case of more than two sources, analytical computations cannot be performed to get the feasible values of the transformation matrix parameters. For instance, in the case of three sources, a possible transformation matrix T bypassing permutation and scaling ambiguities [35] takes the form

$$\boldsymbol{T}(\boldsymbol{\theta}) = \begin{bmatrix} 1 - \theta_1 - \theta_2 & \theta_1 & \theta_2 \\ \theta_3 & 1 - \theta_3 - \theta_4 & \theta_4 \\ \theta_5 & \theta_6 & 1 - \theta_5 - \theta_6 \end{bmatrix},$$
(4.19)

A numerical optimization with several starting points allows to get some realizations of the admissible transformation matrices. Figure 4.5, gives an example of the parameters values in the case of three spectral sources. Once again, one can see the existence of several values of the transformation parameters leading to several admissible solutions. Actually, the feasible NMF solution are mixtures of the original sources. We refer the reader to [36] for a recent paper on this topic.



Figure 4.4 Illustration of feasible NMF solutions in the case of a spectral mixture of two components S1 and S2. The NMF solutions are mixtures of the actual sources.



Figure 4.5 Illustration of feasible NMF solutions in the case of a three component mixture: a-d) simulated sources and mixing coefficients, e-g) joint distribution of the the transformation parameters of the actual sources leading to feasible NMF solutions.

Non-negative matrix factorization algorithms

In this section, we present several widely used algorithms to solve NMF problems. We will mostly focus on the least squares formulation

$$\min_{\boldsymbol{A},\boldsymbol{S}} ||\boldsymbol{X} - \boldsymbol{A}\boldsymbol{S}||_F^2 = \sum_{i,j} (\boldsymbol{X} - \boldsymbol{A}\boldsymbol{S})_{ij}^2 \quad \text{such that} \quad \boldsymbol{A} \succeq \text{ and } \boldsymbol{S} \succeq 0.$$
(4.20)

The Frobenius norm is suitable when assuming Gaussian noise; for example, if X is a dense matrix representing images [8]. However, this is not always the case, in particular for sparse input matrices X which are often encountered in the literature (such as document data sets or matrices arising from large networks). In these cases, other objective functions should be used; we refer the reader to [37, 38] and the references therein.

The problem (4.20) is a difficult non-convex optimization problem and, in fact, is NP-hard [20]. In practice, it is in general solved via standard iterative nonlinear optimization approaches; some of them are described in the next section 4.5.1. Although these approaches do not guarantee to obtaining an optimal solution, they usually generate satisfactory results that are useful for applications. More recently, a new class of NMF methods were introduced that are guaranteed to recover an optimal solution, up to the noise level, given that the input matrix X has a particular structure; this is briefly discussed in section 4.2, and in much more details in Chapter ??.

To add to the complexity, as discussed in the previous section, the optimal solution of (4.20) is in general non-unique (even by considering equivalent solutions due to permutation and scaling); see section 4.4. To alleviate this problem, a standard approach is to incorporate priors, usually via regularizers or additional constraints, into (4.20); see sections 4.5.2 and Chapter **??**. Another problem which we do not discuss in this chapter and which is inherent to most blind source separation problems is the choice of the number R of sources. In the sequel we discuss the NMF algorithm in the case for a fixed value of R.

4.5.1

Iterative factorization methods

Although (4.20) is a difficult nonlinear optimization problem, it has several nice properties. In particular, if we assume a known mixing matrix A, then the problem becomes a convex optimization problem with respect to S,

$$\min_{S \succeq 0} ||\boldsymbol{X} - \boldsymbol{AS}||_F^2. \tag{4.21}$$

This is a particular convex quadratic optimization problem with linear inequality constraints referred to as *non-negative least squares* () [39, 40]. Clearly, the same property holds for A when S is fixed (this structure is sometimes referred to as biconvex). Most iterative methods take advantage of this fact: denoting $(A^{(k)}, S^{(k)})$

4.5

the solution obtained after k iterations (that is, the kth iterate), most methods obey to the following framework:

- 1. Initialize $(A^{(0)}, S^{(0)})$.
- 2. For $k = 1, 2, \ldots$,
- 2.a Compute $S^{(k)}$ such that $||X A^{(k-1)}S^{(k)}||_F^2 \le ||X A^{(k-1)}S^{(k-1)}||_F^2$. 2.b Compute $A^{(k)}$ such that $||X - A^{(k)}S^{(k)}||_F^2 \le ||X - A^{(k-1)}S^{(k)}||_F^2$.

The iterative process is usually stopped according to standard convergence rules, e.g., stabilization of the objective function and/or the iterates; see Chapter **??**.

Before we present some numerical algorithms, it is interesting to note the following issues:

- Because of the symmetry of the problem, since $X \approx AS \iff X^{\mathsf{T}} \approx (AS)^{\mathsf{T}} = S^{\mathsf{T}}A^{\mathsf{T}}$, A and S are updated in the same way for most NMF algorithms (that is, steps 2.a and 2.b above use the same strategy).
- The product AS is the sum of R rank-one factors $a_r s_r$, where a_r is the Rth column of A and s_r the Rth row of S. Therefore, there is always a scaling and permutation degree of freedom (see also section 4.4) since we can permute indistinguishably the rank-one factors and since

$$\boldsymbol{a}_r \boldsymbol{s}_r = (\alpha_r \boldsymbol{a}_r) \left(\frac{1}{\alpha_r} \boldsymbol{s}_r \right)$$

for any $\alpha_r > 0$. Therefore, in practice, the columns of A (or rows of S) are usually normalized to have unit ℓ_2 or ℓ_1 norm.

4.5.1.1 Initializing NMF algorithms

The initial matrices $(\mathbf{A}^{(0)}, \mathbf{S}^{(0)})$ can be selected in many different ways. The most naive approach is to initialize them randomly, using for example the uniform distribution in the interval [0, 1] for each entry of $\mathbf{A}^{(0)}$ and $\mathbf{S}^{(0)}$. This approach is of course very simple and easy to implement, but has the drawback to ignore completely the input matrix \mathbf{X} . In particular, using such a procedure usually leads to a low-rank approximation $\mathbf{A}^{(0)}\mathbf{S}^{(0)}$ which can be rather far from \mathbf{X} . To improve the initial iterate, it is recommended to scale it, that is, to multiply it by a constant $\hat{\alpha}$ such that [41]

$$\hat{\alpha} = \operatorname{argmin}_{\alpha > 0} || \boldsymbol{X} - \alpha \boldsymbol{A}^{(0)} \boldsymbol{S}^{(0)} ||_{F}^{2} = \frac{\langle \boldsymbol{X}, \boldsymbol{A}^{(0)} \boldsymbol{S}^{(0)} \rangle}{\langle \boldsymbol{A}^{(0)} \boldsymbol{S}^{(0)}, \boldsymbol{A}^{(0)} \boldsymbol{S}^{(0)} \rangle} \\ = \frac{\langle \boldsymbol{X} \boldsymbol{S}^{(0)^{\mathsf{T}}}, \boldsymbol{A}^{(0)} \rangle}{\langle \boldsymbol{A}^{(0)^{\mathsf{T}}} \boldsymbol{A}^{(0)}, \boldsymbol{S}^{(0)} \boldsymbol{S}^{(0)^{\mathsf{T}}} \rangle},$$

where argmin denotes the global minimum of an optimization problem.

There exists many more sophisticated initialization approaches for NMF. Quite naturally, the goal of this initialization is to locate a good initial point close to a reasonable factorization in order to (i) avoid bad local minima and (ii) allow the NMF algorithms to converge faster. We list here a few standard approaches:

- The most commonly encountered idea is to use clustering algorithms, such as kmeans or spherical k-means [42]. They are used to initialize the rows of $S^{(0)}$ (using the cluster centroids), while $A^{(0)}$ is obtained either by using the cluster assignment matrix of by solving the corresponding NNLS subproblem; see, e.g., [43, 44] and the references therein.
- · A computationally more expensive but, in general, more effective method is to use the best unconstrained rank-R approximation of matrix X (that can be computed efficiently via the singular value decomposition). Of course, this approximation does not generate non-negative factors and the trick is to (somehow) project them back onto the non-negative orthant [45, 46].
- A procedure that is cheap and effective is to initialize $m{S}^{(0)}$ by selecting a representative subset of the observations, that is, rows of the input matrix X. Selecting this subset can be done in many different ways and many algorithms exist for doing so. This is closely related to the column subset selection problem and to separable NMF which is discussed in section 4.2.

The choice of the initialization procedure depends on the application at hand: which initialization seems to work best? which seems to make more sense? which has a reasonable computational time (this depends on the dimension of the data)?

4.5.1.2 Alternating non-negative least squares, an exact coordinate descent method with 2 blocks of variables

The first algorithm for NMF proposed by Paatero and Tapper in their original paper [7] is based on alternating regression, which is now in general referred to as (ANLS). It is a class of methods that solves the NNLS subproblems exactly, alternatively for A and S:

1. Initialize
$$(A^{(0)}, S^{(0)})$$
.
2. For $k = 1, 2, ...,$
2.a Compute $S^{(k)} = \operatorname{argmin}_{Y \succeq 0} ||X - A^{(k-1)}Y||_F^2$.
2.b Compute $A^{(k)} = \operatorname{argmin}_{Z \succeq 0} ||X - ZS^{(k)}||_F^2$.

ANLS is a so called exact two-block coordinate descent method: there are two blocks of variables (A and S) that are alternatively optimized exactly. This method is guaranteed to converge to a stationary point of (4.20) [47]. ANLS methods differ in the way the NNLS subproblems are solved. In fact, any method from (convex) optimization can potentially be used. Here is a possible classification of methods that can be used to solve NNLS:

- First-order methods such as standard projected gradient [48], (optimal) fast gradient methods [49], and coordinate descent methods (see section 4.5.1.5). These methods only use the first-order information, that is, the value and the gradient of the objective function at each iterate.
- · Higher-order methods such as interior point methods, Quasi-Newton or Newton methods [50]. These method have a faster local convergence rate but each iteration is more expensive.

• Active-set methods that take advantage of the fact that, if one would know the position of the zero entries in A and S, then the NNLS subproblems reduce to unconstrained least-squares problems. The set of zero entries (the active set) is updated in a clever way to guarantee the objective function to decrease at each step and the algorithm to converge [39, 51, 52, 53]. This class of algorithms performs well in practice, but the worst-case complexity is exponential; namely, proportional to the number of active sets which is proportional to 2^R (as the simplex method for linear programming).

We refer the reader to [54] for a survey on NNLS methods.

4.5.1.3 Multiplicative updates

The most popular approach to attack (4.20) is the multiplicative updates (MU) introduced along with the paper of Lee and Seung [8, 55] that really launched the research on NMF. It updates alternatively over A and S using the following updates (dropping the iteration index k for convenience):

$$A \leftarrow A \boxdot \frac{[XS^{\mathsf{T}}]}{[ASS^{\mathsf{T}}]}$$
 and $S \leftarrow S \boxdot \frac{[A^{\mathsf{T}}X]}{[A^{\mathsf{T}}AS]}$, (4.22)

where \Box denotes the component-wise multiplication. Note that similar updates exists for many other objective functions [55, 37]. Note also that the MU were originally introduced in [56] for updating only one factor in order to solve NNLS problems. The MU (4.22) are guaranteed to decrease the objective function of (4.20) while clearly preserving non-negativity of the iterates.

It is interesting to note that the MU can be interpreted as a scaled gradient descent method, that is, a gradient descent method where each entry of the gradient is multiplied by a non-negative constant (this is equivalent to a quasi-Newton method where one would only use a diagonal matrix as an approximation of the Hessian); see [57, Section 1.3.2]. In fact, (4.22) can be equivalently written as

$$A \leftarrow A - \frac{[A]}{[ASS^{\mathsf{T}}]} \boxdot (ASS^{\mathsf{T}} - XS^{\mathsf{T}})$$
 (4.23)

and similarly for S, by symmetry. The term $ASS^{\mathsf{T}} - XS^{\mathsf{T}}$ is the gradient of $\frac{1}{2} ||X - AS||_F^2$ with respect to A.

It is important to note that, when implementing the MU, the term ASS^{T} is computed by first performing the matrix-matrix product SS^{T} . In fact, in that case, the computational cost will be in $\mathcal{O}(TR^2 + PR^2)$ operations while computing AS first requires $\mathcal{O}(PTR)$ operations and $\mathcal{O}(PT)$ space in memory. This is particularly crucial for large and sparse matrices since AS could be dense and could be too large to store in memory.

The main advantage of the MU is the ease of its implementation. However, it has several drawbacks. The main one being that it usually converges rather slowly compared to most other approaches. Another drawback is that the MU are not guaranteed to converge to a stationary point. The main reason being that once an entry is fixed to zero, it can no longer be modified (this is the so called locking phenomenon). However, this can be overcome by unlocking variables at zero using a proper procedure; see the discussion in [38].

It is interesting to note that updating A several times before updating S (and similarly for S) allows a much faster convergence, because the matrix-matrix products XS^{T} and SS^{T} do not need to be recomputed between updates of A when S is fixed [58]. Note also that the multiplicative updates can be used in an ANLS framework, where A and S would be updated until convergence for the NNLS sub-problems (see [59] for a convergence analysis when using the MU for NNLS).

4.5.1.4 Alternating least squares

A naive approach to solve NMF, referred to as alternating least squares (ALS), takes advantage of the fact that unconstrained least-squares problems can be solved very efficiently. To update A (and similarly for S), it first solves the unconstrained least squares problem

$$A \leftarrow \operatorname{argmin}_{Y \in \mathbb{R}^{P \times R}} ||X - YS||_F^2,$$

and then projects the solution back onto the non-negative orthant

 $A \leftarrow \max(\mathbf{0}, A).$

ALS is very easy to implement (e.g., in Matlab, it requires one line of code to update A, namely $A = \max(0, (X*S') / (S*S'))$). However, it is not guaranteed to converge and the objective function of (4.20) usually oscillates under the ALS updates, sometimes drastically. However, it is usually efficient to use ALS as an initialization step, before a convergent algorithm is used, especially for sparse input matrices; see, e.g., [37, 60].

4.5.1.5 Exact coordinate descent method with 2R blocks of variables

A method that works very well in many situations is the so called (HALS) method. It is a block coordinate descent method such as ANLS but has more block of variables. HALS optimizes alternatively over the columns of A and the rows of S, hence there are R blocks of P variables (the columns of A) and R blocks of T variables (the rows of S). HALS, although first suggested in [40], was first implemented an analyzed in [61] (and later in [62]) and independently in [41, 63, 64].

The benefit of considering smaller blocks of variables is that the optimization subproblem for each block is much easier to solve. In fact, the optimal solution for each r-th column of A and each r-th row of S can be written in closed form. Let us derive the formula here for the r-th column of A (again, by symmetry, the same holds for the r-th row of S). Fixing all variables but the R-th row of A, we need to solve the following optimization problem

$$\min_{\boldsymbol{a}_{r} \succeq 0} \left\| \boldsymbol{X} - \boldsymbol{A} \boldsymbol{S} \right\|_{F}^{2} = \min_{\boldsymbol{a}_{r} \succeq 0} \left\| \boldsymbol{X} - \sum_{k \neq r} \boldsymbol{a}_{k} \boldsymbol{s}_{k} - \boldsymbol{a}_{r} \boldsymbol{s}_{r} \right\|_{F}^{2}$$
$$= \min_{\boldsymbol{a}_{r} \succeq 0} \left\| \boldsymbol{Z}^{(r)} - \boldsymbol{a}_{r} \boldsymbol{s}_{r} \right\|_{F}^{2}, \quad (4.24)$$

where a_r is the *r*-th column of A, s_r is the *r*-th row of S, and $Z^{(r)} = X - \sum_{k \neq r} a_k s_k$ is the residual matrix with respect to the *r*-th rank-one factor $a_r s_r$. Interestingly, (4.24) can be decoupled into P independent NNLS problems in one variable:

$$\min_{\boldsymbol{a}_{r} \succeq 0} \left\| \boldsymbol{Z}^{(r)} - \boldsymbol{a}_{r} \boldsymbol{s}_{r} \right\|_{F}^{2} = \sum_{p=1}^{P} \min_{a_{pr} \succeq 0} \left\| \boldsymbol{z}_{p}^{(r)} - a_{pr} \boldsymbol{s}_{r} \right\|_{2}^{2},$$
(4.25)

where vector $\boldsymbol{z}_p^{(r)}$ is composed by the *p*-th row of $\boldsymbol{Z}^{(r)}$. A NNLS problem in one variable is equivalent to a problem of the form

$$\min_{y\ge 0} \alpha y^2 - 2\beta y,\tag{4.26}$$

for some $\alpha > 0$ and $\beta \in \mathbb{R}$. Clearly, if the solution of the unconstrained problem $\min_{y \in \mathbb{R}} \alpha y^2 - 2\beta y$ is non-negative, then it is also the solution of (4.26), otherwise the minimizer is zero. Therefore,

$$\operatorname{argmin}_{y \ge 0} \alpha y^2 - 2\beta y = \max\left(0, \frac{\beta}{\alpha}\right).$$
(4.27)

For the problem $\min_{y \succeq 0} \| \boldsymbol{z}_p^{(r)} - y \boldsymbol{s}_r \|_2^2$, we have

$$\alpha = \boldsymbol{s}_r^\mathsf{T} \boldsymbol{s}_r = ||\boldsymbol{s}_r||_2^2 \quad \text{and} \quad \beta = \boldsymbol{s}_r^\mathsf{T} \boldsymbol{z}_p^{(r)}$$

hence

$$\operatorname{argmin}_{a_{pr} \succeq 0} \left\| \boldsymbol{z}_p^{(r)} - a_{pr} \boldsymbol{s}_r \right\|_2^2 = \max\left(0, \frac{\boldsymbol{s}_r^{\mathsf{T}} \boldsymbol{z}_p^{(r)}}{\|\boldsymbol{s}_r\|_2^2}\right).$$

Finally, in vector form, we have

$$\operatorname{argmin}_{\boldsymbol{a}_{r} \succeq 0} \left\| \boldsymbol{Z}^{(r)} - \boldsymbol{a}_{r} \boldsymbol{s}_{r} \right\|_{F}^{2} = \max \left(0, \frac{\boldsymbol{Z}^{(r)} \boldsymbol{s}_{r}^{\mathsf{T}}}{\|\boldsymbol{s}_{r}\|_{2}^{2}} \right)$$
$$= \max \left(0, \frac{\boldsymbol{X} \boldsymbol{s}_{r}^{\mathsf{T}} - \sum_{k \neq r} \boldsymbol{a}_{k} \boldsymbol{s}_{k} \boldsymbol{s}_{r}^{\mathsf{T}}}{\|\boldsymbol{s}_{r}\|_{2}^{2}} \right). \quad (4.28)$$

HALS updates successively the columns of A and the rows of S using the above closed-form solution:

1. Initialize $(\boldsymbol{A}, \boldsymbol{S})$. 2. For $\ell = 1, 2, ...,$ 2.a For r = 1, 2, ..., R: Update $\boldsymbol{s}_r \leftarrow \max\left(0, \frac{\boldsymbol{a}_r^{\mathsf{T}} \boldsymbol{X} - \sum_{k \neq r} \boldsymbol{a}_r^{\mathsf{T}} \boldsymbol{a}_k \boldsymbol{s}_k}{\|\boldsymbol{a}_r\|_2^2}\right)$. 2.b For r = 1, 2, ..., R: Update $\boldsymbol{a}_r \leftarrow \max\left(0, \frac{\boldsymbol{X} \boldsymbol{s}_r^{\mathsf{T}} - \sum_{k \neq r} \boldsymbol{a}_k \boldsymbol{s}_k \boldsymbol{s}_r^{\mathsf{T}}}{\|\boldsymbol{s}_r\|_2^2}\right)$.

Note that the residuals $Z^{(r)}$ are not computed (as they could be dense while X could be sparse). As for ANLS, HALS is guaranteed to converge to a stationary point of (4.20) [65]. As for the MU, it is possible to accelerate HALS by updating several times the columns of A before updating the rows of S since the terms Xs_r^T and $s_k s_r^T$ do not need to be recomputed [65, 58]. The computational cost of HALS is, up to some negligible factor, the same as for the MU [58]. However, in practice, HALS converges significantly faster. In fact, in most situations, HALS performs the best among ANLS, the MU and ALS; see the references above and also, e.g., [37, 66].

To conclude this section, we refer the reader to the survey [67] on the classification of NMF methods as coordinate descent schemes (exact or approximate), where a more detailed analysis can be found along with some numerical comparisons.

4.5.2

Constrained and penalized factorization methods

As explained in details in sections 4.2 and 4.4, a crucial aspect for practical applications when designing NMF algorithm is to take into account the non-uniqueness issue. The usual way to tackle this is to take into account, in the model, additional prior information depending on the application at hand. Here is a (non-exhaustive) list of additional constraints that can be added to the NMF model:

- Minimum volume. Looking back at the geometric interpretation of NMF from section 4.2, it often makes sense in practice to look for a source matrix S such that H(S) has minimum volume. This enforces the sources to be as close as possible to the data points H(X), while allowing to approximate them well. This is discussed in much details in Chapter ??; see also [33].
- *Sparsity.* In many cases, the activation matrix *A* should be sparse because, for most observations, only a few sources are active. For example, in hyperspectral unmixing, most pixels only contain a few constitutive materials (also called end-members); see Chapter ??. In other cases, the source matrix *S* should be sparse. For example, in document classification, the sources are topics (represented by a set of words) which contain only a few words from the dictionary.
- Orthogonality. NMF can be used as a clustering model, using an additional orthogonality constraints $A^{T}A = I$. In fact, this constraint imposes that each observation is approximated only by a single source since, for any $1 \le r \le R$ and $1 \le p \le P$,

$$a_{pr} > 0 \implies a_{p\ell} = 0 \text{ for all } \ell \neq r.$$

where a_{pr} is the activation coefficient of the *r*-th source for the *p*-th observation. This follows from the fact that $\sum_{p=1}^{P} a_{pr} a_{p\ell} = 0$ for all $\ell \neq r$ and $A \succeq 0$. In practice, this constraint is often too restrictive because most observations result from a combination of several sources. However, adding a penalty term in the objective function of the form $||\mathbf{A}^{\mathsf{T}}\mathbf{A} - \mathbf{I}||_{F}^{2}$ to enforce the mixing matrix \mathbf{A} to be closer to an orthogonal matrix allows to enforce sparsity.

• Spatial information. In imaging applications (see, e.g., section 4.6), it is possible to take into account the spatial information as most neighboring pixels (the observations) will share similarities and their activation coefficients will be similar. For example, if pixels p and p' for $1 \le p, p' \le P$ are neighbors, it is useful to add the following penalty term in the objective function

$$\sum_{r=1}^R |a_{pr} - a_{p'r}|.$$

Note that the absolute is usually used because it allows to preserve the edges in the images, instead of the ℓ_2 norm that smooths the edges; see, e.g., [68, 69].

- Graph regularization. In several situations, it is possible to embed the observations in a graph using some similarity measure. The vertices of the graph are the observations and the (weighted) edges indicate whether two observations are similar. The way this graph is constructed depends on the application. For example, in imaging, the graph could be constructed connecting neighboring pixels. In general, this graph is constructed using some similarity measure between observations and the NMF model takes this information by requiring that similar observations have similar activation coefficients (as shown above for the spatial information). This type of structure can take into account many different priors and is a very flexible way to add prior knowledge in the NMF model [70].
- Sum-to-one constraints. In several applications, the mixing coefficients correspond to a proportions (e.g., in hyperspectral imaging) in which case the entries on each row of A should sum to one.

The additional constraints added into the NMF model are in general taken into account by adding a penalty term in the objective function, with some penalty parameter balancing the importance between the fitting error and the regularization term.

For example, to obtain sparser A and S, it is standard to use the ℓ_1 -norm as a proxy for the ℓ_0 -'norm' that counts the number of nonzero entries [52], and consider a model of the type

$$\min_{oldsymbol{A} \succeq 0, oldsymbol{S} \succeq 0} ||oldsymbol{X} - oldsymbol{A}oldsymbol{S}||_F^2 + \lambda_A ||oldsymbol{A}||_1 + \lambda_S ||oldsymbol{S}||_1$$

where $\|\mathbf{Y}\|_1 = \sum_{i,j} |y_{ij}|$. Doing so, the algorithms described in the previous section 4.5 can usually be adapted to handle these situations. For example, if the penalty term is convex then the NNLS subproblems remain convex hence efficiently solvable. For gradient-based methods, it suffices to account for the additional term(s) in

the objective function and update the way the gradient is computed. Note that, in general, it will be necessary (and non trivial) to properly tune the (penalty) parameters of the model to obtain good solutions, although some general strategies can sometimes be designed; see, e.g., [71] for sparse NMF. We refer the reader to the references above, and the references therein, for more details on these methods.

Another way these constraints can be taken care of is by using a projection. If the projection onto the feasible set can be computed efficiently, then gradient-based method can be adapted easily. For example, non-negativity constraints along with the sum-to-one constraint amounts to optimizing over the unit simplex for which the projection can be computed efficiently [72]. For sparsity, Hoyer [73] introduced a projection onto the set of matrices of a given sparsity level, based on the following measure of sparsity: for a nonzero vector $\boldsymbol{y} \in \mathbb{R}^N$,

$$spar(\boldsymbol{y}) = \frac{\sqrt{N} - \|\boldsymbol{y}\|_1 / \|\boldsymbol{y}\|_2}{\sqrt{N} - 1} \in [0, 1],$$

where spar(y) = 1 when y has a single nonzero entry (hence $||y||_1 = ||y||_2$), and spar(y) = 0 when all entries of y are positive with similar values (hence $||y||_1 = \sqrt{N} ||y||_2$).

To summarize, when using NMF for a particular application, it is crucial to first think carefully about the constraints that the sources and the mixing coefficients should satisfy. This will allow to design a dedicated NMF algorithm that will be able to recover the sought solution.

4.5.3

Geometrical approaches and separability

Another class of NMF problems that has gathered much attention lately is based on intuitions coming from the geometric interpretation of NMF.

This class is referred to as *separable NMF* and requires that the input matrix X satisfies the following separability condition: there exists an index set $\mathcal{K} \subset \{1, 2, \ldots, P\}$ of size R, that is, $Card(\mathcal{K}) = R$, and a non-negative matrix $A \succeq 0$ with R columns such that

$$X = A \underbrace{X(\mathcal{K},:)}_{=S}$$

This means that the sources can be found among the observations. The problem therefore boils down to identifying the rows of X corresponding to the sources.

Geometrically, this means that, if the observations are scaled, we are looking for the vertices of $conv{X}$. This can be done efficiently, even in the presence of noise [74]; see also [75, 76, 77] and the references therein for some recent developments. This assumption makes sense in several applications such as document classification [78], and hyperspectral imaging (this is the so-called pure-pixel assumption). These approaches are described and analyzed in details in Chapter ??.

4.6

Applications of NMF chemical sensing. Two examples of reducing admissible solutions

As mentioned several times before, NMF is a difficult problem in general, being NPhard. However, this does not imply that all NMF instances are difficult to solve. In some applications, the input matrix X can have a specific structure or results from a measurement process that makes the NMF problem less difficult, hence taking this structure into account allows to design much more effective algorithms. For example, we have already seen in section 4.2 that the NMF problem can be solved easily if rank{X} ≤ 2 . This observation has been used to solve the NMF problem in a hierarchical way, and has been shown to perform well in several situations: analyzing Magnetic Resonance Spectroscopy and Imaging (MRSI) data [79], document classification [80] and hyperspectral unmixing [81].

In real-world applications, the given data matrix is not guaranteed to obey the established uniqueness conditions, thereby limiting the practical success of NMF algorithms. Hence, it motivates the development of algorithms to constrain the solution space of a given NMF problem as presented in section 4.5.2). However, in some situations, additional information can be taken into account to reduce the set of admissible solutions. We consider first a data augmentation strategy which consists in coupling multiple data sets. This approach is illustrated on a polarized Raman spectroscopy application. An alternative approach is to perform a data pre-processing stage [82, 83, 84]. To illustrate this later approach, we consider the blind unmixing of spectroscopy imaging data.

4.6.1

Polarized Raman spectroscopy: a data augmentation approach

This section is adapted from [85]. Raman scattering is a light-matter interaction process which reflects the molecular vibration properties of molecules and materials, thus characterizing the chemical composition of the analyzed sample [86, 87]. For materials presenting a regular atomic or molecular structure, a more accurate characterization of the sample can be achieved by using polarizers [88]. In particular, this is the case for crystals as their response to the polarized light excitation will reflect the crystallographic structure of the sample, motivating the development of polarized Raman spectroscopy. This is the type of problem addressed here since we are considering the polarized Raman spectroscopy of a Rutile (TiO2) crystal. It is worth mentioning that polarized Raman spectroscopy has other possible uses. For example, it can also be used for determining the molecular orientation distribution of polymeric materials [89].

4.6.1.1 Raman data description

The Raman measurements presented in this section were carried out in back scattering geometry with the same objective for excitation and collection of light. The confocal Raman spectrometer was equipped with a cooled CCD camera and the laser source was an ionized argon laser emitting at a wavelength $\lambda = 514.5 \ nm$. The analyzed crystal sample is fixed on a rotating stage as shown in Fig.4.6. Two coordinate systems are used, one associated with the laboratory space-fixed coordinates (O, X, Y, Z) and another attached to the analyzed sample (O, x, y, z). The incident light is polarized such that the electric field arriving on the sample is oriented along the Y direction. The scattered light is analyzed by positioning an analyzer in front of the entrance slit of the spectrometer. The analyzer is oriented either along the Y-axis (parallel polarization) or the X-axis (crossed polarization). Thus, the acquisition in one point of the sample yields a pair of spectra, one for the parallel polarization, indexed by Y and another for the crossed polarization, indexed by X. The rotational diversity scheme consists in rotating the sample around the Z-axis (Fig. 4.6) with a fixed angular step (typically 10°) and acquiring two polarized spectra for each step of the rotation. For the rotational diversity acquisition scheme, m polarized spectra are acquired for m different rotation angles $(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m)$ of the analyzed sample. For this angular diversity data, the 'sources" are represented by vibrational modes. Indeed, the vibrational modes are characterized by specific displacements of the atoms from their equilibrium position, which dictate the magnitude of the components of the Raman polarizability tensor. The change of polarized Raman intensity versus rotational angle, for a specific vibrational mode, will therefore be different from another one. Each mode in polarized Raman spectra will thus contribute as one source in the full spectrum.



Figure 4.6 Polarized Raman spectroscopy set-up in backscattering geometry

Under the assumption of instantaneous linear mixture, the acquired data can be

structured as two $m \times n$ matrices, corresponding to the two polarization orientations:

$$X_1 = A_1 S_1 + E_1. (4.29)$$

and

$$X_2 = A_2 S_1 + E_2. (4.30)$$

In (4.29) and (4.30), matrices E_1 and E_2 accounts for the additive noise on the sensors and the model errors.

If we further analyze the underlying physico-chemical phenomenon generating the two data sets, the spectra of pure compounds are the same for the crossed and the parallel polarization [86, 87], since the vibrational modes are imposed by the structure of the crystal. This implies $S_1 = S_2 = S$, which is quite intuitive if we consider a geometrical point of view in which the crossed and parallel polarized spectra are projections of the same signal on two orthogonal axes. By injecting this information into (4.29) and (4.30), we can write:

$$\begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \end{pmatrix} \mathbf{S} + \begin{pmatrix} \mathbf{E}_1 \\ \mathbf{E}_2 \end{pmatrix} \quad . \tag{4.31}$$

Eq. (4.31) points out a mixing model for the polarized spectra with rotational diversity considering both polarized spectra families jointly. Besides the fact that this is a more natural and compact representation of the data, the sample size is doubled in (4.33) compared to (4.29), (4.30) and the number of unknowns is lower; this should improve the accuracy of the estimated source parameters. In order to simplify the presentation we use the following notations:

$$X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}, \quad A = \begin{pmatrix} A_1 \\ A_2 \end{pmatrix}, \quad E = \begin{pmatrix} E_1 \\ E_2 \end{pmatrix}.$$
 (4.32)

Equation (4.31) can thus be re-written in a more concise manner as:

$$\boldsymbol{X} = \boldsymbol{A}\boldsymbol{S} + \boldsymbol{E}.\tag{4.33}$$

4.6.1.2 Raman data processing

Given the nature of the data, the sources and the mixing coefficients are nonnegative, meaning that (4.33) expresses a NMF model. It should be noticed that stacking the data matrices X_1 and X_2 into a bigger matrix X corresponds to a data augmentation strategy. This kind of technique was already proposed for diverse problems such as the analysis of multiple runs of gasoline blending processes [90]. Another example is the joint analysis of UV-visible spectra related to the complexation of the aluminum by caffeic acid and the titration of caffeic acid [91]. Actually, the benefit of matrix augmentation strategy is threefold: it allows to decrease estimation error uncertainties, it may remove rank deficiency and helps in reducing rotational ambiguities.



Figure 4.7 Polarized Raman data versus rotational angle χ for rutile $TiO_2~(110)$ single crystal

The approach was applied to a rutile TiO_2 crystal, as shown in Fig.4.6. The crystallographic face (110) (Hermann-Mauguin international crystallographic symbols) is analyzed. The sample is rotated with respect to Z axis only, meaning $\boldsymbol{\theta} = (0, 0, \chi)$. Fig.4.7 presents the acquired polarized data for the parallel and crossed polarizations (matrices \boldsymbol{X}_1 and \boldsymbol{X}_2). The data was acquired in a spectral range of 100 cm⁻¹ - 800 cm⁻¹, with an angular rotation step of 10 degrees between 0° and 190°.

In the case of TiO_2 , four Raman active modes denoted as A_{1g} , E_g , B_{1g} and B_{2g} (Mulliken symbols for symmetry groups [86]) are expected from theory. However, the B_{2g} mode at 826 cm⁻¹ is out of the spectral window used in the present work (and anyway the B_{2g} has a very low Raman crossed section and is often not detected). The B_{1g} mode with the (110) oriented crystal plane is inactive either in parallel or crossed polarizations. Consequently, one can expect two Raman active modes, *i.e.* two sources, in the data collected here. Nevertheless, three sources are necessary to properly describe the data, as indicated by the magnitude analysis of the singular values of data matrices X_1 , X_2 and X. A theoretical explanation for the presence of this third source is provided in [85].

We illustrate the effect of the joint use of the crossed and parallel polarization data sets on the reduction of the NMF admissible solutions set. The NMF algorithm [92] was used to estimate the three source vectors and the corresponding mixing coefficients. The two data sets were processed separately and jointly and the results are presented on Fig.4.8 for the source spectra and on Fig.4.9 and Fig.4.10 for the mixing coefficients. To evaluate the size of the admissible solutions set, we used 25 independent runs for each plot, with different random initial values for the matrices A and S. As one can see, by processing jointly both polarization data sets (Fig.4.8 (c) and Fig.4.10) the admissible solution domain is largely reduced as compared to the case when only one polarization is used (Fig.4.8 (a),(b) and Fig.4.9). However, the solution is still not unique, which motivated the use of some regularization techniques. In [85], a penalized NMF algorithm derived from a Bayesian source separation approach (called BPSS) [93] was used to further reduce the set of admissible solution. The results and the physical interpretation of the obtained results are not reported here. The interested reader is referred to [85] for the detailed analysis.



Figure 4.8 Source spectra estimated by NMF (25 runs)

4.6.2 Unmixing blurred Raman spectroscopy images

Hyperspectral images may be viewed as a collection of highly resolved spectra. In many cases, the image contains a small number of pure materials - termed endmembers - whose spectral signatures are mixed in each pixel because of limited spatial



Figure 4.9 Estimated coefficients by NMF for each polarization data set separately



Figure 4.10 Estimated coefficients by NMF for both polarization data sets jointly (25 runs)

resolution. Blind spectral unmixing usually refers to the estimation of endmembers and their fractional contribution to each pixel, named abundances. A geometrical framework for spectral unmixing has attracted a lot of attention from researchers in the past two decades. In this approach, each pixel spectrum belongs to a simplex whose vertices are the endmembers we seek. In some cases, the data are known to contain at least one pure pixel per endmember, which are subsequently extracted from the hyperspectral scene; this is equivalent to separable NMF; see Section 4.5.3. When the pure pixel hypothesis does not hold, the geometrical approach to unmixing then consists in finding the Minimum Volume Simplex (MVS) enclosing the data using one of many MV algorithms [94, 95, 96]; see Chapter **??** for more details. Once endmembers have been extracted from the scene, abundances can be estimated using constrained least squares algorithms [97, 98]. However, highly mixed data are beyond the reach of geometrical algorithms because spectral signatures are located near the center of the true endmember [99].

4.6.2.1 Blurring effect modeling

Consider an hyperspectral image measuring radiance on L different spectral bands (*channels*) and N pixels. We gather the data in a $L \times N$ matrix X and use the following notations:

- x^l is the l-th row of X, that is the 2D image at spectral band l after lexicographical ordering into a row vector of length N;
- 2) x_p is the *p*-th column of X, i.e. the $L \times 1$ spectrum of the *p*-th pixel (also termed spectral vector or pixel vector).

Each spectral vector in the image is a linear combination of an known number R of endmembers $\{s_1, \ldots, s_R\}$. When unknown, R can be obtained by some model order estimation method such as *virtual dimensionality* [100]. Ignoring noise for now, the linear mixing model (LMM) writes

$$X = SA$$

where the r-th column of $L \times R$ source matrix S indexes endmember s_r and the p-th column a_p of $P \times R$ abundance matrix A contains the fractional abundance coefficients for x_p :

$$\boldsymbol{x}_p = \boldsymbol{S}\boldsymbol{a}_p = \sum_{r=1}^R a_{pr} \boldsymbol{s}_r.$$
(4.34)

The linear mixing model in spectral imaging is generally based on the following assumptions [101]:

- i) The number of endmembers R is much smaller than the number of bands L, that is $R \ll L$;
- ii) Matrix S is of full column rank, *i.e.* endmembers $\{s_1, \ldots, s_R\}$ are linearly independent;

- iii) Abundance Nonnegativity Constraint (ANC): $a_{pr} \ge 0$ for all p and r;
- iv) Abundance Sum Constraint (ASC): $\sum_{r=1}^{R} a_{pr} = 1$ for all p.

Assumptions (i) and (ii) seem very reasonable in hyperspectral imaging since many bands are collected and the image is made up of a few distinct materials. Assumptions (iii) and (iv) come from the physical interpretation of abundance coefficient a_{pr} as the fractional spatial area occupied by the *r*-th endmember in the *p*-th pixel.

We now account for the fact that the image is degraded during the acquisition process. Under the common linear blur assumption, the 2D image y^{ℓ} observed at a given channel ℓ is obtained as the 2D convolution product of the true image and the channel point-spread function \mathcal{H}^{ℓ} :

$$\boldsymbol{y}^{\ell} = \boldsymbol{x}^{\ell} \boldsymbol{H}^{\ell} \tag{4.35}$$

where the $N \times N$ matrix \mathbf{H}^{ℓ} is a convolution matrix corresponding to \mathcal{H}^{ℓ} . For instance, when the blur is space-invariant for different pixels, $(\mathbf{H}^{\ell})^{\mathsf{T}}$ is a block-Toeplitz matrix where each block is Toeplitz [102]. Each entry of the observed data matrix \mathbf{Y} is given by

$$y_p(\ell) = \sum_{n=1}^{N} h_{pn}^{\ell} x_n(\ell)$$
(4.36)

Using equations (4.34) and (4.36), the overall model combining noise, observation blurring and linear mixing of endmembers writes

$$y_p(\ell) = \sum_{n=1}^{N} \sum_{r=1}^{R} h_{pn}^{\ell} a_{pr} s_r(\ell) + e_p(\ell)$$
(4.37)

where E is the noise term and model (4.37) assumes that the SNR is high enough for the noise to be additive and i.i.d. Gaussian. We observe that the blurred data do not satisfy the linear mixing model since the mixing coefficients $\left(\sum_{n=1}^{N} h_{pn}^{\ell} a_{pr}\right)$ depend on the channel index ℓ . However, in the specific case where the PSF is invariant across channels, the model reduces to the following linear mixing model

$$Y = SAH + E \tag{4.38}$$

with fixed matrix H.

How does the observation process affect the distribution of pixel vectors inside the simplex? The answer to the question obviously depends on the nature of the PSF. Since the entries of H are known to be non-negative, the blurring process tends to average neighboring pixel intensities. This phenomenon causes observed spectral vectors to cluster towards the center of S. This contraction property has important practical consequence: directly applying a NMF algorithms to the observed data may produce incorrect sources, and thus the subsequent estimation of abundances will also be biased. To improve the unmixing performance, it is necessary to deconvolve the data before applying any NMF algorithm.

4.6.2.2 Application to Raman spectroscopy images

In this section, we illustrate the impact of the deblurring process on the performances of the separation of real Raman spectroscopy data. This dataset comprises images of size 98×131 pixels, each pixel being $100 \text{ nm} \times 100 \text{ nm}$, acquired on 337 spectral bands ranging from 800 cm^{-1} to 1200 cm^{-1} . The scene of interest consists in a grain of *sodium acetate* (CH₃COONa) covered with *sodium carbonate* (Na₂CO₃) laying on a *silicon* layer (Si). Part of the sodium carbonate reacts with water vapor to yield hydrated sodium carbonate. These four chemical compounds are the endmembers we seek. A thorough inspection of the data reveals that the silicon compound contributes to all pixels of the image. The extraction of these endmembers is a challenging problem, since silicon is the only compound for which the pure pixel assumption is fulfilled. Because of the inherent high mixing of the data, even minimum volume methods are not supposed to produce good results on this data set.

Given the limited spectral range, the point spread function is considered to be invariant across channels and mixels. It is modeled as a 2D Gaussian function [103] with an experimentally measured full-width at half maximum of 300 nm. We apply our deconvolution algorithm to the data by setting regularization parameters through a trial-and-error process to $\mu = 20$ and $\nu = 5$ using the algorithm proposed in [102]. Non-negative Matrix Factorization with volume constraints (*NMF-vol*) [99] was applied for the data processing. The algorithm operates on both the raw data (where negative pixels have been clipped to zero since the algorithm imposes a nonnegative data matrix) and restored data.

The resulting abundance maps and endmembers are given in figure 4.11. The first extracted endmember corresponds to the silicon layer, which presents a broad band at $910 - 960 \text{ cm}^{-1}$ due to the 2TO harmonic phonon mode of bulk silicon. The deconvolution step allows to denoise its abundance map and more importantly, uncovers structure that was distributed throughout other abundance maps. The second endmember is sodium acetate, with a peak at 930 cm^{-1} due to the intense C-C stretching mode of the acetate molecule [104]. The sodium acetate endmember displays more undesirable contribution from the silicon compound. The spectral shapes of endmembers extracted from the raw data appear noisier, a problem solved by the deconvolution step. Moroevet, the algorithm is able to separate the sodium carbonate compound (third endmember, characterized by the peak at 1080 cm^{-1}) from the hydrated sodium carbonate compound (fourth endmember, peak at 1060 cm^{-1}); both are mixed with the sodium benzoate, as expected. The main gain of deconvolution clearly appears on the third endmember, where the silicon contribution is completely suppressed. Another benefit of the restoration step is to reveal structure hidden in raw abundance maps (first and third endmember) that was not displayed by applying NMF-vol to the raw data.

34



Figure 4.11 Spectral signatures and spatial abundances of the sources obtained by applying NMF with minimum volume contraint on both the raw Raman data and the deblurred data.

4.7 Conclusions

This chapter described the concept of non-negative matrix factorization and its application in the context of source separation in chemical-physical sensing. Although the NMF problem is described in a simple mathematical way, the theoretical formulation of its solution existence and uniqueness remains an active area of study. Moreover, the resolution algorithms are generally application dependent since additional constraints (to non-negativity) should be firstly determined and then accounted for to get an acceptable solutions. Actually, these additional contraints are exploited by developing dedicated factorization algorithms calling to some algebraic tools, optimization techniques, geometrical concepts and statistical inference methods. Active open problems in the NMF field rely on the reduction of the numerical complexity of the algorithms through the development of sequential factorization methods and data compression techniques.

References

- Imbrie, J. and Van Andel, T.H. (1964) Vector analysis of heavy-mineral data. *Geological Society of America Bulletin*, 75, 1131–1156.
- 2 Markham, T.L. (1972) Factorizations of non-negative matrices. *Proceedings of American Mathematical Society*, 32, 45–47.
- Thomas, L.B. (1974) Rank factorizations of nonnegative matrices. *SIAM Review*, 16, 393–394.
- 4 Campbell, S.I. and Poole, G.D. (1981) Computing non-negative rank factorizations. *Linear Algebra and its Applications*, **35**, 175–182.
- 5 Chen, J.C. (1984) Nonnegative rank factorisation of nonnegative matrices. *Linear Algebra and its Applications*, 62, 207–217.
- Cohen, J. and Rothblum, U. (1993) Non-negative ranks decompositions and factorizations of non-negative matrices. *Linear Algebra and its Applications*, 190, 149–168.
- 7 Paatero, P. and Tapper, U. (1994) Positive matrix factorization: a non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics*, 5, 111–126.
- Lee, D. and Seung, H. (1999) Learning the Parts of Objects by Nonnegative Matrix Factorization. *Nature*, 401, 788–791.
- 9 Tauler, R., Kowalski, B., and Fleming, S. (1993) Multivariate curve resolution applied to spectral data from multiple runs of an industrial process. *Analytical Chemistry*, 65, 2040–2047.
- 10 Paatero, P. and Tapper, U. (1994) Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics*, 5, 111–126.
- Paatero, P. (1997) Least squares formulation of robust non-negative factor analysis. *Chemomerics and Intelligent Laboratory Systems*, **37**, 23–35.
- 12 Lee, D. and Seung, H. (1999) Learning the parts of objects by non–negative matrix factorization. *Nature*, 401, 788–791.

- 13 Sajda, P., S.Du, Brown, T., Stoyanova, R., Shungu, D., Mao, X., and Parra, L. (2004) Nonnegative matrix factorization for rapid recovery of constituent spectra in magnetic resonance chemical shift imaging of the brain. *IEEE Transaction* on Medical Imaging, 23 (12), 1453–1465.
- 14 Hoyer, P. (2004) Nonnegative matrix factorization with sparseness constraints. *Journal of Machine Learning Research*, 5, 1457–1469.
- 15 Pascual-Montano, A., Carazo, J., Kochi, K., and Pascual-Marqui, R. (2006) Nonsmooth nonnegative matrix factorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **28** (3), 403–415.
- 16 Moussaoui, S., Brie, D., Mohammad-Djafari, A., and Carteret, C. (2006) Separation of non-negative mixture of non-negative sources using a Bayesian approach and MCMC sampling. *IEEE Transactions on Signal Processing*, 54 (11), 4133–4145.
- 17 Das, G. and Joseph, D. (1990) The Complexity of Minimum Convex Nested Polyhedra, in *Proc. of the 2nd Canadian Conf. on Computational Geometry*, pp. 296–301.
- 18 Aggarwal, A., Booth, H., O'Rourke, J., and Suri, S. (1989) Finding minimal convex nested polygons. *Information and Computation*, 83 (1), 98–110.
- 19 Mond, D., Smith, J., and van Straten, D. (2003) Stochastic factorizations, sandwiched simplices and the topology of the space of explanations. *Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, 459 (2039), 2821–2845.
- 20 Vavasis, S. (2010) On the complexity of nonnegative matrix factorization. *SIAM J.on Optimization*, 20 (3), 1364–1377.
- 21 Gillis, N. and Glineur, F. (2012) On the geometric interpretation of the nonnegative rank. *Linear Algebra and its Applications*, 437 (11), 2685–2712.
- 22 Chistikov, D., Kiefer, S., Marušić, I., Shirmohammadi, M., and Worrell, J. (2016) On restricted nonnegative matrix factorization, in *Proceedings of the 43rd*

International Colloquium on Automata, Languages and Programming (ICALP).

- 23 Goemans, M. (2015) Smallest compact formulation for the permutahedron. *Mathematical Programming*, 153 (1), 5–11.
- 24 Vandaele, A., Gillis, N., and Glineur, F. (2015) On the linear extension complexity of regular n-gons. ArXiv:1505.08031.
- 25 Moitra, A. (2013) An almost optimal algorithm for computing nonnegative rank, in *Proc. of the 24th Annual ACM-SIAM Symp. on Discrete Algorithms* (*SODA '13*), pp. 1454–1464.
- 26 Gillis, N. (2012) Sparse and unique nonnegative matrix factorization through data preprocessing. *Journal of Machine Learning Research*, 13 (Nov), 3349–3386.
- 27 Févotte, C. and Idier, J. (2011) Algorithms for nonnegative matrix factorization with the beta-divergence. *Neural Computation*, 23 (9), 2421–2456.
- 28 Park, E.S., Spiegelman, C.H., and Henry, R.C. (2002) Bilinear extimation of pollution source profiles and amounts by using multivariate receptor models. *Environmetrics*, 13, 775–798.
- 29 Smilde, A., Hoefsloot, H., Kiers, H., Bijlsma, S., and Boelens, H. (2001) Sufficient condition for unique solutions within a certain class of curve resolution models. *Journal of Chemometrics*, 15, 405–411.
- 30 Moussaoui, S., Brie, D., and Idier, J. (2005) Non-negative source separation: Range of admissible solutions and conditions for the uniqueness of the solution, in proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP'2005), Philadelphia, USA.
- 31 Donoho, D. and Stodden, V. (2003) When does non-negative matrix factorization give a correct decomposition into parts?, in Advances in Neural Information Processing Systems 16, MIT Press, Cambridge, United States.
- 32 Laurberg, H., Christensen, M., Plumbley, M.D., Hansen, L.K., and Jensen, S.H. (2008) Theorems on Positive Data: On the Uniqueness of NMF. *Computational Intelligence and Neuroscience*, **ID** 764206, 9 pages.

- 33 Fu, X., Huang, K., Sidiropoulos, N.D., and Ma, W.K. (2019) Nonnegative matrix factorization for signal and data analytics: Identifiability, algorithms, and applications. *IEEE Signal Processing Magazine*, 36 (2), 59–80.
- 34 Sasaki, K., Kawata, S., and Minami, S. (1983) Constrained nonlinear method for estimating component spectra from multicomponent mixtures. *Applied Optics*, 22 (22), 3599–3606.
- 35 Moussaoui, S., Carteret, C., Brie, D., and Mohammad-Djafari, A. (2006) Bayesian analysis of spectral mixture data using Markov chain Monte Carlo methods. *Chemometrics and Intelligent Laboratory Systems*, 81 (2), 137–148.
- 36 Neymeyr, K. and Sawall, M. (2018) On the set of solutions of the nonnegative matrix factorization problem. *SIAM Journal on Matrix Analysis and Applications*, 39 (2), 1049–1069.
- 37 Cichocki, A., Zdunek, R., Phan, A., and Amari, S.I. (2009) Nonnegative Matrix and Tensor Factorizations: Applications to Exploratory Multi-way Data Analysis and Blind Source Separation, Wiley-Blackwell.
- 38 Chi, E. and Kolda, T. (2012) On tensors, sparsity, and nonnegative factorizations. *SIAM J. on Matrix Analysis and Applications*, 33 (4), 1272–1299.
- **39** Lawson, C. and Hanson, R. (1974) Solving Least Squares Problems, Prentice-Hall.
- 40 Bro, R. (1998) Multi-Way Analysis in the Food Industry: Models, Algorithms, and Applications, Ph.D. thesis, University of Copenhagen.
- **41** Ho, N.D. (2008) *Nonnegative Matrix Factorization - Algorithms and Applications*, Ph.D. thesis, Université catholique de Louvain.
- Wild, S., Curry, J., and Dougherty, A. (2004) Improving non-negative matrix factorizations through structured initialization. *Pattern Recognition*, 37 (11), 2217–2232.
- **43** Langville, A., Meyer, C., Albright, R., Cox, J., and Duling, D. (2006) Initializations for the nonnegative matrix factorization, in *Proceedings of the twelfth ACM SIGKDD international conference*

on knowledge discovery and data mining, pp. 23–26.

- 44 Casalino, G., Del Buono, N., and Mencar, C. (2013) Subtractive clustering for seeding non-negative matrix factorizations. *Information Sciences*, (257), 369–387.
- 45 Boutsidis, C. and Gallopoulos, E. (2008) SVD based initialization: A head start for nonnegative matrix factorization. *Pattern Recognition*, 41, 1350–1362.
- 46 Atif, S.M., Qazi, S., and Gillis, N. (2019) Improved svd-based initialization for nonnegative matrix factorization using low-rank correction. *Pattern Recognition Letters*, **122**, 53–59.
- 47 Grippo, L. and Sciandrone, M. (2000) On the convergence of the block nonlinear Gauss-Seidel method under convex constraints. *Operations Research Letters*, 26, 127–136.
- Lin, C.J. (2007) Projected Gradient Methods for Nonnegative Matrix Factorization. *Neural Computation*, 19, 2756–2779.
- 49 Guan, N., Tao, D., Luo, Z., and Yuan, B. (2012) NeNMF: an optimal gradient method for nonnegative matrix factorization. *IEEE Transactions on Signal Processing*, **60** (6), 2882–2898.
- 50 Cichocki, A., Zdunek, R., and Amari, S. (2006) Non-negative Matrix Factorization with Quasi-Newton Optimization, in *Lecture Notes in Artificial Intelligence, Springer*, vol. 4029, vol. 4029, pp. 870–879.
- 51 Bro, R. and De Jong, S. (1997) A fast non-negativity-constrained least squares algorithm. *Journal of chemometrics*, 11 (5), 393–401.
- 52 Kim, H. and Park, H. (2007) Sparse non-negative matrix factorizations via alternating non-negativity-constrained least squares for microarray data analysis. *Bioinformatics*, 23 (12), 1495–1502.
- 53 Kim, J. and Park, H. (2011) Fast nonnegative matrix factorization: An active-set-like method and comparisons. *SIAM J. on Scientific Computing*, 33 (6), 3261–3281.
- 54 Chen, D. and Plemmons, R. (2009) Nonnegativity Constraints in Numerical Analysis, in A. Bultheel and R. Cools

(Eds.), Symposium on the Birth of Numerical Analysis, World Scientific Press.

- 55 Lee, D. and Seung, H. (2001) Algorithms for Non-negative Matrix Factorization. In Advances in Neural Information Processing (NIPS '01), 13.
- 56 Daube-Witherspoon, M. and Muehllehner, G. (1986) An iterative image space reconstruction algorithm suitable for volume ECT. *IEEE Transactions on Medical Imaging*, 5, 61–66.
- 57 Bertsekas, D.P. (2014) Constrained Optimization and Lagrange Multiplier Methods, Academic Press.
- 58 Gillis, N. and Glineur, F. (2012) Accelerated multiplicative updates and hierarchical ALS algorithms for nonnegative matrix factorization. *Neural Computation*, 24 (4), 1085–1105.
- 59 Sha, F., Lin, Y., Saul, L., and Lee, D. (2007) Multiplicative updates for nonnegative quadratic programming. *Neural Computation*, **19** (8), 2004–2031.
- 60 Gillis, N. (2014) The why and how of nonnegative matrix factorization, in *Regularization, Optimization, Kernels, and Support Vector Machines* (eds J. Suykens, M. Signoretto, and A. Argyriou), Chapman & Hall/CRC, Machine Learning and Pattern Recognition Series, pp. 257–291.
- 61 Cichocki, A., Zdunek, R., and Amari, S.I. (2007) Hierarchical ALS Algorithms for Nonnegative Matrix and 3D Tensor Factorization, in *Lecture Notes in Computer Science, Vol. 4666, Springer*, pp. 169–176.
- 62 Cichocki, A. and Phan, A. (2009) Fast local algorithms for large scale Nonnegative Matrix and Tensor Factorizations. *IEICE Transactions on Fundamentals of Electronics*, Vol. E92-A No.3, 708–721.
- 63 Li, L. and Zhang, Y.J. (2009) FastNMF: highly efficient monotonic fixed-point nonnegative matrix factorization algorithm with good applicability. *J. Electron. Imaging*, Vol. 18 (033004).
- **64** Liu, J., Liu, J., Wonka, P., and Ye, J. (2012) Sparse non-negative tensor factorization using columnwise coordinate descent. *Pattern Recognition*,

45 (1), 649-656.

- 65 Hsieh, C.J. and Dhillon, I. (2011) Fast coordinate descent methods with variable selection for non-negative matrix factorization, in *Proc. of the 17th ACM SIGKDD int. conf. on Knowledge discovery and data mining*, pp. 1064–1072.
- **66** Vandaele, A., Gillis, N., Glineur, F., and Tuyttens, D. (2015) Heuristics for exact nonnegative matrix factorization. *Journal of Global Optimization*.
- **67** Kim, J., He, Y., and Park, H. (2013) Algorithms for nonnegative matrix and tensor factorizations: A unified view based on block coordinate descent framework. *Journal of Global Optimization*.
- 68 Zymnis, A., Kim, S.J., Skaf, J., Parente, M., and Boyd, S. (2007) Hyperspectral image unmixing via alternating projected subgradients, in *Signals, Systems and Computers, 2007*, pp. 1164 –1168.
- **69** Iordache, M.D., Bioucas-Dias, J., and Plaza, A. (2011) Total variation regulatization in sparse hyperspectral unmixing, in *Third Worskshop on Hyperspectral Image and Signal Processing: Evolution in Remote Sensing* (WHISPERS), Lisbon.
- 70 Cai, D., He, X., Han, J., and Huang, T. (2011) Graph regularized nonnegative matrix factorization for data representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33 (8), 1548–1560.
- 71 Rapin, J., Bobin, J., Larue, A., and Starck, J.L. (2013) Sparse and non-negative BSS for noisy data. *IEEE Transactions on Signal Processing*, **61** (22), 5620–5632.
- 72 Condat, L. (2015) Fast projection onto the simplex and the ℓ₁ ball. *Mathematical Programming*.
- 73 Hoyer, P. (2004) Nonnegative matrix factorization with sparseness constraints. *Journal of Machine Learning Research*, 5, 1457–1469.
- 74 Arora, S., Ge, R., Kannan, R., and Moitra, A. (2012) Computing a nonnegative matrix factorization – provably, in *Proc.* of the 44th Symp. on Theory of Computing (STOC '12), pp. 145–162.
- 75 Gillis, N. and Vavasis, S. (2015)

Semidefinite programming based preconditioning for more robust near-separable nonnegative matrix factorization. *SIAM J. on Optimization*, **25** (1), 677–698.

- 76 Kumar, A. and Sindhwani, V. (2015) Near-separable non-negative matrix factorization with ℓ₁- and Bregman loss functions, in SIAM Int. Conf. on Data Mining.
- 77 Gillis, N. and Ma, W.K. (2015) Enhancing pure-pixel identification performance via preconditioning. *SIAM J. on Imaging Sciences*, 8 (2), 1161–1186.
- 78 Arora, S., Ge, R., Halpern, Y., Mimno, D., Moitra, A., Sontag, D., Wu, Y., and Zhu, M. (2013) A practical algorithm for topic modeling with provable guarantees, in *Int. Conf. on Machine Learning (ICML '13)*, vol. 28, pp. 280–288.
- **79** Li, Y., Sima, D., Van Cauter, S., Croitor Sava, S., Himmelreich, U., Pi, Y., and Van Huffel, S. (2012) Hierarchical non-negative matrix factorization (hnmf): a tissue pattern differentiation method for glioblastoma multiforme diagnosis using mrsi. *NMR in Biomedicine*, **26** (3), 307–319.
- **80** Kuang, D. and Park, H. (2013) Fast rank-2 nonnegative matrix factorization for hierarchical document clustering, in 19th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '13), pp. 739–747.
- 81 Gillis, N., Kuang, D., and Park, H. (2015) Hierarchical clustering of hyperspectral images using rank-two nonnegative matrix factorization. *IEEE Transactions* on Geoscience and Remote Sensing, 53 (4), 2066–2078.
- 82 Gillis, N. (2012) Sparse and unique nonnegative matrix factorization through data preprocessing. *Journal of Machine Learning Research*, 13 (1), 3349–3386.
- 83 Kumar, N., Moussaoui, S., Idier, J., and Brie, D. (2015) Impact of sparse representation on the admissible solutions of spectral unmixing by non-negative matrix factorization, in *Proc. IEEE Workshop on Hyperspectral Image and Signal Processing (WHISPERS)*, Tokyo, Japan.
- 84 Henrot, S., Soussen, C., Dossot, M., and

Brie, D. (2014) Does deblurring improve geometrical hyperspectral unmixing? *IEEE Transactions on Image Processing*, **23** (3), 1169–1180.

- 85 Miron, S., Dossot, M., Carteret, C., Margueron, S., and Brie, D. (2011) Joint processing of the parallel and crossed polarized raman spectra and uniqueness in blind nonnegative source separation. *Chemometrics and Intelligent Laboratory Systems*, 105 (1), 7–18.
- **86** Long, D.A. (2002) *The Raman effect: a* unified treatment of the theory of Raman scattering by molecules, Wiley.
- **87** Turrell, G. (1972) *Infrared and Raman spectra of crystals*, New York, Academic Press.
- 88 Jiménez, C., Caroff, T., Bartasyte, A., Margueron, S., Abrutis, A., Chaix-Pluchery, O., and Weiss, F. (2009) Raman study of ceo2 texture as a buffer layer in the ceo2/la2zr2o7/ni architecture for coated conductors. *Applied spectroscopy*, 63 (4), 401–406.
- 89 Tanaka, M. and Young, R. (2006) Review polarised raman spectroscopy for the study of molecular orientation distributions in polymers. *Journal of Materials Science*, 41 (3), 963–991.
- **90** Jaumot, J., Menezes, J.C., and Tauler, R. (2006) Quality assessment of the results obtained by multivariate curve resolution analysis of multiple runs of gasoline blending processes. *Journal of chemometrics*, **20** (1-2), 54–67.
- **91** Ruckebusch, C., De Juan, A., Duponchel, L., and Huvenne, J. (2006) Matrix augmentation for breaking rank-deficiency: A case study. *Chemometrics and intelligent laboratory systems*, **80** (2), 209–214.
- 92 Lee, D. and Seung, H. (2000) Algorithms for non-negative matrix factorization, in Advances on Neural Information Processing Systems 13, (NIPS'2000), MIT Press, pp. 556–562.
- 93 Moussaoui, S., Brie, D., Mohammad-Djafari, A., and Carteret, C. (2006) Separation of non-negative mixture of non-negative sources using a bayesian approach and MCMC sampling. *IEEE Transactions on Signal Processing*, 54 (11), 4133–4145.

- 94 Craig, M. (1994) Minimum-volume transforms for remotely sensed data. *IEEE Transactions on Geoscience and Remote Sensing*, 32 (3), 542 –552.
- **95** Winter, M. (1999) N-findr: an algorithm for fast autonomous spectral end-member determination in hyperspectral data, in *Proc. SPIE Conf.on Imaging Spectrometry* V.
- **96** Chang, C.I. (2007) *Hyperspectral Data Exploitation. Theory and Applications*, Wiley Interscience.
- 97 Heinz, D. and Chein-I-Chang (2001) Fully constrained least squares linear spectral mixture analysis method for material quantification in hyperspectral imagery. *Geoscience and Remote Sensing*, *IEEE Transactions on*, 39 (3), 529 –545.
- 98 Chouzenoux, E., Legendre, M., Moussaoui, S., and Idier, J. (2014) Fast constrained least squares spectral unmixing using primal-dual interior-point optimization. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, PP (99), 1–11.
- **99** Miao, L. and Qi, H. (2007) Endmember extraction from highly mixed data using minimum volume constrained nonnegative matrix factorization. *IEEE Transactions on Geoscience and Remote Sensing*, **45** (3), 765–777.
- 100 Chang, C.I. and Du, Q. (2004) Estimation of number of spectrally distinct signal sources in hyperspectral imagery. *Geoscience and Remote Sensing, IEEE Transactions on*, 42 (3), 608 – 619.
- 101 Bioucas-Dias, J., Plaza, A., Dobigeon, N., Parente, M., Du, Q., Gader, P., and Chanussot, J. (2012) Hyperspectral unmixing overview: Geometrical, statistical, and sparse regression-based approaches. Selected Topics in Applied Earth Observations and Remote Sensing, IEEE Journal of, 5 (2), 354 –379.
- 102 Henrot, S., Soussen, C., and Brie, D. (2013) Fast positive deconvolution of hyperspectral images. *IEEE Trans. Image Process*, 22 (2), 828–833.
- 103 De Grauw, C., Sijtsema, N., Otto, C., and Greve, J. (1997) Axial resolution of confocal raman microscopes: Gaussian beam theory and practice. *Journal of Microscopy*, 188 (3), 273–279.

- 40
- 104 Wang, L.Y., Zhang, Y.H., and Zhao, L.J. (2005) Raman spectroscopic studies on single supersaturated droplets of sodium

and magnesium acetate. *The Journal of Physical Chemistry A*, **109** (4), 609–614.

Index

а

admissible NMF 12 alternating non-negative least squares 18

h

hierarchical alternating least squares 20

n NMF uniqueness 11, 12 NNLS 16 non-negative rank 2

р

positive matrix factorization 2