# Decontamination of galaxy spectra using four dispersion directions

Mostafa Bella[1,2], Shahram Hosseini[1], Hicham Saylani[2], Thierry Contini[1], Tristan Grégoire[1],
Andréa Guerrero[1], Yannick Deville[1] on behalf of the Euclid Consortium
[1]IRAP, Université de Toulouse, UPS, CNRS, CNES, 14 Av. Edouard Belin, 31400 Toulouse, France
[2]LETSMP, Faculté des Sciences, Université Ibn Zohr, BP 8106 Cité Dakhla, Agadir, Maroc
Email: {mbella, shosseini, Thierry.Contini, tgregoire, Yannick.Deville}@irap.omp.eu;
h.saylani@uiz.ac.ma; andrea.guerrero@altran.com

*Abstract*—This paper addresses the problem of spectra decontamination in slitless spectroscopy in the context of the *Euclid* space mission. This problem can be treated as a source separation problem which consists of estimating a set of unknown sources from a set of their mixtures. We first present a mixing model linking observed data to source spectra that simultaneously takes into account four light dispersion directions of the grism, then we propose two methods to decontaminate the spectra from the mixed data. Exploitation of all these dispersion directions improves the estimation of the spectrum of an object of interest. Preliminary results obtained using realistic noisy data show the effectiveness of the proposed methods.

*Index Terms*—Source separation, *Euclid* mission, Data concatenation, MPDR beamforming, Slitless spectroscopy.

Fig. 1. Contamination of neighboring object spectra.

## I. INTRODUCTION

*Euclid* is a space telescope of the European Space Agency [1], scheduled for launch in 2023. Its main task is to understand the nature of dark energy and how this energy is responsible for the increasing acceleration of the Universe expansion. *Euclid* will be equipped with a slitless near-infrared spectrograph that will measure the spectra of more than 50 million galaxies. These spectra will then be analyzed to estimate the galaxy redshifts, which should help to better understand how dark energy contributes to this acceleration. The spectroscopy is performed using several grisms which are a combination of prisms and diffraction gratings. A grism differently disperses different wavelengths of the emitted light in a dispersion direction (similar to a rainbow effect). Since an object usually extends over several pixels, the grism output is an image with several rows, which will be hereafter called *spectrogram*[1]. The 1-dimensional spectrum of the object may e.g. be obtained by summing up the rows of its spectrogram. However, as shown in Figure 1, the slitless spectroscopy used in *Euclid* leads to the superposition of spectrograms of neighboring astronomical objects (galaxies and stars), which could lead to redshift measurement errors and uncertainties [2].

[1]Astronomers call it a spectrogram because it is generated by an optical spectrograph. It should not be confused with the spectrogram resulting from the time-frequency analysis of a signal.
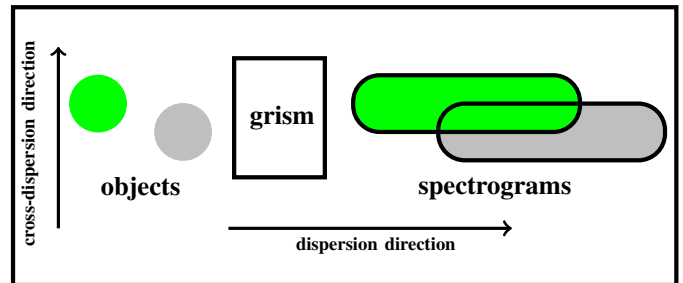
In order to overcome this problem, several spectrograms will be generated in different dispersion directions. As a result, the spectrum of an object of interest is contaminated by different spectra in different directions. Then, the exploitation of all these directions helps to better estimate the spectrum of the object of interest.

Spectra decontamination in slitless spectroscopy can be viewed as a source separation problem that consists of estimating a set of unknown signals, called sources, from a set of their mixtures [3]–[5]. Indeed, we showed in [6]–[8] that under certain assumptions, the contaminated spectrogram of an object of interest in each dispersion direction of the grism can be approximated by a linear instantaneous model. Unfortunately, most spectra are not mutually independent or sparse. Therefore, the blind source separation methods based on independent component analysis or sparse component analysis cannot be used. We proposed in [6], [7] a blind method based on Non-negative Matrix Factorization (NMF) [9]–[11] to decontaminate the spectra, by only considering the spectrograms corresponding to the two dispersion directions 0 and 90 degrees. However, the Euclid consortium recently adopted another observation strategy which excludes the use of a grism in the 90 degree direction. The new strategy consists of generating spectrograms in four dispersion directions, namely 0, 180, 184 and −4 degrees. As a result, our previous proposed methods are no longer applicable to this new strategy. Moreover, it is well known that the NMF-based methods are very sensitive to initialization and do not guarantee the uniqueness of the solution, which is primary in a space project like Euclid.

In this paper, we propose two new methods exploiting spectral data provided in these 4 directions, together with direct images provided by the *Euclid* near infrared photometers. At first, the direct images are used to get an estimate of the mixing coefficients, then the spectrum of the object of interest is estimated from them and observed contaminated spectrograms.

## II. PROPOSED METHODS

### A. Mixing model

We have shown in [6] that for each dispersion direction of the grism, denoted $d_i$, the noiseless mixture containing the contaminated spectrogram of an object of interest (i.e. the target object to be decontaminated) can be approximated by a linear instantaneous model of the form:

$$\mathbf{X}^{(d_i)} = \mathbf{A}^{(i)} \cdot \mathbf{E}^{(i)} + \mathbf{C}^{(i)}, \tag{1}$$

where $\mathbf{X}^{(d_i)}$ is the matrix of size $M_i \times K$ containing the observed data, with $M_i$ the number of rows in the cross-dispersion direction associated to spectrogram of the object of interest and $K$ the number of spectral bands considered identical for all dispersion directions. $\mathbf{A}^{(i)} = \left[ \mathbf{a}_s^{(i)} | \mathbf{A}_c^{(i)} \right]$ is the mixing matrix of size $M_i \times (N_i + 1)$, defined by:

$$\mathbf{A}^{(i)} = \left( \begin{array}{c|ccc} \overbrace{a_s^{(i)}(1)}^{\mathbf{a}_s^{(i)}} & \overbrace{a_{c_1}^{(i)}(1)}^{\mathbf{A}_c^{(i)}} & \cdots & a_{c_{N_i}}^{(i)}(1) \\ \vdots & \vdots & \ddots & \vdots \\ a_s^{(i)}(M_i) & a_{c_1}^{(i)}(M_i) & \cdots & a_{c_{N_i}}^{(i)}(M_i) \end{array} \right) \tag{2}$$

where $N_i$ is the number of the considered contaminant objects in direction $d_i$, $a_s^{(i)}(m)$ and $a_{c_n}^{(i)}(m)$, $(m, n) \in [1, M_i] \times [1, N_i]$, are respectively the mixing coefficients of the object of interest, constituting the vector $\mathbf{a}_s^{(i)} = [a_s^{(i)}(1), ..., a_s^{(i)}(M_i)]^T$, and the contaminants constituting the matrix $\mathbf{A}_c^{(i)}$. $\mathbf{E}^{(i)}$ is the source matrix of size $(N_i + 1) \times K$ whose first row corresponds to the spectrum of the object of interest, whereas its other rows contain 1D spectra of the contaminating objects in direction $d_i$. $\mathbf{C}^{(i)}$ is a known matrix which corresponds to the sky background that we do not consider in this work since it will be subtracted from the observations before estimating the spectra.

As mentioned in Section I, for each object of interest we have four observations for the four dispersion directions (0, 180, 184, and $-4$ degrees). Assuming that the spectrum of the object of interest is the same whatever the dispersion direction, we can merge these four observations in order to improve the estimation of the spectrum of this object. We therefore redefine the observation matrix $\mathbf{X}$ of an object of interest, whose spectrum we want to decontaminate, by

merging its contaminated spectrograms in the four directions in the following manner[2]:

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}^{(d_1)} \\ \mathbf{X}^{(d_2)} \\ \mathbf{X}^{(d_3)} \\ \mathbf{X}^{(d_4)} \end{bmatrix} = \begin{bmatrix} \mathbf{X}^{(0)} \\ \mathbf{X}^{(180)} \\ \mathbf{X}^{(184)} \\ \mathbf{X}^{(-4)} \end{bmatrix} = \mathbf{A} \cdot \mathbf{E}, \tag{3}$$

where $\mathbf{A} = [\mathbf{a}_s | \mathbf{A}_c]$ is the total mixing matrix of size $M \times (N+1)$, with $M = \sum_{i=1}^4 M_i$, $N = \sum_{i=1}^4 N_i$, $\mathbf{a}_s$ is the mixing vector corresponding to the object of interest, defined by :

$$\mathbf{a}_s = [\mathbf{a}_s^{(1)^T}, \mathbf{a}_s^{(2)^T}, \mathbf{a}_s^{(3)^T}, \mathbf{a}_s^{(4)^T}]^T, \tag{4}$$

and $\mathbf{A}_c$ is the contaminant mixing matrix. Assuming that the contaminants in each dispersion direction are distinct, this matrix can then be defined by :

$$\mathbf{A}_c = \begin{bmatrix} \mathbf{A}_c^{(1)} & 0 & 0 & 0 \\ 0 & \mathbf{A}_c^{(2)} & 0 & 0 \\ 0 & 0 & \mathbf{A}_c^{(3)} & 0 \\ 0 & 0 & 0 & \mathbf{A}_c^{(4)} \end{bmatrix}. \tag{5}$$

$\mathbf{E}$ is the matrix of size $(N + 1) \times K$ containing the spectrum of the object of interest in its first row, and the contaminants in each of the four directions in the following rows.

In the next two subsections, we will present two methods that can be used to estimate the spectrum of the object of interest by exploiting the observation matrix $\mathbf{X}$ and the information available on the objects.

### B. First method

The first step of this method is to estimate the total mixing matrix $\mathbf{A}$ using the direct photometric image of the object of interest and those of its contaminants. Indeed, *Euclid* is also equipped with a photometer which provides direct images of all astronomical objects in the field of view [2]. As can be seen in Figure 1, the direct image flux is dispersed by a grism to produce its spectrogram. As a result, for each dispersion direction, the mixing coefficients ($a_s^{(i)}(m)$ or $a_{c_n}^{(i)}(m)$) related to an object (target or contaminant) may be estimated by calculating the sum of pixels on different rows of its direct image, which will allow us to estimate the total mixing matrix $\mathbf{A}$. To estimate these coefficients, we first oversample the photometric image of each object (target or contaminant) in the cross-dispersion direction, then recentre it in its spectrogram, and finally undersample it by the same sampling rate. This allows us to correct the offset between the photometric image and the spectrogram of this object. After this preprocessing, we choose the $M_i \times R$ pixel values of the image centred on the position of the object of interest, where $R$ is the number of columns in the direct image. Then, the sum of the pixels on the m-th row provides the value $a_s^{(i)}(m)$ or $a_{c_n}^{(i)}(m)$ of this object. This process is shown in Figure 2. The mixing matrix $\mathbf{A}$ is then constructed using the calculated coefficients.

---

[2]Note that before combining the spectrograms of the 4 directions, it is necessary to rotate the contaminated spectra $\mathbf{X}^{(180)}$, $\mathbf{X}^{(184)}$ and $\mathbf{X}^{(-4)}$ respectively by rotations of 180, 184 and $-4$ degrees.
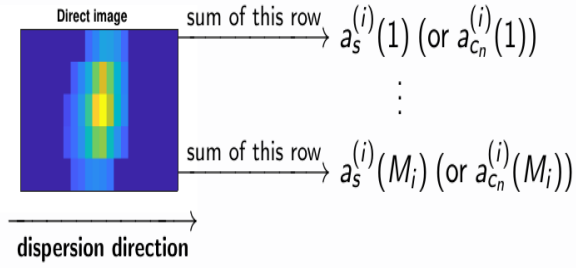
Fig. 2. Mixing coefficient estimation using photometric images.

After constructing the mixing matrix $\mathbf{A}$, a first solution to estimate the source matrix $\mathbf{E}$ consists of minimizing the following criterion, where $||.||_2$ stands for the Frobenius norm:

$$J_1 = ||\mathbf{X} - \mathbf{A}\mathbf{E}||_2^2, \qquad (6)$$

which leads to the following Least SQuares (LSQ) estimate:

$$\mathbf{E} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{X}. \qquad (7)$$

A second solution to estimate $\mathbf{E}$ consists of using the Non-Negative Least SQuares (NNLSQ) method presented in [12]. Indeed, since the spectra to be estimated are by definition non-negative, this method takes into account this non-negativity through the minimization of the criterion.:

$$J_2 = ||\mathbf{X} - \mathbf{A}\mathbf{E}||_2^2 \quad \text{s.t.} \quad \mathbf{E} \geq 0 \qquad (8)$$

where $\mathbf{E} \geq 0$ means that all values of the matrix $\mathbf{E}$ are non-negative. Finally, the spectrum of the object of interest, denoted $\hat{\mathbf{E}}_s$, is the first row of the estimated matrix $\hat{\mathbf{E}}$.

### C. Second method

Unlike the method presented in Subsection II-B which requires the direct photometric image of the object of interest and those of all its contaminants, the method proposed in this Subsection requires only the image of the object of interest in order to estimate its spectrum. Indeed, this image makes it possible to estimate the mixing coefficients $a_s^{(i)}(m)$ of the object of interest in all dispersion directions, which will allow us to estimate the mixing vector $\mathbf{a}_s$ of this object by the same procedure as described in Subsection II-B. This vector will then be exploited by a beamformer to estimate the spectrum of the object of interest. The aim of the beamformer is to enhance the target object while attenuating interferers and noise originating from other positions [13], [14]. In this paper, we are interested in the so-called Minimum Power Distortion-less Response (MPDR) [14] beamformer. This beamformer aims at estimating an optimal filter denoted $\mathbf{w}_{\text{MPDR}}$, whose output minimizes the total output power under the constraint $\mathbf{w}^H\mathbf{a}_s = 1$. The desired filter is the solution of the following minimization problem:

$$\mathbf{w}_{\text{MPDR}} = \underset{\mathbf{w}}{\arg\min} \left\{\mathbf{w}^H\mathbf{R}_\mathbf{X}\mathbf{w}\right\} \quad \text{s.t.} \quad \mathbf{w}^H\mathbf{a}_s = 1, \qquad (9)$$

which leads to the following beamforming coefficients[3] [14]:

$$\mathbf{w}_{\text{MPDR}} = \frac{\mathbf{R}_\mathbf{X}^{-1}\mathbf{a}_s}{\mathbf{a}_s^H\mathbf{R}_\mathbf{X}^{-1}\mathbf{a}_s}, \qquad (10)$$

where $\mathbf{R}_\mathbf{X} = \mathbf{X}\mathbf{X}^H/K$ is the covariance matrix of observations.

Even though MPDR beamforming can be used directly with the total observation matrix $\mathbf{X}$ to estimate the spectrum of the object of interest, we here propose to split the observation matrix $\mathbf{X}$ into $L$ parts, where each part denoted $\mathbf{X}_l$ includes all the rows of the matrix $\mathbf{X}$ and a set of columns of this matrix, and then decontaminate each of these parts separately. Indeed, by adopting this approach, we minimize the risk of falling into the underdetermined case that MPDR beamforming cannot handle well [14]. In fact, each contaminant object usually contaminates only a part of the spectrogram of the target object (see Figure 1), so that the number of contaminants for each of the L parts is lower than their number for the entire spectrum. To do this, we generate $L$ beamformers, each one is given by:

$$\mathbf{w}_{\text{MPDR}}^{(l)} = \frac{\mathbf{R}_{\mathbf{X}_l}^{-1}\mathbf{a}_s}{\mathbf{a}_s^H\mathbf{R}_{\mathbf{X}_l}^{-1}\mathbf{a}_s}, \quad l \in [1, L] \qquad (11)$$

where $\mathbf{R}_{\mathbf{X}_l}$ is the covariance matrix associated to the $l-$th part of the observation matrix $\mathbf{X}$. Finally, the outputs $\hat{\mathbf{E}}_s^{(l)} = \mathbf{w}_{\text{MPDR}}^{(l)H}\mathbf{X}_l$ of each beamformer are concatenated to yield the final estimate of the spectrum of the object of interest as follows:

$$\hat{\mathbf{E}}_s = [\hat{\mathbf{E}}_s^{(1)}, \hat{\mathbf{E}}_s^{(2)}, ..., \hat{\mathbf{E}}_s^{(L)}]. \qquad (12)$$

It should be noted that this second method can attenuate even unmodeled interference such as undetected object spectra or hot pixels, and should lead to better results under these conditions.

### III. TEST RESULTS

To evaluate the performance of our methods, we performed tests on two different scenarios, using realistic simulated observed data, provided by the Euclid consortium. A difference between the two scenarios concerns the presence or not of hot pixels. Note that the observed data are also affected by a strong noise that comes mainly from the acquisition instruments. To measure the decontamination performance we used as a criterion the *Normalized Root-Mean-Square Error* (NRMSE) before and after decontamination, defined by:

$$\text{NRMSE}_{\text{in}} = \frac{1}{4} \sum_{d_i=\{0,180,184,-4\}} \frac{||\mathbf{E}_s - \mathbf{E}_x^{(d_i)}||_2}{||\mathbf{E}_s||_2}, \qquad (13)$$

$$\text{NRMSE}_{\text{out}} = \frac{||\mathbf{E}_s - \hat{\mathbf{E}}_s||_2}{||\mathbf{E}_s||_2}, \qquad (14)$$

where $\mathbf{E}_x^{(d_i)}$ is the 1D spectrum of the observed data in direction $d_i$, $\mathbf{E}_s$ is the true noiseless spectrum of the object of interest and $\hat{\mathbf{E}}_s$ its estimate. In (13) and (14), all spectra are

---

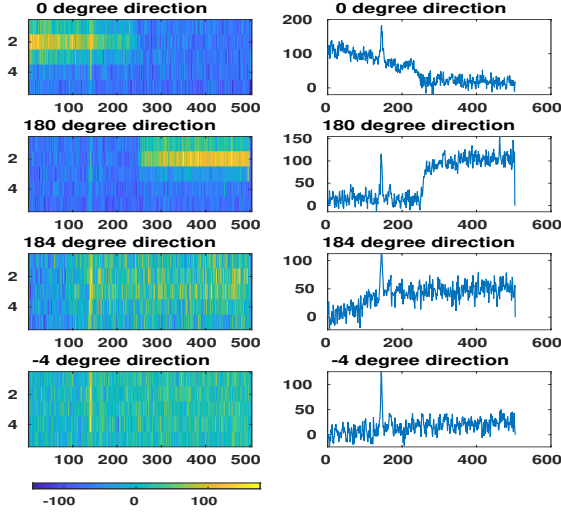[3]Note that it is necessary to centre the observation matrix $\mathbf{X}$ before performing beamforming.

Fig. 3. Observed noisy spectrograms (left) and corresponding 1D spectra (right) in the first scenario.

centered, and $\hat{\mathbf{E}}_s$ and $\mathbf{E}_x^{(d_i)}$ are normalized to have the same variance as $\mathbf{E}_s$. In both scenarios, the proposed methods were compared to the randomly initialized NMF Alternating Least Square (NMF-ALS) method presented in [15]. In the second method, the number of beamformers $L$ was set to 3. We choose observation areas of $5 \times 504$ pixels, 504 corresponding to the number of wavelengths contained in the spectrogram, and 5 corresponding to the spectrogram width in the cross-dispersion direction.

### A. A simple scenario without hot pixels

In our first experiment, we chose a simple scenario where the spectrogram of an object of interest (a galaxy with a redshift of 1.12) is contaminated by the spectra of four objects in the two directions "0" degree and "180" degree, by the spectra of two other objects in the "184"degree direction, and by the spectrum of another object in the "-4" degree direction. The observed spectrograms and corresponding 1D spectra[4] concerning this object are shown in Figure 3.

The true noiseless spectrum of the object of interest and its estimates using our methods presented in Section II and the NMF-ALS method after normalization by their maximum are shown in Figure 4. This figure clearly shows the effectiveness of the proposed methods in decontaminating the spectrum of the object of interest. Indeed, our methods successfully remove all contaminations of other objects in this first scenario. The NMF-ALS method has rather poor performance as we can see in Figure 4. The numerical results are shown in Table I. Note that although the NRMSE$_{out}$ of the proposed methods seems large in this scenario, the NRMSE$_{in}$ of the observed data is equal to 10.67, which is much larger compared to those obtained with our methods. As can be seen, our first method (using NNLSQ) yields the best results in this scenario.

---

[4]The observed 1D spectrum is here defined as the mean of the rows of the spectrogram.

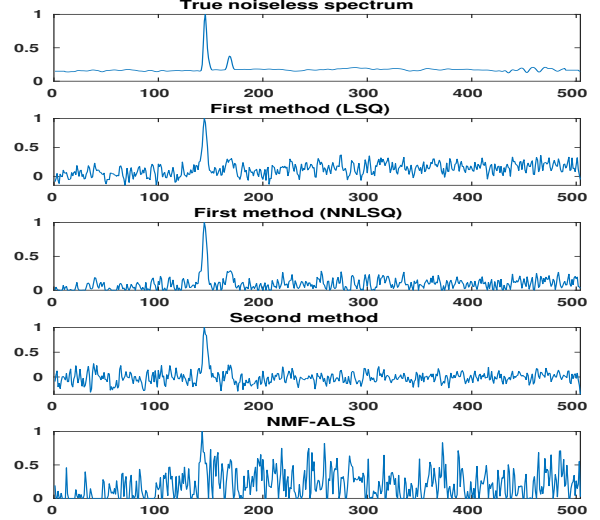| Scenario | NRMSE$_{in}$ | NRMSE$_{out}$ | | | |
|---|---|---|---|---|---|
| | | LSQ | LSQN | Second method | NMF-ALS |
| 1 | 10.67 | 0.90 | **0.73** | 0.92 | 1.19 |
| 2 | 12.74 | 1.26 | 1.09 | **0.97** | 1.38 |



Fig. 4. True noiseless spectrum of the object of interest in the first scenario and its estimates using the three methods.

### B. A complicated scenario with hot pixels

In our second experiment, we chose a complicated scenario where the spectrogram of an object of interest (a galaxy with a redshift of 1.41) is contaminated by the spectra of 8 objects in the "0" degree direction, 7 objects in the "180" degree direction, 3 objects in the "184" degree direction, and one object in the "-4" degree direction. Moreover, there is at least one hot pixel in each dispersion direction, which further complicates this scenario. The observed spectrograms and corresponding 1D spectra concerning this object are shown in Figure 5, where the isolated pics in 1D spectra are due to hot pixels.

Figure 6 shows the true noiseless spectrum of the object of interest and its estimates using the above-mentioned methods after normalization by their maximum. As can be seen in Figure 6, both the first method (using NNLSQ) and the second one successfully remove contamination of other objects as well as all hot pixels. The first method (using LSQ) is also successful in removing contamination from other objects, however, it is not able to remove all hot pixels. As can be seen in Table I, our second method yields the best results in this scenario, which is expected because the second method should work better in the presence of unmodeled interference such as hot pixels. Finally, the NMF-ALS method failed to decontaminate the spectrum of the object of interest in this scenario, which further confirms the effectiveness of our proposed methods.
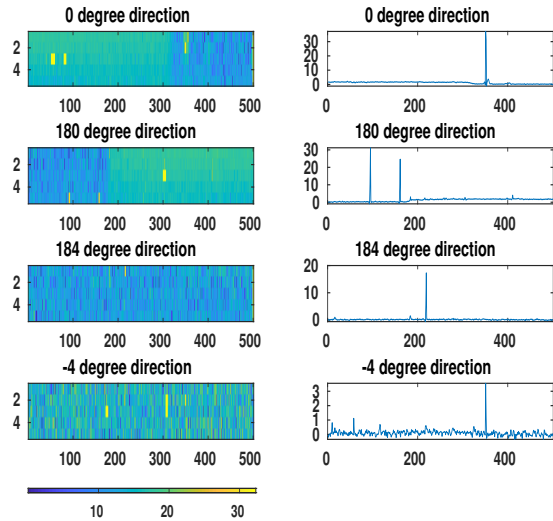
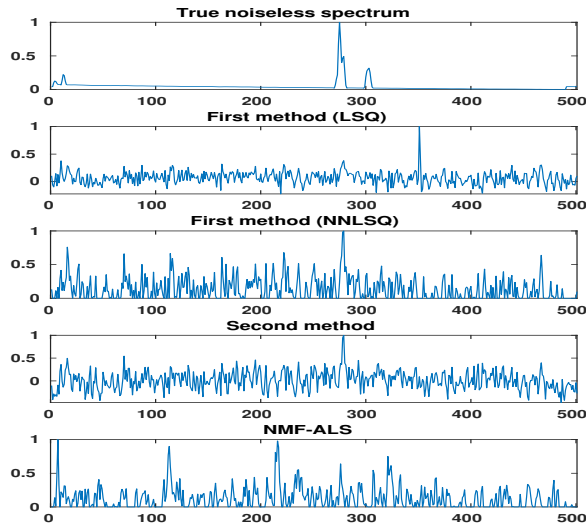Fig. 5. Observed noisy spectrograms (left) and corresponding 1D spectra (right) in the second scenario.



Fig. 6. True noiseless spectrum of the object of interest in the second scenario and its estimates using the three methods.

## IV. CONCLUSION

In this paper, we proposed two new decontamination methods that are based on data concatenation between four dispersion directions of the grism and which are adapted to *Euclid*'s new observation strategy. Indeed, exploiting all these observations allows us to improve the estimation of the spectrum of the object of interest. According to the results of the performed tests, the first method yields the best results (in terms of NRMSE) for objects without hot pixels, while the second method yields the best results for objects with hot pixels. Nevertheless, a more complete performance study would be desirable to confirm this result.

In terms of future work, it would be interesting to replace the local approach used in this paper, where decontamination is performed object by object, by a global approach where all the objects in the field of view are decontaminated at once. Also, the noise covariance matrix has not been taken into account in this paper. If available, this matrix can be used to improve the performance of our two methods. Finally, in our work, we have adopted the linear-instantaneous mixing model which is a simple but approximate model. A more realistic convolutive model can improve the estimation quality.

### REFERENCES

[1] "Euclid consortium page," https://www.euclid-ec.org, 2021.
[2] "Euclid definition study report," https://arxiv.org/pdf/1110.3193.pdf, 2011.
[3] Y. Deville, "Blind source separation and blind mixture identification methods," 2016, pp. 1–33, Wiley Encyclopedia of Electrical and Electronics Engineering.
[4] A. Hyvärinen, J. Karhunen, and E. Oja, "Independent component analysis," 2001, John Wiley and Sons Inc.
[5] C. Jutten and P. Comon, "Handbook of blind source separation: Independent component analysis and applications," 2010, Academic Press, Oxford.
[6] A. Selloum, S. Hosseini, Y. Deville, and T. Contini, "Mixing model in slitless spectroscopy and resulting blind methods for separating galaxy spectra," in 2016 IEEE 26th International Workshop on Machine Learning for Signal Processing (MLSP), 2016, pp. 1–6.
[7] S. Hosseini, A. Selloum, T. Contini, and Y. Deville, "Separation of galaxy spectra measured with slitless spectroscopy," Digital Signal Processing, vol. 106, pp. 102837, 2020.
[8] A. Guerrero, "Méthodes de séparation aveugle de sources et application à l'imagerie hyperspectrale en astrophysique et observation de la terre," 2019, PhD thesis.
[9] A. Cichocki, R. Zdunek, A. Phan, and S. Amari, "Nonnegative matrix and tensor factorizations: Applications to exploratory multi-way data analysis and blind source separation," October 2009.
[10] X. Fu, K. Huang, N.D. Sidiropoulos, and W. Ma, "Nonnegative matrix factorization for signal and data analytics: Identifiability, algorithms, and applications," IEEE Signal Processing Magazine, vol. 36, no. 2, pp. 59–80, 2019.
[11] Lee D.D. and H.S. Seung, "Learning the parts of objects by non-negative matrix factorization," Nature, vol. 401, pp. 788–791, 1999.
[12] R.J. Lawson, C.L. Hanson, "Solving least squares problems," 1974.
[13] J. Capon, "High-resolution frequency-wavenumber spectrum analysis," Proceedings of the IEEE, vol. 57, pp. 1408–1418, August 1969.
[14] H.L. Van Trees, "Adaptive beamformers," in Optimum Array Processing. John Wiley and Sons, Ltd, 2002, pp. 710–916.
[15] P. Paatero and U. Tapper, "Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values," Environmetrics, vol. 5, pp. 111–126, 1994.