

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.

1
2
3 **Assessment of the statistical significance of classifications in infrared spectroscopy based**
4 **diagnostic models**

5
6 David Pérez-Guaita¹, Julia Kuligowski², Salvador Garrigues³, Guillermo Quintás⁴ and Bayden.R.
7 Wood^{1*}
8

9
10 ¹Centre for Biospectroscopy, School of Chemistry, Monash University, Clayton, 3800 VIC,
11 Australia.
12

13
14 ² Neonatal Research Group, Health Research Institute La Fe, 46026 Valencia, Spain
15

16
17 ³Analytical Chemistry Department, University of Valencia, Edifici Jeroni Muñoz, 46100 Burjassot,
18 Spain ⁴Leitat Technological Center, Bio In Vitro Division, 46026 Valencia, Spain
19

20
21 *To whom correspondence should be addressed
22

23 Fourier transform infrared (IR) spectroscopy in combination with multivariate data analysis is a
24 versatile tool that can be applied for the diagnosis. However, a rigorous validation of the obtained
25 models is necessary in order to obtain robust results. This work evaluates the advantages of the use
26 of permutation testing for determining the statistical significance of the misclassification errors
27 obtained from IR based diagnostic models through cross validation (CV). The model performance,
28 estimated by CV, is compared to a distribution of CV-performance values obtained using randomly
29 permuted class labels. The distribution of ‘random CV-values’ is considered as a null distribution
30 and used to establish the significance of the model estimators obtained using real class labels. ATR-
31 FTIR spectra of serum samples were classified using random forest (RF) classifiers according to
32 two criteria, the tag number (a randomly assigned pseudo class membership) and the level of urea
33 (real class). CV errors obtained were compared to the null distribution of CV errors from a
34 permutation test and an independent validation set. The procedure was evaluated testing typical
35 conditions leading to overoptimistic estimations provided by the CV like e.g. the size of subsamples
36 used during CV, variable selection and the use of replicates. Results show that for the tag number
37 (pseudo class), CV indicated classification errors between 23 and 33 % depending on the subsample
38 size employed. Those values were even lower when variable selection or replicates were used.
39 However, permutation testing indicated that those CV errors were non-significant. In contrast, for
40 sample classification according to their levels of urea, all cross validation errors were found to be
41 significant. Although the proposed method is computationally intensive, it provides a simple way of
42 calculating an empirical p-value of the CV-estimator, thus establishing the statistical significance
43 and providing a feasibility indicator especially useful for studies where the number of samples is
44 limited.
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1. Introduction

In the last decade, experimental studies revealed the potentials of vibrational spectroscopy for the diagnosis of disease in a wide variety of biological fluids and tissues¹. Specifically, the combination of multivariate analysis and Raman or infrared (IR) spectroscopy has been repeatedly proposed for medical diagnosis based on the high amount of information encrypted in IR and Raman spectra, which provides a bimolecular fingerprint of the metabolic state of the patient². A literature review between 2012 to 2014 reports its application for the diagnosis of several diseases such as asthma³ or diverse types of cancer including glioma⁴, breast⁵, ovarian^{6,7}, oral⁸ or bladder cancer⁹. In addition, in the case of blood and related samples such as serum and plasma, attenuated total reflectance (ATR) IR measurements offer a minimum-invasive, fast and cost-effective diagnostic tool¹⁰. The access to well characterized samples is fundamental for exploiting the potentials of vibrational spectroscopy. Selection of the sample size required to derive a statistically significant estimation is critical to produce robust results and it depends on a number of factors including the size of the effect and intrinsic biological and instrumental variation. However, difficulties in obtaining well characterized samples (e.g. blood, tissues) may hinder the development of this type of technology and, frequently, reported studies include a limited number of samples. In these situations, external independent test sets for the evaluation of the prediction capabilities of the models are not always available and, alternatively, a number of cross validation (CV) procedures are employed. In spite of its usefulness, CV-based approaches might lead to an underestimation of the classification error¹¹ leading to overly optimistic descriptions of the model discriminant performance. In the fields of metabolomics and proteomics it is not unusual to generate mega-variate datasets with a high variable to sample ratio¹². Because of the need to establish reliable indicators of the performance of the classifications when biomarker identification is involved, the aforementioned problems related with CV have been already addressed. One approach is based on non-parametric permutation tests, which imply the random reallocation of class labels in order to establish the statistical significance

1
2 of a figure of merit obtained using CV^{13,14}. In spite of being a quite common approach in other
3
4 fields such as proteomics¹³ or metabolomics¹², this strategy is not yet common practice in IR
5
6 biospectroscopy and besides a conference proceedings of Lloyd et al¹⁵, it is difficult to find any
7
8 application of this procedure in the biospectroscopy field .
9

10
11 The aim of this paper is to introduce permutation testing as a feasible indicator of the statistical
12
13 significance of CV-figures of merit of IR based diagnostic models. For this purpose, ATR-Fourier
14
15 transform infrared (FTIR) spectra of serum samples were recorded and discriminant models by
16
17 means of random forest (RF) classifiers were developed. Reproducing preliminary studies, a limited
18
19 number of samples were used for calculating and cross validating the RF classifiers. Results of a
20
21 permutation test performed over the CV error were compared with the actual prediction capability
22
23 of those models, established by validating them with a large and representative set of samples.
24
25 Serum samples were classified in two groups according to: i) the urea concentrations in sample, and
26
27 ii) to an arbitrary sample number un-related to the sera composition. The importance of the access
28
29 to statistically reliable performance parameters was further underlined in common situations of IR
30
31 based diagnostic models, which are known to enhance overoptimistic CV errors, i.e. the use of
32
33 measurement replicates and variable selection.
34
35
36
37
38

39 **2. Experimental**

40 **2.1 Samples and Spectral Acquisition**

41
42 ATR-FTIR Spectra of 122 serum samples from different patients were acquired as described
43
44 elsewhere¹⁶. Briefly, spectra containing 235 data points recorded in the range 1800-850 cm⁻¹ with a
45
46 spectral resolution of 4 cm⁻¹ were recorded using a Tensor 27 spectrophotometer from Bruker
47
48 (Bremen, Germany), equipped with temperature-stabilized deuterated lanthanum tryglycine sulfate
49
50 (DLATGS) detector and a nine reflection diamond/ZnSe DuraDisk from Smiths Detection Inc.
51
52 (Warrington, UK). 150 microliters of serum were deposited on the ATR cell and the spectrum was
53
54
55
56
57
58
59
60

1
2 measured using the empty clean cell as a background. Both, spectrum and background were
3
4 obtained averaging 100 scans, and the water contribution to the spectrum was eliminated by
5
6 subtracting a spectrum of water measured in the same conditions.
7
8

9
10 Measurements for each sample were recorded in triplicate and the average of the replicates was
11
12 used in sections 3.1 and 3.2 and the actual individual replicates were used in section 3.3.
13

14 15 **2.2 Samples Sets and Class Assignment**

16
17 As we have stressed in the introduction, in most of the cases the number of samples are limited on
18
19 the studies for the diagnosis of illnesses based on IR spectroscopy. Since our investigation aimed to
20
21 simulate those preliminary studies, only 30 samples were employed for calibration. Besides,
22
23 reproducing the conditions of a representative blind test, 92 samples were used only as an
24
25 independent test set in order to establish the actual prediction capability of the models against a
26
27 real-world external set of samples not included in any step of both the calibration and selection of
28
29 RF parameters. Sample classes were assigned according to two criteria:
30
31

- 32
33
- 34 1) First, samples were classified based on the concentration of urea measured by
35
36 employing an enzymatic reference method used at the hospital for the
37
38 conventional routine analysis. 61 samples (15 for calibration and 46 for validation)
39
40 contained urea concentration values above 60 mg/dL and the other 61 (15 for
41
42 calibration and 46 for validation) below 40 mg/dL. Urea concentration in serum
43
44 can be measured using FTIR spectroscopy¹⁶ and so, a discriminative model should
45
46 be a priori to provide good predictive performance.
47
48
 - 49 2) Samples were distributed based on an arbitrary hospital tag number: 56 samples
50
51 (15 for calibration and 41 for validation) with a tag finished by 1, 2, 3 and 4 were
52
53 assigned to 'class A' and 66 samples (15 for calibration and 51 for validation)
54
55 with a tag finished by 0, 5, 6, 7, 8 and 9 were assigned to 'class B'. This
56
57
58
59
60

1
2 classification was therefore arbitrary and thus, any classification obtained
3
4 considering this pseudo-class could be attributed to chance, model overfitting or to
5
6 an optimistic estimation of the CV-figures of merit.
7
8

9 10 **2.3 Random Forest Classification**

11
12 Random forest classification was performed in Matlab 2012b (The Mathworks, Natick, USA) using
13 the routine available at <https://code.google.com/p/randomforest-matlab/> and further calculations
14
15 were carried out using in-house written scripts.
16
17

18
19 In order to increase the number of variables, an augmented data matrix was obtained by
20 concatenating the raw and the first derivate spectra calculated using a Savitzky-Golay routine,
21
22 resulting in a calibration dataset with 30 samples and 986 variables. The tuning parameters of the
23
24 RF classifiers were set following the recommendations of Ollesch et al⁹. The number of trees
25
26 employed for each random classifier was three times the number of variables and the number of
27
28 wavenumber-features used in each tree was a third of the number of variables). All other parameters
29
30 of the RF classifiers were set to the default values established in the routines.
31
32
33
34
35

36 **2.4 Cross validation, external validation and permutation test**

37
38 For performing the Montecarlo CV only the calibration samples were used. A set of M/n samples,
39
40 where “M” is the number of samples in the calibration set and “n” the number of data splits used
41
42 during cross validation, were randomly selected and predicted by the RF classifiers obtained using
43
44 the remaining samples. This step was repeated 60 times, changing the calibration and internal
45
46 validation subsets selected and registering the CV-predicted classes in order to obtain a
47
48 representative value of the CV from the all possible subsets. Based on CV-predicted classes, the
49
50 samples were classified using a majority rule: each sample was classified in the group of the
51
52 majority of CV-predicted values. The error rate was then obtained by dividing the number of the
53
54 bad classified samples by the numbers of samples of the calibration set. Then, the external test set
55
56
57
58
59
60

1
2 of 92 samples, was used for establishing the actual classification performance of the models. In this
3
4 case the whole calibration dataset was employed for computing the RF classifiers and the same
5
6 parameters employed for the CV in terms of tuning parameters of the RF classifiers. In the case of
7
8 the external validation after variable selection, RF classifiers were performed including only the
9
10 variables selected in the routine.
11

12
13
14 Permutation testing included a random class permutation followed by an ensemble RF using the
15
16 dataset with the permuted class labels. CV was carried out as described above. Those two steps
17
18 were repeated 200 times and the CV error was obtained. The distribution of the so-called ‘random
19
20 CV-values’ was used as a null distribution for establishing the statistical significance of the model
21
22 CV estimators obtained using real class labels. This statistical significance was calculated as the
23
24 number of permuted values lower than the CV error obtained using the real classes divided by the
25
26 number of permuted values.
27
28

29 30 **2.5 Variable selection** 31

32
33 The variable selection (selection of the wavenumbers which provide the best classification)
34
35 procedure was based on a Montecarlo procedure as described in a previous reference⁹, setting the
36
37 tunable parameters to less restrictive values to facilitate the selection of correlated but
38
39 uninformative variables. Montecarlo feature selection utilizes the repetition of a cyclic routine
40
41 randomly selecting 23 samples in each cycle. In the cyclic procedure, an ensemble RF classifier was
42
43 created and validated using the Montecarlo CV as described in the section 2.3, recording the mean
44
45 Gini values for each variable and the CV error. Then, 20% of variables with the lowest Gini values
46
47 were eliminated and the CV was repeated. This cycle of CV and elimination of variables was
48
49 repeated 12 times until only 45 values were retained. From the variables selected from the 10 RF
50
51 models, those selected by the model with the best classification error rate were finally chosen as
52
53 variables selected for the cyclic routine. This cyclic routine was repeated 10 times and the final
54
55 variables chosen were those that were selected in 3 or more out of the 10 cyclic routines.
56
57
58
59
60

3. Results

3.1 Statistical significance of cross validation errors

The selection of the employed n-folds (i.e. number of sub-samples used during CV) is critical when CV errors are used as the selection of a high n-fold could lead to an overoptimistic evaluation of the model. RF models for the classification by sample tag number and urea levels were performed and validated using an independent validation set (external validation) and CV with different n-folds (15, 10 and 4). Results (see Table 1) evidenced a strong association between the CV error and the n-fold. In all cases, the CV error was lower than the one obtained using the independent validation set and increasing the n-fold, the CV error decreased. More importantly, in the case of the classification according to the sample tag number, classification errors between 26 and 44% were obtained. Those values are significantly lower than the one obtained by the independent validation, which was 52±1% and hence corresponded to the expected value for a non-statistically significant class difference. This demonstrates the capability of RF for recognizing apparent discriminant patterns even in random data, which only was evidenced after a proper external validation with an independent sample set and might be not detected using CV-figures of merit.

Table 1: Error parameters obtained from the validation of the RF classifiers for different n-fold values.

Class	Parameter	n-fold		
		15	10	4
TAG NUMBER	Independent validation error (%)	51.8±0.5%	51.8±0.5%	55.4%±1.0%
	CV Error (%)	26±4%	38±4%	44%±3%
	p-value	0.36	0.33	0.35
UREA LEVEL	Independent validation error (%)	32±2%	32±3%	32±2%
	CV Error (%)	8.2±0.9%	13.3±0.7%	15.7±0.7%
	p-value	<0.005	0.005	0.005

Permutation testing was applied for calculating the statistical significance of the aforementioned CV errors. The values of the CV errors obtained for the classification of the permuted classes are represented in Figure 1. It is remarkable that in the case of the 15-fold CV (leave-2-out), the distribution of CV errors obtained for 200 permutations are always lower than the theoretically expected 50% error value. This shows the ability of RF to provide low CV errors from IR based classification/diagnostic methods no matter which class/illness is under study. However, comparing the CV errors to the distribution of the CV errors obtained from permuted classes enables the evaluation of the statistical significance of the model under study.

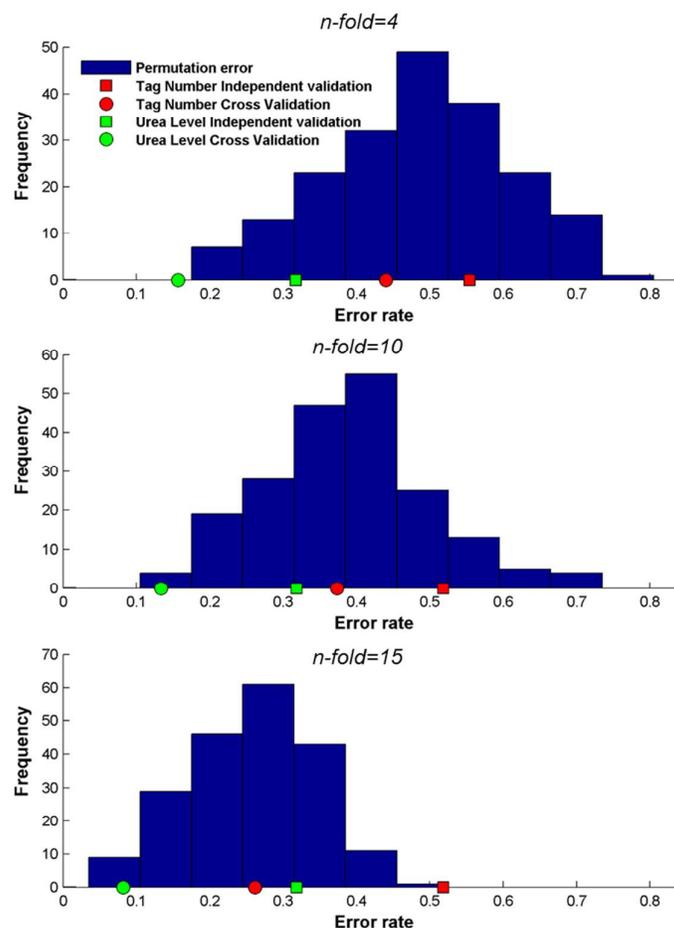


Figure 1: Distribution of errors obtained from the independent and cross validation of the models using RF classifiers.

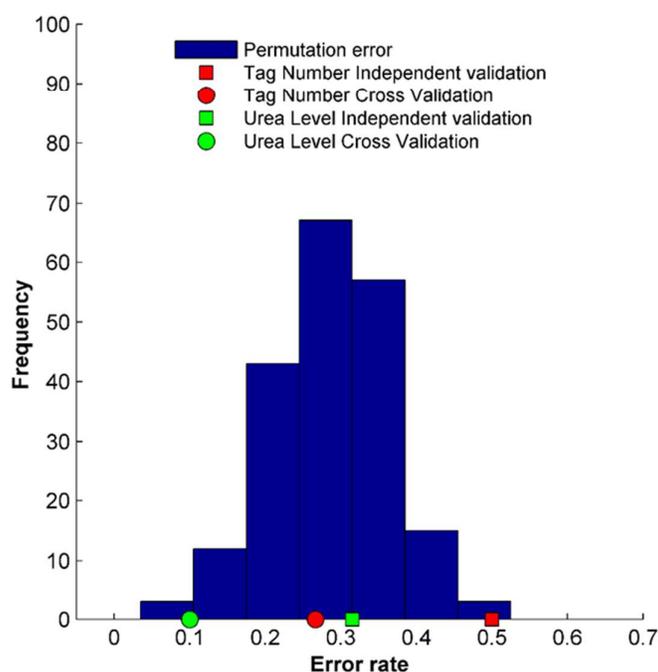
1
2 The mean of the Gaussian distribution of the prediction moves to lower values as the n-fold
3 increases. Therefore, although the CV error also decreases with increasing n-fold, their statistical
4 significance does not show substantial changes throughout the different n-folds under study. In the
5 case of the arbitrary classes based on the sample tag number, the obtained CV was in the range of
6 the mean CV error obtained from permuted classes as shown in Figure 1 and the calculated
7 statistical significance was $p > 0.05$ as shown in Table 1. Thus, the ability of permutation testing for
8 detecting overoptimistic estimations of the CV error at different n-folds is demonstrated. In
9 contrast, for the classification between samples with low and high levels of urea, the CV errors
10 obtained were significantly lower than those obtained by the permuted classes ($p \leq 0.005$). It has to
11 be also considered that the prediction error obtained using the independent validation dataset was
12 two times higher than the one obtained using CV, which indicates, as previously shown¹², that
13 permutation testing focuses on assessing the significance of the classification without taking into
14 account the real predictability of the population of samples to predict. An n-fold value of 4 (leave-7-
15 out) was selected for the further calculations.

3.2 Statistical significance of CV errors after feature selection

36 The development of diagnostic tools based on the modelling of IR spectra faces a critical challenge
37 regarding the sample to variable ratio, especially considering early studies where the number of
38 samples is limited. Therefore, feature selection is normally used for “producing a smaller number of
39 variables (“feature = input variable”) that are more informative than the original whole set of
40 wavenumber-variables”¹⁷. The main aim of this procedure is the improvement of the diagnostic
41 model performance and the detection of spectral biomarkers that can help to improve the
42 understanding of the biological mechanisms of the illness, also providing a scientific justification of
43 the IR based classification. However, the introduction of this feature selection could lead to
44 overoptimistic results and the apparent improvement in the CV error obtained after variable
45 selection can also be originated from the susceptibility of variable selection towards chance

1
2 correlations¹⁸. This effect is a well-known source of complications in the metabolomics field and
3
4 the likelihood of obtaining a model based on chance correlation depends on the sample to variable
5
6 ratio and the correlation structure of the data¹⁹. In addition, the presence of chance correlation on
7
8 regression models for the prediction of glucose concentration has been widely studied^{20,21}.
9

10
11 The aim of this section was to assess the use of the permutation testing for evaluating the statistical
12
13 significance of IR models after variable selection. Variable selection was carried out for the
14
15 classification of samples according to the sample tag number and low/high urea levels followed by
16
17 cross validation and permutation testing. From results shown in Figure 2 it can be appreciated that,
18
19 as compared to the 4-fold distribution before variable selection (see Figure 1), the variable selection
20
21 procedure improved the CV error of the classification according to the sample tag number from
22
23 44% to 27%. However, the CV distribution of errors obtained from the permuted values also moves
24
25 to lower errors and consequently no significant difference to the CV error was observed ($p=0.46$).
26
27 In contrast, for the classification according to the urea levels, the CV error obtained after variable
28
29 selection was still significantly lower than errors obtained from permutation testing ($p=0.015$).
30
31
32
33
34
35



1
2
3
4
5 *Figure 2: Distribution of errors obtained from the independent and cross validation of the models*
6 *performed for classification of the samples after variable selection using the RF classifiers.*
7
8
9

10 Figure 3 shows the frequency of the selected variables in each of the 10 cycles performed during the
11 variable selection process, where red dots correspond to the variables that were finally selected. The
12 distribution of the variables selected differs from the classification considered. In the case of the
13 arbitrary determination of the sample tag number (see Figure 3b), several wavenumber values
14 arbitrarily distributed across the IR bands were selected. It should be noted that the use of these
15 variables improved the CV error in comparison to the use of the whole spectrum, and if the
16 statistical significance was not considered, they could have been mistakenly assigned as
17 “biomarkers” for the sample tag number. These results demonstrate that the variable selection
18 procedure could find correlations on trivial wavenumber values which apparently improved the
19 prediction capability of the model, but which did not have any relationship with the studied
20 parameter. Nevertheless, the use of the permutation test and a thorough investigation of the
21 wavenumber values “feature selected” demonstrates the irrelevance of these markers and models,
22 which are likely associated with chance correlations without statistically significant results.
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39

40 In the case of the classification according to urea levels, the variables selected correspond to
41 characteristic urea bands, such as a shoulder of the band at 1040 (ρNH_2)²² cm^{-1} in the untreated and
42 the first derivative spectrum and the bands at 1552 ($\nu\text{C-N}$)²² and at 1155 (ρNH_2)²² cm^{-1} in the first
43 derivative spectrum. It has to be noted that the bands selected are not the most intense bands of the
44 urea spectrum. However, this might be due to strong interference in the 1600-1400 cm^{-1} region
45 caused by amide I and II bands. The selected wavenumber values are in good agreement with the
46 selectivity ratio²³ obtained from a PLS regression performed on the data set for the prediction of the
47 concentration of urea (see black spectra)¹⁶.
48
49
50
51
52
53
54
55
56
57
58
59
60

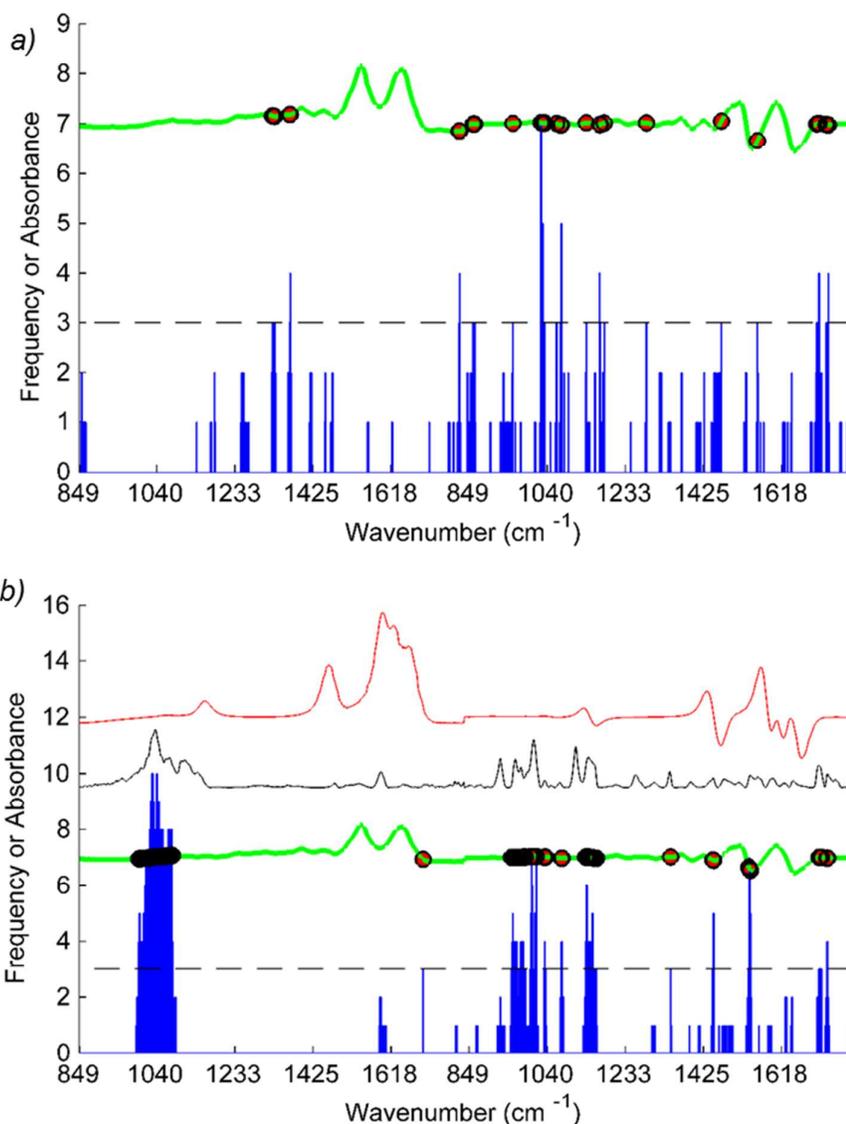


Figure 3: Wavenumber values chosen by the feature selection procedure for the classification according to the sample tag number (a) and urea levels (b). Blue bars indicate the variable selection frequency in the 10 cyclic routines performed. Variables with a frequency higher than the threshold (dashed line) were finally retained and are indicated as red points with black circles in the mean sample spectrum (green line). Panel (b) also shows the ATR spectrum of a urea standard solution at 1000 mg/dL (red line) and the selectivity ratio values for a PLS regression model (black

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

line). The Y axis label “Frequency” corresponds to blue bars and the Y axis label “Absorbance” corresponds to the spectra. Selectivity ratio is unitless.

3.3 The effect of replicates on the statistical significance on RF models

The high variable to sample ratio can be also compensated by including measurement replicates of the sample in the model. The use of measurement replicates is usually justified by the need of taking into account the variability of different spots of the samples. The ‘sample replicate trap’ occurs when replicates of the same physical samples are introduced in both the calibration and the validation subsets employed during model CV and external validation²⁴. This effect can be corrected by using a continuous block cross validation employing sub models of m-blocks being m the number of replicates. However this process can be troublesome if the number of replicates is not the same for all samples.

The effect of the use of replicates during cross validation of RF models was evaluated employing permutation testing. Figure 4 represents the CV and independent validation errors obtained for the two class labels considered as well as the distribution of the CV errors for 200 permutations when two and three replicates are used. Once again, the distribution of the errors obtained using the permuted classes was dramatically shifted to lower errors as the number of replicates increased. According to the CV error, the RF was able to correctly classify 82 and 93% of the samples according to the sample tag number using 2 and 3 replicates, respectively. As expected, the permutation testing indicated that those classifications were not statistically significant ($p=0.6$ and 0.5). Regarding the classification of the samples according to urea levels, the prediction errors were found to be significantly different from those obtained from randomly permuted classes in all cases ($p<0.05$). It has to be remarked that, although the use of replicates improved the CV errors, it

decreased the real prediction capability of the model, as evidenced by the increase of the independent prediction errors obtained using 1, 2 and 3 replicates.

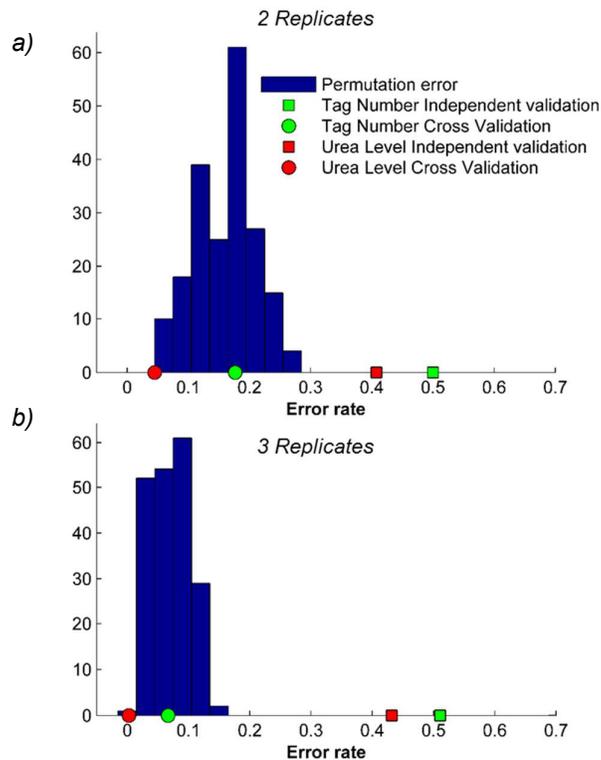


Figure 4: Distribution of the classification errors of the RF classifiers obtained using 2 (a) and 3 replicates (b).

4. Conclusions

Results obtained show that the application of permutation testing is suitable for calculating the statistical significance of IR based diagnostic models. In the case of the sample tag number, the manifested CV errors obtained were not statistically significant after permutation testing. In the case of the classification according to urea levels, CV errors were found to be significant. It has to be stipulated that the method did not provide an estimation of the predictive capacity of the model, which can only be achieved by an external validation set. However, it provides a significance value of the classification, which contains information about the significance of the differences between the spectra of different classes. Therefore this method can be used as an early estimator of the

1
2 validity of a hypothesis or for making a decision about the utility of incorporating further samples to
3
4 the system. Although there are several repetitions of the procedures involved in the modelling such
5
6 as the CV and variable selection, which can become computationally intensive, it is nonetheless a
7
8 versatile tool that can be used in systems with a high variable to sample ratio or with an unbalanced
9
10 number of samples in each class and hence we recommend the use of permutation testing especially
11
12 in the case of preliminary studies.
13
14
15
16
17
18

19 Acknowledgements

20
21
22 JK acknowledges the Sara Borrell CD12/00667 grant (*Instituto Carlos III*, Ministry of Economy
23 and Competitiveness, Spain). B.R.W. is supported by an Australian Research Council (ARC) Future
24
25 Fellowship (FT120100926) and the project is supported by an ARC Discovery grant DP120100431.
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

- 41 1. C. Kendall, M. Isabelle, F. Bazant-Hegemark, J. Hutchings, L. Orr, J. Babrah, R. Baker, and N.
42 Stone, *Analyst*, 2009, **134**, 1029–1045.
- 43 2. D. I. Ellis and R. Goodacre, *The Analyst*, 2006, **131**, 875–885.
- 44 3. A. Sahu, K. Dalal, S. Naglot, P. Aggarwal, and C. Murali Krishna, *PLoS ONE*, 2013, **8**, e78921.
- 45 4. J. R. Hands, P. Abel, K. Ashton, T. Dawson, C. Davis, R. W. Lea, A. J. S. McIntosh, and M. J.
46 Baker, *Anal. Bioanal. Chem.*, 2013, **405**, 7347–7355.
- 47 5. M. Khanmohammadi and A. B. Garmarudi, *TrAC Trends Anal. Chem.*, 2011, **30**, 864–874.
- 48 6. K. Gajjar, J. Trevisan, G. Owens, P. J. Keating, N. J. Wood, H. F. Stringfellow, P. L. Martin-
49 Hirsch, and F. L. Martin, *The Analyst*, 2013, **138**, 3917–3926.
- 50 7. G. L. Owens, K. Gajjar, J. Trevisan, S. W. Fogarty, S. E. Taylor, B. Da Gama-Rose, P. L. Martin-
51 Hirsch, and F. L. Martin, *J. Biophotonics*, 2014, **7**, 200–209.
- 52 8. A. Sahu, S. Sawant, H. Mamgain, and C. M. Krishna, *Analyst*, 2013, **138**, 4161–4174.
- 53 9. J. Ollesch, S. L. Drees, H. M. Heise, T. Behrens, T. Brüning, and K. Gerwert, *Analyst*, 2013, **138**,
54 4092–4102.
- 55 10. G. Bellisola and C. Sorio, *Am. J. Cancer Res.*, 2011, **2**, 1–21.
- 56 11. K. H. Esbensen and P. Geladi, *J. Chemom.*, 2010, **24**, 168–187.
- 57
58
59
60

12. C. M. Rubingh, S. Bijlsma, E. P. P. A. Derks, I. Bobeldijk, E. R. Verheij, S. Kochhar, and A. K. Smilde, *Metabolomics*, 2006, **2**, 53–61.
13. S. Smit, M. J. van Breemen, H. C. J. Hoefsloot, A. K. Smilde, J. M. F. G. Aerts, and C. G. de Koster, *Anal. Chim. Acta*, 2007, **592**, 210–217.
14. J. A. Westerhuis, H. C. J. Hoefsloot, S. Smit, D. J. Vis, A. K. Smilde, E. J. J. van Velzen, J. P. M. van Duijnhoven, and F. A. van Dorsten, *Metabolomics*, 2008, **4**, 81–89.
15. G. R. Lloyd, J. Hutchings, L. M. Almond, H. Barr, C. Kendall, and N. Stone, 2012, vol. 8219, p. 82190C–82190C–6.
16. D. Perez-Guaita, J. Ventura-Gayete, C. Pérez-Rambla, M. Sancho-Andreu, S. Garrigues, and M. de la Guardia, *Microchem. J.*, 2013, **106**, 202–211.
17. J. Trevisan, P. P. Angelov, P. L. Carmichael, A. D. Scott, and F. L. Martin, *Analyst*, 2012, **137**, 3202–3215.
18. J. Kuligowski, D. Pérez-Guaita, J. Escobar, M. de la Guardia, M. Vento, A. Ferrer, and G. Quintás, *Talanta*, 2013, **116**, 835–840.
19. K. Baumann, *QSAR Comb. Sci.*, 2005, **24**, 1033–1046.
20. B. Zhao, Y. Cao, R. Liu, and K. Xu, *Guang Pu Xue Yu Guang Pu Fen Xi Guang Pu*, 2012, **32**, 934–938.
21. R. Liu, W. Chen, X. Gu, R. K. Wang, and K. Xu, *J. Phys. Appl. Phys.*, 2005, **38**, 2675–2681.
22. J. Grdadolnik and Y. Maréchal, *J. Mol. Struct.*, 2002, **615**, 177–189.
23. T. Rajalahti, R. Arneberg, F. S. Berven, K.-M. Myhr, R. J. Ulvik, and O. M. Kvalheim, *Chemom. Intell. Lab. Syst.*, 2009, **95**, 35–48.
24. K. A. Bakeev, *Process Analytical Technology: Spectroscopic Tools and Implementation Strategies for the Chemical and Pharmaceutical Industries*, John Wiley & Sons, 2008.