



HAL
open science

Machine Learning for automatic identification of new minor species

Frédéric Schmidt, Guillaume Cruz Mermey, Justin Erwin, Séverine Robert, Lori Neary, Ian Thomas, Frank Daerden, Bojan Ristic, Manish Patel, Giancarlo Bellucci, et al.

► **To cite this version:**

Frédéric Schmidt, Guillaume Cruz Mermey, Justin Erwin, Séverine Robert, Lori Neary, et al.. Machine Learning for automatic identification of new minor species. *Journal of Quantitative Spectroscopy and Radiative Transfer*, 2021, 259, pp.107361. 10.1016/j.jqsrt.2020.107361 . hal-02965570

HAL Id: hal-02965570

<https://hal.science/hal-02965570v1>

Submitted on 21 Nov 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

Machine Learning for automatic identification of new minor species

Frédéric Schmidt¹, Guillaume Cruz Mermy¹, Justin Erwin², Séverine Robert²,
Lori Neary², Ian R. Thomas², Frank Daerden², Bojan Ristic², Manish R.
Patel³, Giancarlo Bellucci⁴, Jose-Juan Lopez-Moreno⁵, Ann-Carine Vandaele²

¹ *Université Paris-Saclay, CNRS, GEOPS, 91405, Orsay, France,* ² *Belgian Institute for Space Aeronomy (BIRA-IASB), Avenue Circulaire, 3 B-1180 Brussels Belgium,* ³ *School of Physical Sciences, The Open University, Milton Keynes, MK7 6AA, U.K.,* ⁴ *INAF-Istituto di Astrofisica e Planetologia Spaziali, Rome, ITALY,* ⁵ *Instituto de Astrofisica de Andalucía CSIC*

Abstract

One of the main difficulties to analyze modern spectroscopic datasets is due to the extremely large amount of data. For example, in atmospheric transmittance spectroscopy, the solar occultation channel (SO) of the NOMAD instrument onboard the ESA ExoMars2016 satellite called Trace Gas Orbiter (TGO) had produced ~ 10 millions of spectra in ~ 20000 acquisition sequences since the beginning of the mission in April 2018 until 15 January 2020. Other datasets are even larger with \sim billions of spectra for OMEGA onboard Mars Express or CRISM onboard Mars Reconnaissance Orbiter. Usually, new lines are discovered after a long iterative process of model fitting and manual residual analysis. Here we propose a new method based on unsupervised machine learning, to automatically detect new minor species. Although precise quantification is out of scope, this tool can also be used to quickly summarize the dataset, by giving few endmembers ("source") and their abundances.

The methodology is the following: we proposed a way to approximate the dataset non-linearity by a linear mixture of abundance and source spectra (endmembers). We used unsupervised source separation in form of non-negative matrix factorization to estimate those quantities. Several methods are tested on synthetic and simulation data. Our approach is dedicated to detect minor species spectra rather than precisely quantifying them. On synthetic example, this approach is able to detect chemical compounds present in form of 100 hidden spectra out of 10^4 , at 1.5 times the noise level. Results on simulated spectra of NOMAD-SO targeting CH_4 show that detection limits goes in the range of 100-500 ppt in favorable conditions. Results on real martian data from NOMAD-SO show that CO_2 and H_2O are present, as expected, but CH_4 is absent. Nevertheless, we confirm a set of new unexpected lines in the database, attributed by ACS instrument Team to the CO_2 magnetic dipole.

Keywords: spectroscopy, atmosphere, data mining, machine learning, unsupervised, source separation, non-negative matrix factorization

1. Introduction

In modern exploration science, one has to face a major challenge : how to learn something new from analyzing a large dataset collection while taking into account what we already know. If the current knowledge outweighs the analysis, the discovery of new elements may be difficult. Usually, in the field of spectroscopy, one can compare laboratory spectra, model and observation spectra. Going back and forth leads to discovery of new lines by identifying unexpected residuals in the observation data (not expected by the model). Sometimes, initial identification of lines can be wrong. As an example, spectroscopic evidence of atmospheric CO₂ ice cloud was reported after the discovery of an emission spike at a wavelength of 4.3 μm from Mariner 6 and 7 infrared probings of the bright martian limb (Herr and Pimentel, 1970), but this spectral feature was mistaken for a resonant scattering band of CO₂ fluorescence (López-Valverde et al., 2005).

For one single spectrum, one can use simulation algorithm (see for instance Faisal et al., 2020). For large datasets, simplest ideas would be to scrutinize average spectra, or potential band depth distribution. Unfortunately, in the case of low signal-to-noise ratio (SNR, defined as signal / standard deviation of noise), such methods fail (as will be illustrated in the toy example). Analyzing residuals after modeling is a good method but it requires a lot of work.

Several statistical tools with various approaches have been proposed, such as the Principal Component Analysis (PCA) (Penttilä et al., 2018; Geminale et al., 2015), or Independent Component Analysis (ICA) (Shashilov et al., 2006; Erard et al., 2009), but most of them require a human operator to pick endmembers and trends since those methods are nothing more than a change of representation. Furthermore, none of these methods guarantees positivity of the *component* (which are sometimes also called *source*), which can be problematic during the interpretation. Recently, advanced machine learning methods based on non-negative matrix factorization have been proposed (Lee and Seung, 1999; Moussaoui et al., 2008; Dobigeon et al., 2009; Schmidt et al., 2010; Gillis and Glineur, 2012; Hinrich and Mørup, 2018). This approach is completely different from PCA/ICA: each source is positive and represents an endmember / a trend. A source is not one spectrum extracted from the dataset but a statistical reconstruction. By using this approach, the human operator doesn't have to identify endmembers/trends anymore, since they are automatically picked by the algorithm in form of source. Furthermore when there are statistical / spectral correlations between sources PCA/ICA fails because it assumes orthogonality / independence, which is not the case for non-negative matrix factorization.

Based on this new approach, we propose a tool:

- to give an overview and quickly summarize a large and complex spectroscopic dataset with simple variables

- 42 • to detect potential new spectroscopic features (unexpected minor species,
43 new absorption lines,...)
- 44 • to be performed in a fully blind way (without prior information on neither
45 the spectra, nor the abundances).

46 The target observation type of this study is solar occultation. This measure-
47 ment principle has been proposed as early as 1900, an interesting review was
48 published by Smith and Hunten (1990). Several recent instruments used this
49 technique to investigate the composition of the Earth's (SCIAMACHY/ENVISAT
50 Bovensmann et al. 1999), Mars' (SPICAM Bertaux et al., 2000) or Venus' at-
51 mospheres (SPICAV Bertaux et al., 2007). Here we will focus on the recent
52 NOMAD instrument (Vandaele et al., 2015), and especially the SO channel,
53 designed to study the Martian atmosphere and its trace gases, such as methane.
54 Indeed the presence of CH₄ on Mars is a very hot topic for the planetary science
55 community (Giuranna et al., 2019; Korabev et al., 2019; Moores et al., 2019).
56 In the present article, we propose to apply the tool for potential CH₄ detection.
57 Nevertheless, the approach can be extended to other types of spectroscopic
58 measurements.

59 2. Dataset

60 We propose here to focus on the Nadir and Occultation for MARS Discovery
61 (NOMAD) instrument onboard ESA's ExoMars Trace Gas Orbiter and espe-
62 cially the Solar Occultation (SO) channel (Vandaele et al., 2015). NOMAD is
63 a compact, high-resolution, dual channel IR spectrometer (SO and LNO) cou-
64 pled with a highly miniaturized UV-visible spectrometer (UVIS), capable of
65 operating in different observation modes: solar occultation, nadir and limb.

66 The SO channel operates at wavenumbers from 2320 cm⁻¹ to 4550 cm⁻¹
67 (wavelength 2.2 to 4.3 μm), using an echelle grating with a groove density of
68 4 lines/mm in a Littrow configuration in combination with an Acousto-Optic
69 Tunable Filter (AOTF) for spectral order selection. The width of the selected
70 spectral ranges varies from 20 to 35 cm⁻¹ depending on the selected diffrac-
71 tion order. The detector is an actively cooled HgCdTe Focal Plane Array. SO
72 achieves an instrument line profile resolution of 0.15 cm⁻¹, corresponding to a
73 resolving power $\lambda/\Delta\lambda$ of approximately 25000. All details of the instrument are
74 available in Neefs et al. 2015 and Vandaele et al. 2018. The orders with the
75 maximum sensitivity to CH₄ are: 119, 134 and 136. We will use the data from
76 the beginning of the mission in April 2018 until 15 January 2020, in calibration
77 version 1p0a. Due to temperature change, the spectral registration varies, pro-
78 ducing a shift up to ~10 spectels. We corrected it by aligning the full dataset to
79 a reference spectra (arbitrarily chosen with the maximum band depth of water)
80 by cross-correlation. No sub-spectel resampling has been performed but a sim-
81 ple shift. When the calibration will be improved, this step will most probably be
82 replaced by a routine correction. The data are available on the ESA/Planetary
83 Science Archive after a 6 months embargo period.

84 **3. Method**

85 In this section, we first describe the data pretreatment required for non-
86 negative matrix factorization purpose followed by the data mining method.

87 *3.1. Data pretreatment*

88 After calibration, the NOMAD SO spectra are in transmittance $T = I/I_0$,
89 depending on wavenumber ν , with I the observed light intensity through the
90 atmosphere and I_0 the solar spectra measured outside the atmosphere.

91 Assuming that the atmosphere is homogeneous, and that multiple scattering
92 and refraction are negligible (Smith and Hunten, 1990; Bovensmann et al., 1999),
93 the optical depth τ is a linear combination of $E(\nu)$ the total extinction, and ϵ
94 the slant column density, for each chemical species i :

$$\tau(\nu) = -\log T(\nu) \approx \sum_{i=1}^{N_S} E_i(\nu) \cdot \epsilon_i + MC(\nu) \quad (1)$$

95 with N_S , the total number of species and $MC(\nu)$ a modeled continuum
96 described below.

97 The slant column density ϵ is directly related to the total number of particles
98 $N(s)$ along the line of sight s :

$$\epsilon = \int N(s) ds \quad (2)$$

99 While the extinction by gas is usually highly structured, absorption by parti-
100 cles, scattering by molecules and particles, and also reflection at the surface are
101 broadband features. Such large features are modeled by a continuum $MC(\nu)$,
102 often taken as a polynomial, that is filtered out.

103 The problem with this continuum removal rationale is that when the optical
104 depth is large, the SNR is decreased and the noise effect on continuum removal
105 amplified (see Sup. Mat.).

106 Instead of using this rationale, we propose to first correct for the continuum
107 $C(\nu)$ in the transmittance space:

$$T^*(\nu) = T(\nu) - C(\nu) \quad (3)$$

108 Then convert the spectra into absorbance:

$$X(\nu) = 1 - T^*(\nu) \quad (4)$$

109 The final step is the linear mixture :

$$X(\nu) \approx \sum_{i=1}^{N_S} S_i(\nu) \cdot A_i \quad (5)$$

110 with $S_i(\nu)$ the source spectra and A_i the spectral abundance. In this de-
111 scription, the physical meaning of $S_i(\nu)$ and A_i is lost but the apparent SNR

112 is dramatically increased, which is much more important for our analysis. Nev-
 113 ertheless, assumptions required in eq. 1 are usually not relevant. Radiative
 114 transfer model used for precise quantification is highly non-linear.

115 One has to consider that this unsupervised linear unmixing problem is al-
 116 ready very difficult for machine learning. Solving non-linear model in a unsu-
 117 pervised way is a research area that is clearly not solved yet. In addition, we
 118 would like to focus on spectral detection, rather than quantification. Thus, we
 119 will focus on $S(\nu)$ much more than A . We will show that for linear, but also
 120 non-linear simulation and real data, meaningful $S(\nu)$ can be retrieved. Due
 121 to non-linearity, A may differ significantly from truth, but the big tendencies
 122 should be respected. After the quick-look analysis, estimating S_i and A , one
 123 must go back to the real data. The most trivial strategy is to pick the spectra
 124 X out of the collection, with the highest abundance of a selected source S_i .

125 In the following, we will use the continuum estimation $C(\nu)$ using asymmet-
 126 ric least square (Eilers and Boelens, 2005), with parameters : $\nu_{smooth} = 10^3$ and
 127 $p = 1 - 10^{-2}$, 10 number of iterations.

128 3.2. Non negative matrix factorization

129 For a collection of spectra, eq. 5 can be written in matrix form $\mathbf{X}_{kj} \approx \mathbf{S}_{ki} \cdot \mathbf{A}_{ij}$,
 130 with i the source index (from 1 to N_S), j the observation index (from 1 to N_O)
 131 and k the wavenumber index (from 1 to N_ν). Thus, one have to estimate \mathbf{S} and
 132 \mathbf{A} , by minimizing the objective function:

$$F = \|\mathbf{X} - \mathbf{S} \cdot \mathbf{A}\|^2 \quad (6)$$

133 with $\|\cdot\|$, the Frobenius norm (usual L_2 norm).

134 Several algorithms have been proposed to solve this problem, subject to
 135 positivity (both \mathbf{S} and \mathbf{A} are non-negative). Such problem is called Non negative
 136 Matrix Factorization (NMF). This constraint is important to keep the physical
 137 meaning, but also to promote sparsity of \mathbf{S} (a signal is sparse when most of the
 138 values are close to zero except several non-zero values). Let $\hat{\mathbf{S}}$ and $\hat{\mathbf{A}}$ be the
 139 estimation of those quantities.

140 *MU*. We propose to use the Multiplicative Updates (MU) of Lee and Seung
 141 (1999) accelerated by Gillis and Glineur (2012). We used the convergence pa-
 142 rameter $\alpha_{MU} = 1$. Other alternative algorithms are possible but give equiv-
 143 alent results since they minimize the same cost function. This algorithm has
 144 the advantage of very fast computation time but the result may depend on
 145 initialization.

146 *BPSS2*. We propose to test another kind of algorithm: the Bayesian Prior
 147 Source Separation (Moussaoui et al., 2006; Dobigeon et al., 2009), that has been
 148 optimized (Schmidt et al., 2010), hereafter called BPSS2. This algorithm has the
 149 main advantage to account for extra constraint : the sum-to-one or sum-lower-
 150 than-one on the abundances ($\sum_i A_{ij} = 1$) that also promotes sparsity of \mathbf{S} . This
 151 algorithm, based on Monte Carlo approach is much more time consuming. One

152 approach to reduce the computation time is to select only relevant spectra out
 153 of the dataset (Moussaoui et al., 2008), but then the statistics may be biased
 154 (Schmidt et al., 2010). Thanks to the advances of computer capabilities, we
 155 propose to treat the full dataset. This kind of algorithm is very slow but since
 156 the formulation is Bayesian, it converge toward an unique solution.

157 *psNMF*. In order to regularize the problem of eq. 6, one can add an extra
 158 penalization term to enforce sparsity on \mathbf{A} (only few non zeros elements in \mathbf{A})
 159 (Kim and Park, 2007) :

$$F = \|\mathbf{X} - \mathbf{S}\cdot\mathbf{A}\|^2 + \lambda \|\mathbf{A}\|_1 \quad (7)$$

160 With $\|\cdot\|_1$, the L_1 norm. The first term is called data attachment term (the
 161 usual squared difference). The second is called regularization term. The prob-
 162 lem with this approach, is that hyperparameter λ is not known and has to be
 163 tuned manually. A recent approach has been proposed to solve this problem in
 164 the Bayesian framework (Hinrich and Mørup, 2018). The main idea is to encom-
 165 pass all variables and hyperparameters in a unique problem that is estimated
 166 with variational update principle. We will refer this algorithm to probability
 167 sparse NMF (psNMF). This algorithm has the advantage to have a reduced
 168 computation time and no hyperparameter tuning. It also has a regularization
 169 term to avoid strong dependence of the initialization on the final solution.

170 In order to estimate the precision of the reconstruction, we used the Root
 171 Mean Square Difference *RMSE*:

$$RMSE = \frac{\sqrt{\langle (\mathbf{X} - \dot{\mathbf{S}}\cdot\dot{\mathbf{A}})^2 \rangle}}{\langle \mathbf{X} \rangle} \quad (8)$$

172 With $\langle \cdot \rangle$, the mean.

173 Once the sources are estimated, we quantify their relevance for the global
 174 dataset. From the total reconstruction $\dot{\mathbf{X}}_{kj} = \dot{\mathbf{S}}_{ki} \cdot \dot{\mathbf{A}}_{ij}$, for all i , we can estimate
 175 the contribution of source i' , that is to say: $\dot{\mathbf{X}}_{kj}^i = \dot{\mathbf{S}}_{ki'} \cdot \dot{\mathbf{A}}_{i'j}$. Thus, the relevance
 176 of source i is defined as:

$$R^i = \frac{\langle |\dot{\mathbf{X}}^i - \dot{\mathbf{X}}| \rangle}{\langle \dot{\mathbf{X}} \rangle} \quad (9)$$

177 This definition is convenient since the sum of all R^i is one (this property
 178 is only present when sources and abundances are positive) and we can easily
 179 estimate the % contribution of each source in the final reconstruction. One has
 180 to note that relevance is not a measure of presence or not of a minor specie (for
 181 instance CH_4) but a measure of how important is the source over the dataset.
 182 Major species, should always have a larger relevance than minor species. In the
 183 following, we plot all sources results by decreasing order of relevance.

184 *3.3. Band depth (BD)*

185 We used the following band depth definition, difference of the geometric
186 mean of two reference wavenumbers in the continuum, compared to the band:

$$BD = X(\nu_l)^{\frac{\nu_c - \nu_l}{\nu_r - \nu_l}} \cdot X(\nu_r)^{\frac{\nu_r - \nu_c}{\nu_r - \nu_l}} - X(\nu_c) \quad (10)$$

187 with X the observed spectra in transmittance, ν_c the wavenumber of the
188 center of band, ν_l the wavenumber of the reference level on the left (smaller
189 wavenumber), ν_r the wavenumber of the reference level on the right (larger
190 wavenumber).

191 **4. Synthetic tests**

192 We simulated several synthetic observations in different conditions, to mimic
193 the case of NOMAD-SO. The first section describes a simple toy model example
194 and the second one presents extensive tests of this toy model with various cases.
195 By *hidden spectra*, *hidden compounds* and *hidden CH₄*, we always refer to a
196 spectral dataset with a dominant major component (here water) and a minor
197 specie (here CH₄). The goal of the proposed approach is to pick up a source,
198 containing CH₄ only.

199 *4.1. Toy example*

200 *4.1.1. Synthetic dataset*

201 In order to demonstrate the usefulness of our method, we propose here a toy
202 example in a very difficult case. We will see that usual method fails detection
203 but our method is able to detect the hidden compounds.

204 For this toy example, we simulate a linear mixture of $N_O = 10^4$ observations
205 spanned over $N_\nu = 320$ spectels (see fig. 1) similar to order 136 of NOMAD-
206 SO. Each spectrum is a mixture of a spectra of water vapor S_{H_2O} (coming
207 from one actual source estimated from real data using psNMF) and theoretical
208 methane S_{CH_4} from Villanueva et al. (2018), with corresponding abundances
209 A_{H_2O} , A_{CH_4} :

$$X = S_{H_2O} \cdot A_{H_2O} + S_{CH_4} \cdot A_{CH_4} + n \quad (11)$$

210 The noise n is assumed to be a Gaussian process with a standard deviation
211 of $\sigma=0.001$ and no bias: $n = \mathcal{G}(0, \sigma)$. All spectra contain pure water vapor with
212 a coefficient following $A_{H_2O} = 5/6 \cdot \beta(1, 10) + 1/6 \cdot \mathcal{U}(0, 1)$, a mixture of beta (β)
213 distribution for 5/6 of the sample and an uniform (\mathcal{U}) distribution for 1/6 of
214 the sample. This process mimics well the water vapor band depth distribution
215 (BD, see definition in section 3.3) of the real dataset (see Fig.2). As the baseline
216 of S_{H_2O} is not zero, we also mimic baseline correction errors. In addition 100
217 spectra out of 10000 contain methane with $A_{CH_4} = 1$, such that the band depth
218 of S_{CH_4} is at $3\text{-}\sigma$ level. Please note that the model to generate the data is not
219 fulfilling the sum-to-one constraint, but fully fulfilling the positivity constraint.
220 Given the defined noise and signal level, the noise *RMSD* is 0.16.

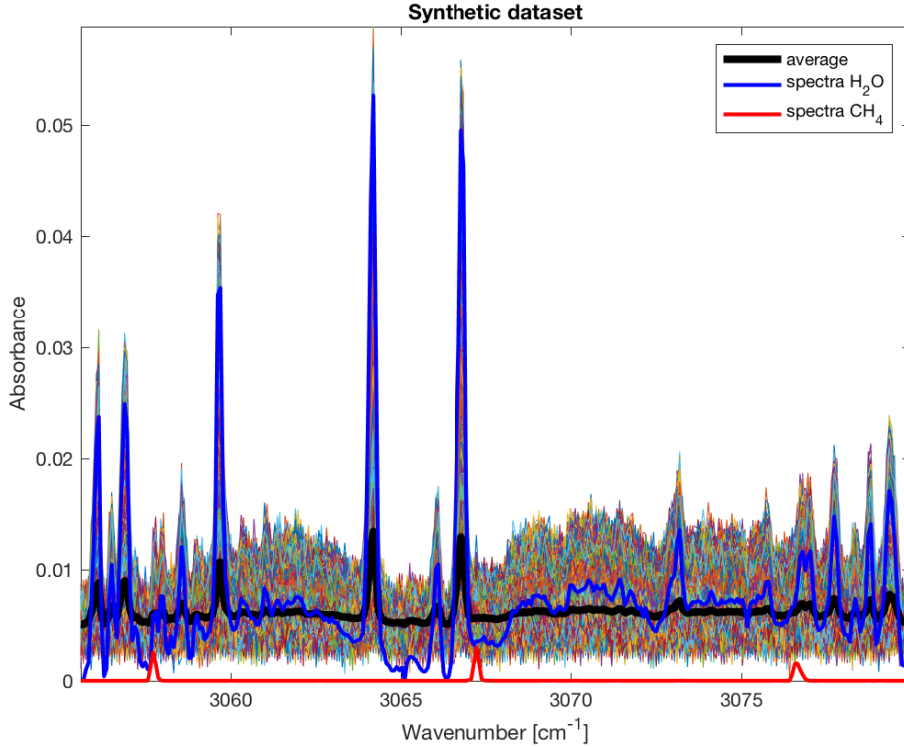


Figure 1: Synthetic dataset containing 10^4 spectra with various abundances of H_2O and 100 containing CH_4 at $3\text{-}\sigma$ level of the noise. In blue the reference spectra S_{H_2O} of H_2O (coming from actual data analysis). In red the reference spectra S_{CH_4} of CH_4 (from theoretical data).

221 The final synthetic dataset is represented in Fig. 1.

222 In order to check the quality of the estimation, we simply compute the
 223 correlation coefficient between S_{CH_4} and the estimated N_S sources $\hat{\mathbf{S}}_i$, using:

$$Q = \text{corr} \left\{ S_{CH_4}, \hat{\mathbf{S}}_i \right\} \quad (12)$$

224 The i th source with the maximum correlation is identified to CH_4 contri-
 225 bution. The value to the maximum correlation is used as metric to assess the
 226 quality of the retrieval.

227 4.1.2. Results

228 By plotting the 10000 samples of the dataset, one is able to identify easily
 229 the H_2O bands. Nevertheless, we cannot observe the target CH_4 in the average
 230 spectrum, even at $3\text{-}\sigma$ level, because it is lost in the baseline changes.

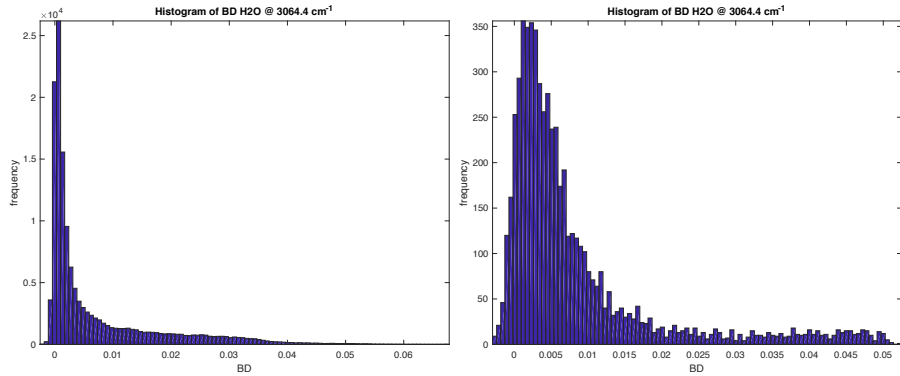


Figure 2: Water vapor Band Depth distribution (left) in the real observation (right) modeled by the toy example.

231 The second simple tool for detection would be the analysis of the band depth.
 232 Figure 3 (left) shows the histogram of the main CH_4 band that exhibits no sign
 233 of the presence of CH_4 (no asymmetry in the positive part). Figure 3 (right)
 234 represents the 100 spectra with the maximum CH_4 BD at 3067.2 cm^{-1} . Again,
 235 no particular elements can be used to argue for detection.

236 Figure 4 represents the results from the non-negative matrix factorization
 237 using psNMF algorithm. One can clearly identify both H_2O and CH_4 sources.
 238 Since those 2 chemical compounds are not correlated in abundance, ($A_{\text{H}_2\text{O}}$ and
 239 A_{CH_4} are independent), two different source spectra are identified. Please note
 240 that the relevance of source 4 is very low (0.4%), meaning that only 0.4% of the
 241 variability in the dataset is due to CH_4 , a very low value, as expected for minor
 242 species.

243 In this case, the correlation coefficient between estimated abundances \hat{A}_4
 244 and true ones A_{CH_4} is 0.73. Since the quantification of abundance is a more
 245 difficult problem, we will not pay excessive attention on this parameter.

246 4.1.3. Convergence and computation time

247 We set the MU algorithm convergence to relative difference of the cost func-
 248 tion $< 10^{-8}$ and a maximum running time of 1000 seconds. For psNMF, we set
 249 the relative difference of the cost function to $< 10^{-7}$ and a maximum iteration
 250 to 2000. For BPSS2, we compute a minimum burn in of 1000 iterations and after
 251 that when the long term statistics (1000 last iterations) of the Markov Chain is
 252 close to the short term statistics (100 last iterations), convergence is considered
 253 to be reached. Then another 1000 iterations are computed to estimate the final
 254 solution statistics.

255 We run the 3 identified tools 10 times on the same dataset with different
 256 noise realization, and compute mean and standard deviation from these 10 ex-
 257 periments. Results are presented in Table 1. One can clearly see that the even
 258 if the convergence is set, there is a high variability in MU results, due to the

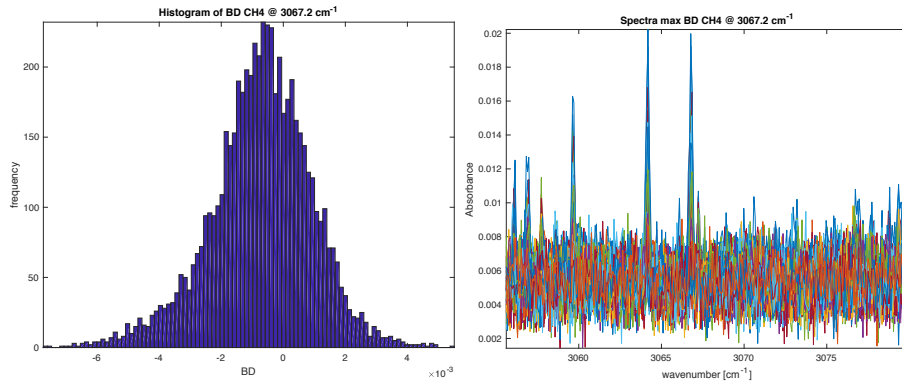


Figure 3: (left) Histogram of Band Depth at 3067.2 cm^{-1} from the dataset containing 100 CH_4 at $3\text{-}\sigma$ level out of 10^4 spectra. (right) 100 spectra with the maximum Band Depth at 3067.2 cm^{-1} specific of CH_4 . Signal is dominated by water and by noise. No specific signature of CH_4 is visible.

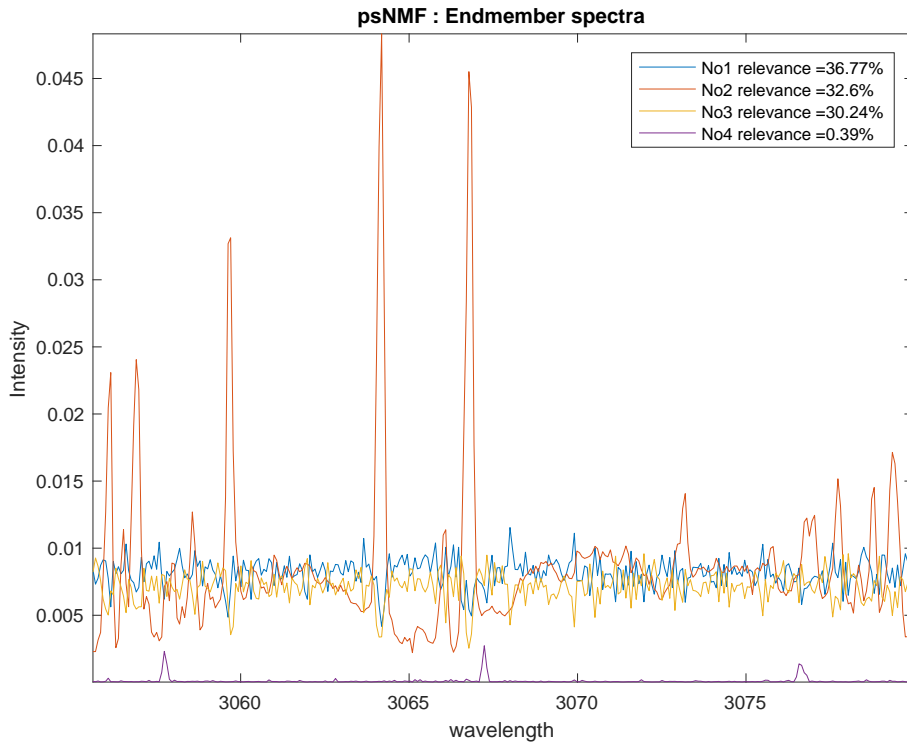


Figure 4: Results of the psNMF algorithm for $N_S = 4$. Sources 1 and 3 are identified to the level with significant noise contribution, source 2 is identified to H_2O (correlation coef. with groundtruth 0.99), and source 4 is CH_4 (correlation coef. with groundtruth 0.98). Relevance is computed from Eq. 9.

	MU	psNMF	BPSS2
Quality Q	0.35 ± 0.12	0.822 ± 0.005	0.41 ± 0.06
$RMSD$ relative error	0.1455 ± 2.10^{-6}	0.1461 ± 5.10^{-6}	0.1468 ± 3.10^{-4}
Computation time (s)	13 ± 8	46 ± 9	413 ± 21

Table 1: Results (mean and standard deviation) from 10 realizations of a toy synthetic example with $N_S = 5$ (in agreement with next section on synthetic tests), $N_O = 10000$, $N_\nu = 320$ and 300 CH_4 spectra hidden at a level of 1 std of the noise. Quality is computed as a correlation coefficient (see Eq. 12). $RMSD$ is computed from Eq. 8. Computation time is expressed in second.

259 lack of regularization. On this particular example, the best is clearly psNMF
 260 algorithm.

261 The $RMSD$ is computed for all cases and shown in Table 1. We can observe
 262 that the value is almost equivalent, around 0.146, for all method but MU is
 263 slightly better, due to the fact that the cost function has no other term. MU
 264 algorithm is just minimizing the reconstruction. As a comparison, the $RMSD$
 265 expected for a perfect reconstruction of the signal (and not the noise) of this
 266 toy example is 0.16. With 5 sources (significantly more than the 3 sources we
 267 define in this toy example), the noise is also fitted, as expected.

268 The quality Q is the only parameter to assess the quality of the algorithm to
 269 detect minor specie (here CH_4). In this particular toy example, psNMF seems
 270 to be the best algorithm, providing a source correlated with groundtruth CH_4
 271 with a correlation coefficient up to 0.8. We will extensively test this performance
 272 in the next section.

273 We also estimate the computation time on a 2.9 GHz Intel Core i7 with 16 Go
 274 DDR3 RAM as an example. All algorithms are implemented in ©Matlab using
 275 parallelized matrix computation. Results, presented in Table 1, demonstrate
 276 that MU is faster than psNMF but both are clearly less resources consuming
 277 than BPSS2. From the computation time and efficiency, we excluded BPSS2
 278 from the next tests.

279 4.2. Extended synthetic tests

280 For the first set of tests, we used the same toy model described in section
 281 4.1, except with 100 CH_4 spectra hidden at a level of 2 and 3 standard deviation
 282 of the noise (this number is called “factor above noise level”). In order to have
 283 robust results, we made 10 realizations and averaged the results.

284 Figure 5 represents the results as a function of the number of sources N_S . It
 285 presents two quality indicators of the results: the average correlation coefficient
 286 Q (see Eq. 12) and the fraction of realization with acceptable results (with
 287 $Q > 0.5$). We can observe that the psNMF is always better than MU on
 288 average at cost of an higher variability (higher standard deviation). Adding
 289 sources seems to always increase the detection until reaching a plateau around
 290 $N_S = 5$. Adding more sources will not drastically increase/decrease the source
 291 estimation. Nevertheless, it requires more computation time for a larger number

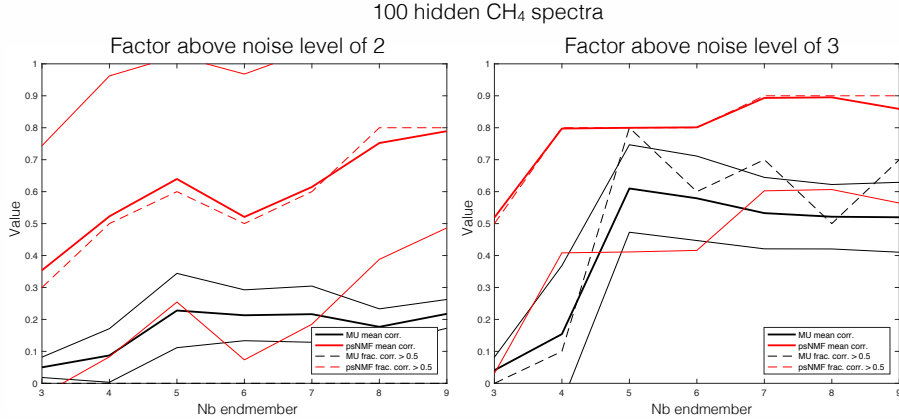


Figure 5: Results of the MU and psNMF algorithm for $N_S = 3$ to 9, $N_O = 10000$, $N_\nu = 320$, as a function of the number of source. The average Q of 10 realizations of the best estimated source (thick lines and standard deviation in thin lines) and the fraction of acceptable results (with $Q > 0.5$). (left) with a factor above noise level of 2 (right) with factor above noise level of 3.

292 of source (approximately x2 between 3 and 9 sources but the computation time
 293 always stays below 200 seconds).

294 For the second set of tests, we used the same toy model, except with 50 and
 295 100 CH_4 spectra hidden at a level of 0.7, 1, 1.2, 1.5, 2.0, 2.5 and 3 standard
 296 deviation of the noise (this number is called “factor above noise level”). In
 297 order to have robust results, we made 10 realizations and averaged the results.
 298 Results are always with $RMSD < 0.18$ with an average ~ 0.16 . $RMSD$ from
 299 the noise level is 0.16 whatever the experiment (the CH_4 is low enough so that
 300 it’s contribution to $RMSD$ is negligible), so the reconstruction is in average as
 301 expected.

302 Figure 6 presents two quality indicators of the results: the average correla-
 303 tion coefficient Q (see Eq. 12) and the fraction of realization with acceptable
 304 results (with $Q > 0.5$). Both indicators indicate that the method psNMF clearly
 305 outperforms MU at high factor above noise level. From our visual inspection of
 306 the results, we define the detection limit when at least 50% of the results are
 307 with $Q > 0.5$ (correlation coefficient > 0.5). This definition is debatable but
 308 there is no absolute way of defining it. Figure 6 shows that the detection limit
 309 is at 1.5 factor above noise level for 100 hidden spectra case, around 2 for 50
 310 hidden spectra. Below this limit, none of the method is able to detect the CH_4
 311 spectra from the noise. For 20 hidden spectra, even at a factor above noise level
 312 of 3, none of the methods is able to detect the CH_4 spectra. One can also note
 313 that the psNMF is less stable since the standard deviation is much larger.

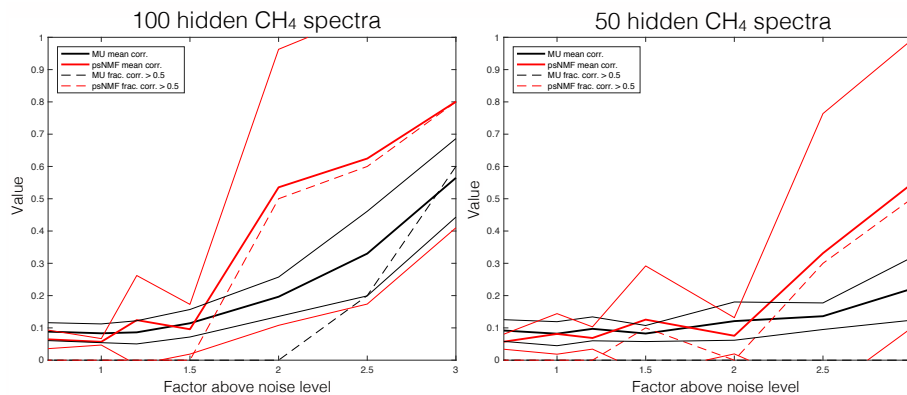


Figure 6: Results of the MU and psNMF algorithm for $N_S = 5$, $N_O = 10000$, $N_v = 320$, as a function of the factor above noise level. The average Q of 10 realizations of the best estimated source (thick lines and standard deviation in thin lines) and the fraction of acceptable results (with $Q > 0.5$). (left) with 100 hidden CH_4 spectra (right) with 50 hidden CH_4 spectra

314 5. Simulation of NOMAD-SO

315 5.1. Simulation dataset

316 This second dataset has been generated with the most precise direct model,
 317 taking into account the full non-linear radiative transfer and instrumental effects
 318 to produce synthetic transmittance, highly comparable with actual observations.
 319 Synthetic transmittances were made for real NOMAD-SO observation files using
 320 the relevant geometry and instrument parameters to attempt to include the
 321 variability inherent in the true measurements.

322 Model atmospheres for each occultation were developed from the GEM-Mars
 323 general circulation model (Neary and Daerden, 2018; Daerden et al., 2019). The
 324 output of the model were provided for 1 Martian day every 10 solar longitude,
 325 and 48 timesteps per Martian day. Atmospheric profiles were developed for each
 326 occultation by interpolating the model temperature and pressure to the solar
 327 longitude, local solar time, latitude, longitude, and tangent altitude relative to
 328 the areoid.

329 To construct the simulated transmittance spectra, the high resolution irra-
 330 diances were computed for each occultation assuming a spherically symmetry
 331 and the tangent atmosphere developed from GEM-Mars for several different
 332 abundance of methane and water, which were simulated as constant volume
 333 mixing ratios. The spectroscopic data for methane and water were taking from
 334 HITRAN 2016 using CO_2 broadening (Gordon et al., 2017; Gamache et al.,
 335 2016; Fissiaux et al., 2014). The instrument forward model was then applied
 336 to each simulation by considering the AOTF bandpass, instrument Instrument
 337 Line Shape (ILS), blaze function, spectel to wavenumber calibration, and the
 338 contribution of light coming from the main order and nearby orders. The final
 339 synthetic transmittance spectra is the ratio of this low-resolution irradiance to
 340 the top-of-atmosphere low resolution irradiance.

	CH ₄ [ppt]	H ₂ O [ppm]	fraction of CH ₄ [%]	noise level
Value	0; 100; 500; 1000	0; 10; 100	1; 5; 10; 50; 100	0.001; 0.0001

Table 2: Simulation parameters. Fraction of CH₄ is fraction of spectra containing methane hidden in the simulation dataset.

341 The AOTF/echelle instrument was modeled using the latest available cali-
342 bration (Liuzzi et al., 2019; Aoki et al., 2019), considering order addition from
343 $+/- 2$ nearby orders (5 total). The spectral calibration of NOMAD-SO varies
344 because it is affected by the instrument temperature, and is provided for each
345 individual NOMAD spectra. The 320 spectels cover the range 3056.1 cm^{-1} to
346 3080.4 cm^{-1} with a wavenumber step of 0.0763 cm^{-1} .

347 No simulation of dust has been performed. Due to the limited spectral range
348 on a single order, about 25 cm^{-1} , the major effect of dust and other aerosols
349 is relatively flat baseline, which we remove at the pre-treatment of the spectra.
350 When dust is optically thick, then non-linearity may appear that are out of the
351 scope of this simulation.

352 The simulation dataset consist of 12486 spectra, simulating observations
353 of order 136 in the same configuration as the 106 solar occultations actually
354 observed from May to December 2018.

355 We add to the dataset a random noise with standard deviation of 0.001 and
356 0.0001 in order to simulate the instrumental noise (corresponding to SNR of 100
357 and 1000 approximately).

358 We hide spectra containing CH₄ in a fraction of the total number of spectra
359 from 1% to 100% in a random manner. In real observation, CH₄ may be spa-
360 tially / temporally coherent but the number of scenarios is infinite. We feel that
361 the random case is interesting enough to be tested. One has to note that con-
362 trarily to the previous toy model of section 4, here abundance are quantitative
363 abundance in the atmosphere.

364 The simulation parameters are summed up in table 2.

365 5.2. Detection limits

366 We applied the psNMF method with $N_S = 5$, which is the most promising
367 one from the previous analysis. We compute the analysis 10 times for 10 different
368 random noise realizations and average the results in order to present robust
369 conclusion. We select a pure CH₄ and a pure H₂O spectra (noted P_{CH_4} and
370 P_{H_2O}) from the simulation as reference spectra.

371 5.2.1. Methods to analyze the results

372 The main difference with the toy model section in 4 is that H₂O and CH₄
373 may be highly mixed in the sources. Simple correlation coefficient to pick the
374 best source is thus not efficient enough. We propose here another approach to
375 estimate the best source.

376 For each estimated source $\hat{S}_{:i}$, we analyze it as a linear mixture of P_{H_2O} and
377 P_{CH_4} :

$$\dot{S}_{:i} = P_{H_2O} \cdot \alpha_{H_2O,i} + P_{CH_4} \cdot \alpha_{CH_4,i} \quad (13)$$

378 This problem is called supervised detection algorithm since P_{H_2O} and P_{CH_4}
 379 are known, contrary to the general one, presented in Eq. 5, where source spectra
 380 are not known. The source i^* with the maximum α_{CH_4,i^*} is selected as the best
 381 target CH_4 source, called *best source* hereafter.

382 We then propose to use three indicators of good detection :

- 383 • *Fraction of the 4 main CH_4 peaks detected* (at 3057.7, 3063.4, 3067.2 and
 384 3076.6 cm^{-1}). This is computed using the peak detection algorithm from
 385 ©Matlab on both simulation and best source with a tolerance of 2 spec-
 386 tels, i.e. detected peaks can be 2 spectels off the expected one. The peak
 387 must be with a maximum amplitude larger than 1/1000 the maximum of
 388 $\dot{S}_{:i^*}$ to be considered significant. Please note that even there are only 5
 389 possible fraction (0, 0.25, 0.5, 0.75 and 1), since we average on 10 realiza-
 390 tions, any number can appear.
- 391 • *Mean distance to the expected center.* Mean distance in spectel between
 392 the CH_4 peaks detected in the best source and the reference one.
- 393 • *Abundance of CH_4 in the source.* α_{CH_4} (from Eq. 13), which describes
 394 the amplitude of the CH_4 peaks in the best source.

395 5.2.2. Analysis of the results

396 Figure 7 summarizes all the results. Fraction of the 4 main CH_4 peaks
 397 detected in the most relevant source has always a standard deviation < 0.43
 398 and a mean value of 0.06 over the 10 realizations. The Mean distance to the
 399 expected center has always a standard deviation < 0.40 and a mean value of
 400 0.07 over the 10 realizations. The abundance of CH_4 in the source has always a
 401 standard deviation < 0.05 and a mean value of 0.005 over the 10 realizations.

402 This figure shows that the detection limits clearly depend on CH_4 density,
 403 but also on the fraction of hidden CH_4 and noise level, as expected. Abundance
 404 of CH_4 in the source α_{CH_4} maximum is 25%, meaning that in any cases H_2O is
 405 dominating the best source and so both CH_4 and H_2O are present in each best
 406 source. This is because CH_4 is a minor specie (as expected from the conditions
 407 of our simulation), its absorption band generally follows the air-mass, as H_2O
 408 does. So there is no particular source for CH_4 only.

409 When more than two lines are detected, we can consider it as a detection.
 410 This limit is reached for $CH_4 \geq 500$ ppt for 10 and 100 ppm of H_2O . Never-
 411 theless, the detection limits lies between 100 and 500 ppt in the case of 10 ppm
 412 of H_2O vapor since the detection is perfect (100% of the 4 main CH_4 peaks
 413 detected) occurs for a fraction of CH_4 5 to 50%. Interestingly, the optimum
 414 detection is not when 100% of the spectra contains CH_4 , but more between
 415 5-50 %. This behavior is due to the statistics that is richer when also CH_4 is
 416 lacking in certain spectra. When 100% of spectra contain CH_4 , the statistical
 417 variability of the dataset is mainly due to airmass (atmosphere is assumed to

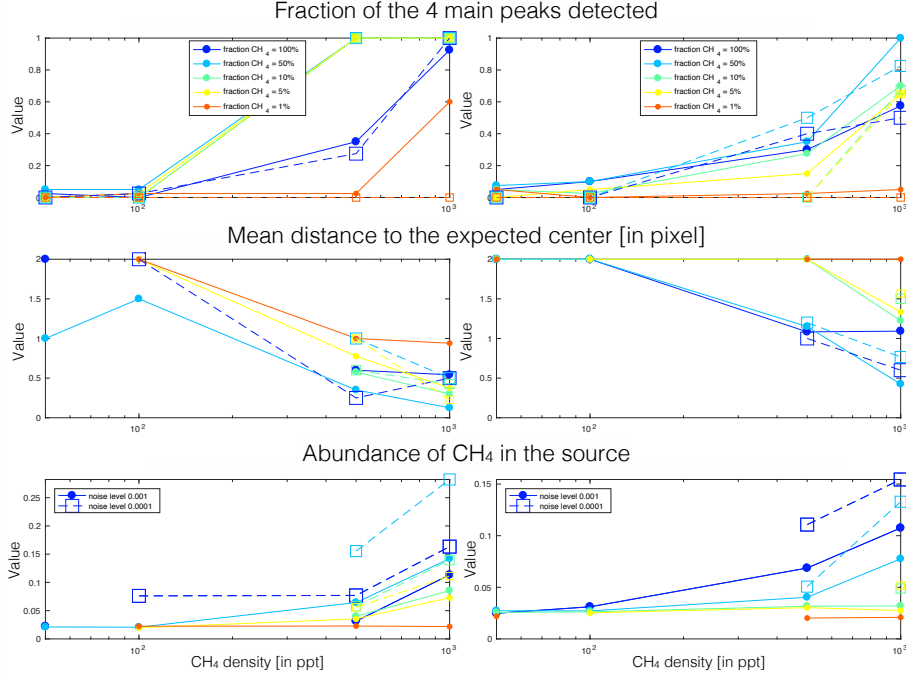


Figure 7: Results of the psNMF algorithm for $N_S = 5$ on simulation dataset, averaged over 10 noise realizations, for different noise levels (0.001 and 0.0001) and different fractions of hidden CH_4 (1%, 5%, 10%, %, 100%). Hidden CH_4 are taken within the same orbital sequences. The left panels represent results for 10 ppm of water vapor and the right ones for 100 ppm of H_2O . From top to bottom, we show: a) Fraction of the 4 main CH_4 peaks detected in the *best source*; b) Mean distance to the expected center in spectel and c) Abundance of CH_4 in the source α_{CH_4} . Please note that the absence of plotted data means that no source was successfully detected.

418 be well mixed). So both CH_4 and H_2O are varying together and there is less
 419 statistics to base the detection on.

420 Noise level does not affect first the fraction of the 4 main CH_4 peaks but
 421 increases the spectral shift of the band center. In addition, it clearly affects the
 422 abundance and thus the band depth.

423 In conclusion, from this simulation analysis, one could expect detection limits
 424 of CH_4 in the range 100-500 ppt when operating in favorable conditions.

425 6. Real data analysis

426 In this section, we report the results of actual NOMAD data, focusing on
 427 diffraction orders with potential CH_4 lines: 119, 134 and 136, are shown respec-
 428 tively on Fig. 8, 9 and 10. We used the 821 ingress and egress transit orbits
 429 for order 119, 2358 orbits for order 134 and 703 for order 136. We filter spectra
 430 with $\text{SNR} > 100$. Results are compared with NOMAD simulations (Villanueva

	119	134	136
N_O	134045	365985	140064
$N_S = 4$	0.476	0.575	0.634
$N_S = 5$	0.456	0.553	0.609
$N_S = 6$	0.442	0.553	0.585
$N_S = 10$	0.410	0.484	0.544

Table 3: Number of spectra N_O and $RMSE$ relative errors for 4 to 10 number of sources N_S resulting from the analysis of all observations of NOMAD data up to 15 January 2020, using the psNMF algorithm. $RMSE$ is computed from Eq. 8.

et al., 2018) using the calibration pipeline. This process adds ghost lines from adjacent orders, as in real data. Table 3 summarizes the relative error and the number of spectra. The approach here is to compute the analysis with psNMF using $N_S = 5$ in agreement with the previous section.

Please remind that our approach is fully blind: no spectral information has been included in the analysis (nothing about H_2O , CO_2 or CH_4).

For all orders, sources of H_2O are estimated, as expected. Also a source presenting a residual of the continuum is always present. Due to non-linearities of the radiative transfer, the acquisition process (temperature dependence) and the wavenumber shift, the molecular species appears sometimes in different sources.

Order 136 gives the 1 source related to the background and 4 sources related to H_2O . All 4 sources of water have the peaks but with different relative intensities and wavenumber shift.

For order 119, both CO_2 and H_2O lines are identified (see Fig. 8). Since those two components are uncorrelated, separated sources are found by the algorithm.

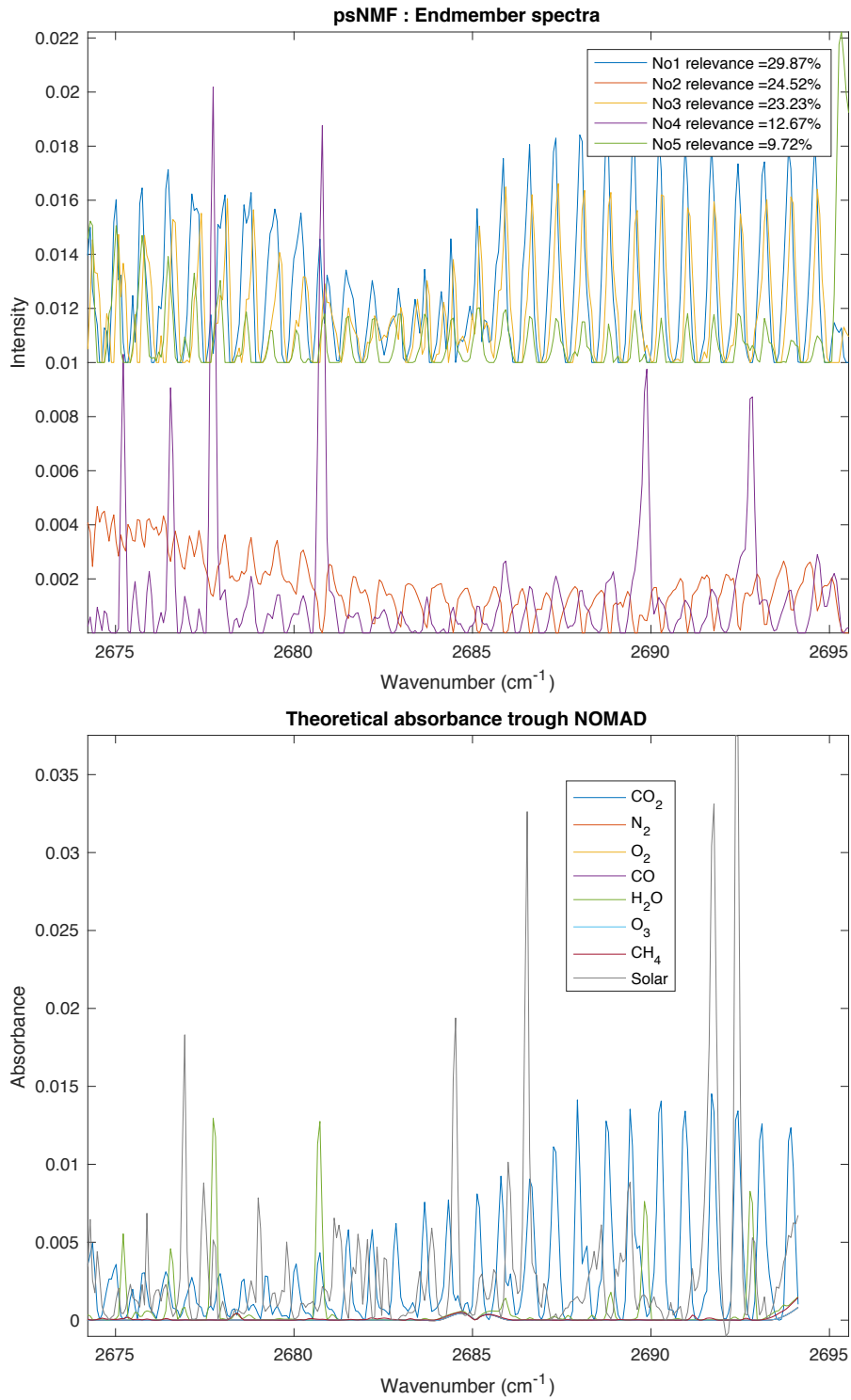
Interestingly, order 134 presents a source with unexpected lines. The main lines are at positions : 3016.70, 3017.07, 3018.12, 3019.54, 3020.90, 3022.25, 3023.60, 3024.96, and 3027.29 cm^{-1} . These lines has been also detected in the ACS instrument data and attributed to CO_2 magnetic dipole transition (Trokhimovskiy et al., 2020). Further analysis shall be done to compare both NOMAD AND ACS data.

Solar lines are never appearing in the sources. They are self-corrected by the calibration since we don't use a reference solar spectra but the solar observation during the transit when the tangent altitude is so high that there is no martian atmosphere (typically > 200 km).

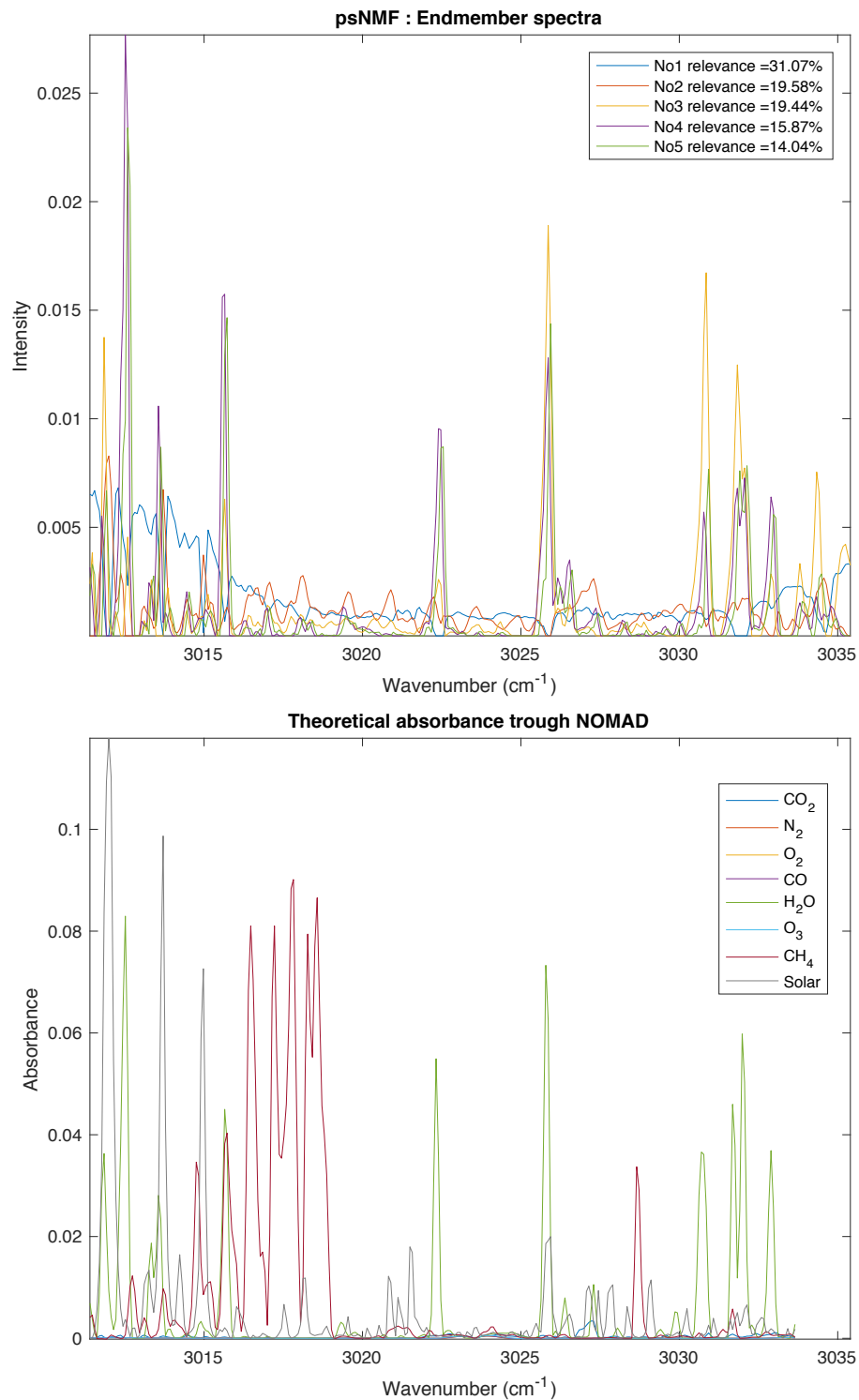
None of the analyzed orders presents sources related to CH_4 .

7. Discussions and Conclusion

We implemented a new strategy to analyze spectroscopic datasets. This strategy is fully unsupervised, so that any kind of absorption bands can be discovered. The amount of prior information required is thus very low. The computation can be done on a regular hardware for the most common database and within reasonable amount of time (~ 100000 spectra).



18
 Figure 8: Results of the psNMF algorithm for the diffraction order 119 for $N_S = 5$. The sources 1, 3 and 5 are identified to CO_2 (shift of 0.01 for clarity). The source 2 is identified to the background level (continuum misestimation). The source 4 is identified to H_2O . No source seems to be related to CH_4 .



19
 Figure 9: Results of the psNMF algorithm for the diffraction order 134 for $N_S = 5$. The source 1 is identified to the background level (continuum misestimation), the sources 3, 4 and 5 are identified to H_2O . The sources 2 present unmodeled lines that are not present in the spectroscopic database. These lines has been first detected in the ACS instrument data and attributed to CO_2 magnetic dipole transition (Trokhimovskiy et al., 2020). No source seems to be related to CH_4 .

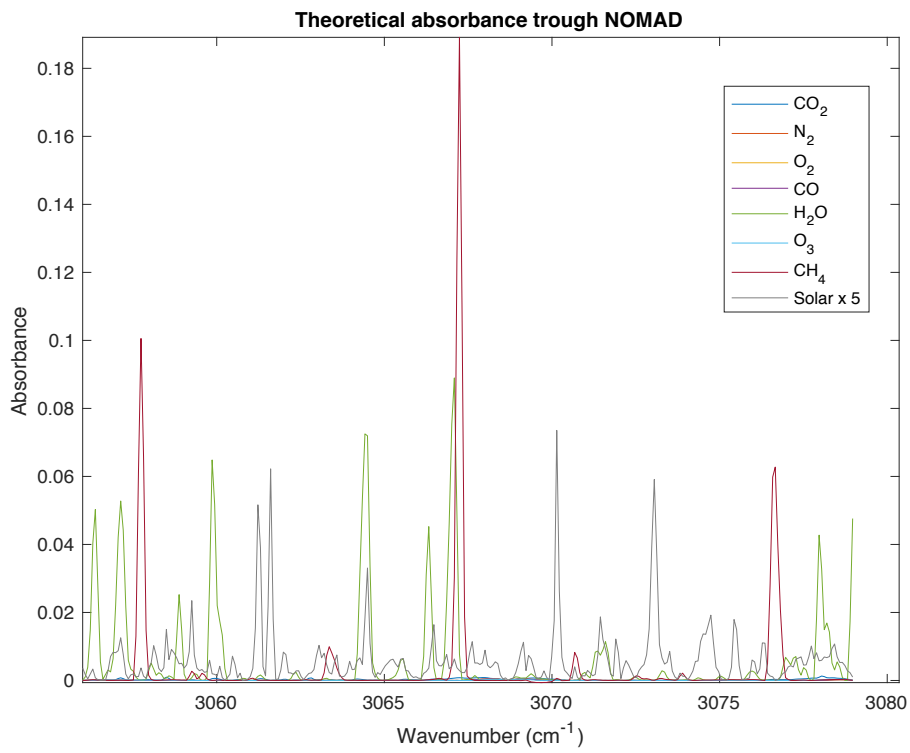
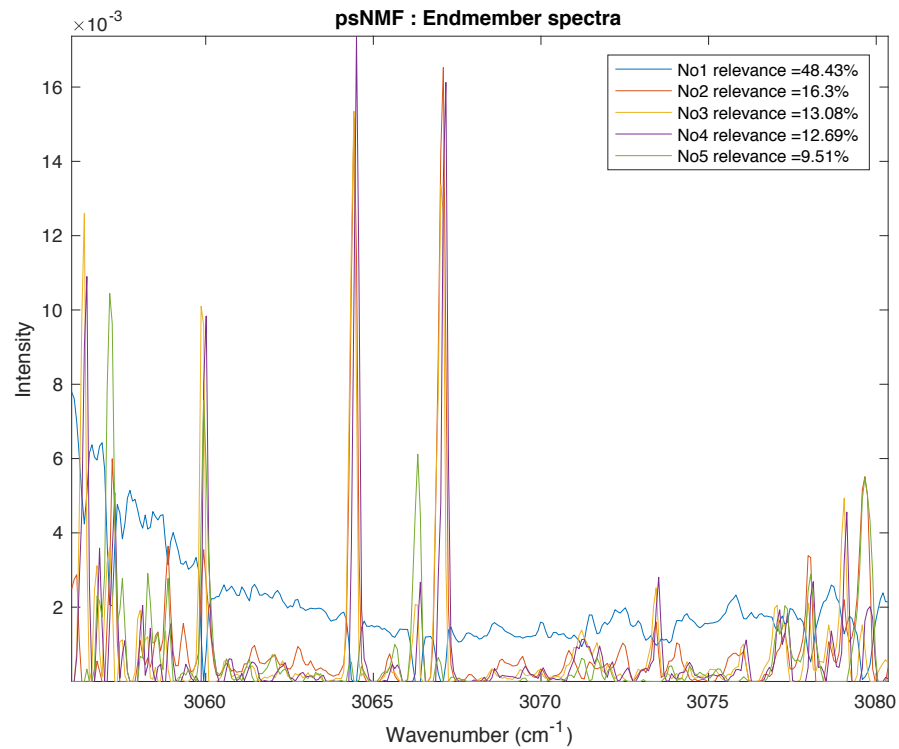


Figure 10: Results of the psNMF algorithm for the order 136 for $N_S = 5$. The source 1 is identified to the level background (continuum misestimation), the sources 2, 3, 4 and 5 are identified to H_2O , either directly either from the adjacent orders. No source seems to be related to CH_4 .

464 We illustrate the approach for typical atmospheric spectroscopy. We first put
465 forward a synthetic test, based on simple linear mixing to give a toy example
466 and to identify the best promising algorithm. The psNMF clearly outperformed
467 MU and BPSS2.

468 Then we proposed a simulation, based on realistic radiative transfer and
469 instrumental effects, applied on NOMAD-SO spectra. The detection limits goes
470 below 500 ppt in favorable conditions, with reduced H₂O and low noise level.
471 The same range of detection limits is reach with usual approach of model fitting
472 at a much higher computation cost and analysis effort. Given the simplicity
473 of use, this tool may be relevant to handle large and complex datasets at first
474 glance. As a perspective, analysis of residuals after the non-linear retrieval of
475 the data may lower the detection limits. One can then test if the residuals are
476 simply Gaussian noise, or if they may contain interesting features.

477 Interestingly, a molecular specie not well mixed in the atmosphere can be
478 most easily detected with our approach.

479 The last section presented the results of the application on real NOMAD-SO
480 data, using orders 119, 134 and 136, selected as they are representative of the
481 baseline strategy of measurements in NOMAD, allowing characterization of H₂O
482 and potential detection of CH₄. The outcome is that no CH₄ has been identified,
483 but H₂O and CO₂ are detected. Interestingly a new set of spectral lines has
484 been discovered in the NOMAD data. These lines has been first detected in
485 the ACS instrument data and attributed to CO₂ magnetic dipole transition
486 (Trokhimovskiy et al., 2020). We thus confirm their presence with our current
487 analysis.

488 One way to go back to the data is to pick the real data with the highest
489 source contribution $\hat{\mathbf{A}}$. Our quicklook analysis is thus only a starting point of
490 a more complete scientific analysis. This second step will require much more
491 prior information (chemical compounds, fundamental spectroscopic constants,
492 radiative transfer model, ...).

493 Future work should apply the proposed approach to other datasets, such
494 as other NOMAD-SO orders, or other spectroscopic datasets (including hy-
495 perspectral images) from laboratory measurements, ground based telescopes or
496 space-born spectrometers. The approach is generic enough to treat datasets
497 that can be at first order approximated to a linear mixture.

498 *Acknowledgements*

499 We acknowledge support from the "Institut National des Sciences de l'Univers"
500 (INSU), the "Centre National de la Recherche Scientifique" (CNRS) and "Centre
501 National d'Etudes Spatiales" (CNES) through the "Programme National
502 de Planétologie" and the ExoMars TGO programs. The NOMAD experiment
503 is led by the Royal Belgian Institute for Space Aeronomy (BIRA-IASB), as-
504 sisted by Co-PI teams from Spain (IAA-CSIC), Italy (INAF-IAPS), and the
505 United Kingdom (Open University). This project acknowledges funding by the
506 Belgian Science Policy Office (BELSPO), with the financial and contractual co-
507 ordination by the ESA Prodex Office (PEA 4000103401, 4000121493), by Span-
508 ish Ministry of Science and Innovation (MCIU) and by European funds under

509 grants PGC2018-101836-B-I00 and ESP2017-87143-R (MINECO/FEDER), as
510 well as by UK Space Agency through grants ST/R005761/1, ST/P001262/1,
511 ST/R001405/1 and ST/R001405/1 and Italian Space Agency through grant
512 2018-2-HH.0. This work was supported by the Belgian Fonds de la Recherche
513 Scientifique - FNRS under grant number 30442502 (ET-HOME). The IAA/CSIC
514 team acknowledges financial support from the State Agency for Research of the
515 Spanish MCIU through the Center of Excellence Severo Ochoa award for the
516 Instituto de Astrofísica de Andalucía (SEV-2017-0709). US investigators were
517 supported by the National Aeronautics and Space Administration. Canadian
518 investigators were supported by the Canadian Space Agency.

519 Aoki, S., Vandaele, A. C., Daerden, F., Villanueva, G. L., Liuzzi, G., Thomas,
520 I. R., Erwin, J. T., Trompet, L., Robert, S., Neary, L., Viscardy, S., Clancy,
521 R. T., Smith, M. D., Lopez-Valverde, M. A., Hill, B., Ristic, B., Patel, M. R.,
522 Bellucci, G., Lopez-Moreno, J.-J., the NOMAD team, 2019. Water vapor
523 vertical profiles on mars in dust storms observed by tgo/nomad. *Journal of*
524 *Geophysical Research: Planets* 124 (12), 3482–3497.

525 Bertaux, J.-L., Fonteyn, D., Korablev, O., Chassefière, E., Dimarellis, E.,
526 Dubois, J., Hauchecorne, A., Cabane, M., Rannou, P., Levasseur-Regourd,
527 A., Cernogora, G., Quemerais, E., Hermans, C., Kockarts, G., Lippens, C.,
528 Maziere, M., Moreau, D., Muller, C., Neefs, B., Simon, P., Forget, F., Hour-
529 din, F., Talagrand, O., Moroz, V., Rodin, A., Sandel, B., Stern, A., oct 2000.
530 The study of the martian atmosphere from top to bottom with SPICAM light
531 on mars express. *Planetary and Space Science* 48 (12-14), 1303–1320.

532 Bertaux, J.-L., Nevejans, D., Korablev, O., Villard, E., Quémerais, E., Neefs,
533 E., Montmessin, F., Leblanc, F., Dubois, J., Dimarellis, E., Hauchecorne, A.,
534 Lefèvre, F., Rannou, P., Chaufray, J., Cabane, M., Cernogora, G., Souchon,
535 G., Semelin, F., Reberac, A., Ransbeek, E. V., Berkenbosch, S., Clairquin,
536 R., Muller, C., Forget, F., Hourdin, F., Talagrand, O., Rodin, A., Fedorova,
537 A., Stepanov, A., Vinogradov, I., Kiselev, A., Kalinnikov, Y., Durry, G.,
538 Sandel, B., Stern, A., Gérard, J., oct 2007. SPICAV on venus express: Three
539 spectrometers to study the global structure and composition of the venus
540 atmosphere. *Planetary and Space Science* 55 (12), 1673–1700.

541 Bovensmann, H., Burrows, J. P., Buchwitz, M., Frerick, J., Noël, S., Rozanov,
542 V. V., Chance, K. V., Goede, A. P. H., jan 1999. SCIAMACHY: Mission ob-
543 jectives and measurement modes. *Journal of the Atmospheric Sciences* 56 (2),
544 127–150.

545 Daerden, F., Neary, L., Viscardy, S., Muñoz, A. G., Clancy, R., Smith, M.,
546 Encrenaz, T., Fedorova, A., 2019. Mars atmospheric chemistry simulations
547 with the gem-mars general circulation model. *Icarus* 326, 197–224.

548 Dobigeon, N., Moussaoui, S., Tourneret, J.-Y., Carteret, C., Dec. 2009.
549 Bayesian separation of spectral sources under non-negativity and full addi-
550 tivity constraints. *Signal Processing* 89 (12), 2657–2669.

- 551 URL [http://www.sciencedirect.com/science/article/
552 B6V18-4W9XDSW-2/2/f3d4b6f457b91e5ccfcce8ffcf41bb18](http://www.sciencedirect.com/science/article/B6V18-4W9XDSW-2/2/f3d4b6f457b91e5ccfcce8ffcf41bb18)
- 553 Eilers, P. H., Boelens, H. F., 2005. Baseline correction with asymmetric least
554 squares smoothing.
- 555 Erard, S., Drossart, P., Piccioni, G., Jan. 2009. Multivariate analysis of visible
556 and infrared thermal imaging spectrometer (virtis) venus express nightside
557 and limb observations. *J. Geophys. Res.* 114, –.
558 URL <http://dx.doi.org/10.1029/2008JE003116>
- 559 Faisal, M., Windholz, L., Kröger, S., apr 2020. Systematic investigations of the
560 hyperfine structure constants of niobium i levels. part i: Constants of upper
561 odd parity energy levels between 16,672 and 31,025 cm⁻¹ and discovery of a
562 new level. *Journal of Quantitative Spectroscopy and Radiative Transfer* 245,
563 106873.
- 564 Fissiaux, L., Delière, Q., Blanquet, G., Robert, S., Vandaele, A. C., Lepère,
565 M., mar 2014. CO₂-broadening coefficients in the ν_4 fundamental band of
566 methane at room temperature and application to CO₂-rich planetary atmo-
567 spheres. *Journal of Molecular Spectroscopy* 297, 35–40.
- 568 Gamache, R. R., Faresé, M., Renaud, C. L., aug 2016. A spectral line list for
569 water isotopologues in the 1100–4100 cm⁻¹ region for application to CO₂-rich
570 planetary atmospheres. *Journal of Molecular Spectroscopy* 326, 144–150.
- 571 Geminalé, A., Grassi, D., Altieri, F., Serventi, G., Carli, C., Carrozzo, F.,
572 Sgavetti, M., Orosei, R., D'Aversa, E., Bellucci, G., Frigeri, A., 2015. Re-
573 moval of atmospheric features in near infrared spectra by means of principal
574 component analysis and target transformation on mars: I. method. *Icarus*
575 253 (0), 51 – 65.
576 URL [http://www.sciencedirect.com/science/article/pii/
577 S0019103515000640](http://www.sciencedirect.com/science/article/pii/S0019103515000640)
- 578 Gillis, N., Glineur, F., apr 2012. Accelerated multiplicative updates and hierar-
579 chical ALS algorithms for nonnegative matrix factorization. *Neural Compu-
580 tation* 24 (4), 1085–1105.
- 581 Giuranna, M., Viscardy, S., Daerden, F., Neary, L., Etiopé, G., Oehler, D.,
582 Formisano, V., Aronica, A., Wolkenberg, P., Aoki, S., Cardesin-Moinelo, A.,
583 de la Parra, J. M.-Y., Merritt, D., Amoroso, M., apr 2019. Independent con-
584 firmation of a methane spike on mars and a source region east of gale crater.
585 *Nature Geoscience* 12 (5), 326–332.
- 586 Gordon, I., Rothman, L., Hill, C., Kochanov, R., Tan, Y., Bernath, P., Birk, M.,
587 Boudon, V., Campargue, A., Chance, K., Drouin, B., Flaud, J.-M., Gamache,
588 R., Hodges, J., Jacquemart, D., Perevalov, V., Perrin, A., Shine, K., Smith,
589 M.-A., Tennyson, J., Toon, G., Tran, H., Tyuterev, V., Barbe, A., Császár, A.,
590 Devi, V., Furtenbacher, T., Harrison, J., Hartmann, J.-M., Jolly, A., Johnson,

- 591 T., Karman, T., Kleiner, I., Kyuberis, A., Loos, J., Lyulin, O., Massie, S.,
592 Mikhailenko, S., Moazzen-Ahmadi, N., Müller, H., Naumenko, O., Nikitin,
593 A., Polyansky, O., Rey, M., Rotger, M., Sharpe, S., Sung, K., Starikova, E.,
594 Tashkun, S., Auwera, J. V., Wagner, G., Wilzewski, J., Wcislo, P., Yu, S.,
595 Zak, E., 2017. The hitran2016 molecular spectroscopic database. *Journal of*
596 *Quantitative Spectroscopy and Radiative Transfer* 203, 3 – 69.
- 597 Herr, K. C., Pimentel, G. C., Jan. 1970. Evidence for Solid Carbon Dioxide in
598 the Upper Atmosphere of Mars. *Science* 167, 47–49.
- 599 Hinrich, J. L., Mørup, M., 2018. Probabilistic sparse non-negative matrix fac-
600 torization. In: *Latent Variable Analysis and Signal Separation*. Springer In-
601 ternational Publishing, pp. 488–498.
- 602 Kim, H., Park, H., may 2007. Sparse non-negative matrix factorizations via al-
603 ternating non-negativity-constrained least squares for microarray data anal-
604 ysis. *Bioinformatics* 23 (12), 1495–1502.
- 605 Korablev, O., Vandaele, A. C., Montmessin, F., Fedorova, A. A., Trokhi-
606 movskiy, A., Forget, F., Lefèvre, F., Daerden, F., Thomas, I. R., Trompet,
607 L., Erwin, J. T., Aoki, S., Robert, S., Neary, L., Viscardy, S., Grigoriev,
608 A. V., Ignatiev, N. I., Shakun, A., Patrakeevev, A., Belyaev, D. A., Bertaux,
609 J.-L., Olsen, K. S., Baggio, L., Alday, J., Ivanov, Y. S., Ristic, B., Mason, J.,
610 Willame, Y., Depiesse, C., Hetey, L., Berkenbosch, S., Clairquin, R., Queirolo,
611 C., Beeckman, B., Neefs, E., Patel, M. R., Bellucci, G., López-Moreno, J.-J.,
612 Wilson, C. F., Etiope, G., Zelenyi, L., Svedhem, H., Vago, J. L., The ACS
613 and NOMAD Team, apr 2019. No detection of methane on mars from early
614 ExoMars trace gas orbiter observations. *Nature* 568 (7753), 517–520.
- 615 Lee, D. D., Seung, H. S., Oct. 1999. Learning the parts of objects by non-
616 negative matrix factorization. *Nature* 401 (6755), 788–791.
617 URL <http://dx.doi.org/10.1038/44565>
- 618 Liuzzi, G., Villanueva, G., Mumma, M., Smith, M., Daerden, F., Ristic, B.,
619 Thomas, I., Vandaele, A., Patel, M., López-Moreno, J., Bellucci, G., 2019.
620 Methane on mars: New insights into the sensitivity of ch4 with the NO-
621 MAD/ExoMars spectrometer through its first in-flight calibration. *Icarus* 321,
622 671–690.
- 623 López-Valverde, M., López-Puertas, M., López-Moreno, J., Formisano, V.,
624 Grassi, D., Maturilli, A., Lellouch, E., Drossart, P., aug 2005. Analysis of
625 non-LTE emissions at in the martian atmosphere as observed by PFS/mars
626 express and SWS/ISO. *Planetary and Space Science* 53 (10), 1079–1087.
- 627 Moores, J. E., Gough, R. V., Martinez, G. M., Meslin, P.-Y., Smith, C. L.,
628 Atreya, S. K., Mahaffy, P. R., Newman, C. E., Webster, C. R., mar 2019.
629 Methane seasonal cycle at gale crater on mars consistent with regolith ad-
630 sorption and diffusion. *Nature Geoscience* 12 (5), 321–325.

- 631 Moussaoui, S., Brie, D., Mohammad-Djafari, A., Carteret, C., 2006. Separation of non-negative mixture of non-negative sources using a bayesian approach and mcmc sampling. *Signal Processing, IEEE Transactions on* [see also *Acoustics, Speech, and Signal Processing, IEEE Transactions on*] 54 (11), 4133–4145.
- 636 Moussaoui, S., Hauksdóttir, H., Schmidt, F., Jutten, C., Chanussot, J., Brie, D., Douté, S., Benediktsson, J., Jun. 2008. On the decomposition of mars hyperspectral data by ica and bayesian positive source separation. *Neurocomputing* 71 (10-12), 2194–2208.
- 640 URL <http://www.sciencedirect.com/science/article/B6V10-4RV17HX-4/1/739950d227add850ec0720718c1c2362>
- 642 Neary, L., Daerden, F., 2018. The gem-mars general circulation model for mars: Description and evaluation. *Icarus* 300, 458–476.
- 644 Neefs, E., Vandaele, A. C., Drummond, R., Thomas, I. R., Berkenbosch, S., Clairquin, R., Delanoye, S., Ristic, B., Maes, J., Bonnewijn, S., Pieck, G., Equeter, E., Depiesse, C., Daerden, F., Ransbeeck, E. V., Nevejans, D., Rodriguez-Gómez, J., López-Moreno, J.-J., Sanz, R., Morales, R., Candini, G. P., Pastor-Morales, M. C., del Moral, B. A., Jeronimo-Zafra, J.-M., Gómez-López, J. M., Alonso-Rodrigo, G., Pérez-Grande, I., Cubas, J., Gomez-Sanjuan, A. M., Navarro-Medina, F., Thibert, T., Patel, M. R., Bellucci, G., Vos, L. D., Lesschaeve, S., Vooren, N. V., Moelans, W., Aballea, L., Glorieux, S., Baeke, A., Kendall, D., Neef, J. D., Soenen, A., Puech, P.-Y., Ward, J., Jamoye, J.-F., Diez, D., Vicario-Arroyo, A., Jankowski, M., sep 2015. NOMAD spectrometer on the ExoMars trace gas orbiter mission: part 1—design, manufacturing and testing of the infrared channels. *Applied Optics* 54 (28), 8494.
- 657 Penttilä, A., Martikainen, J., Gritsevich, M., Muinonen, K., feb 2018. Laboratory spectroscopy of meteorite samples at UV-vis-NIR wavelengths: Analysis and discrimination by principal components analysis. *Journal of Quantitative Spectroscopy and Radiative Transfer* 206, 189–197.
- 661 Schmidt, F., Schmidt, A., Treguier, E., Guiheneuf, M., Moussaoui, S., Dobi-geon, N., 2010. Implementation strategies for hyperspectral unmixing using bayesian source separation. *Geoscience and Remote Sensing, IEEE Transactions* 48 (11), 4003–4013.
- 665 Shashilov, V. A., Xu, M., Ermolenkov, V. V., Lednev, I. K., nov 2006. Latent variable analysis of raman spectra for structural characterization of proteins. *Journal of Quantitative Spectroscopy and Radiative Transfer* 102 (1), 46–61.
- 668 Smith, G. R., Hunten, D. M., 1990. Study of planetary atmospheres by absorp-tive occultations. *Reviews of Geophysics* 28 (2), 117.

- 670 Trokhimovskiy, A., Perevalov, V., Korablev, O., Fedorova, A. F., Olsen, K. S.,
671 Bertaux, J.-L., Patrakeev, A., Shakun, A., Montmessin, F., Lefèvre, F., Luka-
672 shevskaya, A., jul 2020. First observation of the magnetic dipole CO₂ absorp-
673 tion band at 3.3 μm in the atmosphere of mars by the ExoMars trace gas
674 orbiter ACS instrument. *Astronomy & Astrophysics* 639, A142.
- 675 Vandaele, A., Neefs, E., Drummond, R., Thomas, I., Daerden, F., Lopez-
676 Moreno, J.-J., Rodriguez, J., Patel, M., Bellucci, G., Allen, M., Altieri, F.,
677 Bolsée, D., Clancy, T., Delanoye, S., Depiesse, C., Cloutis, E., Fedorova, A.,
678 Formisano, V., Funke, B., Fussen, D., Geminale, A., Gérard, J.-C., Giuranna,
679 M., Ignatiev, N., Kaminski, J., Karatekin, O., Lefèvre, F., López-Puertas,
680 M., López-Valverde, M., Mahieux, A., McConnell, J., Mumma, M., Neary,
681 L., Renotte, E., Ristic, B., Robert, S., Smith, M., Trokhimovsky, S., Auwera,
682 J. V., Villanueva, G., Whiteway, J., Wilquet, V., Wolff, M., dec 2015. *Science*
683 objectives and performances of NOMAD, a spectrometer suite for the
684 ExoMars TGO mission. *Planetary and Space Science* 119, 233–249.
- 685 Vandaele, A. C., , Lopez-Moreno, J.-J., Patel, M. R., Bellucci, G., Daerden,
686 F., Ristic, B., Robert, S., Thomas, I. R., Wilquet, V., Allen, M., Alonso-
687 Rodrigo, G., Altieri, F., Aoki, S., Bolsée, D., Clancy, T., Cloutis, E., De-
688 piesse, C., Drummond, R., Fedorova, A., Formisano, V., Funke, B., González-
689 Galindo, F., Geminale, A., Gérard, J.-C., Giuranna, M., Hetey, L., Ignatiev,
690 N., Kaminski, J., Karatekin, O., Kasaba, Y., Leese, M., Lefèvre, F., Lewis,
691 S. R., López-Puertas, M., López-Valverde, M., Mahieux, A., Mason, J., Mc-
692 Connell, J., Mumma, M., Neary, L., Neefs, E., Renotte, E., Rodriguez-Gomez,
693 J., Sindoni, G., Smith, M., Stiepen, A., Trokhimovsky, A., Auwera, J. V.,
694 Villanueva, G., Viscardy, S., Whiteway, J., Willame, Y., Wolff, M., jun 2018.
695 NOMAD, an integrated suite of three spectrometers for the ExoMars trace gas
696 mission: Technical description, science objectives and expected performance.
697 *Space Science Reviews* 214 (5).
- 698 Villanueva, G., Smith, M., Protopapa, S., Faggi, S., Mandell, A., sep 2018.
699 Planetary spectrum generator: An accurate online radiative transfer suite for
700 atmospheres, comets, small bodies and exoplanets. *Journal of Quantitative*
701 *Spectroscopy and Radiative Transfer* 217, 86–104.