

Analyst

Accepted Manuscript



This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.

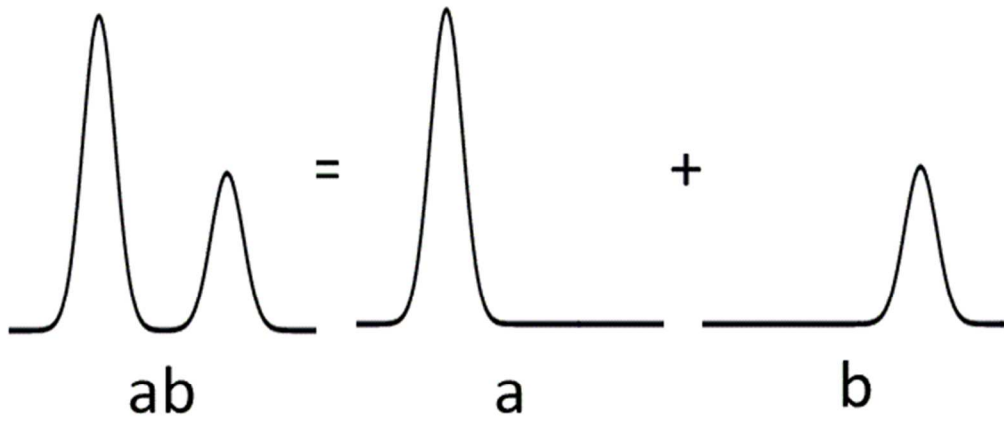


Table of Contents.

Identification by spectral analysis using Canberra distance as a novel metric for the comparison of spectra.

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

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

New algorithm for identification of components in a mixture: application to Raman spectra of solid amino acids

*Tomasz Roliński, * Sylwester Gawinkowski, Agnieszka Kamińska, and Jacek Waluk*

Institute of Physical Chemistry, Polish Academy of Sciences,

Kasprzaka 44/52, 01-224 Warsaw, Poland

*Corresponding author: E-mail: rolinski@ichf.edu.pl; Fax: +48 22 343 3333

ABSTRACT: The procedure of identifying components in a mixture was developed and tested on Raman spectra of mixtures of solid amino acids, using the spectra of single amino acids as templates. The method is based on finding the optimum scaling coefficients of the linear combination of template spectra that minimize the Canberra distance between measured and reconstructed spectra. The Canberra distance, used here as a measure of dissimilarity between spectra, defines non-convex objective function in the related optimization process. In view of the possibility of presence of local minima, *differential evolution*, which is a non-gradient stochastic method for finding global minimum, was chosen for optimization. The method was tested on twenty measured spectra of mixtures of solid powders containing one to eight amino acids taken

1
2
3 from the collection of twenty that are coded in living organisms. The results show that the
4
5 procedure can successfully identify several amino acids, and, in general, several components in a
6
7 mixture. The method was shown to compare favorably against the *least squares* and *partial least*
8
9 *squares* methods, the procedures used in commercially available chemometrics packages.
10
11

14 1. INTRODUCTION

16
17 Many papers have been devoted to identification of substances by their spectra and the specific
18
19 issues they address are diverse. For example, the problem may concern looking for a single
20
21 substance by comparing its measured spectrum with successive entries in a library of spectra.
22
23 Tanabe and Saeki¹ examined the possibility of identifying single substances by their IR spectra
24
25 and the Pearson correlation coefficient. Several factors were investigated influencing the
26
27 efficiency of the procedure, such as wavenumber range and the spacing between adjacent data
28
29 points, both related to the number of sampling points, as well as wavenumber accuracy and
30
31 sample purity. Another problem concerns the case of assigning reference spectra in a library as
32
33 components of a measured spectrum representing a mixture. Mallick et al.² compared several
34
35 methods of calculating coefficients of components of a mixture spectrum assumed to be a linear
36
37 combination of reference Raman spectra. They included all library reference spectra into the
38
39 combination, which implied solving one problem of high dimension. Thus the efficiency of a
40
41 numerical procedure was very important in this case. The methods were tested with simulated
42
43 measurements obtained from a statistical model with the most important error sources. The work
44
45 by Drake et al.³ dealt with the case of linear dependence of some reference spectra in the library,
46
47 which must take place in case the number of spectra exceeds the number of points in the
48
49 spectrum, and found that *non-negative least squares with active set* method described by Lawson
50
51 and Hanson⁴ was suitable for this task. Another possible scenario is looking for a particular
52
53
54
55
56
57
58
59
60

1
2
3 substance in a mixture containing excipients or contaminants. O'Connell et al.⁵ first pre-
4 processed a large number of spectra by normalization to the strongest peak and calculating the
5 first derivative. Both steps were justified by the Principal Component Analysis (PCA).^{6,7} Then
6
7
8 they used several classification methods to discriminate the target analyte. These included
9
10 Principal Component Regression (PCR),⁶ Support Vector Machine,⁸ K-Nearest Neighbors,⁹ a
11
12
13 decision tree,¹⁰ and others.
14
15
16

17
18 Beyond correlation coefficient there are several other similarity or dissimilarity measures for
19 matching spectra, such as Euclidean distance, city-block distance,¹¹ Tanimoto coefficient
20 (Jaccard index),¹² cosine of an angle between spectra. Li et al.¹³ dealt with the general analysis of
21 correlation coefficient, Euclidean cosine and their first-difference counterparts applied to
22 simulated spectra of one peak and ten peaks. The authors studied the influence of changing peak
23 width and peak position on the similarity (dissimilarity) indices. They recommended that such
24 indices should be used locally in predefined windows of significant intensities to increase the
25 reliability of the results. Varmuza et al.¹⁴ studied correspondence between spectral similarity,
26 measured by correlation coefficient, mean of the absolute and squared differences or Euclidean
27 cosine, and structural similarity measured by the Tanimoto index. The authors performed random
28 queries to a compound database, retrieving hit lists of compounds with similar IR spectra and
29 found the average for the Tanimoto coefficient between query and hit list compounds. The
30 method was used to characterize the performance of a spectral similarity search.
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

48 In this work we propose a nonstandard dissimilarity measure between spectra, the so-called
49 Canberra distance,¹⁵ which is the sum of relative errors in intensities for successive
50 wavenumbers. This index has been studied theoretically,¹⁶ and it is presently applied in genomics
51 as a measure of similarity.¹⁷ In this procedure, the mixture is assumed to be a linear combination
52
53
54
55
56
57
58
59
60

1
2
3 of reference spectra. Then one tries to minimize the objective function, defined by the Canberra
4 distance between the reconstructed and measured spectra, by varying the coefficients in the
5 combination. Since the objective function is not a convex function of its coefficients, the
6 uniqueness of the minimum is not guaranteed. Thus, some optimization procedure is required
7 that is capable of finding the global minimum in the presence of local minima. Non-gradient
8 stochastic optimization methods are suitable for this task. One method from this class,
9 *differential evolution*, has recently gained popularity.^{18,19} It is a representative of a wider class of
10 genetic algorithms that finds the global solution with high degree of probability, which proved to
11 be the case in our present analysis. An example of the application of a genetic algorithm was
12 presented by Forshed et al.²⁰ for peak alignment procedure for NMR metabonomic data. The
13 authors divided two spectra into common segments, and tried to shift sideways and stretch or
14 shrink one of them by linear interpolation to fit the other one. The optimum values for this
15 segment transformation are found by means of a genetic algorithm. The dissimilarity function
16 (Canberra distance) of our paper can be viewed as a weighted city-block distance (sum of
17 absolute differences), the weights being the inverses of the sums of absolute values of second
18 derivatives of compared spectra for successive wavenumbers. The idea of using weights to cope
19 with the problem when the range of values is wide was proposed by Liu et al.²¹ The authors used
20 weighted Pearson product-moment correlation coefficient to compare high-performance liquid
21 chromatograms and obtained better results than in the case of the non-weighted coefficient. It is
22 methodologically appropriate to precede any identification process by the correlation analysis
23 between template spectra. In our previous work²² we performed comparison of the same set of
24 spectra of twenty amino acids as in this paper using the intensities of strongest peaks and their
25 positions, as well as Pearson correlation coefficient as measures of similarity.
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

2. METHODOLOGY

The method was analyzed using as templates twenty measured Raman spectra of proteinogenic amino acids enumerated in Table 1. For a detailed description of amino acid samples, as well as the measurement conditions and results we refer the reader to the work by Roliński et al.²² The samples were purchased from Aldrich and Merck and used without any additional purification. Raman spectra of single amino acids and their mixtures were recorded with a Renishaw InVia Raman microscope using 632.8 nm line of the HeNe laser and ×20 objective. The laser power at the sample was 50 mW or less. The microscope was equipped with 1200 grooves/mm grating, cutoff optical filters, and 1024×256 pixels Peltier-cooled RenCam CCD detector, which allowed registering the Stokes part of Raman spectra with 5-6 cm⁻¹ spectral resolution and 2 cm⁻¹ wavenumber accuracy. To exclude the possibility of the orientational dependence of the signal on the polarization of the laser beam, the samples were finely pulverized and at least 100 spectra were recorded for each sample using automatic translation stage and then averaged.

The measured spectra of solid amino acids are shown in Figures 1 and 2. The high wavenumber range, 2500-3500 cm⁻¹, was not included in the identification process, providing a bigger challenge for the analytical algorithm, since a region was neglected for which one observes substantial differences between the spectra of mixture components.

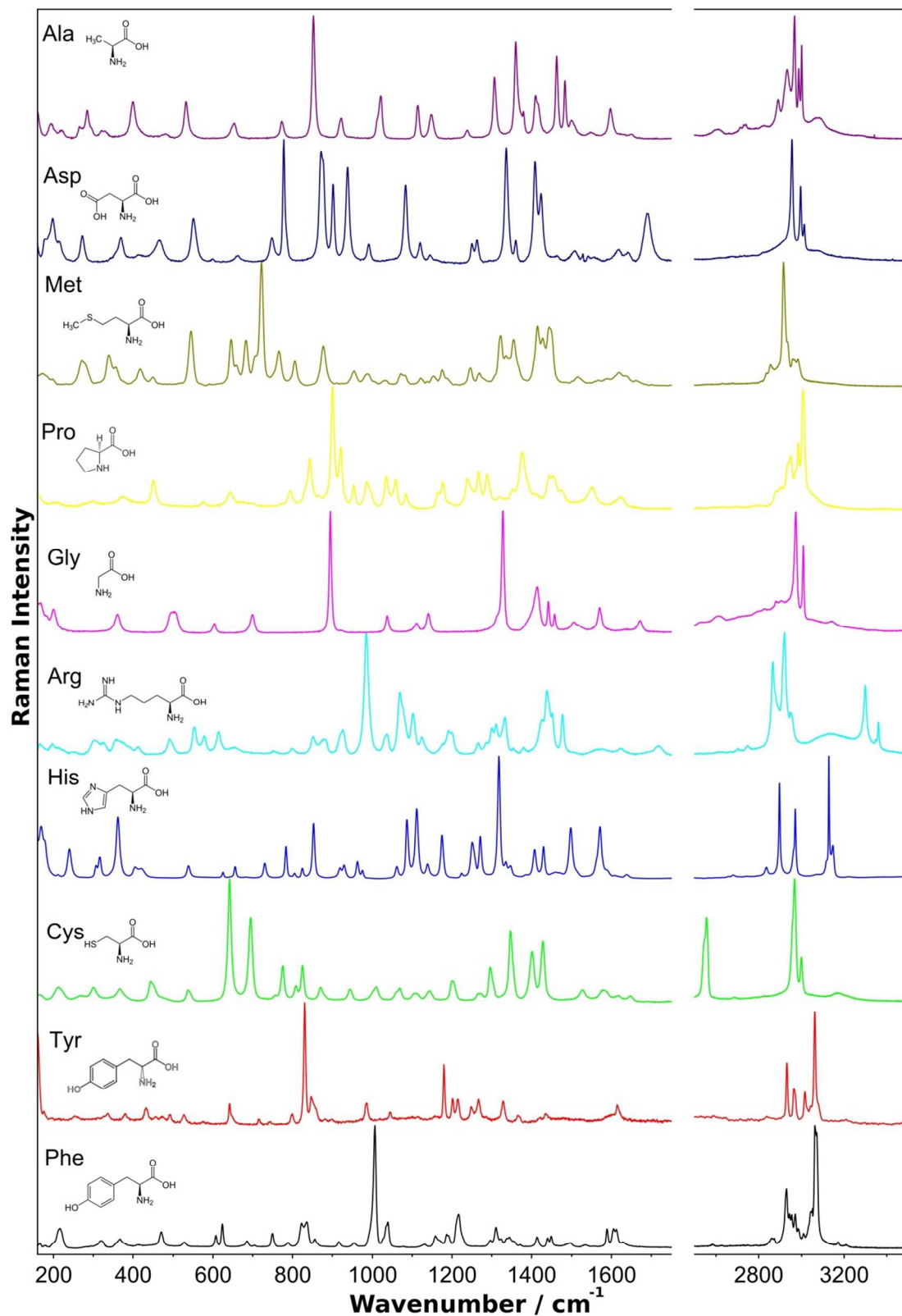


Figure 1. Experimental spectra of solid amino acids measured at 293 K using 632.8 nm laser line.

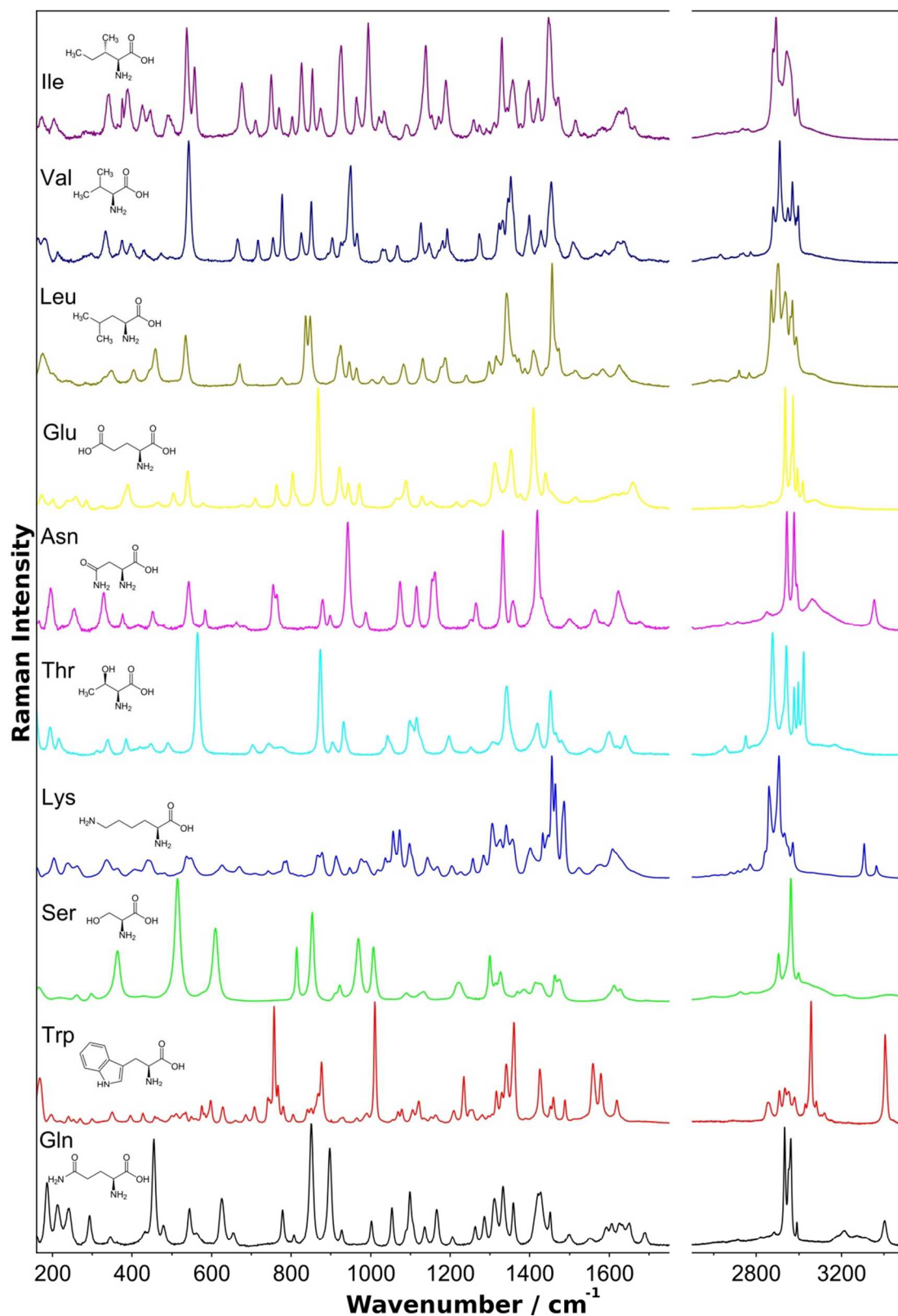


Figure 2. Experimental spectra of solid amino acids measured at 293 K using 632.8 nm laser line.

The spectra of mixtures and templates were first limited to the range 300-1700 cm^{-1} and then scaled so that the strongest peak within each spectrum equaled 100. This can be viewed as a simple preprocessing step.

Visual comparison of the spectra of mixtures with those of successive templates showed that the shifts of the corresponding peaks were very small, which greatly simplified the analysis since otherwise one would have to devise a measure of similarity that could compensate for this.²³ Since mixture intensities are sums of component intensities, one can expect good correspondence for intensities in some spectral regions only where a given component spectrum is dominant. If the fit were good for all considered energies (wavenumbers), this would mean that the mixture spectrum is trivial, i.e. containing one component. Moreover, even in the case of comparison of one component mixture and the corresponding template, the differences in intensities can be attributed to the measurement bias and the error in the definition of baselines for different spectra. It is also known that even under ideal measurement conditions the error in the value of intensity is theoretically proportional to the square root of the value. With the above in view we defined a dissimilarity function that is more responsive to the differences in position between peaks and their differences in widths than the differences in intensity:

$$\sum_{i=1}^N \frac{|R_L^{(2)}(w_i) - M_k^{(2)}(w_i)|}{|R_L^{(2)}(w_i)| + |M_k^{(2)}(w_i)|} \quad (1)$$

where:

$R_L^{(2)} = \sum_{l \in L} c_l a_l^{(2)}$, $L \subset \{1, 2, 3, \dots, M\}$ is the reconstructed second derivative of mixture spectrum;

$a_l^{(2)}$, $l = 1, \dots, M$ is the second derivative of template spectrum number l (see Table 1 and Figures 1 and 2);

c_l , $l = 1, \dots, M$ is the l -th scaling coefficient;

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

$M_k^{(2)}$, $k \in \{1,2,3, \dots, P\}$ is the second derivative of the measured mixture spectrum (see Table 4);

w_i , $i = 1,2,3, \dots, N$ is the wavenumber corresponding to the i -th point in the spectrum.

In our analysis $M = 20$, $N = 1401$,

$$w_i = 299 + i \text{ cm}^{-1}, \quad i = 1,2,3, \dots, 1401 \quad (2)$$

because we considered the spectral region of 300-1700 cm^{-1} , and the set of measured mixtures $\{1,2,3, \dots, P\}$ is replaced by $\{a, b, c, \dots, t\}$ (see Table 4). The function (1) is the sum of local distances for the wavenumbers w_i . The local distance is the absolute value of the difference between $M_k^{(2)}$ and $R_L^{(2)}$ for a given wavenumber, divided by the sum of their absolute values. The definition of (1) can be viewed as the Canberra distance¹⁵ between objects $R_L^{(2)}$ and $M_k^{(2)}$. In our case this function is minimized with respect to each c_l under restrictions $c_l \geq 0$, $l \in L \subset \{1,2,3, \dots, 20\}$.

Since there is no analytical formula for the spectra, the second derivative has to be evaluated numerically. To this end we first obtained a vector of intensities for the wavenumbers w_i from (2) by interpolating initial data, and then applied Savitzky-Golay filtering by trying to fit a second order polynomial locally to the data window of 21 points. The order of the polynomial and the width of the window were chosen by trial and error by comparing the original and the fitted data. An example of fitting is presented in Figures 3 and 4.

One should note first that if there is a difference in sign of the compared second derivatives of spectra for a given energy(wavenumber) then the corresponding term in (1) reaches its maximum value of one. On the other hand, if the signs are the same then this term approaches the maximum value only for a big difference between the second derivative spectra (see Figure 5). Obviously, if the values of the compared second derivatives are the same the term equals zero.

1
2
3 Second, if we scale both second derivative spectra by some constant then the value of (1) does
4 not change. This means that those parts of second derivative spectra that are small in value have
5 the same influence on the dissimilarity function as those which are large in value. Third, under
6 assumption that the background varies slowly with respect to the curve representing spectrum,
7 the use of second derivatives has additional benefit, because it minimizes the error arising from
8 the subtraction of the background, which is not defined precisely.
9

10 The expression (1) as a function of the coefficients c_l is non-convex, which means that it might
11 have local minima, contrary to the case of the *least squares* problem, where one has one global
12 minimum. In practice we find a global minimum for the function (1) by
13

- 14 • choosing a suitable minimization method, e.g. *differential evolution* stochastic
15 optimization method,^{18,19}
16
- 17 • trying to avoid solving problems for large L , i.e. for many templates in the linear
18 combination, by analyzing templates successively (one at a time), which will be
19 explained in detail later on,
20
- 21 • monitoring the results of successive optimizations.
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

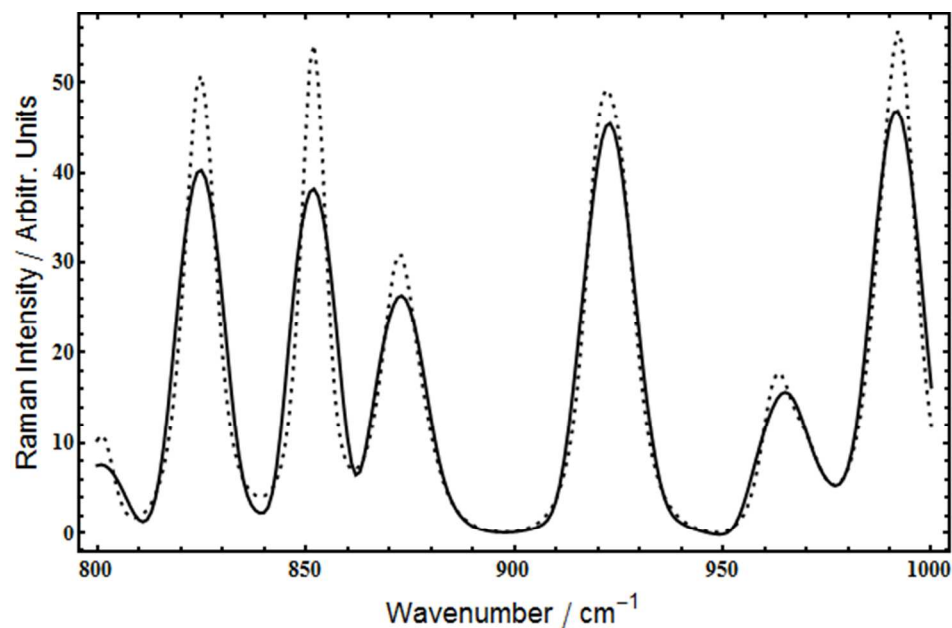


Figure 3. The result of Savitzky-Golay filtering for the interpolated data (second order polynomial, 21 points window). Dashed line, interpolated data; solid line, result of filtering. For clarity, the data was restricted to the range 800-1000 cm^{-1} .

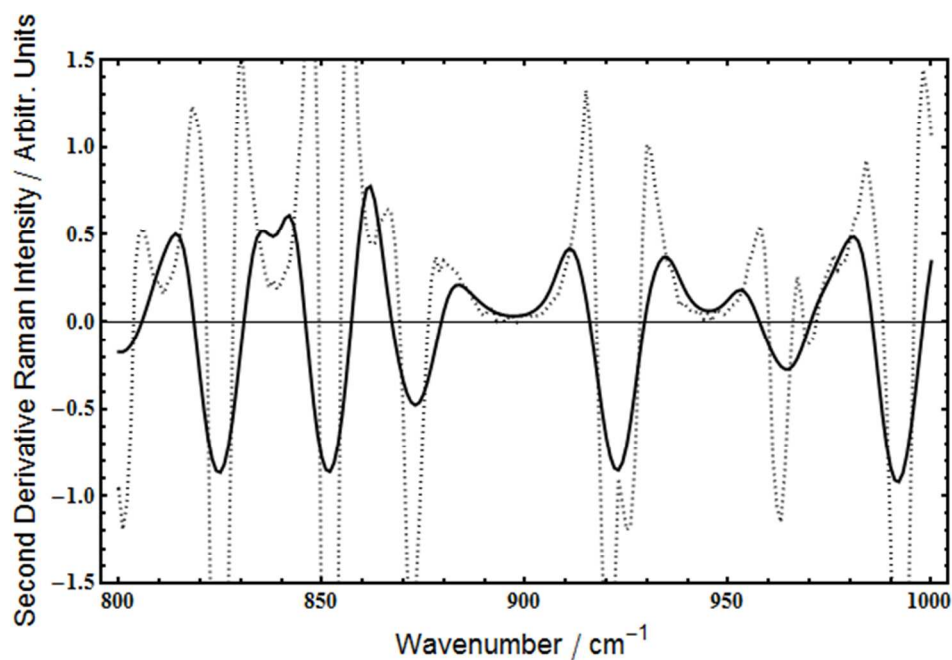


Figure 4. Comparison of numerical calculation of the second derivative of interpolated data. Dashed line, second derivative of data obtained by finite differences for adjacent points; solid

line, second derivative of data obtained by Savitzky-Golay filtering of the interpolated data (second order polynomial, 21 point window). For clarity, data was restricted to the range of 800-1000 cm^{-1} .

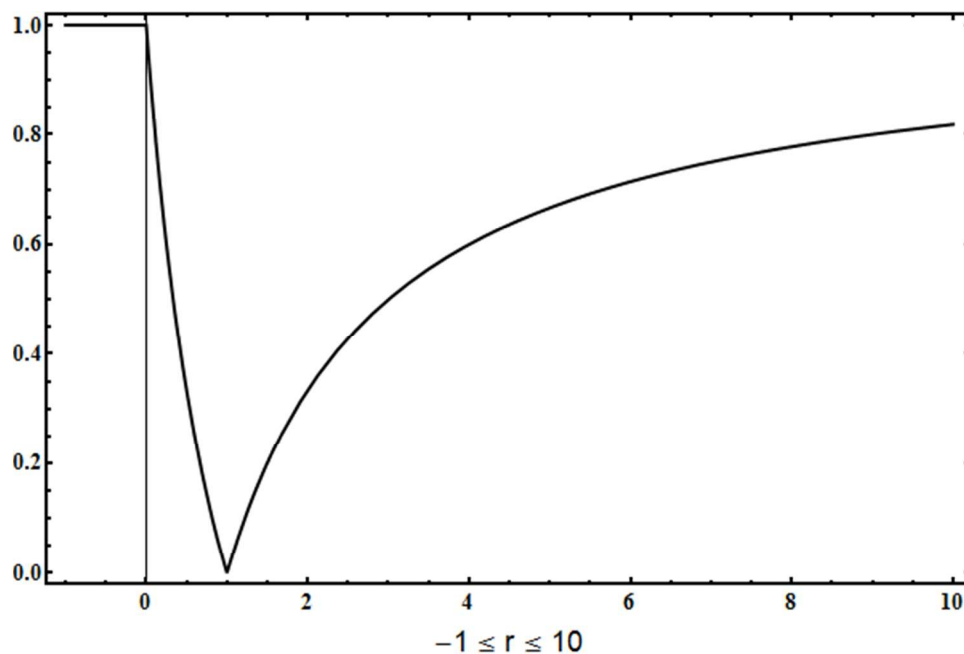


Figure 5. A plot of the function $f(r) = \frac{|r-1|}{|r|+1}$ for $-1 \leq r \leq 10$ showing the behaviour of local distance for a particular energy(wavenumber), i.e. one term of the dissimilarity function (1). We assume that the second derivative of intensity of the measured mixture spectrum equals one and the second derivative of intensity of the reconstructed spectrum equals r .

In what follows we shall be using acronyms for the corresponding amino acids taken into analysis. The correspondence is given in Table 1.

Now we describe the method in detail.

1st iteration. We try to find a single template matching the mixture spectrum the best, i.e. the one for which the optimum scaling coefficient gives the smallest value of function (1); see (1) for the case $L = \{m\}$, $1 \leq m \leq 20$ (see also Table 2 and Figures 6 and 7).

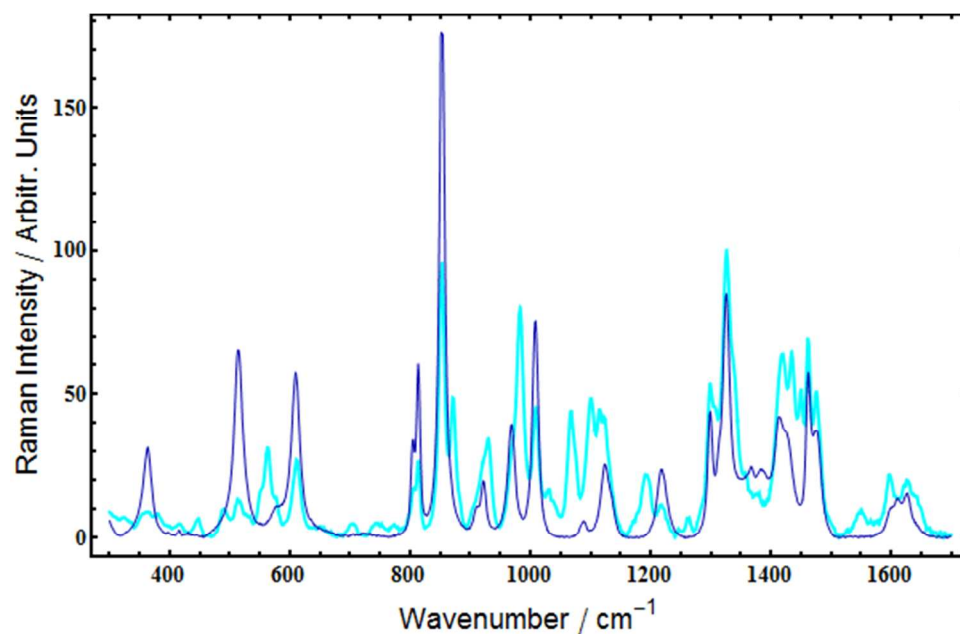


Figure 6. Spectra of mixture f (cyan thick line) and serine (dark blue thin line) multiplied by the corresponding optimum coefficient (see Tables 1, 2 and 4).

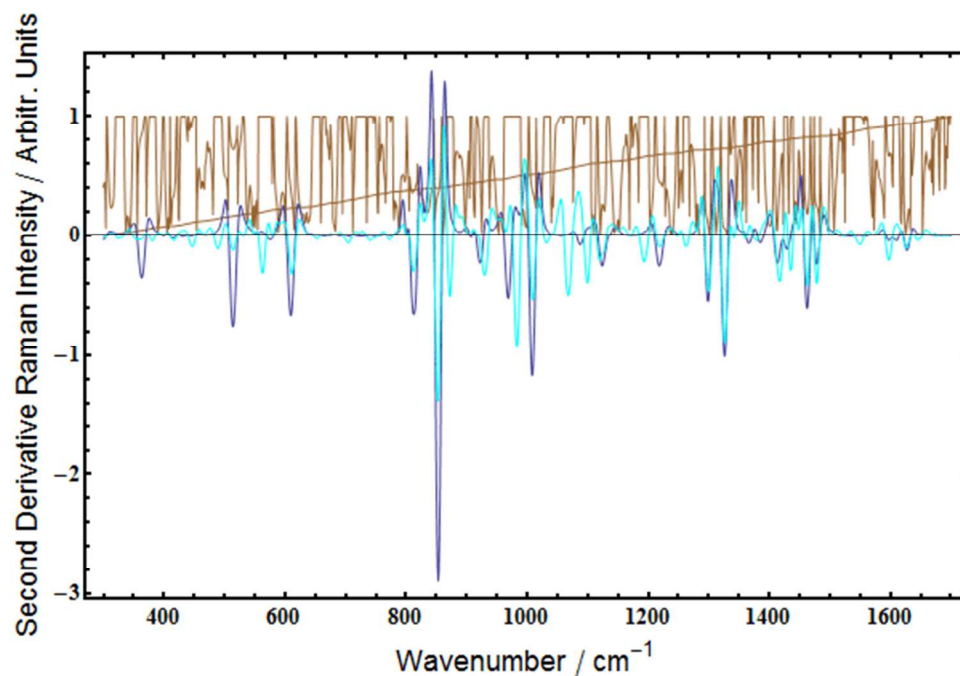


Figure 7. Second derivatives of the spectra from Figure 6. Local distances (components of (1)) are shown as brown line. The line increasing monotonically is the normalized cumulative value

1
2
3 of the local distances. Note that this line increases uniformly, which means that intervals of large
4 value of second derivative do not dominate the value of the global distance, i.e., the value of (1).
5
6

7
8 We claim that this contributes to the quality of the method.
9

10
11 **n-th iteration.** Let us assume we have already accepted some number of amino acid spectra
12 forming, with their coefficients, a linear combination. The combination defines the reconstructed
13 spectrum at the $(n-1)$ -th iteration. Now we want to check if the next template spectrum should be
14 included in it. We try to accept a single template from the collection of templates that have not
15 yet been included in the reconstructed spectrum. This template, together with the reconstructed
16 spectrum, forms a linear combination for the optimization process. It should be underlined that
17 now we try to find two scaling coefficients: one for the template and second for the whole
18 reconstructed spectrum from $(n-1)$ -th iteration. The chosen combination should yield the smallest
19 value for the dissimilarity function (1). If the stop condition is not fulfilled then we accept the
20 best new template and perform additional optimization for all already accepted templates, thus
21 defining the new reconstructed spectrum (see Table 3 and Figure 8). The stop condition means
22 that the absolute value of the difference in values of the function (1) for the best and the worst
23 match divided by the value for the worst match is less than some threshold value or the value of
24 the scaling coefficient corresponding to the best match falls below another threshold.
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

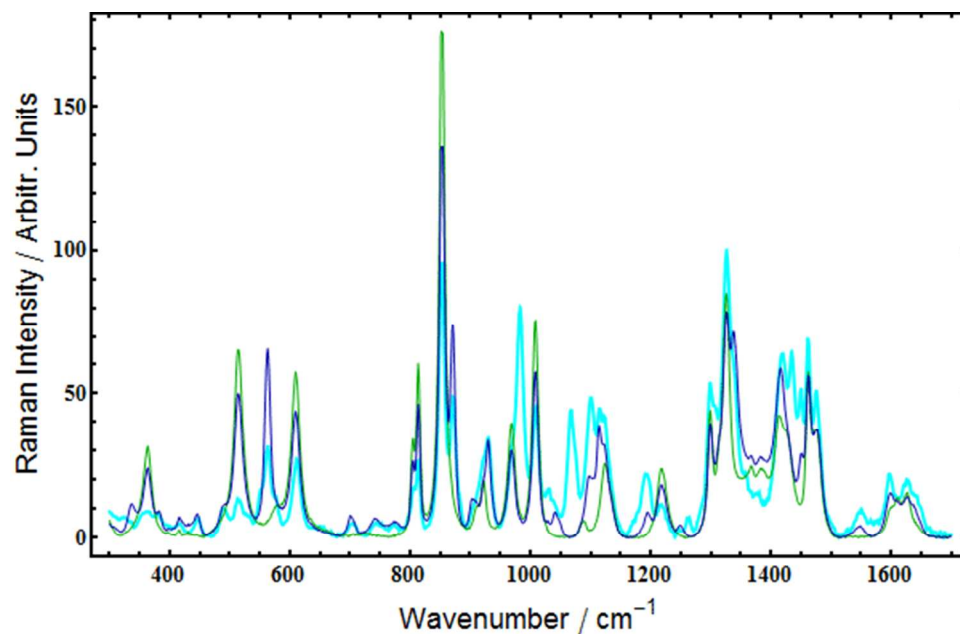


Figure 8. Measured Raman spectrum of mixture f (cyan thick line), spectrum of serine (dark green thin line) with optimum coefficient (see Table 2) and the linear combination of serine and threonine (dark blue thin line) with optimum coefficients (see Table 3). It can be seen that the combination fits the mixture spectrum better than single template of serine if we give more weight to the distribution of peaks than their intensities according to the construction of function (1) (see description below), see also Tables 1 and 4.

3. RESULTS

The method was verified using twenty measured spectra of mixtures of solid powders containing from one to eight amino acids (see Table 4) taken from the set of twenty presented in Table 1. The mixtures contained approximately equal volumes of components, which does not mean that their contributions to the measured spectrum were equal, as different substances yield weaker or stronger Raman signals depending on their polarizability. We tested the power of the method in qualitative analysis of samples, i.e. in identifying the components. The results are presented in Table 5. There may be two kinds of errors in the analysis: identification of

1
2
3 substance not present in the mixture (false positive) and failure of detecting a substance present
4
5 in the mixture (false negative). We assumed that both errors are equally serious and, accordingly,
6
7 tried to minimize their sum by adjusting the stop condition (see the description of the n -th
8
9 iteration of the method), which in the case of the analyzed mixtures (see Table 4) means that the
10
11 difference in values of function (1) for the best and the worst match divided by the value for the
12
13 worst match must be less than 2.4% or the value of the scaling coefficient of a potential
14
15 component corresponding to the best match must fall below 0.04. The number of
16
17 misclassifications for the optimum stop condition can serve as a measure of the quality of the
18
19 identification algorithm. There are no false positive and only two false negative cases: they
20
21 concern the mixtures with a high number of components (five and eight).
22
23
24
25
26
27
28

29 4. COMPARISON WITH THE *NON-NEGATIVE LEAST SQUARES* (NLS) METHOD

30
31 The objective function (1) is nonstandard for component identification in mixtures. The
32
33 standard one is the Euclidean distance corresponding to the *least squares* (LS) method, and
34
35 consequently LS can serve as a benchmark. The LS method has been used extensively in linear
36
37 mixture analysis (see the work by Heinz and Chang²⁴ and references therein). Two options are
38
39 possible, namely we can use zero or second derivative spectra. We performed calculations for the
40
41 second derivative spectra and obtained vastly different numbers of misclassifications for spectra
42
43 normalized as described in the Section 2 (normalization to the strongest peak both for templates
44
45 and mixture spectra) and for spectra normalized as described in Section 5. On the contrary, for
46
47 zero derivative spectra the numbers of misclassifications were quite close for both
48
49 normalizations and this is the reason why we present here the results for zero derivative spectra
50
51 with normalization of Section 2. Now the following function is minimized:
52
53
54
55
56
57
58
59
60

$$\sum_{i=1}^N |R_L(w_i) - M_k(w_i)|^2 \quad (3)$$

where:

R_L, M_k are defined analogously to $R_L^{(2)}, M_k^{(2)}$ in (1), except that we now consider zero derivative of spectra, and consequently the superscript ⁽²⁾ is omitted in the definition;

L, k, w_i are defined in (1) and (2).

This function is the sum of squared differences of intensities in the mixture spectrum and the linear combination of template spectra. The function is minimized with respect to all the c_l coefficients under restrictions $c_l \geq 0, l \in L \subset \{1,2,3, \dots, 20\}$.

We repeated the algorithm of Section 2 by calculating optimum coefficients of linear combinations for the objective function (3). In this case they can be found more efficiently by the Lawson and Hanson algorithm,⁴ but it was sufficient to use the *FindMinimum* procedure from the *Mathematica* package,²⁵ as time was not here a parameter for optimization. The results are presented in Table 6. The stop condition in this case means that the difference in values of function (2) for the best and the worst match divided by the value for the worst match must be less than 11% or the value of the scaling coefficient of a potential component corresponding to the best match must fall below 0.04. As for the Canberra distance case, the parameters in the stop condition were adjusted to obtain the least number of misclassifications. This number equals eleven, and it is much higher than in the case of function (1) (see Table 6).

5. COMPARISON WITH THE *PARTIAL LEAST SQUARES* (PLS) METHOD

Another standard procedure used widely in chemometrics for calibration is the *partial least squares* (PLS) method. It was introduced for the first time several decades ago in econometrics and then it gained popularity in chemistry for modelling the relationship between some

1
2
3 explanatory (easily obtainable) variables and the difficult or expensive to obtain response
4 variables.^{6,26} This is a method for solving the *least squares* (LS) problem approximately. The
5 matrix of intensities from LS is replaced by a matrix of much simpler structure, and usually of
6 lesser rank, that can be represented as a sum of some number of outer products of vectors of
7 scores and loadings. The number is equal to the rank of the matrix and is referred to as a number
8 of factors. The replacement is particularly useful if the columns of the intensity matrix are
9 strongly correlated, which means that the explanatory (independent) variables are correlated, as
10 well as in the case of noisy data.^{7,26} The colinearity of variables is unavoidable if the number of
11 explanatory variables is greater than the number of observations, which is usually the case if we
12 seek the signal contributions of substances using spectral intensities. It must be added that for the
13 PLS method the calculated approximate matrix is dependent on the values of the response
14 variable. In a simpler technique called principal component regression (PCR) the approximate
15 matrix is defined independently of the response variable. This technique is related to the singular
16 value decomposition procedure from the linear algebra and relies on choosing only the singular
17 vectors from this decomposition related to the largest singular values. It was shown that the PLS
18 technique leads often to a faster reduction of the residuals than PCR.⁷ The technique has been
19 implemented in many software chemometrics packages, e.g. *Grams*,²⁷ *Unscrambler*[®] *X*,²⁸ and this
20 section can be treated as a comparison of the identification capabilities of the commercially
21 available software with the identification power of the Canberra distance (1). Generally, there are
22 the training and the testing steps in the analysis. First, we treated the set of templates as the
23 training set and found the scores and loadings to model signal contributions. We applied the
24 PLS1 version of the algorithm, which means that we independently calibrated the signal
25 contribution of each amino acid. For a given amino acid the contribution related to its
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 corresponding spectrum was set equal to one, and for the rest of the templates the contributions
4 were set to zero. Second, we tested the model on twenty mixtures (see Table 4), predicting the
5 signal contributions with the calculated scores and loadings. The sequential character of the
6 algorithm was maintained by successive spectral subtractions of the identified templates
7 multiplied by the found contribution coefficients (scaling coefficients) from the analyzed mixture
8 spectrum. So first we subtracted the template with the highest calculated contribution in the
9 mixture, then repeated the procedure to find the next highest contribution for the difference
10 spectrum and the rest of potential component spectra. The procedure stopped when the calculated
11 contribution dropped below a preset threshold level. It should be underlined here that for the
12 procedure from Section 2 we did not perform the subtraction operation, but tried to find the best
13 fit for the combination of the reconstructed spectrum and a potential component spectrum.
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

29 Since we decided to model the real signal contribution ratios of constituents in the mixture, the
30 measurement and scaling of templates and mixtures were different than in Sections 2, 3 and 4.
31 The spectra of templates (components) were measured for the same time so that their intensity
32 ratios would reflect the contribution ratios in a mixture. Then we scaled both templates and
33 mixtures so that the average of the strongest peaks in all templates equaled 100 and for each
34 mixture spectrum the integral over the whole wavenumber range equaled the average of the
35 integrals for templates. The above can be viewed as a simple preprocessing step.
36
37
38
39
40
41
42
43
44
45

46 We performed some number of simulations by varying the number of factors in the PLS
47 method and analyzing centered (after the subtraction of the intensity average) or non-centered
48 (original, positive) intensity vectors (spectra). The PLS procedure is essentially quantitative and
49 its quality can be assessed by the *predictive residual sum of squares* (PRESS).⁶ Here we use PLS
50 for identification and therefore we must define the contribution threshold for confirmation of the
51
52
53
54
55
56
57
58
59
60

1
2
3 presence of a component in the mixture. If we find the optimum value for this threshold
4
5 corresponding to the smallest number of misclassifications, i.e. the sum of false positives and
6
7 false negatives for a given test set of mixtures, then this number can serve as a measure of
8
9 quality of identification corresponding to PRESS in the basic quantitative case.
10
11

12
13 The best results were obtained for the case of four factors and non-centered data together with
14
15 the optimum contribution threshold of 0.04, which yielded ten misclassifications: four false
16
17 negatives and six false positives (see Table 7). Since the mixtures were prepared by mixing
18
19 approximately equal volumes of powders the threshold of 0.04 seems to be rather small, which
20
21 means that if the volumes of components in mixtures were very small this method of
22
23 identification would probably fail. The number of four factors is relatively small if we compare it
24
25 with the number of explanatory variables (intensities for 1401 wavenumbers), which means that
26
27 many variables contain similar information on the signal contributions. We also obtained twelve
28
29 misclassifications for two factors only in the PLS method. Interestingly, the number of
30
31 misclassifications increased to sixteen for nine factors, which probably means that more detailed
32
33 data in spectra was treated as noise. Comparing the results presented in Table 7 with those in
34
35 Table 6 leads to the conclusion that PLS does not have more identification power than the non-
36
37 negative LS method.
38
39
40
41
42
43
44
45

46 6. CONCLUSIONS

47
48 A method for identifying components in a mixture was developed and tested on powder
49
50 mixtures of amino acids. The procedure is based on the linear model of mixture and involves
51
52 searching for scaling coefficients of the linear combination of template spectra minimizing a
53
54 function of dissimilarity referred to in the literature as Canberra distance.^{15,17} The Canberra
55
56 distance is related to a non-convex objective function in the optimization process and
57
58
59
60

1
2
3 consequently the method requires a stochastic optimization algorithm in view of the possibility
4 of existence of local minima. The method was tested using twenty measured spectra of mixtures.
5
6 Each mixture contained approximately equal volumes of powders of several amino acids taken
7
8 from the collection of twenty. The number of amino acids varied between one and eight. The
9
10 method does not attempt to find coefficients of the combination of all twenty template spectra
11
12 simultaneously, but accepts them into the reconstructed spectrum of the mixture successively,
13
14 starting from those most similar to the measured spectrum of mixture. Most components were
15
16 identified correctly: there were only two false negative cases for mixtures of five and eight
17
18 components and zero false positives (see Table 5). These results were achieved for the optimum
19
20 values of two threshold parameters (see Section 3) defined for all considered mixtures. The
21
22 results compare favorably with those obtained using the *non-negative least squares* (NLS)
23
24 method, which, for the two optimum parameters, gave eleven misclassifications (see Table 6),
25
26 and those provided by the *partial least squares* (PLS) method, which, for one optimum
27
28 parameter, gave ten misclassifications (see Table 7).
29
30
31
32
33
34
35

36 It should be mentioned that PLS is much faster than the remaining two methods, especially the
37
38 method based on the Canberra distance, which, however, is superior in identification power. The
39
40 speed of the PLS method is due to the two-step process mentioned in Section 5. The model
41
42 parameters calculated in the training step serve for prediction of amino acid signal contributions
43
44 in all analyzed mixtures, which includes few vector multiplications only, without solving any
45
46 equations. On the contrary, for the two previous methods we performed a series of optimizations
47
48 for each mixture to find scaling coefficients, though in the NLS method these optimizations were
49
50 relatively fast, because they involved minimizing simple quadratic functions. Moreover, one
51
52
53
54
55
56
57
58
59
60

1
2
3 should be aware that only one contribution threshold parameter is required for PLS, whereas for
4
5 the other two methods there are two of them, which may somewhat bias the results.
6
7

8 The total time for the analysis of all 20 mixtures and 20 templates for the Canberra distance
9
10 method amounted to 210 minutes (approximately 5-14 minutes for two component mixtures up
11
12 to 17-26 minutes for more than five components in a mixture), whereas for the PLS method it
13
14 took only 3 seconds if we do not count the training step. Of course the present method can be
15
16 made faster if we consider less data points (wavenumbers) in spectra. Moreover, the optimization
17
18 method (*differential evolution*) is time-consuming, as it requires many calls to the objective
19
20 function. Therefore, substituting it with a simple gradient method could greatly accelerate the
21
22 identification process, but at the expense of the possibility of falling into a local (not global)
23
24 minimum of the objective function and, consequently, increasing the number of
25
26 misclassifications. Both ways of acceleration, however, were not verified in practice. We think
27
28 that by the already obtained results we could combine the two methods together. First, we could
29
30 use the PLS method with parameters defined so as to make the false negative cases (almost)
31
32 absent, and then use the present method to additionally verify the already identified components
33
34 treated as a limited set of templates.
35
36
37
38
39
40

41 In practice we often face the situation where there are some unknown components in the
42
43 mixture spectrum, i.e. the components that cannot be spanned by the templates, and the
44
45 applicability of the method in such cases is important. Obviously, the stronger the unknown
46
47 components with respect to templates the less identification power of the method. The method
48
49 works sequentially, i.e. the first identified components are more dominant in the Raman signal, at
50
51 least with respect to the chosen Canberra metric. Since the method performed well in the case of
52
53 5-8 components in the mixture, this suggests that if the sought components are reasonably strong
54
55
56
57
58
59
60

1
2
3 they should be identified. However there is the problem of the threshold condition in this case,
4
5 i.e. how close the template spectrum should be to the mixture spectrum to be considered as its
6
7 part. In general, the problem can hardly be solved. In practice, one is interested in 1-3
8
9 components. They can be chosen from the set of templates by the devised algorithm, and then it
10
11 can be checked visually how many peaks in the compared spectra coincide, taking also into
12
13 account how strong they are in the templates.
14
15
16

17 ACKNOWLEDGEMENTS

18
19 The research was supported by the European Union within European Regional Development
20
21 Fund, through grant Innovative Economy (POIG.01.01.02-00-008/08).
22
23
24
25

26 REFERENCES

- 27
28
29
30 1 K. Tanabe, S. Saeki, *Anal. Chem.*, 1975, 47, 118.
31
32
33 2 M. Mallick, B. Drake, H. Park, A. D. Register, P. West, R. Palkki, A. Lanterman, D. Emge,
34
35 in *12th International Conference on Information Fusion*, Seattle, WA, USA, July 6-9, 2009.
36
37
38 3 B. Drake J. Kim, M. Mallick, H. Park, in *Proceedings of the Thirteenth International*
39
40 *Conference on Information Fusion*, Edinburgh, UK, 2010.
41
42
43 4 C. L. Lawson, R. J. Hanson, *Solving Least Squares Problems*, Prentice-Hall, 1974.
44
45
46 5 M. L. O'Connell, T. Howley, A. G. Ryder, M. N. Leger, M. G. Madden, in *Proceeding of:*
47
48 *Opto-Ireland 2005: Optical Sensing and Spectroscopy*, Ireland, 2005, p 340.
49
50
51
52 6 P. Geladi, B. R. Kowalski, *Anal. Chim. Acta* 1986, 185, 1.
53
54
55 7 L. Elden, *Comput. Stat. Data An.*, 2004, 46, 11.
56
57
58
59
60

1
2
3 8 N. Cristianini, J. Shawe-Taylor, *An Introduction to Support Vector Machines and Other*
4
5 *Kernel-based Learning Methods*, Cambridge University Press, 2000.
6
7

8
9 9 B. S. Everitt, S. Landau, M. Leese, D. Stahl, *Miscellaneous Clustering Methods*, in *Cluster*
10
11 *Analysis*, 5th edn., John Wiley & Sons, Ltd., Chichester, UK, 2011.
12
13

14 10 L. Breiman, J. Friedman, R. Olshen, C. Stone, *Classification and Regression Trees*,
15
16 Wadsworth, Belmont, CA, 1984.
17
18

19
20 11 M. Barile, *Taxicab Metric*, From MathWorld--A Wolfram Web Resource, created by Eric
21
22 W. Weisstein, <http://mathworld.wolfram.com/TaxicabMetric.html> .
23
24

25 12 A. H. Lipkus, *J. Math. Chem.*, 1999, 26, 263.
26
27

28 13 J. Li, D. B. Hibbert, S. Fuller, G. Vaughn, *Chemometr. Intell. Lab.*, 2006, 82(1-2), 50.
29
30

31 14 K. Varmuza, M. Karlovits, W. Demuth, *Anal. Chim. Acta*, 2003, 490, 313.
32
33

34 15 G. N. Lance, W. T. Williams, *Comput. J.*, 1966, 9, 60.
35
36

37
38 16 G. Jurman, S. Riccadonna, R. Visintainer, C. Furlanello, in *Advances in Ranking, NIPS 09*
39
40 *Workshop*, ed. S. Agrawal, C. Burges, K. Crammer, 2009, p 22.
41
42

43 17 J. Wu, M. Gan, W. Zhang, R. Jiang, *Int. J. Biosci. Biochem. Bioinform.*, 2011, 1, 102.
44
45

46 18 R. Storn, K. J. Price, *Global Optim.*, 1997, 11, 341.
47
48

49
50 19 V. Plagianakos, E. W. Weisstein, *Differential Evolution*, From MathWorld—A Wolfram
51
52 Web Resource, <http://mathworld.wolfram.com/differentialevolution.html>.
53
54

55
56 20 J. Forshed, I. Schuppe-Koistinen, S. P. Jacobsson, *Anal. Chim. Acta*, 2003, 487(2), 189.
57
58
59
60

1
2
3 21 Y. Liu, Q. Meng, R. Chen, J. Wang, S. Jiang, Y. Hu, *J. Chromatogr. Sci.*, 2004, 42(10),
4
5 545.
6

7
8
9 22 T. Roliński, S. Gawinkowski, A. Kamińska, J. Waluk, in *Optical Spectroscopy and*
10
11 *Computational Methods in Biology and Medicine*, ed. M. Baranska, Springer, 2014, 14, 329.
12

13
14 23 R. de Gelder, R. Wehrens, J. A. Hageman, *J. Comput. Chem.*, 2001, 22(3), 273.
15

16
17 24 D. C. Heinz, C-I. Chang, *IEEE T. Geosci. Remote*, 2001, 39(3), 529.
18

19
20 25 Wolfram Mathematica 8; <http://www.wolfram.com/> .
21
22

23
24 26 S. Wold, M. Sjostrom, L. Eriksson, *Chemometr. Intell. Lab.*, 2001, 58, 109.
25

26
27 27 GRAMS Suite 9.1, <http://gramssuite.com/> .
28
29

30
31 28 The Unscrambler[®]X 10.3, <http://www.camo.com/> .
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 1. Collection of twenty amino acids taken into spectral analysis with their acronyms.

	Amino acid	Acronym
1	Arginine	Arg
2	Proline	Pro
3	Alanine	Ala
4	Phenylalanine	Phe
5	Cysteine	Cys
6	Asparagine	Asn
7	Glutamine	Gln
8	Leucine	Leu
9	Threonine	Thr
10	Valine	Val
11	Isoleucine	Ile
12	Glutamic acid	Glu
13	Glycine	Gly
14	Aspartic acid	Asp
15	Lysine	Lys
16	Methionine	Met
17	Histidine	His
18	Serine	Ser
19	Tryptophan	Trp
20	Tyrosine	Tyr

Table 2. The first iteration of the method for mixture f (see Table 4). The results are sorted according to the values in the second column. First column, acronym of the amino acid (see Table 1); second column, values of the function (1) for optimum coefficients; third column, difference between the current value and the value of function (1) for the worst match expressed as percent of the value for the worst match; fourth column, the optimum scaling coefficient for the corresponding amino acid. In the present case the best match corresponds to serine.

Amino acid	Values of function (1)	Difference in the values of (1) [%]	Optimum coefficient
Ser	865.	-22.0	1.76
Thr	872.	-21.0	1.14
Arg	910.	-18.0	1.95
Gln	999.	-10.0	0.76
Pro	1009.	-9.1	1.77
Ala	1010.	-9.1	0.99
Ile	1016.	-8.5	0.58
Lys	1029.	-7.3	1.23
Asn	1031.	-7.2	0.73
Cys	1033.	-7.0	2.27
Trp	1043.	-6.0	0.93
Tyr	1046.	-5.8	2.96
Asp	1049.	-5.5	0.93
Met	1051.	-5.4	1.78
Phe	1054.	-5.1	2.30
Leu	1070.	-3.7	0.67

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

Val	1070.	-3.6	0.67
Glu	1077.	-3.0	0.99
Gly	1097.	-1.2	2.54
His	1110.	0.0	1.84

Table 3. The second iteration of the method for mixture f (see Table 4). The results are sorted according to the second column. First column, acronym of an amino acid (see Table 1); second column, value of function (1) for optimum coefficients; third column, difference between the current value and the value of function (1) for the worst match as percent of the value for the worst match; fourth column, optimum scaling coefficient for the spectrum reconstructed in the first iteration (serine); fifth column, optimum scaling coefficient for the corresponding amino acid spectrum. In this case the best match is the combination of spectra of serine and threonine (see Tables 1 and 2).

Amino acid	Value of function (1)	Difference in the values of (1) [%]	Optimum coefficient for reconstructed spectrum	Optimum coefficient for corresponding amino acid
Thr	653.	-25.0	0.76	0.63
Arg	774.	-11.0	0.51	1.00
Ala	810.	-6.3	0.84	0.36
Gln	811.	-6.2	0.85	0.14
Asn	817.	-5.5	0.69	0.27
Ile	822.	-4.9	0.73	0.16
Asp	825.	-4.6	0.70	0.27
Gly	831.	-3.9	0.75	0.51
Cys	841.	-2.8	0.92	0.14
Pro	848.	-1.9	0.93	0.37
Glu	851.	-1.6	0.76	0.33
Phe	852.	-1.4	0.97	0.08
Trp	854.	-1.3	0.97	0.10
Leu	857.	-0.9	0.93	0.10

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

His	859.	-0.7	1.00	0.12
Val	861.	-0.5	0.87	0.04
Tyr	862.	-0.4	0.83	0.65
Met	863.	-0.2	1.06	0.00
Lys	865.	0.0	1.00	0.00

Table 4. Sample identification letters and qualitative composition of measured mixtures.

Mixture	Composition
a	Phe Ala
b	His Arg Pro
c	Tyr Asn
d	Gly Thr Ser Gln Leu
e	Trp Glu Gln Ile
f	Ser Thr Arg Ala
g	Met Ala His Gly Leu
h	Glu Leu Ile
i	His Gly Leu
j	Glu Met Lys
k	His Gly Tyr Pro
l	Thr Ala
m	Asn
n	Met Val Leu Ile
o	Phe Gln
p	Tyr Thr Pro Asn Asp
q	Gly Glu Ala
r	Tyr Trp His Arg Ser
s	Trp Asp
t	Tyr Cys Phe His Ser Leu Thr Ile

Table 5. Identification results of measured mixture spectra (denoted by letters; see Table 4). All the amino acids shown in the second column were identified correctly using adjusted mechanical stop condition of the iterative algorithm for all mixtures. The third column shows two false negative cases involving serine and isoleucine from the mixtures r and t of five and eight components, correspondingly. The last column shows zero false positive cases.

Mixture	Identified correctly	False negatives	False positives
a	Phe Ala		
b	His Arg Pro		
c	Tyr Asn		
d	Gly Thr Ser Gln Leu		
e	Trp Glu Gln Ile		
f	Ser Thr Arg Ala		
g	Met Ala His Gly Leu		
h	Glu Leu Ile		
i	His Gly Leu		
j	Glu Met Lys		
k	His Gly Tyr Pro		
l	Thr Ala		
m	Asn		
n	Met Val Leu Ile		
o	Phe Gln		
p	Tyr Thr Pro Asn Asp		
q	Gly Glu Ala		
r	Tyr Trp His Arg	Ser	
s	Trp Asp		

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

Table 6. Identification results of measured mixture spectra (see Table 4) in the case of benchmark objective function (3). All amino acids shown in the second column were identified correctly using mechanical stop condition from the iterative algorithm (see Section 4); the third column shows false negative cases, fourth column shows false positive cases.

Mixture	Identified correctly	False negatives	False positives
a	Phe Ala		
b	His Arg	Pro	
c	Tyr Asn		
d	Gly Thr Ser Gln Leu		
e	Trp Glu Gln	Ile	
f	Ser Thr Arg	Ala	
g	Met Ala His Gly Leu		
h	Glu Leu Ile		
i	His Gly Leu		
j	Glu Met Lys		
k	His Gly Tyr Pro		
l	Thr Ala		
m	Asn		
n	Met Val Leu Ile		
o	Phe	Gln	Trp
p	Tyr Thr Pro Asn	Asp	
q	Gly Glu	Ala	Lys
r	Tyr Trp His Arg	Ser	
s	Trp Asp		
t	Tyr Cys Phe His Ser Ile	Leu Thr	

Table 7. Identification results of measured mixture spectra denoted by letters (see Table 4) in the case of PLS method. All amino acids shown in the second column were identified correctly using the optimum mechanical stop condition from the iterative algorithm (see Section 5); the third column shows false negative cases, the fourth column shows false positive cases.

Mixture	Identified correctly	False negatives	False positives
a	Phe Ala		
b	His Arg Pro		
c	Tyr Asn		
d	Gly Thr Ser Gln Leu		
e	Trp Glu Gln Ile		
f	Ser Thr Arg	Ala	
g	Met Ala His Gly Leu		
h	Glu Leu Ile		
i	His Gly Leu		
j	Glu Met Lys		
k	His Gly Tyr Pro		
l	Thr Ala		
m	Asn		Gly
n	Met Val Leu Ile		
o	Phe Gln		Ser Trp
p	Tyr Thr Pro Asn Asp		
q	Gly Glu Ala		His Lys
r	Tyr Trp His Arg	Ser	
s	Trp Asp		Ile
t	Tyr Cys His Ser Leu Ile	Phe Thr	