

1 **REVISION 2 – Correction date Feb. 26, 2014**

2 **Multivariate analysis of Raman spectra for the identification of sulfates - Implications**  
3 **for ExoMars**

4 Guillermo Lopez-Reyes<sup>1</sup>, Pablo Sobron<sup>2,3,4</sup>, Catherine Lefebvre<sup>2</sup> and Fernando Rull<sup>1</sup>

5 <sup>1</sup>Unidad Asociada UVA-Centro de Astrobiología. Edificio INDITI, Av. Francisco Valles 8,  
6 Parque Tecnológico de Boecillo, Parcela 203, Boecillo 47151, Spain.

7 <sup>2</sup>Space Science & Technology, Canadian Space Agency, 6767. Rte. de l'Aéroport, St.  
8 Hubert, Quebec J3Y 8Y9, Canada.

9 <sup>3</sup>MalaUva Labs, 822 Allen Ave #A, St. Louis, Missouri, 63104, U.S.A.

10 <sup>4</sup>SETI Institute, 189 Bernardo Ave #100, Mountain View, California 94043, U.S.A.

11 **ABSTRACT**

12 We have built three multivariate analysis mathematical models based on principal  
13 component analysis (PCA), partial least squares (PLS), and artificial neural networks  
14 (ANN) to detect sulfate minerals in geological samples from laser Raman spectral data. We  
15 have critically assessed the potential of the models to automatically detect and quantify the  
16 abundance of selected Ca-, Fe-, Na-, and Mg-sulfates in binary mixtures. Samples were  
17 analyzed using a laboratory version of the Raman Laser Spectrometer (RLS) instrument  
18 onboard the European Space Agency 2018 ExoMars mission. Our results show that PCA  
19 and PLS, can be used to quantify to some extent the abundance of mineral phases. PCA  
20 separated hydrated from dehydrated mixtures and classified mixtures depending on the  
21 phases abundances. PLS provided relatively good calibration curves for these mixtures.  
22 Upon spectral pre-processing, ANN provided the most precise qualitative and quantitative  
23 results. The detection of mineral phases was 100% accurate for pure samples, as was for  
24 binary mixtures where the abundance of mineral phases was > 10%. The outputs of the

25 ANN were proportional to the phase abundance of the mixture, thus demonstrating the  
26 ability of ANN to quantify the abundance of different phases without the need for  
27 calibration. Taken together, our findings demonstrate that multivariate analysis provides  
28 critical qualitative and quantitative information about the studied sulfate minerals.

29 **Keywords:** Sulfates, ExoMars, Raman spectroscopy, Multivariate Analysis,  
30 Qualitative, Quantitative

### 31 **INTRODUCTION**

32 Laser Raman spectroscopy has been proposed as a powerful tool for the identification  
33 of minerals in the context of planetary exploration, including Mars (Sharma et al., 2003;  
34 Sobron et al., 2013a; Sobron et al., 2008; Wang et al., 2003; Wiens et al., 2007), Europa  
35 (Angel et al., 2012; Sobron et al., 2013b), Venus (Lambert et al., 2010), the Moon (Ling et  
36 al., 2009), and asteroids (Kong and Wang, 2010). In addition, the feasibility of using laser  
37 Raman spectroscopy for the detection of biosignatures in terrestrial analogues to Mars has  
38 been demonstrated, e.g. (Bower et al., 2013; Dickensheets et al., 2000; Edwards et al.,  
39 2012; Edwards et al., 2011; Edwards et al., 2003; Ellery and Wynn-Williams, 2003; Steele  
40 et al., 2010; Wynn-Williams and Edwards, 2000). A Raman Laser Spectrometer (RLS) is  
41 part of the science instrument payload of the European Space Agency 2018 ExoMars  
42 mission; the RLS instrument will target mineralogical and astrobiological investigations on  
43 the surface and subsurface of Mars (Rull et al., 2011a; Rull et al., 2011b).

44 The current concept of operation of the RLS instrument is a raster analysis of crushed  
45 drill-core materials (Rull et al., 2011a). In this configuration, the geological and  
46 morphological context of the spots analyzed by RLS will be lost, as the crushing stage will  
47 preclude correlation between RLS spectra and the imagery acquired by the rover Close-Up  
48 Imager CLUPI (Josset et al., 2012). While the synergy between these instruments in the

49 current ExoMars payload configuration has been demonstrated (Lopez-Reyes, 2013a),  
50 identification of the mineral phases present in the geological targets and quantification of  
51 their abundance with RLS will mostly rely on spectral data and not morphology or texture.  
52 Therefore, in order to enable unambiguous identification and quantification of phase  
53 abundance, robust spectral processing methods are needed.

54 Mineral identification using Raman spectroscopy is often performed by comparing  
55 acquired spectra to reference spectra available from the literature and different databases  
56 e.g., the RRUFF project database (Downs, 2006). A number of algorithms have been  
57 developed that enable an automated identification of Raman spectra using traditional  
58 univariate analysis, *i.e.*, the description of individual variables in a given spectrum  
59 (Hermosilla, 2013; Kriesten et al., 2008; Perez-Pueyo et al., 2004; Sobron et al., 2008).  
60 These algorithms, however, fail to accurately estimate mineral abundance in complex  
61 geological samples, though applications for the quantitative analysis of relatively simple  
62 mixtures have been proposed (Lopez-Reyes, 2013a; Schumacher et al., 2011; Vagenas,  
63 2003).

64 Multivariate Analysis Techniques (MVAT) are statistical techniques which deal with  
65 simultaneous measurements on many variables, and aim at understanding the relationships  
66 between these many variables to predict the values of important properties not directly  
67 measurable (Johnson and Wichern, 2002). Some examples of MVAT applied to the  
68 analysis of Raman spectra in the literature show that principal component analysis (PCA) is  
69 capable of differentiating mineral species such as carbonates, sulfates, oxides and silicates  
70 in geological samples (Lafuente, 2012). Partial least squares (PLS) has been used to  
71 determine the quality of biodiesel fuels (Ghesti et al., 2007). Artificial neural networks  
72 (ANN) have been designed for the identification and quantification of inorganic salts in

73 water solutions with Raman spectra (Dolenko et al., 2005). Also, combinations and  
74 comparisons of these techniques for chemometrical analysis from Raman spectra (mostly  
75 qualitative) have been reported (Dorfer et al., 2010; Ishikawa and Gulick, 2013; Özbalci et  
76 al., 2013).

77 In this work we evaluate the feasibility of using PCA, PLS, and ANN for the  
78 identification and quantification of sulfate salts in binary mixtures through analysis of  
79 spectra recorded using a breadboard of the flight Raman Laser Spectrometer. The use of  
80 binary mixtures is a first step to evaluate these techniques, prior to being tested with more  
81 complex samples. We utilized a set of pure sulfates synthesized in the laboratory in order to  
82 prepare the binary mixtures. The sulfates we have considered have been proposed as  
83 priority targets for astrobiological investigation of Mars (Chou et al., 2013; King and  
84 McLennan, 2010; Knoll et al., 2005; Wang et al., 2011) because of their association with  
85 liquid water on certain locations on Mars and sulfate reducing bacteria in terrestrial  
86 analogue environments (Fernandez-Remolar et al., 2012; Nixon et al., 2013; Sanchez-  
87 Andrea et al., 2012). Such environments are known to preserve chemical and  
88 morphological fossils (Bonny and Jones, 2003), thus testifying to the importance of the  
89 detection and identification of sulfates and the detailed study of their degree of hydration in  
90 order to address Mars' hydrologic history and potential for habitability.

91 In terms of identification, Raman spectroscopy is a very powerful technique for the  
92 analysis of the mineralogy of terrestrial sulfate samples analogue to Mars (Sobron et al.,  
93 2014; Sobron and Alpers, 2013; Sobron et al., 2009), as well as for the detailed  
94 characterization of the hydration states of this kind of salts, which is critical for the rigorous  
95 interpretation of the hydrologic history of Mars (Chou et al., 2013; Ling et al., 2008; Wang  
96 et al., 2006).

97

## MATERIALS AND METHODS

98

### Samples

99 A set of 17 Raman spectra of sulfates were used as input for training the PCA, PLS,  
100 and ANN MVAT models; hereinafter, we will refer to this set of spectra as training set.  
101 Table 1 lists the materials used to record our set of spectra. The materials were synthesized  
102 in the laboratory using standard techniques (Ling et al., 2008) -- they were characterized  
103 using X-ray diffraction in order to certify their mineralogical composition. The Raman  
104 spectra of these materials were acquired using Raman instrumentation described in (Ling et  
105 al., 2008; Wang et al., 2006).

106 Apart from the training set, two more sets of spectra have been defined: the validation  
107 and the test sets. A random selection of the three sets of spectra spectra is shown in Figure  
108 1. The training set was divided in four subsets: Ca-, Mg-, Fe-, and Na-sulfates. Each of  
109 these four groups included sulfates with several degrees of hydration. The validation set  
110 consists on a set of binary mixed spectra with different proportions (0.1:0.9, 0.25:0.75,  
111 0.5:0.5, 0.75:0.25, 0.9:0.1) of all mineral pairs of the samples. It was synthetically  
112 generated by computing linear combinations of these spectra, parameterized with the  
113 expected proportion of the mixture and the cross-section of the mixed materials.  
114 Randomized noise with realistic amplitude was added to the synthetic spectra to guarantee  
115 differentiation. Though probably not totally accurate, this set of spectra is expected to  
116 behave similarly to weight-proportion mixtures, providing an easy and convenient way to  
117 execute the models with hundreds of spectra without the actual need of preparing mixtures  
118 and acquiring spectra.

119 The third set of spectra, the test set, included the Raman spectra of powdered  
120 mixtures of anhydrite (CaSO<sub>4</sub>), thenardite (NaSO<sub>4</sub>), and MgSO<sub>4</sub>. This set of spectra is

121 used to test the models as well as to assess the goodness of using the computed spectra of  
122 the validation set as part of the model training/validation procedure. The Raman spectra of  
123 anhydrite and MgSO<sub>4</sub> show non-overlapping peaks with thenardite, thus facilitating the  
124 spectral processing described below. Two subsets of samples (anhydrite + thenardite, and  
125 thenardite + MgSO<sub>4</sub>) with a total of 7 samples each were prepared by mixing these  
126 materials in the following weight proportions: 0.01:0.99, 0.1:0.9, 0.25:0.75, 0.5:0.5,  
127 0.75:0.25, 0.9:0.1 and 0.99:0.01. These mixtures were acquired with the RLS instrument,  
128 contrary to the training and validation spectra. This way, the models independence with  
129 respect to the instruments responses could be assessed. Thirty spectra of each mixture were  
130 acquired and averaged, providing one final spectrum of each mixture. These spectra are  
131 representative of the overall behavior of the samples, as the average of different  
132 acquisitions of 30 spectra have proved to be almost equal when obtained from an ExoMars-  
133 type powdered sample with the RLS instrument (Lopez-Reyes. 2013a).

#### 134 **Instrument and Raman spectroscopy**

135 In its current configuration, the ExoMars rover will crush the subsurface drilled  
136 samples and provide the instruments a flattened surface of the powdered sample for  
137 analysis. The RLS instrument, in its baseline mode operation, will analyze from 20 to 40  
138 points of each sample, with a 50 μm spot size and an irradiance level of 0.6 – 1.2 kW/cm<sup>2</sup>  
139 of a 532 nm continuous wave laser (Rull et al., 2011a). In order to test the analytical  
140 capabilities of the instrument in an operation-like environment, including fully automated  
141 analysis on powdered samples, a RLS ExoMars Simulator has been developed at the  
142 University of Valladolid – CSIC – Center of Astrobiology Associated Unit ERICA  
143 (Foucher, 2012; Lopez-Reyes, 2013a; Rull, 2011). The Simulator was used to acquire the  
144 test set spectra from the mixtures with different phase abundances.

145 In order to improve the accuracy of the analytical models described below, and given  
146 that the aim of the models is to distinguish among different types of sulfates, only spectral  
147 regions that are relevant to sulfates were considered. These spectral regions are: (1) the  
148 sulfate symmetric stretching  $\nu_1$  (950 to 1100  $\text{cm}^{-1}$ ); (2) the sulfate asymmetric stretching  $\nu_3$   
149 (1100 to 1220  $\text{cm}^{-1}$ ); (3) the sulfate symmetric and asymmetric bending  $\nu_2$  and  $\nu_4$ ,  
150 respectively (100-750  $\text{cm}^{-1}$ ); (4) the water bending (1600-1700  $\text{cm}^{-1}$ ); and (5) the water and  
151 OH stretching (2800-3800  $\text{cm}^{-1}$ ). The spectra were baseline-corrected to remove  
152 background contributions (e.g., fluorescence), and normalized in intensity in a way such the  
153 maximum peak intensity is 1. In our case, the separation between consecutive points of the  
154 spectrum is 0.5  $\text{cm}^{-1}$ , which provides a total number of variables of about 4000 (each  
155 variable corresponding to the intensity at a determined wavenumber of the selected spectral  
156 regions), which is the size of the input data to the different models.

## 157 **MVAT**

158 This section describes the MVAT that have been used in this work: Principal  
159 Component Analysis (PCA), Partial Least-Squares regression (PLS) and Artificial Neural  
160 Networks (ANNs). PCA and PLS extract latent variables from the system in order to  
161 represent the system in a complexity-reduced variable system. Ideally, the extracted latent  
162 variables will respond to physical properties of the model. The difference between PCA and  
163 PLS is that, while PCA extracts the variables, PLS also performs a regression on the  
164 expected responses of the system for a determined set of inputs. ANNs, on the other hand,  
165 is a technique that can model any non-linear function by example-based training of  
166 computational networks. While ANNs philosophy makes it possible to have direct outputs  
167 of the sample presence and abundance, PCA and PLS provide responses based on the latent  
168 variables of the system that need to be classified and/or calibrated to extract the mineral

169 phases presence/abundance values. The scope of this work does not include the  
170 classification step for PCA and PLS responses, but it aims at evaluating the techniques  
171 ability to differentiate among such samples. All these techniques are described below in  
172 more detail, and were implemented for this work using Mathworks MATLAB.

### 173 **Principal Component Analysis**

174 PCA is a multivariate technique that computes the variance-covariance structure of a  
175 set of variables through a few linear combinations of them, with the objective of reducing  
176 the number of data variables (Johnson and Wichern, 2002). PCA calculates new variables  
177 called Principal Components (PCs) as linear combinations of the original variables. All the  
178 principal components are orthogonal to each other, so the variables in the principal  
179 component space do not provide redundant information. The principal components as a  
180 whole form an orthogonal basis for the space of the data, which can thus be represented in  
181 this new space. There is an infinite number of ways to construct an orthogonal basis to  
182 represent the data, so PCA calculates the PCs taking into account that they are uncorrelated  
183 among them, while maximizing their variance with coefficient vectors of unit length (as the  
184 variance can easily be increased by multiplying by a constant (Johnson and Wichern,  
185 2002)). In other words, the first PC is calculated, among all the possibilities, as the single  
186 axis (linear combination of the original variables) in the space that provides the greatest  
187 variance by any projection of the data on that axis. The second PC is another single axis  
188 that provides the second greatest variance by any projection of the data, and that is  
189 orthogonal to the first PC. The third will be also perpendicular to the previous two, and so  
190 on. This way, the set of PCs conforms a new set of coordinates which can represent the  
191 original data, where most of the variance of the system can be explained with only a few of  
192 the PCs. This is due to the fact that, usually, many of the variables of a system are highly



193 correlated, making it possible to remove some of them and still have a variable set which  
194 can account for most of the variability of the system.

195 The PCA model was trained with the 17 spectra from the pure sulfates (training set),  
196 and then applied to the validation and test spectra sets.

### 197 **Partial Least Squares**

198 Partial Least Squares (PLS) is a common term for a family of multivariate modeling  
199 methods that appeared in the 1980's to solve problems in social sciences (e.g. (Wold et al.,  
200 1983)), but that can be applied to a large variety of modeling problems (Martens and Naes,  
201 1992), including Raman spectroscopy, as discussed above. The underlying assumption of  
202 all PLS methods is, as for PCA, that the observed data is generated by a system which is  
203 driven by a small number of latent (not directly observed or measured) variables, called  
204 components. In addition, PLS performs a regression of the expected responses for each set  
205 of input variables (observations). Thus, this technique is some kind of combination of PCA  
206 and linear regression: PLS creates orthogonal score vectors (the latent vectors or  
207 components) by maximizing the covariance between two different blocks of variables  
208 (Rosipal and Krämer 2006) which correspond to the input variables (predictor variables or  
209 observations) and the expected responses (predicted variables or predictions). The higher  
210 the number of computed components, the better the model fits the system. In our work, the  
211 SIMPLS algorithm (de Jong, 1993) has been used, which calculates the PLS factors directly  
212 as linear combinations of the original variables. With SIMPLS, the PLS factors are  
213 determined such as to maximize the covariance between the predictor variables and the  
214 expected responses, while obeying certain orthogonality and normalization restrictions (de  
215 Jong, 1993).

216 The responses chosen for the analysis of Raman spectra of sulfates were the weight  
217 atomic fraction of the sulfate cation (calcium, magnesium, iron and sodium), as well as the  
218 weight ratio of the water bound to the sulfate molecule, providing a set of five responses.  
219 These responses were chosen to reflect the composition of the samples, thus providing the  
220 model with a physical basis for the regression of the input variables. Given that the shifts in  
221 the Raman bands are produced as a consequence of the frequency shifts associated to the  
222 cation-sulfate group oscillators, the cation ratio is the subjacent physical property that we  
223 choose as a response in our model. The spectra were baseline-removed, normalized with  
224 respect to the maximum peak height, mean centered and scaled in order to avoid biasing the  
225 model.

226 The model was trained with the spectra from the 17 pure salts (training set). As in  
227 every model where input variables are regressed to expected responses, there is a risk of  
228 over-fitting. To minimize this effect, we optimized our model using a leave-one-out cross-  
229 validation method with the validation and test spectra sets, as outlined in the Results and  
230 Discussion section.

### 231 **Artificial Neural Network**

232 An Artificial Neural Network (ANN) is a mathematical procedure for transforming  
233 inputs into desired outputs using highly connected networks of relatively simple processing  
234 units called neurons (Johnson and Wichern, 2002). Each neuron performs a mathematical  
235 operation that produces an output which is a function (usually non-linear) of a series of  
236 biased and weighted inputs coming from other neurons. The underlying principle is that  
237 neural networks be modeled after the neural activity in the human brain, in which the  
238 interconnection of very simple functional units (the neurons) can solve many complex and  
239 non-linear problems in a very fast way. Thus, ANNs consist on parallel computational

240 models comprised of densely interconnected adaptive processing units (the neurons). These  
241 networks provide a tool for modelling nonlinear static or dynamic systems, making use of  
242 their adaptive nature, based on “learning by example”. This feature makes this kind of  
243 computational models very appealing in applications in which the understanding of the  
244 problem is little or incomplete, but where training data is readily available (Hassoun, 1995).

245 ANNs can provide very fast outputs, as it only has to compute a limited number of  
246 very simple operations that are easily parallelizable. However, the design and training of  
247 the network can be a hard task. ANNs are set in layers of interconnected networks, in which  
248 the first layer has as many neurons as inputs in the system, and the output layer has as many  
249 neurons as required outputs. All intermediate layers are called hidden layers, and can have  
250 any number of neurons (Johnson and Wichern, 2002).

251 For our design, many different architectures were trained and evaluated to obtain the  
252 network with the best performance. In the end, a three-layer network with 33 neurons on  
253 the hidden layer was chosen. The input layer was configured with 33 neurons each  
254 corresponding to determined spectral positions. The neurons were configured with log-  
255 sigmoid transfer functions. The output layer consisted on 17 outputs, each corresponding to  
256 one of the sulfates (where each output should take values between 0 and 1, proportional to  
257 the abundance of the sample).

258 As suggested by (Koujelev et al., 2010), and in order to improve the ANN model  
259 performance, only a selection of spectral positions is fed into the network as input; the  $\nu_1$   
260 peak positions of all the sulfates, as well as the most intense non-overlapping secondary  
261 peaks of the pure samples were selected as inputs. This allows the definition of a predefined  
262 set of spectral positions that will be extracted from the input spectra and fed to the network.  
263 In order to determine these peaks, the input spectra to the ANN have to be pre-processed

264 so that only the most intense peaks have non-zero intensities, as depicted in Figure 2: from  
265 all the inputs to the model, those below a determined threshold will be set to 0. The  
266 definition of this threshold depends on the sample noise and can be decided for each  
267 spectrum. In addition, a deconvolution of peaks is needed for mixtures where the peaks are  
268 partially overlapped.

269 The training process was performed using a Levenberg-Marquardt back-propagation  
270 algorithm (Levenberg, 1944; Marquardt, 1963) that used spectra of pure sulfates (training  
271 set) plus some spectra of mixtures with 0.25:0.75, 0.5:0.5 and 0.75:0.25 proportions. The  
272 network was trained to provide outputs proportional to the abundance of each sulfate. To  
273 avoid over-fitting, the early-stopping technique with the validation set was used. This  
274 consists in stopping the iterative training process when the output errors for the validation  
275 set increase. Finally, to test the network in a more representative scenario, the test set from  
276 the RLS instrument was fed to the network.

## 277 **RESULTS AND DISCUSSION**

### 278 **PCA**

279 The training of the PCA model showed that the first three components PC1, PC2 and  
280 PC3 explain more than 80% of the variance of the training data set, and 90% if the first five  
281 components are considered (Fig. 3). This means that 90% of the variance of the system is  
282 explained with only 5 of the calculated orthogonal variables. A dendrogram showing the  
283 interconnections for the 3 PCs model is depicted in Figure 4. Models with higher number of  
284 components (up to ten) do not present better separation among different cations, and  
285 worsen the discrimination between low- and high-hydration sulfates. In addition, the lower  
286 the number of components, the more general the model can be, so the 3-PC model was  
287 selected as the optimal one. The representation of the scores of the training samples in the

288 new variable system (with only the first two components, for representation convenience) is  
289 displayed in Figure 5 as circles. From the validation set, only the calculated scores for  
290 mixed spectra of sulfates of the same cation in 50:50 proportions (worst case) are plotted in  
291 Figure 5 as triangles. The mixed spectra contain sulfates with different degree of hydration.  
292 The scores in Figure 5 show how PCA succeeds in separating, mostly along PC1, the low  
293 hydration mixtures from the high hydrated ones, though it fails to distinguish among  
294 different types of cations. For example, the model scores the mixture of Fe-sulfates with 4  
295 and 7 water molecules between those two elements, and the same happens with the Ca-  
296 sulfates. However, the model fails to correctly separate the Mg-sulfates by hydration level.  
297 Therefore, the general conclusion is that PC1 can only be used to separate low hydration  
298 from high hydration sulfates, but not to distinguish among different hydration states of  
299 same-cation sulfates.

300       While no additional direct associations between the principal components and the  
301 physical properties of the different molecules (as for example the cation ratio) can be  
302 inferred, the PCA analysis of sulfates can provide useful information when representing the  
303 PC1-PC2 scores. For example, the representation of the scores of the test spectra set is  
304 depicted in Figure 6 (for graphical simplicity, only two PCs are represented). This figure  
305 shows how the mixtures are placed between the pure components depending on their  
306 relative abundance: the higher the abundance of a sulfate of the mixture, the closer to the  
307 corresponding pure sulfate score. This would mean that some kind of quantification could  
308 be possible based on the PCA model, even when only trained with pure Raman spectra.

### 309       **PLS**

310       The PLS regression of the pure sulfates spectra to their expected weight cation atomic  
311 fractions and water presence was performed for different numbers of components. To

312 decide the optimal number of components, the Mean Squared Error (MSE) of the  
313 predictions for the training, validation and tests sets was calculated. Figure 7 shows some  
314 spectra from the training set and the corresponding loadings, where the energies that have  
315 the greatest effect on the PLS predictions can be observed. The application of the model to  
316 the spectra sets yielded the prediction errors shown in Figure 8 – a. As expected, the  
317 prediction error for the training samples (blue line) tends to 0 with increasing number of  
318 components, as the model is fitted better. The values of the error for the regression of the  
319 test set of spectra (black lines) indicates that the best fitting (minimum prediction error)  
320 occurs for a 12-component PLS model. As a general rule, lower numbers of components  
321 imply a more general response of the model. Since the spectra of the natural samples were  
322 acquired with a different experimental setup than the spectra of the training and validation  
323 sets, we interpret the 13-component model (where the minimum prediction error for the  
324 validation set is found – see red line) as one in which the 13<sup>th</sup> component accounts for the  
325 spectrometer response. Thus, the 12-component model is considered as the optimum for our  
326 set of sulfates spectra; more than 98% of the variance of the system can be explained with  
327 this model (Fig. 8 – b), apparently with no over-fitting.

328         The average absolute prediction error values for the training, validation and test sets  
329 are presented in Table 2 for the 12-component model. This table shows how the average  
330 error is close to 0 for the training and validation sets, while the model prediction is biased  
331 in the prediction values for the test set spectra. The RMSEP (Root Mean Square Error of  
332 Prediction) for each test is presented in Table 3. This value represents the prediction error  
333 deviation, to compare the prediction accuracy between the different spectra sets. In  
334 addition, the predicted vs. expected responses correlations are shown in Table 4, and the  
335 responses represented in Figure 9. This figure presents the expected vs. calculated cation

336 ratios for each spectrum. As each spectrum is represented as one individual point in the  
337 graph, the estimation error is defined by the y-axis displacement from the expected value.  
338 Thus, the estimation error is 0 when the point is placed on the line with unitary slope ( $y =$   
339  $x$ ). These data show that a correlation between the predicted and the expected responses is  
340 present with this model for all the training, validation and test spectra sets.

341 The RMSEP results for the training set in Table 3 show that the model has the lower  
342 prediction capabilities for the hydration response, which can be explained by the influence  
343 of the many close-to-zero values of these responses, as shown in the plot of the training  
344 responses in Figure 9 – a. It is important to note how the results of the test samples provide  
345 too high values of the iron cation response, especially in Figure 9 – d, when it should  
346 always be 0.

347 For the validation spectra, the prediction accuracy is lower (higher RMSEP) than for  
348 the training spectra, as expected, which is also reflected in its lower correlation and higher  
349 prediction bias. The model behavior is worse for well-balanced mixtures (0.5:0.5) than for  
350 mixtures with unbalanced proportions, as the model was trained with pure samples only.  
351 Most of the outlier points observed in Figure 9 – b belong to the mixtures in this proportion.  
352 The values in Table 4 for the validation set correspond to the averaged correlation for all  
353 the proportions, but these improve between 0.5% and 1% when the 0.5:0.5 mixture is not  
354 considered.

355 The results for the test spectra show better correlations than the training set in the Na-  
356 sulfates case, though the RMSEP value is much higher than for the training and validation  
357 sets. This can be explained by the fact that the correlation value for the training samples is  
358 biased with all the spectra from the rest of the sulfates, which are not represented with the  
359 test samples. This can be readily observed in Figure 9 – a, b vs. Figure 9 – c, d; the slopes

360 of the curves are more or less unitary for the training and validation tests, which indicates a  
361 certain accuracy of the model (and which is reflected in similar RMSEP values and low  
362 bias), while the test sets show a non-unitary slope linear correlation between the expected  
363 and calculated sodium response (with a much higher RMSEP value and prediction bias).  
364 This implies that the 12-component PLS model seems to fail to directly predict hydration  
365 states and cation abundances for this case. However, it still provides linear calibration  
366 curves that could be used to compute these values.

### 367 ANN

368 The training process of our ANN consisted on providing outputs proportional to the  
369 abundance of the materials. As ANN can model any non-linear function, this seemed to be  
370 the most convenient way to do it, contrary to PLS, where the underlying physical principles  
371 tried to be modeled.

372 In order to evaluate the identification accuracy of the ANN, we established the  
373 criteria to consider that a sulfate is detected when the corresponding output is higher than a  
374 determined threshold. For this network, this threshold was set to 0.06. As the output of the  
375 network ranges from 0 to 1, this value will be the theoretically lowest detection threshold  
376 for this model (i.e. no sulfates will be detected below concentrations of 6%).

377 Under these premises, the training set spectra were detected with 100% accuracy with  
378 this ANN. Furthermore, the major phase present in all the samples of validation and test  
379 sets was also detected in 100% of the cases. Both phases of the mixtures were detected with  
380 100% accuracy only when the minor phase was present with at least 10% abundance. In  
381 other words, the ANN detects 100% of the minerals present in binary mixtures with  
382 proportions ranging from 10:90 to 90:10. This implies that the ANN model provides a  
383 robust system for the qualitative detection of sulfates in this kind of mixtures, with a



384 detection threshold for minor phases of around 10%, even for spectra acquired with a  
385 different hardware setup than the training spectra.

386         The representation of the ANN outputs for spectra of mixtures with different  
387 proportions (validation and test sets) shows that the network outputs also provide  
388 information on the relative abundance of the materials. As part of the validation set was  
389 used to train the model, the results for this set proved very good also in terms of  
390 quantification of mineral abundance, as expected. More interesting are the results for the  
391 test spectra, which are shown in Figures 10 and 11, where the estimated concentration  
392 values from the ANN with respect to the expected ones are represented. These results show  
393 that the model accuracy might somehow depend on the samples (e.g., the results for the  
394 mixture in Fig. 11 show that the Mg-sulfate concentration tends to be underestimated for  
395 mixtures where it is the major component). However, a certain degree of correlation  
396 between the modeled values and the actual abundances of the mixtures exists. This is an  
397 interesting result, especially bearing in mind that this is true for the test set spectra, while  
398 the model was trained with pure spectra obtained with a different spectrometer (training  
399 set) and computed spectra of mixtures calculated from those samples.

400         The representation of the maximum and minimum values of the outputs of the ANN  
401 which do not correspond to the minerals present in the spectra can be seen in Figures 10  
402 and 11 – b. The representation of these values, which should always be 0, is of relevance to  
403 show that, with a threshold of 6%, no false identifications are obtained. The conclusion is  
404 that the use of ANN looks promising for providing robust and reliable results under the  
405 described premises, not only for identifying the phases present in binary mixtures of  
406 sulfates, but also to provide some kind of quantification of their abundance.

407 **MVAT comparative**

408 PCA, PLS and ANN models have been trained based-only in 17 spectra of pure  
409 sulfates and mixed spectra computed as linear combinations of those. This procedure has  
410 allowed evaluating these analytical techniques without the need to actually prepare all  
411 possible combinations of samples for the calibration of the models.

412 The analysis of selected regions from the Raman spectrum implies the analysis of  
413 several thousands of variables at the same time. To deal with this amount of information,  
414 PCA and PLS calculate new sets of orthogonal variables as linear combinations of the  
415 original ones. PLS then regresses these variables to expected responses, while PCA doesn't.  
416 These new variables correspond to latent variables which ideally should be directly related  
417 to physical properties of the system (e.g., the degree of hydration of a sulfate). However,  
418 this is not always the case. This is probably due to the non-linear nature of the Raman  
419 emission, which PCA and PLS try to model with linear processing. ANNs, on the other  
420 hand, can provide non-linear transfer functions. Thus, they might be more adequate for the  
421 modeling of non-linear effects (as the Raman emission). However, ANNs require a  
422 relatively low number of input variables to provide any relevant results. This technique  
423 does not perform a reduction to latent variables on its own, so it has been made by only  
424 inputting the most relevant spectral positions (corresponding to the wavenumbers of the  
425 most representative peaks of the training spectra).

426 We have shown the ability of these MVAT models of providing useful qualitative and  
427 even quantitative information for simple binary mixtures, even when the training and  
428 testing were performed with spectra recorded with different hardware setups. As discussed  
429 in the previous sections, PCA separated low from high hydrated sulfates, and also  
430 somewhat classified samples of mixtures depending on their relative abundance. PLS

431 model outputs presented good correlations to the expected responses. However, though  
432 well correlated, the responses in some cases were relatively far from the expected values.  
433 The conclusion is that PCA and PLS provided a classification method which needs a  
434 previous calibration for the sample under analysis, and a classification method. On the  
435 other hand, the ANN model outputs directly provided the abundance of the corresponding  
436 salt, in addition to a 100% qualitative detection for mixtures with abundances as low as  
437 10%.

438 To overcome the various limitations of these MVATs, a synergy between them might  
439 be interesting to improve the overall performance of the models. Some classifiers for the  
440 qualitative analysis of minerals have been proposed based on integrated PCA and ANN  
441 models (Dorfer et al., 2010; Ishikawa and Gulick, 2013). Future research will thus focus in  
442 developing models for the quantification of mineral abundances which integrate different  
443 MVATs.

#### 444 **IMPLICATIONS FOR EXOMARS**

445 ExoMars' RLS instrument will determine the structural and compositional features of  
446 materials in rocks and soils at the surface and subsurface of Mars. The ExoMars samples  
447 will be collected by a drill, then crushed and delivered to a suite of instruments located in  
448 the rover's analytical laboratory, where the RLS instrument sits. A crushing station will  
449 provide homogenized powdered samples that will likely feature complex mixtures of  
450 mineral phases.

451 Fast, robust, unsupervised RLS data processing tools able to interpret the intricate  
452 spectra that will be obtained from Martian samples would benefit ExoMars mission  
453 operations in that they may directly support the daily tactical operations of the rover. The  
454 different multivariate analysis techniques methodologies – PCA, PLS, and ANN – we have

455 discussed here promise to provide an efficient way to process RLS data during the ExoMars  
456 mission.

457 Future work in this direction will focus on exploring the capabilities of these  
458 methodologies to evaluate more complex mixtures; these will include additional synthetic  
459 sulfates, oxides, clays, phyllosilicates, carbonates, perchlorates, as well as natural samples.

460 We will carry out this research using the RLS ExoMars Simulator we have developed  
461 at the University of Valladolid-Centro de Astrobiologia and the different RLS models and  
462 prototypes that we are developing with the Spanish National Institute for Aerospace  
463 Technology (INTA).

#### 464 **ACKNOWLEDGMENTS**

465 We kindly thank A. Wang for providing the Raman spectra of the pure sulfate  
466 samples used in this work and A. Koujelev for the helpful insights on the implementation of  
467 ANN. R. Leveille and the Canadian Space Agency facilitated this project. Preliminary  
468 results of this work have been reported at the 15th Annual Conference of the International  
469 Association for Mathematical Geosciences (Lopez-Reyes, 2013b). GLR acknowledges the  
470 University of Valladolid (Spain) for providing funding for the project. PS and CL  
471 acknowledge support from the Natural Sciences and Engineering Research Council of  
472 Canada (NSERC) and the Canadian Space Agency.

#### 473 **REFERENCES**

474 Angel, S.M., Gomer, N.R., Sharma, S.K., and McKay, C. (2012) Remote Raman  
475 spectroscopy for planetary exploration: A review. *Applied Spectroscopy*, 66(2), 137-150.

- 476 Bonny, S., and Jones, B. (2003) Microbes and mineral precipitation, Miette Hot Springs,  
477 Jasper National Park, Alberta, Canada. *Canadian Journal of Earth Sciences*, 40(11), 1483-  
478 1500.
- 479 Bower, D.M., Steele, A., Fries, M.D., and Kater, L. (2013) Micro Raman spectroscopy of  
480 carbonaceous material in microfossils and meteorites: improving a method for life  
481 detection. *Astrobiology*, 13(1), 103-113.
- 482 Chou, I.M., Seal, R.R., and Wang, A. (2013) The stability of sulfate and hydrated sulfate  
483 minerals near ambient conditions and their significance in environmental and planetary  
484 sciences. *Journal of Asian Earth Sciences*, 62, 734-758.
- 485 de Jong, S. (1993) SIMPLS: An alternative approach to partial least squares regression.  
486 *Chemometrics and Intelligent Laboratory Systems*, 18(3), 251-263.
- 487 Dickensheets, D.L., Wynn-Williams, D.D., Edwards, H.G.M., Schoen, C., Crowder, C.,  
488 and Newton, E.M. (2000) A novel miniature confocal microscope/Raman spectrometer  
489 system for biomolecular analysis on future Mars missions after Antarctic trials. *Journal of*  
490 *Raman Spectroscopy*, 31(7), 633-635.
- 491 Dolenko, T.A., Burikov, S.A., and Sugonjaev, A.V. (2005) Neural network technologies in  
492 Raman spectroscopy of water solutions of inorganic salts. In H.J. Byrne, E. Lewis, B.D.  
493 MacCraith, E. McGlynn, J.A. McLaughlin, G.D. Osullivan, A.G. Ryder, and J.E. Walsh,  
494 Eds. *Opto-Ireland 2005: Optical Sensing and Spectroscopy*, 5826, p. 298-305.

- 495 Dorfer, T., Schumacher, W., Tarcea, N., Schmitt, M., and Popp, J. (2010) Quantitative  
496 mineral analysis using Raman spectroscopy and chemometric techniques. *Journal of Raman*  
497 *Spectroscopy*, 41(6), 684-689.
- 498 Downs, R.T. (2006) The RRUFF Project: an integrated study of the chemistry,  
499 crystallography, Raman and infrared spectroscopy of minerals. Program and Abstracts of  
500 the 19th General Meeting of the International Mineralogical Association in Kobe, Japan.,  
501 O03-13.
- 502 Edwards, H.G.M., Hutchinson, I., and Ingley, R. (2012) The ExoMars Raman spectrometer  
503 and the identification of biogeological spectroscopic signatures using a flight-like  
504 prototype. *Analytical and Bioanalytical Chemistry*, 404(6-7), 1723-1731.
- 505 Edwards, H.G.M., Hutchinson, I.B., Ingley, R., Waltham, N.R., Beardsley, S., Dowson, S.,  
506 and Woodward, S. (2011) The search for signatures of early life on Mars: Raman  
507 spectroscopy and the Exomars mission. *Spectroscopy Europe*, 23(1).
- 508 Edwards, H.G.M., Newton, E.M., Dickensheets, D.L., and Wynn-Williams, D.D. (2003)  
509 Raman spectroscopic detection of biomolecular markers from Antarctic materials:  
510 evaluation for putative Martian habitats. *Spectrochimica Acta Part A: Molecular and*  
511 *Biomolecular Spectroscopy*, 59(10), 2277-2290.
- 512 Ellery, A., and Wynn-Williams, D. (2003) Why Raman spectroscopy on Mars?--a case of  
513 the right tool for the right job. *Astrobiology*, 3(3), 565-79.
- 514 Fernandez-Remolar, D.C., Preston, L.J., Sanchez-Roman, M., Izawa, M.R.M., Huang, L.,  
515 Southam, G., Banerjee, N.R., Osinski, G.R., Flemming, R., Gomez-Ortiz, D., Prieto

- 516 Ballesteros, O., Rodriguez, N., Amils, R., and Dyar, M.D. (2012) Carbonate precipitation  
517 under bulk acidic conditions as a potential biosignature for searching life on Mars. *Earth  
518 and Planetary Science Letters*, 351, 13-26.
- 519 Foucher, F., Lopez-Reyes, G., Bost, N., Rull, F., Rößmann, P. and Westall, F. (2012) Effect  
520 of grain size distribution on Raman analyses and the consequences for in situ planetary  
521 missions. *Journal of Raman Spectroscopy*.
- 522 Ghesti, G.F., de Macedo, J.L., Resck, I.S., Dias, J.A., and Dias, S.C.L. (2007) FT-Raman  
523 spectroscopy quantification of biodiesel in a progressive soybean oil transesterification  
524 reaction and its correlation with H-1 NMR spectroscopy methods. *Energy & Fuels*, 21(5),  
525 2475-2480.
- 526 Hassoun, M. H., (1995) *Fundamentals of Artificial Neural Networks*. MIT Press.
- 527 Hermosilla Rodriguez, I., Lopez-Reyes, G., Llanos, D.R., and Rull Perez, F. (2014)  
528 Automatic raman spectra processing for Exomars. In E. Pardo-Igúzquiza, C. Guardiola-  
529 Albert, J. Heredia, L. Moreno-Merino, J.J. Durán, and J.A. Vargas-Guzmán, Eds.  
530 *Mathematics of Planet Earth*, p. 127-130. Springer Berlin Heidelberg.
- 531 Ishikawa, S.T., and Gulick, V.C. (2013) An automated mineral classifier using raman  
532 spectra. *Computers & Geosciences*, 54(0), 259-268.
- 533 Johnson, R.A., and Wichern, D.W. (2002) *Applied multivariate statistical analysis*. Prentice  
534 hall Upper Saddle River, NJ.
- 535 Josset, J.-L., Westall, F., Hofmann, B.A., Spray, J.G., Cockell, C., Kempe, S., Griffiths,  
536 A.D., De Sanctis, M.C., Colangeli, L., Koschny, D., Pullan, D., Föllmi, K., Diamond, L.,

- 537 Josset, M., Javaux, E., Esposito, F., and Barnes, D. (2012) CLUPI, a high-performance  
538 imaging system on the ESA-NASA rover of the 2018 ExoMars mission to discover  
539 biofabrics on Mars. EGU General Assembly Conference Abstracts, 14, p. 13616.
- 540 King, P.L., and McLennan, S.M. (2010) Sulfur on Mars. *Elements*, 6(2), 107-112.
- 541 Knoll, A.H., Carr, M., Clark, B., Farmer, J.D., Fischer, W.W., Grotzinger, J.P., McLennan,  
542 S.M., Malin, M., Schroder, C., Squyres, S., Tosca, N.J., and Wdowiak, T. (2005) An  
543 astrobiological perspective on Meridiani Planum. *Earth and Planetary Science Letters*,  
544 240(1), 11-11.
- 545 Kong, W.G., and Wang, A. (2010) Planetary Laser Raman Spectroscopy for surface  
546 exploration on C/D-type asteroids: A case study. 41st Lunar and Planetary Science  
547 Conference, p. Abstract #2730. Lunar and Planetary Institute, Houston.
- 548 Koujelev, A., Sabsabi, M., Motto-Ros, V., Laville, S., and Lui, S.L. (2010) Laser-induced  
549 breakdown spectroscopy with artificial neural network processing for material  
550 identification. *Planetary and Space Science*, 58(4), 682-690.
- 551 Kriesten, E., Alsmeyer, F., Bardow, A., and Marquardt, W. (2008) Fully automated indirect  
552 hard modeling of mixture spectra. *Chemometrics and Intelligent Laboratory Systems*,  
553 91(2), 181-193.
- 554 Lafuente, B., Sansano, A., Navarro, R., Rull, F., Martínez Frías, J., Medina, J., Lopez, G.,  
555 Sobron, P. (2012) Multivariate analysis of Raman spectra for geological classification and  
556 identification: Application to Exomars mission. *GeoRaman 2012*, Nancy.



- 557 Lambert, J.L., Morookian, J., Roberts, T., Polk, J., Smrekar, S., Clegg, S.M., Weins, R.C.,  
558 Dyar, M.D., and Treiman, A. (2010) Standoff LIBS and Raman spectroscopy under venus  
559 conditions. 41st Lunar and Planetary Science Conference, p. Abstract #2608. Lunar and  
560 Planetary Institute, Houston.
- 561 Levenberg, K. (1944) A method for the solution of certain problems in least squares.  
562 Quarterly of Applied Mathematics, 2, 164-168.
- 563 Ling, Z.C., Wang, A., Jollif, B.L., Arvidson, R.E., and Xia, H.R. (2008) A systematic  
564 Raman, Mid-IR, and Vis-NIR spectroscopic study of ferric sulfates and implications for  
565 sulfates on Mars. 39th Lunar and Planetary Science Conference, p. Abstract #1463. Lunar  
566 and Planetary Institute, Houston.
- 567 Ling, Z.C., Wang, A., Jolliff, B.L., Li, C., Liu, J., Bian, W., Ren, X., Mu, L.L., and Su, Y.  
568 (2009) Raman spectroscopic study of quartz in lunar soils from Apollo 14 and 15 missions.  
569 40th Lunar and Planetary Science Conference, p. Abstract #1823. Lunar and Planetary  
570 Institute, Houston.
- 571 Lopez-Reyes, G., Rull, F., Venegas, G., Westall, F., Foucher, F., Bost, N., Sanz, A., Catalá-  
572 Espí, A., Vegas, A., Hermosilla, I., Sansano, A., and Medina, J. (2013a) Analysis of the  
573 scientific capabilities of the ExoMars Raman Laser Spectrometer instrument. European  
574 Journal of Mineralogy, 25(5), 721-733.
- 575 Lopez-Reyes, G., Sobron, P., Lefevbre, C. and Rull, F. (2013b) Application of multivariate  
576 analysis techniques for the identification of sulfates from Raman spectra. 15th International  
577 Association for Mathematical Geosciences Conference, Madrid. Marquardt, D.W. (1963)

- 578 An algorithm for least-squares estimation of nonlinear parameters. *Journal of the Society*  
579 *for Industrial & Applied Mathematics*, 11(2), 431-441.
- 580 Martens, H., and Naes, T. (1992) *Multivariate Calibration*. Wiley.
- 581 Nixon, S.L., Cockell, C.S., and Cousins, C.R. (2013) Plausible microbial metabolisms on  
582 Mars. *Astronomy & Geophysics*, 54(1), 13-16.
- 583 Özbalci, B., Boyacı, İ.H., Topcu, A., Kadılar, C., and Tamer, U. (2013) Rapid analysis of  
584 sugars in honey by processing Raman spectrum using chemometric methods and artificial  
585 neural networks. *Food Chemistry*, 136(3–4), 1444-1452.
- 586 Perez-Pueyo, R., Soneira, M.J., and Ruiz-Moreno, S. (2004) A fuzzy logic system for band  
587 detection in Raman spectroscopy. *Journal of Raman Spectroscopy*, 35(8-9), 808-812.
- 588 Rull, F., Lopez, G., Catala, A., Medina, J., Sansano, A., Sanz, A., Sobron, F. (2011) Raman  
589 spectroscopy analysis on powdered samples inside the Exomars mission. *ESA Congrex*  
590 *Lisbon 2011*.
- 591 Rull, F., Maurice, S., Diaz, E., Tato, C., Pacros, A., and Team., t.R. (2011a) The Raman  
592 Laser Spectrometer (RLS) on the ExoMars 2018 Rover Mission. *42nd Lunar and Planetary*  
593 *Science Conference LPSC*, p. Abstract #2400. Lunar and Planetary Institute, Houston.
- 594 Rull, F., Sansano, A., Díaz, E., Canora, C.P., Moral, A.G., Tato, C., Colombo, M.,  
595 Belenguer, T., Fernández, M., Manfredi, J.A.R., Canchal, R., Dávila, B., Jiménez, A.,  
596 Gallego, P., Ibarria, S., Prieto, J.A.R., Santiago, A., Pla, J., Ramos, G., Díaz, C., and  
597 González, C. (2011b) ExoMars Raman laser spectrometer for Exomars. *Society of Photo-*  
598 *Optical Instrumentation Engineers (SPIE), Conference Series*, 8152, 12.

- 599 Sanchez-Andrea, I., Rojas-Ojeda, P., Amils, R., and Luis Sanz, J. (2012) Screening of  
600 anaerobic activities in sediments of an acidic environment: Tinto River. *Extremophiles*,  
601 16(6), 829-839.
- 602 Schumacher, W., Kuehnert, M., Roesch, P., and Popp, J. (2011) Identification and  
603 classification of organic and inorganic components of particulate matter via Raman  
604 spectroscopy and chemometric approaches. *Journal of Raman Spectroscopy*, 42(3), 383-  
605 392.
- 606 Sharma, S.K., Lucey, P.G., Ghosh, M., Hubble, H.W., and Horton, K.A. (2003) Stand-off  
607 Raman spectroscopic detection of minerals on planetary surfaces. *Spectrochimica Acta Part*  
608 *A: Molecular and Biomolecular Spectroscopy*, 59(10), 2391-2407.
- 609 Sobron, P., Sobron, F., Sanz, A., and Rull, F. (2008) Raman signal processing software for  
610 automated identification of mineral phases and biosignatures on Mars. *Applied*  
611 *Spectroscopy*, 62(4), 364-370.
- 612 Sobron, P., Sanz, A., Acosta, T., and Rull, F. (2009) A Raman spectral study of stream  
613 waters and efflorescent salts in Rio Tinto, Spain. *Spectrochimica Acta part A: Molecular*  
614 *and Biomolecular Spectroscopy*, 71(5), 1678-82.
- 615 Sobron, P., and Alpers, C.N. (2013a) Raman spectroscopy of efflorescent sulfate salts from  
616 iron mountain mine superfund site, California. *Astrobiology*, 13(3), 270-8.
- 617 Sobron, P., Bishop, J., Blake, D., Chen, B. and Rull, F. (2014) Natural Fe-bearing oxides  
618 and sulfates from the Rio Tinto Mars analogue – Critical assessment of VNIR reflectance

- 619 spectroscopy, laser Raman spectroscopy, and XRD as mineral identification tools.  
620 American Mineralogist, This volume.
- 621 Sobron, P., Lefebvre, C., Koujelev, A., and Wang, A. (2013b) Why Raman and LIBS for  
622 exploring icy moons? 44th Lunar and Planetary Science Conference, p. Abstract #2381.  
623 Lunar and Planetary Institute, Houston.
- 624 Steele, A., McCubbin, F.M., Fries, M., Glamoclija, M., Kater, L., and Nekvasil, H. (2010)  
625 Graphite in an Apollo 17 impact melt breccia. *Science*, 329(5987), 51-51.
- 626 Vagenas, N.V., Kontoyannis, C. G. (2003) A methodology for quantitative determination of  
627 minor components in minerals based on FT-Raman spectroscopy - The case of calcite in  
628 dolomitic marble. *Vibrational Spectroscopy*, 32(2), 261-264.
- 629 Wang, A., Freeman, J.J., Jolliff, B.L., and Chou, I.M. (2006) Sulfates on Mars: A  
630 systematic Raman spectroscopic study of hydration states of magnesium sulfates.  
631 *Geochimica Et Cosmochimica Acta*, 70(24), 6118-6135.
- 632 Wang, A., Haskin, L.A., Lane, A.L., Wdowiak, T.J., Squyres, S.W., Wilson, R.J., Hovland,  
633 L.E., Manatt, K.S., Raouf, N., and Smith, C.D. (2003) Development of the Mars  
634 microbeam Raman spectrometer (MMRS). *Journal of Geophysical Research*, 108(E1),  
635 5005.
- 636 Wang, A., Zheng, M.P., Kong, F.J., Ling, Z.C., Kong, W.G., Sobron, P., and Jolliff, B.L.  
637 (2011) A Low T, High RH, and potentially life-friendly environment within the martian  
638 salt-rich subsurface in equatorial regions. 42nd Lunar and Planetary Science Conference, p.  
639 Abstract #2049. Lunar and Planetary Institute, Houston.

640 Wiens, R.C., Sharma, S.K., Clegg, S.M., Misra, A.K., and Lucey, P.G. (2007) Combined  
641 remote Raman spectroscopy and LIBS instrumentation for Mars astrobiology exploration.  
642 Seventh International Conference on Mars, p. Abstract #3092. Lunar and Planetary  
643 Institute, Houston.

644 Wold, S., Martens, H., and Wold, H. (1983) The multivariate calibration problem in  
645 chemistry solved by the PLS method. Lecture Notes in Mathematics, 973, 286-293.

646 Wynn-Williams, D.D., and Edwards, H.G.M. (2000) Proximal analysis of regolith habitats  
647 and protective biomolecules in situ by Laser Raman Spectroscopy: Overview of terrestrial  
648 antarctic habitats and Mars analogs. *Icarus*, 144(2), 486-503.

649

650

651

Hydration state	Mg	Ca	Fe	Na
Anhydrous	Anhydrous Mg-Sulfate	Anhydrite	--	Thenardite
1/2 H <sub>2</sub> O	--	Bassanite	--	--
1 H <sub>2</sub> O	Kieserite	--	Szomolnokite	--
2 H <sub>2</sub> O	Sanderite	Gypsum	--	--
3 H <sub>2</sub> O	Mg-sulfate tri-hydrate	--	--	--
4 H <sub>2</sub> O	Starkeyite	--	Rozenite	--
5 H <sub>2</sub> O	Pentahydrate	--	--	--
6 H <sub>2</sub> O	Hexahydrate	--	--	--
7 H <sub>2</sub> O	Epsomite	--	Melanterite	--
10 H <sub>2</sub> O	--	--	--	Glauber's salt
11 H <sub>2</sub> O	Meridianiite	--	--	--

652 Table 1. Sulfates used for the analysis with multivariate techniques

653

Response	Training set (pure sulfates)	Validation set (Mixed spectra)	Test set (Anhydrite + Thenardite)	Test set (Thenardite + MgSO <sub>4</sub> )
Hydration ratio	-4.5 <sup>-17</sup>	-0.0087	0.0814	0.0509
Ca- ratio	3.3 <sup>-17</sup>	0.0027	-0.0249	0.0046
Mg- ratio	-1.9 <sup>-17</sup>	-0.0061	-0.0288	-0.0625
Fe- ratio	-1.3 <sup>-17</sup>	0.0021	-0.0737	-0.0211
Na- ratio	-8.1 <sup>-18</sup>	0.0075	0.1123	0.0962
Average	-1.1 <sup>-17</sup>	-0.0004	0.0132	0.0136

654 Table 2. Average prediction error values with the 12 component PLS model

655

656

657

658

659

<b>Response</b>	<b>Training set (pure sulfates)</b>	<b>Validation set (Mixed spectra)</b>	<b>Test set (Anhydrite + Thenardite)</b>	<b>Test set (Thenardite + MgSO<sub>4</sub>)</b>
<b>Hydration ratio</b>	0.0788	0.0821	0.0901	0.0838
<b>Ca- ratio</b>	0.0072	0.0212	0.0327	0.0178
<b>Mg- ratio</b>	0.0109	0.0193	0.0356	0.0639
<b>Fe- ratio</b>	0.0134	0.0227	0.0897	0.0925
<b>Na- ratio</b>	0.0269	0.0236	0.1368	0.1280
<b>Average</b>	0.0274	0.0338	0.0770	0.0772

660 Table 3. RMSEP values (error variance) with the 12 component PLS model

661

662

663

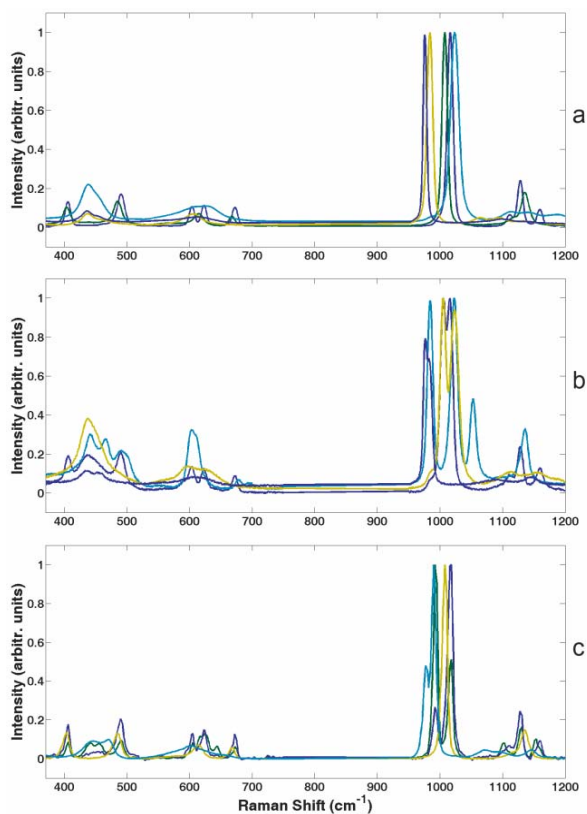
<b>Response</b>	<b>Training set (pure sulfates)</b>	<b>Validation set (Mixed spectra)</b>	<b>Test set (Anhydrite + Thenardite)</b>	<b>Test set (Thenardite + MgSO<sub>4</sub>)</b>
<b>Hydration ratio</b>	92.9%	89.5%	--	--
<b>Ca- ratio</b>	99.7%	97.9%	98.1%	--
<b>Mg- ratio</b>	98.7%	96.2%	--	98.4%
<b>Fe- ratio</b>	98.9%	97.1%	--	--
<b>Na- ratio</b>	93.9%	93.3%	98.9%	98.8%
<b>Average</b>	96.8%	94.8%	98.5%	98.6%

664 Table 4. Correlation values (in %) for the Calculated vs. Expected responses with the

665 12 component PLS model

666

667

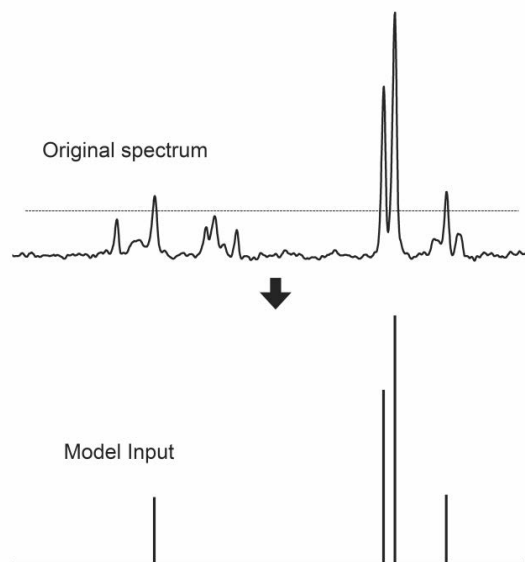


668

FIGURE 1

669

Figure 1. Example spectra from the training (a), validation (b) and test sets (c)



670

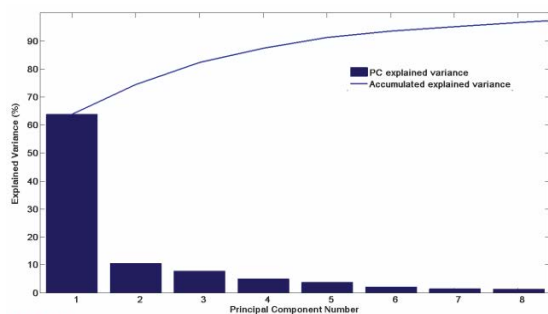
FIGURE 2

671

Figure 2. Spectrum pre-processing example for the ANN model

672

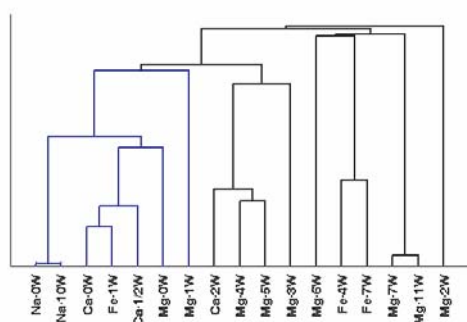




673

FIGURE 3

674 Figure 3. Explained variance of the data set (pure spectra of sulfates) after PCA  
675 analysis

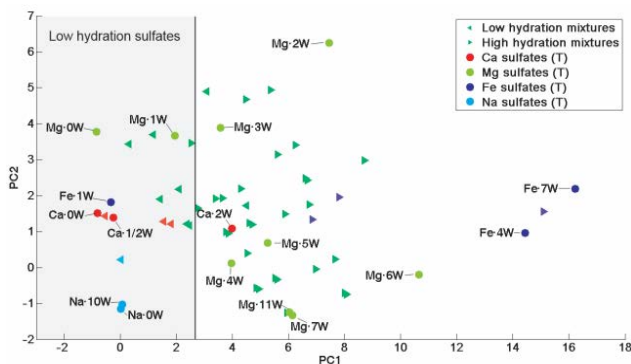


676

FIGURE 4

677 Figure 4. Dendrogram showing the classification of the PCA model with 3 PCs. As it  
678 can be seen, low-hydrated sulfates are contained in the blue branch of the dendrogram.

679

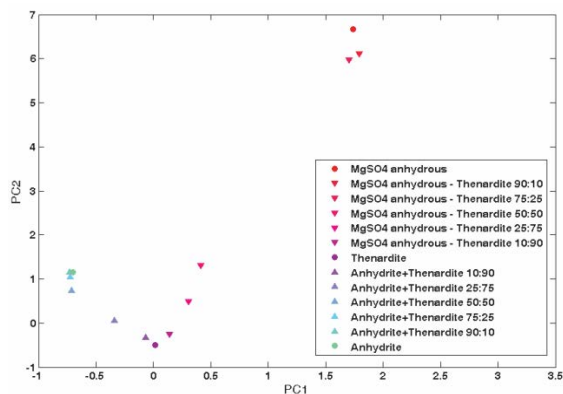


680

FIGURE 5

681 Figure 5. Representation of the scores for the training set (pure sulfates, circle) and  
682 50:50 mixtures of the sulfates of the same cation from the validation set (triangles).

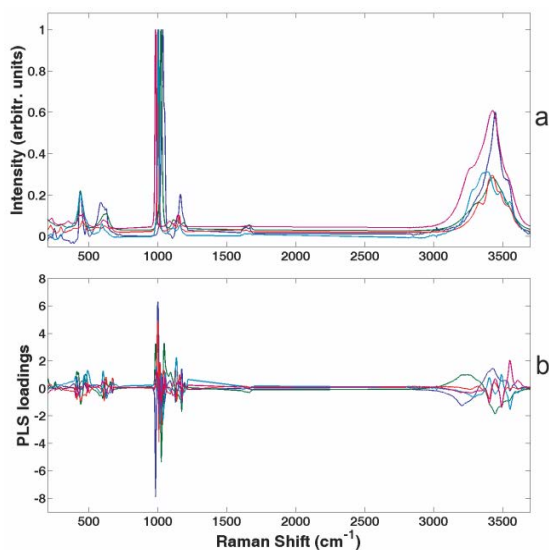
683 Different colors represent different cations. It can be seen how the low-hydration samples  
684 and mixtures are found with lower values of PC1.



685 **FIGURE 6**

686 Figure 6. Scores for the averaged spectra of the different mixtures acquired with the  
687 ExoMars Raman instrument. The scores are placed along lines between the endmembers at  
688 a distance dependent on the mixture concentration.

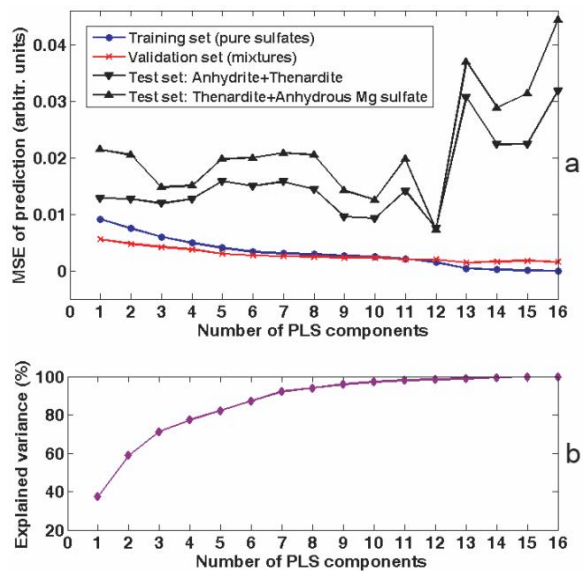
689



690 **FIGURE 7**

691 Figure 7. Some spectra from the training set (a) and corresponding PLS loadings (b)

692



693

FIGURE 8

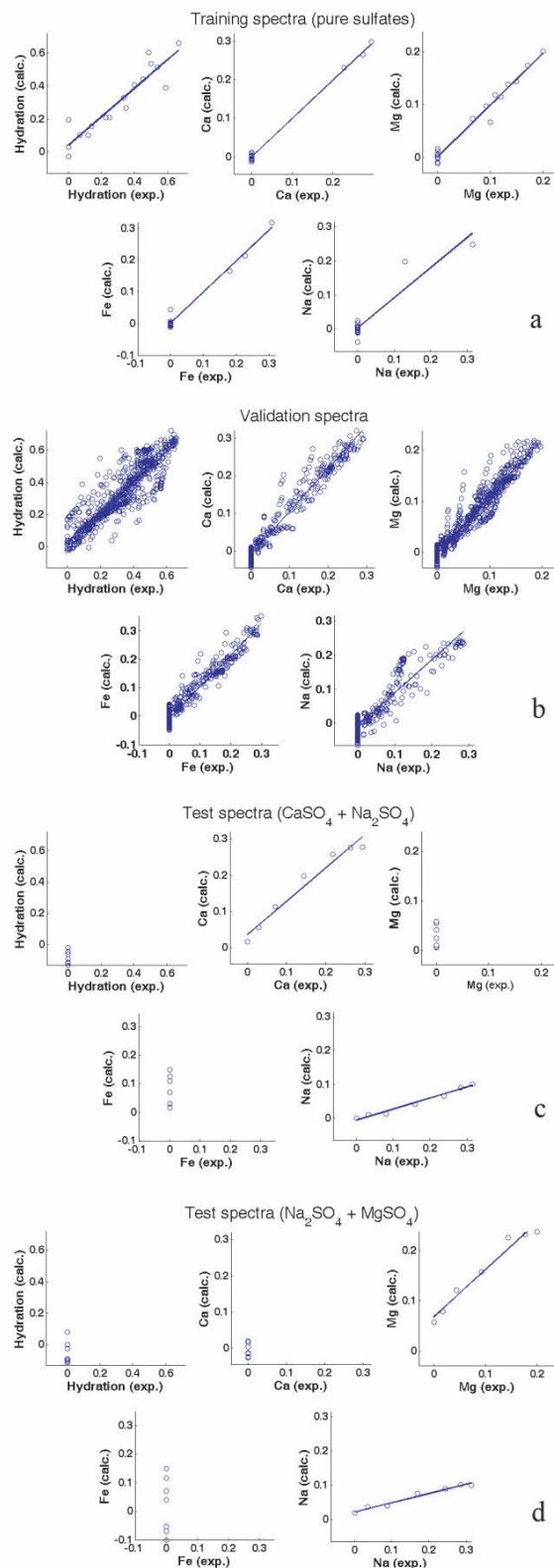
694

Figure 8. MSE of prediction (a), and accumulated explained variance (b) obtained by

695

the models with different numbers of components.

696

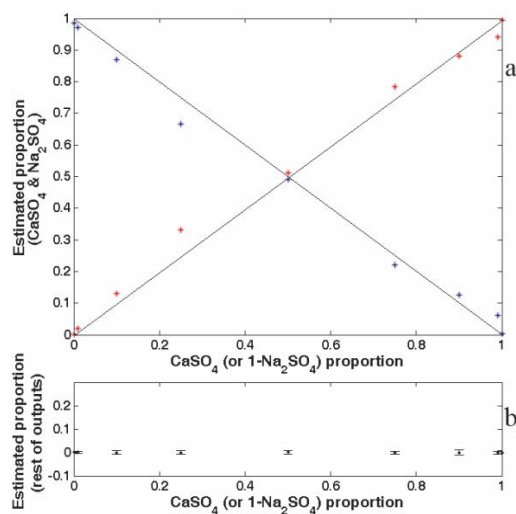


697

FIGURE 9

698 Figure 9. Calculated vs. Expected responses (cation weight ratio) for the training (a),  
699 validation (b) and test (c,d) spectra for the 12 component PLS model

700



701

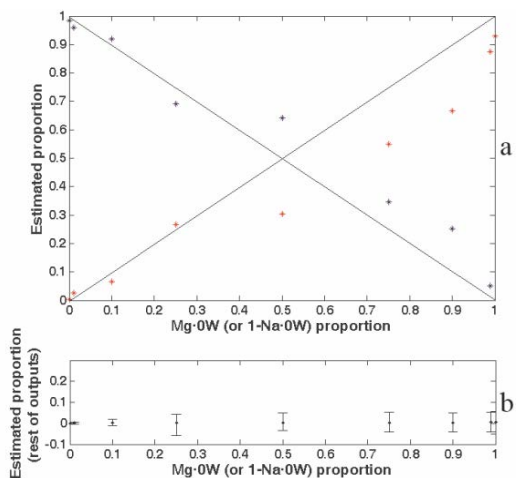
**FIGURE 10**

702 Figure 10. ANN outputs for the test set. Outputs corresponding to Anhydrite (CaSO<sub>4</sub>  
703 – red) and Thenardite (Na<sub>2</sub>SO<sub>4</sub> – blue) (a). Outputs corresponding to the rest of sulfates (b).

704 Vertical axis represents the proportion estimated by the ANN, while horizontal axis

705 represents the expected (known) proportion.

706



707

FIGURE 11

708

Figure 11. ANN outputs for the test set (mixtures in several proportions from RLS).

709

Outputs corresponding to Thenardite ( $\text{Na}_2\text{SO}_4$  – blue) and Mg-sulfate ( $\text{MgSO}_4$  – red) (a).

710

Outputs corresponding to the rest of sulfates (b). Vertical axis represents the proportion

711

estimated by the ANN, while horizontal axis represents the expected (known) proportion.