

Central Lancashire Online Knowledge (CLoK)

Title	Aluminium foil as an alternative substrate for the spectroscopic interrogation of endometrial cancer
Type	Article
URL	https://clock.uclan.ac.uk/21857/
DOI	https://doi.org/10.1002/jbio.201700372
Date	2018
Citation	Paraskevaidi, Maria, Medeiros-De-morais, Camilo De lelis orcid iconORCID: 0000-0003-2573-787X, Raglan, Olivia, Lima, Kássio M.G., Paraskevaidis, Evangelos, Martin-Hirsch, Pierre L., Kyrgiou, Maria and Martin, Francis L (2018) Aluminium foil as an alternative substrate for the spectroscopic interrogation of endometrial cancer. <i>Journal of Biophotonics</i> , 11 (7). e201700372. ISSN 1864-063X
Creators	Paraskevaidi, Maria, Medeiros-De-morais, Camilo De lelis, Raglan, Olivia, Lima, Kássio M.G., Paraskevaidis, Evangelos, Martin-Hirsch, Pierre L., Kyrgiou, Maria and Martin, Francis L

It is advisable to refer to the publisher's version if you intend to cite from the work.
<https://doi.org/10.1002/jbio.201700372>

For information about Research at UCLan please go to <http://www.uclan.ac.uk/research/>

All outputs in CLoK are protected by Intellectual Property Rights law, including Copyright law. Copyright, IPR and Moral Rights for the works on this site are retained by the individual authors and/or other copyright owners. Terms and conditions for use of this material are defined in the <http://clock.uclan.ac.uk/policies/>

1 **Aluminium foil as an alternative substrate for the**
2 **spectroscopic interrogation of endometrial cancer**

3 Maria Paraskevaidi^{a,1}, Camilo L.M. Morais^{a,b}, Olivia Raglan^c, Kássio M.G. Lima^b, Evangelos
4 Paraskevaidis^d, Pierre L. Martin-Hirsch^e, Maria Kyrgiou^c and Francis L. Martin^{a,1}

5 ^a*School of Pharmacy and Biomedical Sciences, University of Central Lancashire, Preston PR1*
6 ²*HE, UK*

7 ^b*Institute of Chemistry, Biological Chemistry and Chemometrics, Federal University of Rio*
8 ⁸*Grande do Norte, Natal 59072-970, Brazil*

9 ^c*Institute of Reproductive and Developmental Biology, Faculty of Medicine, Imperial College,*
10 ¹⁰*London W12 0HS, UK*

11 ^d*Department of Obstetrics and Gynaecology, University Hospital of Ioannina, Ioannina 45500,*
12 ¹²*Greece*

13 ^e*Department of Obstetrics and Gynaecology, Lancashire Teaching Hospitals NHS Foundation,*
14 ¹⁴*Sharoe Green Unit, Fullwood, Preston PR2 9HT, UK*

15
16 ¹To whom correspondence should be addressed. Email: mparaskevaidi@uclan.ac.uk or
17 flmartin@uclan.ac.uk

22 **Abstract**

23 Biospectroscopy has the potential to investigate and characterise biological samples and could,
24 therefore, be utilised to diagnose various diseases in a clinical environment. An important
25 consideration in spectrochemical studies is the cost-effectiveness of the substrate used to
26 support the sample, as high expense would limit their translation into clinic. In this paper, the
27 performance of low-cost aluminium (Al) foil substrates was compared with the commonly used
28 low-emissivity (low-E) slides. Attenuated total reflection-Fourier transform infrared (ATR-
29 FTIR) spectroscopy was used to analyse blood plasma and serum samples from women with
30 endometrial cancer and healthy controls. The two populations were differentiated using
31 principal component analysis with support vector machines (PCA-SVM) with 100% sensitivity
32 in plasma samples (endometrial cancer=70; healthy controls=15) using both Al foil and low-E
33 slides as substrates. The same sensitivity results (100%) were achieved for serum samples
34 (endometrial cancer=60; healthy controls=15). Specificity was found higher using Al foil
35 (90%) in comparison to low-E slides (85%) and lower using Al foil (70%) in comparison to
36 low-E slides in serum samples. The establishment of Al foil as low-cost and highly-performing
37 substrate would pave the way for large-scale, multi-centre studies and potentially for routine
38 clinical use.

39

40

41

42

43

44 **Introduction**

45 Vibrational spectroscopy is increasingly utilised in biomedical research as a valuable tool in
46 disease investigation. Allowing the analysis of a variety of biological samples, such as cells,
47 tissues and biofluids, this spectrochemical analysis has a bright future ahead, not only in
48 scientific/laboratory research but also in clinical practice. The key factor that renders this
49 analytical method a perfect diagnostic tool, in comparison to other molecular methods, is its
50 non-destructive, cost-effective and label-free nature. Over the years, infrared (IR) and Raman
51 spectroscopic techniques have been employed to study a number of different diseases like
52 cancer, neurological diseases, prenatal disorders and many others ¹⁻¹⁰. Within the field of
53 disease investigation, spectroscopy has the potential to diagnose and monitor a disease, while
54 at the same time assessment of surgical margins of a tumour or determination of the subtype
55 of a disease is also feasible.

56 Most spectroscopic studies so far, with only a few exceptions ^{8, 10, 11}, have included a
57 limited number of subjects which appears to be an important limitation for the establishment
58 of the method and its migration into clinics ¹²⁻¹⁴. Standardisation and validation of methods
59 should be performed in large clinical trials for more robust and trustworthy results. A further
60 issue that limits the ability for clinical implementation relates to experimental methodology.
61 Specifically, inconsistencies in the pre-analytical stages of sample collection and preparation
62 to spectral collection and data analysis. A fundamental factor of the analytical procedure is the
63 use of the correct substrate in order to avoid non-biological interference from the substrate in
64 use. Unfortunately, the majority of the available, “featureless” substrates are high-cost ^{15, 16},
65 something which prevents their use in large scale studies and routine analysis. Previous studies
66 have even developed data correction algorithms to remove the substrate’s signal after the
67 collection of the raw spectra ¹⁶⁻¹⁸.

68 Different types of substrates are selected depending on the spectroscopic technique used
69 each time (*e.g.*, IR or Raman spectroscopy), as well as on the chosen sampling mode [*e.g.*,
70 transmission IR, transfection IR or attenuated total reflection (ATR)]. Namely, some of the
71 substrates that have been used for IR and Raman spectroscopy over the years include barium
72 fluoride (BaF₂), calcium fluoride (CaF₂), zinc selenide (ZnSe), gold-coated (Au), silver or
73 silver-coated (Ag), fused silica (SiO₂) and low-emissivity (low-E) slides^{19,20}. However, due to
74 their expense, efforts are being made to introduce novel, low-cost substrates that would
75 facilitate the analysis of hundreds, even thousands, of samples cost-effectively. Glass substrates
76 are routinely used in medical laboratories and hospitals for preparation of analysis of various
77 types of biological samples; however, glass has been found unsuitable for spectroscopy as it
78 generates background signal and distorts the biological information coming from the samples
79¹⁹. Therefore, an ideal approach would be to take advantage of the extremely cost-effective
80 glass slides by covering them with a metallic surface that would eliminate any background
81 noise. Previous proof-of-concept studies have been conducted showing aluminium (Al) foil's
82 potential as a suitable substrate^{11, 16, 21}. A robust and inexpensive substrate for both IR and
83 Raman spectroscopic methods would be extremely beneficial. However, there has been no
84 conclusive study comparing its effect on diagnostic accuracy with other, widely used
85 substrates.

86 In this study, we used ATR-FTIR spectroscopy to explore whether Al foil could be an
87 appropriate substrate for spectroscopic investigations. ATR-FTIR uses an internal reflection
88 element (IRE) with a high refractive index to direct the beam to the sample; an evanescent
89 wave is created, penetrating the sample by a few microns in order to derive its chemical
90 information²². A commonly used substrate for ATR-FTIR measurements is the low-E slide,
91 which has been effectively used in numerous biological studies in the past²³⁻²⁵. Therefore, we
92 compared our results from the low-E slides with those from Al foil slides to assess the

93 performance of the latter with regard to the diagnostic accuracy. For the purpose of this piece
94 of work, we analysed blood samples from women with endometrial cancer, as well as from
95 benign cases used as controls. Endometrial cancer develops in the endometrium (*i.e.*, inner
96 lining of the uterus) and is the fourth most common gynaecological cancer in the developing
97 world, with an increasing incidence in postmenopausal women; in 2012 alone, 319,000 new
98 cases were diagnosed worldwide ²⁶. Although symptoms of endometrial cancer develop
99 relatively early, which allows “timely” diagnosis and early intervention, a more objective, less
100 expensive and non-invasive method of diagnosing this type of cancer is highly desirable and
101 clinically indicated. Currently, a diagnosis is based on microscopic histological examination of
102 endometrial tissue, which is dependent on subjective interpretation, therefore allowing human
103 error.

104 **Materials and Methods**

105 **Blood plasma and serum analysis**

106 The collection of all samples for this study was approved by the institutional review board at
107 Imperial College Healthcare NHS Trust (tissue bank sub-collection number GYN/HG/13-020).
108 All patients provided informed consent for use of their samples in this study. This study
109 included age-matched cohorts; plasma samples were available for 70 endometrial cancer
110 patients and 15 non-cancer individuals used as controls; serum samples were available for 60
111 endometrial cancer patients and 15 controls. At time of diagnosis, patients were not receiving
112 any medications such as Tamoxifen treatments which might affect the outcomes. Also women
113 who had hyperplasia or hypertension have been excluded. Both blood plasma and serum
114 samples were collected and stored at -80°C until analysis; prior to spectroscopic interrogation,
115 the samples were left to defrost at room temperature before 50 µL of each were deposited on a
116 substrate and left to air-dry for approximately 30 min. All of the samples were analysed in

117 duplicates using two different substrates: the IR-reflective glass slides (MirrIR Low-E slides,
118 Kevley Technologies, USA) and cheap, microscope glass slides covered with Al foil. The latter
119 were carefully flattened with the shiny side of the foil being exposed to achieve a greater level
120 of reflectivity. Covering the slide with Al foil required ~30-45 seconds with one slide taking
121 up to 3 different samples, rendering the slide preparation time insignificant.

122 **Spectrochemical Analysis**

123 All blood samples were analysed using a Tensor 27 FTIR spectrometer with Helios ATR
124 attachment (Bruker Optics Ltd, Coventry, UK). The sampling area, defined by the internal
125 reflection element (IRE), which was a diamond crystal, was approximately $250\ \mu\text{m} \times 250\ \mu\text{m}$.
126 The slide with the sample is placed onto a moving platform with the sample facing up; the
127 platform is then moved upward to achieve good contact with the diamond crystal. Spectral
128 resolution was $8\ \text{cm}^{-1}$ with two times zero-filling, giving a data-spacing of $4\ \text{cm}^{-1}$ over the range
129 $4000\text{-}400\ \text{cm}^{-1}$; 32 co-additions and a mirror velocity of 2.2 kHz were used for optimum signal
130 to noise ratio. A CCTV camera attachment was used to locate the area of interest and spectra
131 were acquired from ten different locations to minimize bias. Also, in order to take into
132 consideration the natural phenomenon of “coffee ring” effect, spectra were mainly collected
133 from the periphery of each drop where the absorbance intensity was higher, as important
134 components, such as proteins and nucleic acids, migrate towards the edge of the drop after
135 drying²⁷. The ATR crystal was cleaned with distilled water before moving to a different sample
136 and a background spectrum was acquired to take into account any atmospheric changes.

137 **Spectral data handling and analysis**

138 All spectral information was converted to suitable files (.txt) before input to MATLAB
139 (Mathworks, Natick, USA). Pre-processing and computational analysis of the data was
140 performed using PLS Toolbox version 7.9.3 (Eigenvector Research, Inc., Manson, USA) and

141 an in-house developed IRootLab toolbox (<http://trevisanj.github.io/irootlab/>). Pre-processing
142 of the acquired spectra is an essential step of all spectroscopic experiments and is used to
143 correct problems associated with spectral acquisition, instrumentation or even sample handling
144 before further multivariate analysis ²⁸. In this study, spectra were cut at the biochemical
145 fingerprint region (1800-900 cm⁻¹), rubberband baseline corrected and vector normalised.

146 The samples were divided into training (~70%), validation (~15%) and test (~15%) sets
147 on a patient basis before chemometric analysis, using the Kennard-Stone sample selection
148 algorithm ²⁹; all spectra collected for each individual were used for model construction. In total,
149 60 samples were used for training ($n = 600$ spectra), 12 for validation ($n = 120$ spectra) and 13
150 for test ($n = 130$ spectra) with plasma samples; and 53 for training ($n = 530$ spectra), 11 for
151 validation ($n = 110$ spectra) and 11 for test ($n = 110$ spectra) with serum samples. The training
152 set was used for model construction, the validation set for optimization of the number of
153 principal components and latent variables used, and the test set for final model evaluation.
154 Cross-validation venetian blinds (10 splits with 1 sample per split) was used for optimization
155 of support vector machines (SVM) parameters (cost, epsilon, gamma and number of support
156 vectors) in principal component analysis with support vector machines (PCA-SVM).

157 For the classification of endometrial cancer and non-cancer cases a number of
158 chemometric techniques was used, such as partial least squares discriminant analysis (PLS-
159 DA); and principal component analysis followed by linear discriminant analysis (PCA-LDA),
160 quadratic discriminant analysis (PCA-QDA) and support vector machines (PCA-SVM).

161 PLS-DA is one of the most known chemometric technique of supervised classification.
162 It is based on a linear classification model for which the classification criterion is obtained by
163 partial least squares (PLS) analysis ³⁰. For this, PLS is applied to the data reducing the original
164 variables (*e.g.*, wavenumbers) to a few number of latent variables (LVs) in an interactive

165 process, in which the category variables for each class in the training set (*e.g.*, ± 1) is used to
166 optimise the model. A straight line that divides the classes' regions is then found ³¹.

167 Similarly to PLS, PCA also reduces the original data into a few set of variables called
168 principal components (PCs). These variables are orthogonal to each other and account most of
169 the explained variance from the original data set. They are composed of scores and loadings
170 that are used to identify similarities/dissimilarities among the samples and the weight that each
171 variable contributes for the PCA model, respectively ³². However, differently from PLS, the
172 category variables are not used for this reduction. To perform a supervised classification model,
173 the PCA scores are employed as input variables for discriminant algorithms. This procedure
174 avoids collinearity problems and also speeds up computational analysis.

175 LDA and QDA are discriminant algorithms that create a classification rule between the
176 classes based on a Mahalanobis distance. The main difference between these techniques is that
177 LDA uses a pooled covariance matrix to calculate the discriminant function between the
178 classes, whereas QDA uses the variance-covariance matrix of each class separately ³³.
179 Therefore, QDA usually achieves better performance than LDA when analysing complex data
180 sets where the variance structures between the classes are very different. The LDA (L_{ik}) and
181 QDA (Q_{ik}) classification scores are calculated following the equations ³⁴:

$$182 \quad L_{ik} = (\mathbf{x}_i - \bar{\mathbf{x}}_k)^T \Sigma_{\text{pooled}}^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}_k) - 2 \log_e \pi_k \quad (1)$$

$$183 \quad Q_{ik} = (\mathbf{x}_i - \bar{\mathbf{x}}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}_k) + \log_e |\Sigma_k| - 2 \log_e \pi_k \quad (2)$$

184 in which \mathbf{x}_i is the vector containing the classification variables for sample i (*e.g.*, PCA scores
185 for A components); $\bar{\mathbf{x}}_k$ is the mean vector of class k ; Σ_k is the variance-covariance matrix of
186 class k ; Σ_{pooled} is the pooled covariance matrix; and π_k is the prior probability of class k .
187 These last three terms are calculated by ³⁴:

188
$$\Sigma_k = \frac{1}{n_k-1} \sum_{i=1}^{n_k} (\mathbf{x}_i - \bar{\mathbf{x}}_k) (\mathbf{x}_i - \bar{\mathbf{x}}_k)^T \quad (3)$$

189
$$\Sigma_{\text{pooled}} = \frac{1}{n} \sum_{k=1}^K n_k \Sigma_k \quad (4)$$

190
$$\pi_k = \frac{n_k}{n} \quad (5)$$

191 where n_k is the number of samples of class k ; n is the total number of samples in the training
 192 set; and K is the number of classes.

193 On the other hand, SVM is a technique that classifies data sets in a completely non-
 194 linear fashion. For this, SMVs classifiers work by finding a classification hyperplane that
 195 separates the data clusters providing the largest margin of separation ³⁵. During model
 196 construction, the data is transformed into a different feature space by means of a kernel function
 197 that is responsible for the SVM classification ability ³³. The most common kernel function is
 198 the radial basis function (RBF). The SVM classifier takes the form of ³⁶:

199
$$f(x) = \text{sign}\left(\sum_{i=1}^{N_{SV}} \alpha_i y_i K(\mathbf{x}_i, \mathbf{z}_j) + b\right) \quad (6)$$

200 where N_{SV} is the number of support vectors; α_i is the Lagrange multiplier; y_i is the class
 201 membership (*e.g.*, ± 1); b is the bias parameter; and $K(\mathbf{x}_i, \mathbf{z}_j)$ is the RBF kernel function,
 202 calculated by:

203
$$K(\mathbf{x}_i, \mathbf{z}_j) = \exp\left(-\gamma \|\mathbf{x}_i - \mathbf{z}_j\|^2\right) \quad (7)$$

204 in which \mathbf{x}_i and \mathbf{z}_j are samples measurement vectors; and γ is the parameter that determines
 205 the RBF width.

206 **Results and Discussion**

207 By employing the above-mentioned multivariate techniques (PCA-LDA, PLS-DA, PCA-QDA
 208 and PCA-SVM), it was demonstrated that some provided superior performance than others.

209 The techniques were very different from each other and were used following a parsimonious
210 order (PCA-LDA < PLS-DA < PCA-QDA < PCA-SVM). It is natural to expect an
211 improvement of the results when more robust algorithms are applied, as the classification
212 methods varied from a linear (PCA-LDA and PLS-DA) to a completely non-linear
213 classification algorithm (PCA-SVM). Analysis of the plasma samples deposited on Al foil
214 showed classification to be: 68% sensitivity and 70% specificity (68% accuracy) after PLS-
215 DA; 47% sensitivity and 75% specificity after PCA-LDA (51% accuracy); 83% sensitivity and
216 45% specificity after PCA-QDA (78% accuracy); 100% sensitivity and 90% specificity (98%
217 accuracy) after PCA-SVM. For plasma samples that were deposited on low-E slides the results
218 were: 65% sensitivity and 65% specificity (65% accuracy) after PLS-DA; 46% sensitivity and
219 85% specificity (52% accuracy) after PCA-LDA; 96% sensitivity and 15% specificity (84%
220 accuracy) after PCA-QDA; 100% sensitivity and 85% specificity (98% accuracy) after PCA-
221 SVM (Table 1). All PCA-based models for plasma samples were built with 10 PCs, accounting
222 99.11% and 96.84% of explained variance for Al and low-E substrates, respectively. PLS-DA
223 was built with 10 LVs, accounting 98.97% and 95.28% of explained variance for Al and low-
224 E substrates, respectively.

225 After applying classification algorithms for the blood serum samples, the results using
226 Al foil as a substrate were: 82% sensitivity and 75% specificity (81% accuracy) after PLS-DA;
227 90% sensitivity and 40% specificity (81% accuracy) after PCA-LDA; 94% sensitivity and 50%
228 specificity (86% accuracy) after PCA-QDA; 100% sensitivity and 70% specificity (94%
229 accuracy) after PCA-SVM. When using serum samples on low-E slides the results were: 78%
230 sensitivity and 90% specificity (80% accuracy) after PLS-DA; 63% sensitivity and 50%
231 specificity (61% accuracy) after PCA-LDA; 97% sensitivity and 20% specificity (83%
232 accuracy) after PCA-QDA; 100% sensitivity and 75% specificity (95% accuracy) after PCA-
233 SVM (Table 2). All PCA-based models for serum samples were built with 10 PCs, accounting

234 for 98.78% and 97.50% of explained variance for Al and low-E substrates, respectively. PLS-
235 DA was built with 10 LVs, accounting for 98.43% and 90.24% of explained variance for Al
236 and low-E substrates, respectively.

237 Overall, PCA-SVM was found to provide optimal results for both plasma and serum
238 samples regardless of the substrate that was used (Fig. 1 and 2). This was due to the fact that
239 PCA-SVM can create a more complex decision boundary between the classes, classifying even
240 non-linearly separable data^{33,35}. In addition, SVM creates large margins of separation between
241 the classes, which provides more stability to the classifier. In this sense, small disturbances or
242 noises do not cause misclassification³⁵. Standard deviation (SD) was higher for Al foil in
243 comparison to low-E slides (Fig. 1 and 2). This improved the Al foil classification models as
244 more sources of variation were contemplated during model construction, thus creating well-
245 distributed boundary functions and increasing the robustness of the classification. The SD in
246 the training set decreases the degree of overfitting and provides better predictive capacity³⁷.
247 The PCA-SVM cost function and optimization to estimate RBF parameters are shown in Fig.
248 3, where the red 'X' mark represents the optimal value. This optimization was performed in
249 order to avoid overfitting and to ensure classification stability. Fig. 4 shows the reference and
250 predicted class labels (1 for control; and 2 for cancer) using PCA-SVM with the samples from
251 the test set; if the yellow (predicted) and blue (reference) lines are superposed, then the values
252 are equal (*i.e.*, no misclassification). For all substrates and type of samples (plasma and serum),
253 there was no misclassification in the cancer set, reflecting the 100% sensitivity of PCA-SVM
254 models. A degree of misclassification was observed in the control set, particularly when using
255 serum samples. More specifically, specificity was higher in Al foil (90%) in contrast to low-E
256 (85%); this has provided the slightly higher accuracy in Al foil (98.5%) in contrast to low-E
257 (97.7%), in the plasma dataset. This can be seen in Fig. 4A and 4B as there are two and three
258 misclassified spectra, respectively ("continuous" peaks represent more than one spectrum). The

259 specificity differed slightly in the serum dataset too, when Al foil (70%) and low-E (75%) were
260 used. This contributed to the slightly lower accuracy in Al foil (94.5%) in contrast to low-E
261 (95.5%). In this case, there were six misclassified spectra for Al foil and five for low-E slides.
262 Although both PCA-LDA and PCA-QDA were regularized to correct classes having different
263 sizes (prior probability term in eq. 1 and 2), the number of errors is larger on the smaller class
264 (healthy control) due to the influence of the unequal class sizes to the classifiers. To summarise,
265 Al foil has been seen to perform better than low-E in the plasma dataset, while in the serum
266 dataset it achieved slightly lower specificity, but still high enough and comparable to low-E.

267 PCA-SVM models (Fig. 5) have different loadings profiles according to the type of
268 sample and substrate. The loadings are dependent on the nature of the dataset used for the PCA
269 model and they can differ depending on the input. Even though the same sample type is used,
270 the change of the substrate has subsequently changed the spectral profile as well. Any variation
271 above the instrumental noise can cause variation in the loading profiles. For instance,
272 differentiation was also observed at specific spectral peaks between Al foil and low-E
273 substrates (Fig. S1). Even though some spectral regions were visually similar, the reasoning of
274 using multivariate analysis is to overcome visual interpretation which can be inaccurate.
275 Therefore, a statistical *t*-test (95% confidence level) has been performed to calculate p-values
276 for each spectral point between Al foil and low-E as well as between plasma and serum. The
277 results showed that many wavenumbers were statistically significant ($p < 0.05$, 95% confidence
278 level) irrespectively of the visual similarities (Fig. S2). Additionally, the fact that PC1
279 accounted for low values of explained variance (70.09% for plasma-Al; 38.98% for plasma
280 low-E; 69.48% for serum-Al; and 28.69% for serum low-E) due to the high complexity of the
281 biological dataset, makes the loadings interpretation even harder.

282 Using aluminium as substrate, larger coefficients were found between ~ 1000 - 1150 cm^{-1}
283 ¹ for both plasma and serum samples, indicating possible glycogen and phosphate absorptions;

284 between ~ 1400 - 1480 cm^{-1} , indicating possible stretching vibrations of COO^- groups in fatty
285 acids and amino acids; at $\sim 1504\text{ cm}^{-1}$ for serum, signalling amide II absorption; and at ~ 1744
286 cm^{-1} for plasma, indicating C=O stretching of lipids³⁸. Using low-E slides as substrate for
287 plasma samples, larger coefficients were found at $\sim 1628\text{ cm}^{-1}$ (amide I), $\sim 1655\text{ cm}^{-1}$ (amide I)
288 and $\sim 1744\text{ cm}^{-1}$ (C=O stretching of lipids); whereas for serum samples, the coefficients were
289 greater at $\sim 1504\text{ cm}^{-1}$ (amide II), $\sim 1620\text{ cm}^{-1}$ (base carbonyl stretching and ring breathing mode
290 in nucleic acids) and 1655 cm^{-1} (amide I)³⁸. Such absorptions are known for signalling
291 biological changes using mid-IR spectroscopy¹⁹.

292 The classification accuracies achieved for the segregation between endometrial cancer
293 patients and controls are remarkably high (~ 95 - 98%), suggesting that blood-based ATR-FTIR
294 spectroscopy could potentially be an accurate and objective diagnostic tool for endometrial
295 cancer. Investigation of a panel of multiple biomolecules could be the reason for the achieved
296 accuracies. Several molecular biomarkers have been suggested over the years, such as
297 carcinoembryonic antigen (CEA), cancer antigen 125 (CA125), cancer antigen 15-3 (CA15-3),
298 immunosuppressive acidic protein (IAP), human epididymis protein-4 (HE4), apolipoprotein-
299 1 (ApoA-1), prealbumin (TTR) and transferrin (TF); a combination of CA125 and HE4 has
300 also been implied to improve diagnosis and classification of the disease³⁹⁻⁴². However, the
301 resulting sensitivities and specificities of the above-mentioned biomarkers are low, rendering
302 them clinically unusable. Therefore, spectroscopic methods are ideal, as they can
303 simultaneously extract information from a range of molecules. Another possible rationale
304 behind the diagnostic results could be the existence of circulating tumour DNA (ctDNA)
305 fragments in the bloodstream of cancer patients, which would make them distinct from the
306 normal population^{43,44}. Nowadays, ctDNA is increasingly investigated and is considered to be
307 useful as a biomarker for malignancy cases⁴⁵. Nevertheless, for an accurate and specific
308 biomarker detection, vibrational spectroscopy would need to be complemented with other

309 techniques as well, or maybe make use of labels or antibodies that would be molecule specific.
310 IR spectroscopy alone indicates some molecular fragments which are indicative of
311 biomolecules, such as proteins, lipids or carbohydrates. However, each spectral peak may
312 'hide' more than one molecules and thus, it is not preferred to assign specific biomarkers to
313 specific peaks.

314 In this study, plasma samples resulted in slightly higher diagnostic accuracies (~98%)
315 in contrast to serum samples (~95%). Current studies are unclear on whether serum or plasma
316 is a better source for ctDNA ⁴⁴. However, plasma has been previously found superior and the
317 specificity obtained using serum has been related to a higher concentration of normal cell-free
318 DNA (cfDNA), produced by the lysis of white blood cells during clotting ^{46, 47}. This could
319 potentially justify the lower classification rates found when using serum.

320 Careful consideration of the substrate, onto which the biological sample is placed, is
321 critical in order to collect reproducible and high-quality spectra. When comparing the
322 classification results coming from Al foil and low-E slides (Fig. 1 and 2), it is clear that Al foil
323 not only achieved equally high results with low-E but, in the plasma dataset, it even provided
324 slightly higher sensitivities and specificities (Fig. 1). Previous work has indicated that Al foil
325 generates no background noise, leaving the quality of the biological spectra unaffected; our
326 study used a larger number of subjects, which was needed to verify these preliminary results
327 and also study the impact on the sensitivity and specificity. Studies have also demonstrated the
328 enhancement of the IR signal in ATR mode when the sample is deposited onto metal surfaces
329 creating a similar effect to surface enhanced Raman spectroscopy (SERS), which has been
330 given the name surface enhanced IR absorption spectroscopy (SEIRAS). Molecules on metal
331 surfaces show 10-100 times stronger signal than without the metal ⁴⁸⁻⁵¹ and on the basis of this
332 we have hypothesized that Al foil slides may also promote this effect. However, this requires
333 further and more detailed investigation that will be the focus of a future study. The economic

334 cost of low-E slides has been estimated before and is not extremely high, especially when
335 compared with substrates like CaF₂ and Au-coated slides ²¹. Nonetheless, when it becomes a
336 matter of routine use, in a clinical setting for instance, the annual cost becomes considerably
337 high and this could render biospectroscopy prohibitive for translation into clinical practice. The
338 fact that Al foil slides are suitable for both IR and Raman studies is also an important advantage
339 as it would ease clinical implementation. The results of our study have shown that Al foil
340 slides could make an ideal, cost-effective substrate for biomedical studies employing
341 vibrational spectroscopy.

342 **Conclusion**

343 To summarise, biospectroscopy could potentially be used as a screening tool for endometrial
344 cancer in postmenopausal women as it provides exceptionally high sensitivities and
345 specificities with a simple blood test. This could automatically enable a large number of women
346 to be assessed on a daily basis. Using disposable, low-cost and, at the same time, high-
347 performance substrates would allow for universal studies with thousands of participants; this
348 would probably also generate an interest for multi-centre studies which could further validate
349 the pre-analytical, analytical and post-analytical phases of biospectroscopy.

350 **Acknowledgements**

351 MP would like to thank Rosemere Cancer Foundation for funding and Giouli Paraskevaidou
352 for providing effective feedback. CLMM would like to acknowledge CAPES/Doutorado Pleno
353 no Exterior/No. 88881.128982/2016-01 for financial support. KMGL acknowledges CNPq
354 (grant 305962/2014-0) for financial support.

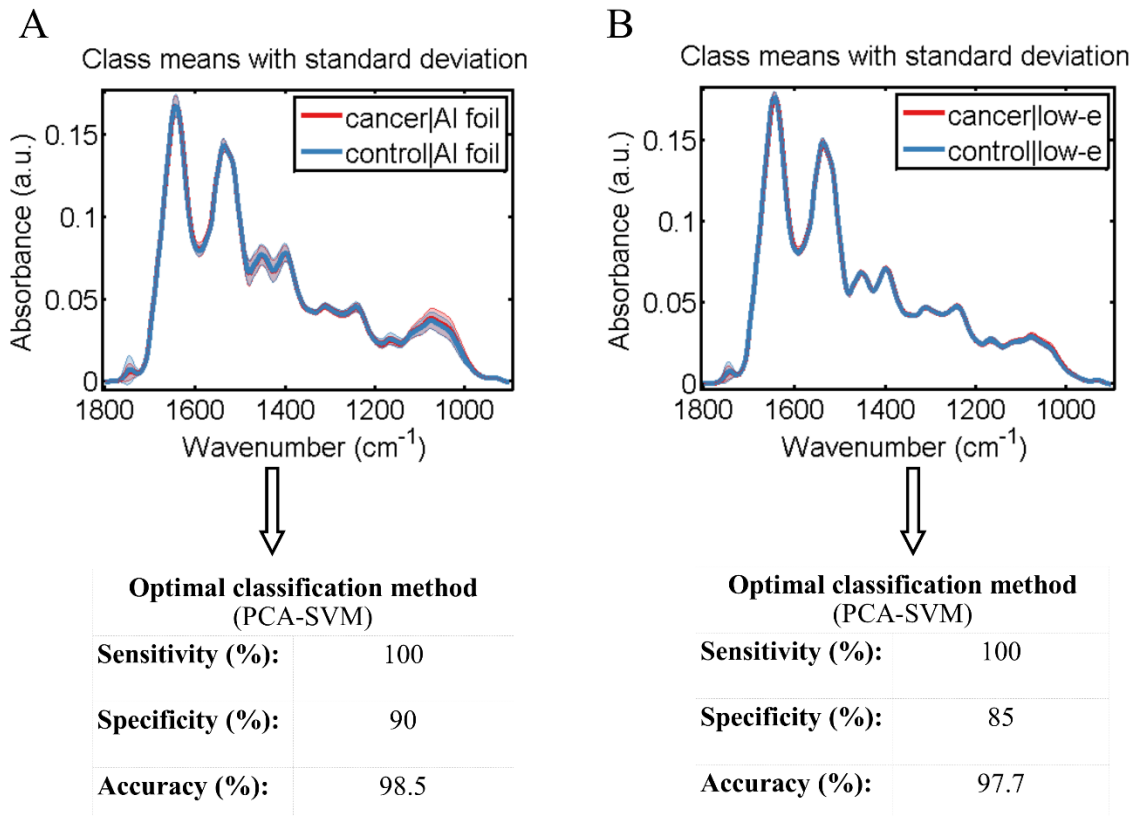
355

356

357

358 **Figures**

PLASMA

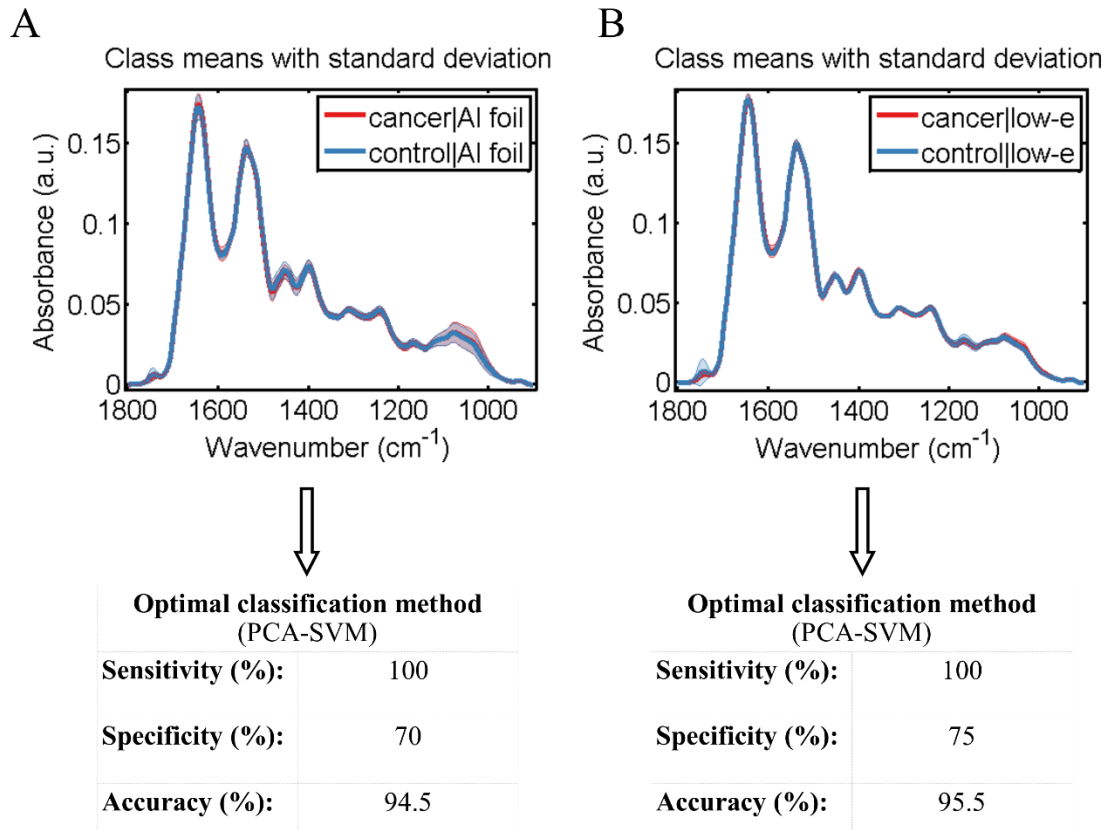


359

360 **Figure 1: Pre-processed spectra of plasma comparing endometrial cancer (n=70) with**
361 **controls (n=15).** (A) Endometrial cancer versus healthy controls; samples were analysed on
362 aluminium (Al) foil. Sensitivity and specificity were 100% and 90%, respectively. (B)
363 Endometrial cancer versus healthy controls; samples analysed on low-E slides. Sensitivity and
364 specificity were 100% and 85%, respectively.

365

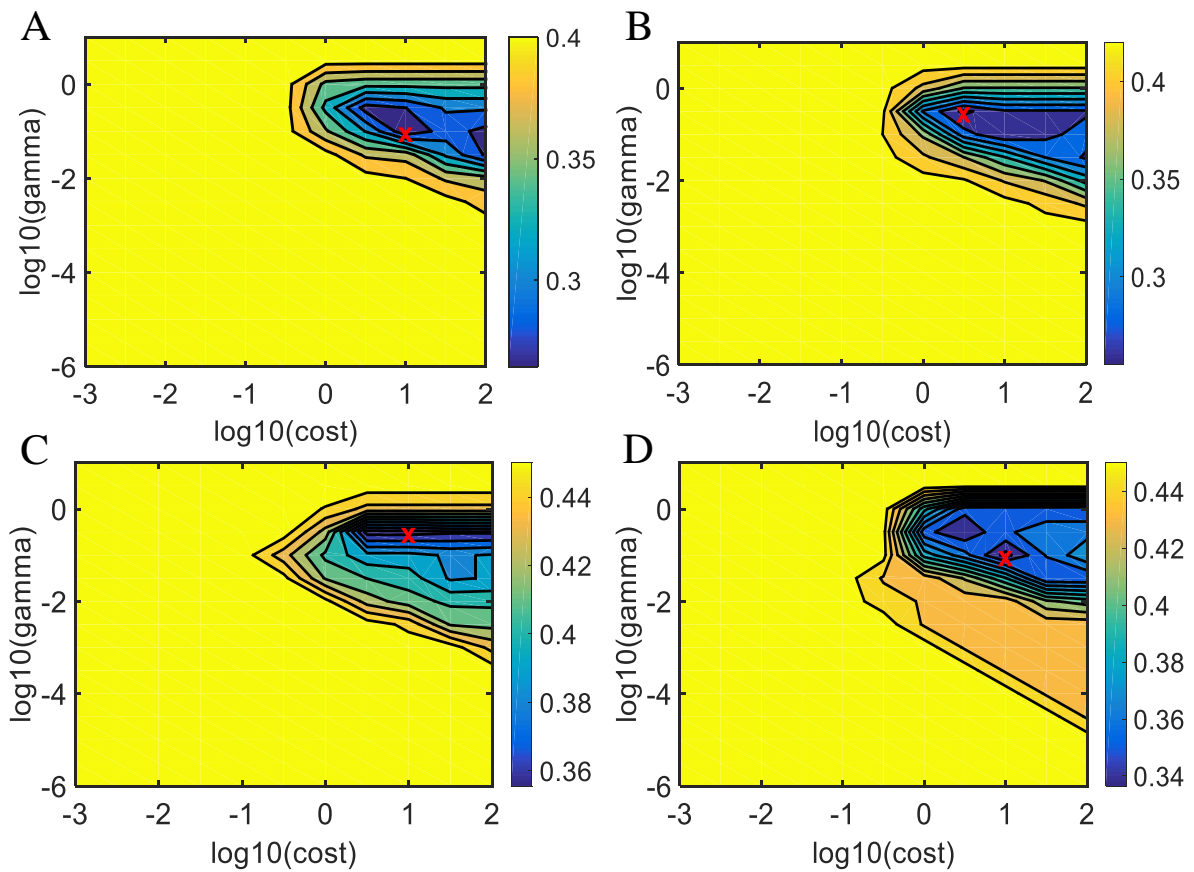
SERUM



366

367 **Figure 2: Pre-processed spectra of serum comparing endometrial cancer (n=60) with**
368 **controls (n=15).** (A) Endometrial cancer versus healthy controls; samples were analysed on
369 aluminium (Al) foil. Sensitivity and specificity were 100% and 70%, respectively. (B)
370 Endometrial cancer versus healthy controls; samples analysed on low-E slides. Sensitivity and
371 specificity were 100% and 75%, respectively.

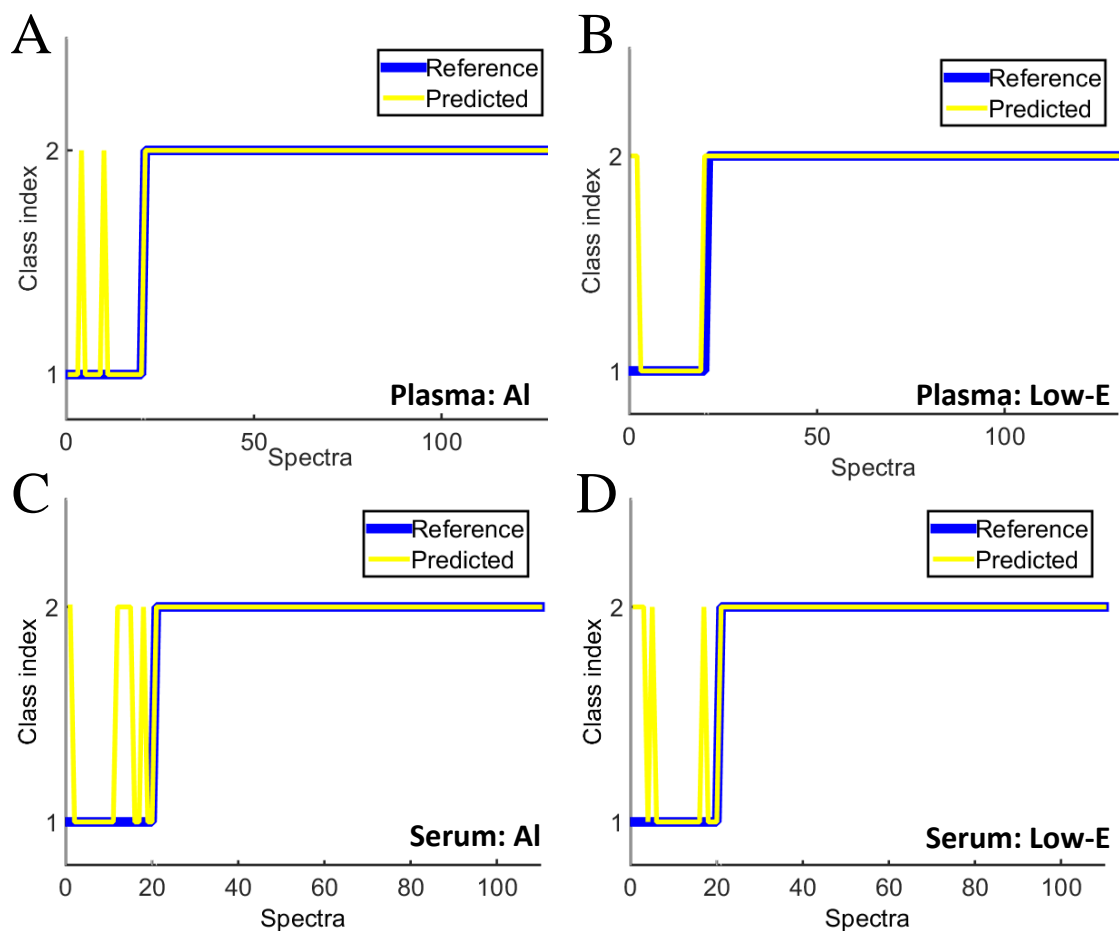
372



373

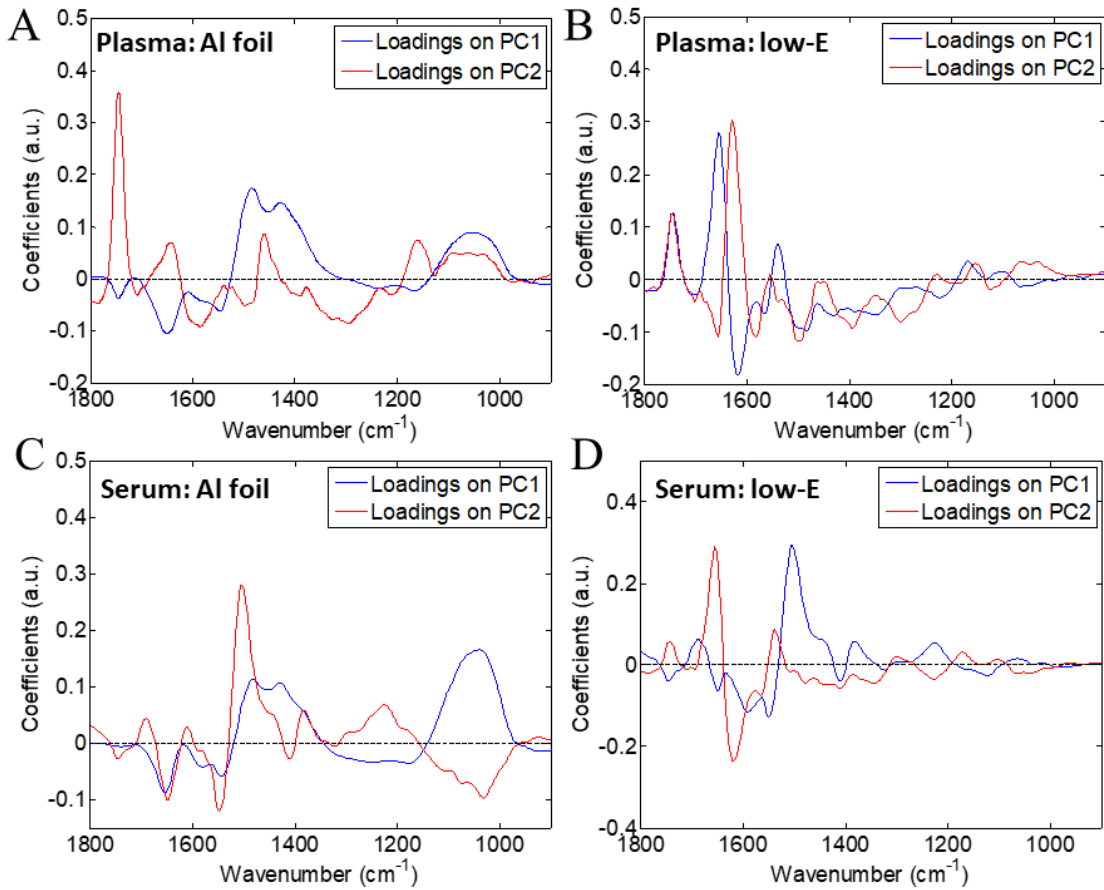
374 **Figure 3: PCA-SVM cost function and radial basis function (RBF) parameter**
 375 **optimization.** (A) Plasma samples with aluminium (Al) foil as a substrate. (B) Plasma samples
 376 with low-E slides as a substrate. (C) Serum samples with aluminium (Al) foil as substrate. (D)
 377 Serum samples with low-E slides as substrate. Gamma: RBF parameter (γ). Colour bar:
 378 misclassification rate using cross-validation.

379



380

381 **Figure 4: Reference and predicted class labels using PCA-SVM in the test set.** (A) Plasma
 382 samples with aluminium (Al) foil as a substrate; sensitivity was 100% and specificity 90% (two
 383 misclassified spectra). (B) Plasma samples with low-E slides as a substrate; 100% sensitivity
 384 and 85% specificity (three misclassified spectra). (C) Serum samples with aluminium (Al) foil
 385 as substrate; 100% sensitivity and 70% specificity (six misclassified spectra). (D) Serum
 386 samples with low-E slides as substrate; 100% sensitivity and 75% specificity (five
 387 misclassified spectra). Class 1 = control; and class 2 = cancer.



388

389 **Figure 5: Loading plots generated after PCA analysis.** (A) Loadings on PC1, PC2 for
 390 plasma samples deposited on aluminium (Al) foil slides. (B) Loadings on PC1, PC2 for plasma
 391 samples deposited on low-E slides. (C) Loadings on PC1, PC2 for serum samples deposited on
 392 aluminium (Al) foil slides. (D) Loadings on PC1, PC2 for serum samples deposited on low-E
 393 slides.

394

395

396

397

398

399

400

401

402

403

404

405 **Tables**

406

407 **Table 1:** Classification algorithms applied after the analysis of blood plasma samples.
 408 Results for both substrates, aluminium foil and low-E slide, are shown below.

409 **Correct classification rate (%):**

	Training (%)	Validation (%)	Test (%)
Aluminium foil			
PLS-DA	69.1	64.5	68.5
PCA-LDA	67.8	65.0	51.5
PCA-QDA	85.2	80.0	77.7
PCA-SVM	99.0	93.3	98.5
Low-E			
PLS-DA	71.1	71.8	65.4
PCA-LDA	62.7	54.2	52.3
PCA-QDA	85.2	82.5	83.8
PCA-SVM	99.8	97.5	97.7

410

411 **Quality parameters (%):**

	Accuracy (%)	Sensitivity (%)	Specificity (%)
Aluminium foil			
PLS-DA	68.5	68.2	70.0
PCA-LDA	51.5	47.3	75.0
PCA-QDA	77.7	83.6	45.0
PCA-SVM	98.5	100	90.0
Low-E			
PLS-DA	65.4	65.5	65.0
PCA-LDA	52.3	46.4	85.0
PCA-QDA	83.8	96.4	15.0
PCA-SVM	97.7	100	85.0

412

413

414 **Table 2:** Classification algorithms applied after the analysis of blood serum samples.
 415 Results for both substrates, aluminium foil and low-E slide, are shown below.

416 **SERUM**

417 **Correct classification rate (%):**

	Training (%)	Validation (%)	Test (%)
Aluminium foil			
PLS-DA	80.0	79.1	80.9
PCA-LDA	72.1	79.1	80.9
PCA-QDA	84.3	79.1	86.4
PCA-SVM	98.3	93.6	94.5
Low-E			
PLS-DA	85.7	71.8	80.0
PCA-LDA	70.2	65.5	60.9
PCA-QDA	84.2	88.2	82.7
PCA-SVM	99.1	98.2	95.5

418

419 **Quality parameters (%):**

	Accuracy (%)	Sensitivity (%)	Specificity (%)
Aluminium foil			
PLS-DA	80.9	82.2	75.0
PCA-LDA	80.9	90.0	40.0
PCA-QDA	86.4	94.4	50.0
PCA-SVM	94.5	100	70.0
Low-E			
PLS-DA	80.0	77.8	90.0
PCA-LDA	60.9	63.3	50.0
PCA-QDA	82.7	96.7	20.0
PCA-SVM	95.5	100	75.0

420

421

422

423

424 **References**

425 [1] G. Graça, A. S. Moreira, A. J. V. Correia, B. J. Goodfellow, A. S. Barros, I. F. Duarte, I. M. Carreira, E.
 426 Galhano, C. Pita, M. d. C. Almeida, A. M. Gil *Anal Chim Acta*. **2013**, *764*, 24-31.
 427 [2] A. Khoshmanesh, M. W. A. Dixon, S. Kenny, L. Tilley, D. McNaughton, B. R. Wood *Anal Chem*.
 428 **2014**, *86*, 4379-4386.
 429 [3] S. E. Taylor, K. T. Cheung, I. I. Patel, J. Trevisan, H. F. Stringfellow, K. M. Ashton, N. J. Wood, P. J.
 430 Keating, P. L. Martin-Hirsch, F. L. Martin *Br J Cancer*. **2011**, *104*, 790-797.
 431 [4] E. Staniszevska, K. Malek, M. Baranska *Spectrochim Acta Mol Biomol Spectrosc*. **2014**, *118*, 981-
 432 986.
 433 [5] K. Gajjar, L. D. Heppenstall, W. Pang, K. M. Ashton, J. Trevisan, I. I. Patel, V. Llabjani, H. F.
 434 Stringfellow, P. L. Martin-Hirsch, T. Dawson, F. L. Martin *Anal Methods*. **2013**, *5*, 89-102.
 435 [6] N. Stone, R. Baker, K. Rogers, A. W. Parker, P. Matousek *Analyst*. **2007**, *132*, 899-905.
 436 [7] F. M. Lyng, E. Ó. Faoláin, J. Conroy, A. Meade, P. Knief, B. Duffy, M. Hunter, J. Byrne, P. Kelehan,
 437 H. Byrne *Exp Mol Pathol*. **2007**, *82*, 121-129.
 438 [8] J. R. Hands, G. Clemens, R. Stables, K. Ashton, A. Brodbelt, C. Davis, T. P. Dawson, M. D.
 439 Jenkinson, R. W. Lea, C. Walker, M. J. Baker *J Neurooncol*. **2016**, *127*, 463-472.
 440 [9] P. Carmona, M. Molina, M. Calero, F. Bermejo-Pareja, P. Martinez-Martin, A. Toledano *J*
 441 *Alzheimers Dis*. **2013**, *34*, 911-920.
 442 [10] M. Paraskevaidi, C. L. M. Morais, K. M. G. Lima, J. S. Snowden, J. A. Saxon, A. M. T. Richardson,
 443 M. Jones, D. M. A. Mann, D. Allsop, P. L. Martin-Hirsch, F. L. Martin *Proc Natl Acad Sci USA*. **2017**,
 444 *114*, E7929-e7938.
 445 [11] A. Shapiro, O. N. Gofrit, G. Pizov, J. K. Cohen, J. Maier *Eur Urol*. **2011**, *59*, 106-112.
 446 [12] K. Maheedhar, R. A. Bhat, R. Malini, N. B. Prathima, P. Keerthi, P. Kushtagi, C. M. Krishna
 447 *Photomed Laser Surg*. **2008**, *26*, 83-90.
 448 [13] K. Guze, H. C. Pawluk, M. Short, H. Zeng, J. Lorch, C. Norris, S. Sonis *Head Neck*. **2015**, *37*, 511-
 449 517.
 450 [14] A. Sahu, K. Dalal, S. Naglot, P. Aggarwal, C. M. Krishna *Plos One*. **2013**, *8*, e78921.
 451 [15] J. Zheng, L. He *Compr Rev Food Sci Food Saf*. **2014**, *13*, 317-328.
 452 [16] L. T. Kerr, H. J. Byrne, B. M. Hennelly *Anal Methods*. **2015**, *7*, 5041-5052.
 453 [17] B. D. Beier, A. J. Berger *Analyst*. **2009**, *134*, 1198-1202.
 454 [18] H. J. Byrne, P. Knief, M. E. Keating, F. Bonnier *Chem Soc Rev*. **2016**, *45*, 1865-1878.
 455 [19] M. J. Baker, J. Trevisan, P. Bassan, R. Bhargava, H. J. Butler, K. M. Dorling, P. R. Fielden, S. W.
 456 Fogarty, N. J. Fullwood, K. A. Heys, C. Hughes, P. Lasch, P. L. Martin-Hirsch, B. Obinaju, G. D.
 457 Sockalingum, J. Sulé-Suso, R. J. Strong, M. J. Walsh, B. R. Wood, P. Gardner, F. L. Martin *Nat Protoc*.
 458 **2014**, *9*, 1771-1791.
 459 [20] H. J. Butler, L. Ashton, B. Bird, G. Cinque, K. Curtis, J. Dorney, K. Esmonde-White, N. J. Fullwood,
 460 B. Gardner, P. L. Martin-Hirsch, M. J. Walsh, M. R. McAinsh, N. Stone, F. L. Martin *Nat. Protoc*. **2016**,
 461 *11*, 664-687.
 462 [21] L. Cui, H. J. Butler, P. L. Martin-Hirsch, F. L. Martin *Anal Methods*. **2016**, *8*, 481-487.
 463 [22] M. Paraskevaidi, P. L. Martin-Hirsch and F. L. Martin, in: *Characterization Tools for Nanoscience*
 464 *and Nanomaterials*, Springer, **2017** (in press).
 465 [23] K. Gajjar, J. Trevisan, G. Owens, P. J. Keating, N. J. Wood, H. F. Stringfellow, P. L. Martin-Hirsch,
 466 F. L. Martin *Analyst*. **2013**, *138*, 3917-3926.
 467 [24] K. Wehbe, J. Filik, M. D. Frogley, G. Cinque *Anal Bioanal Chem*. **2013**, *405*, 1311-1324.
 468 [25] K. A. Heys, R. F. Shore, M. G. Pereira, F. L. Martin *Environ Sci Technol*. **2017**, *51*, 8672-8681.
 469 [26] Cancer Research UK, [http://www.cancerresearchuk.org/health-professional/cancer-](http://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-cancer-type/uterine-cancer/incidence#heading-Ten)
 470 [statistics/statistics-by-cancer-type/uterine-cancer/incidence#heading-Ten](http://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-cancer-type/uterine-cancer/incidence#heading-Ten), Accessed February 2018.
 471 [27] M. J. Baker, S. R. Hussain, L. Lovergne, V. Untereiner, C. Hughes, R. A. Lukaszewski, G. Thieffn, G.
 472 D. Sockalingum *Chem Soc Rev*. **2016**, *45*, 1803-1818.

- 473 [28] J. Trevisan, P. P. Angelov, P. L. Carmichael, A. D. Scott, F. L. Martin *Analyst*. **2012**, *137*, 3202-
474 3215.
- 475 [29] R. W. Kennard, L. A. Stone *Technometrics*. **1969**, *11*, 137-148.
- 476 [30] D. B. Hibbert *Pure Appl Chem*. **2016**, *88*, 407-443.
- 477 [31] R. G. Brereton, G. R. Lloyd *J Chemometrics*. **2014**, *28*, 213-225.
- 478 [32] R. Bro, A. K. Smilde *Anal Methods*. **2014**, *6*, 2812-2831.
- 479 [33] S. J. Dixon, R. G. Brereton *Chemom Intell Lab Syst*. **2009**, *95*, 1-17.
- 480 [34] F. S. Costa, P. Silva, C. Lelis, R. Theodoro, T. Arantes, K. de Lima *Anal Methods*. **2017**.
- 481 [35] P. de Boves Harrington *Anal Chem*. **2015**, *87*, 11065-11071.
- 482 [36] C. L. Morais, F. S. Costa, K. M. Lima *Anal Methods*. **2017**, *9*, 2964-2970.
- 483 [37] H. Martens, P. Geladi, *Multivariate calibration*, Wiley Online Library, **1989**.
- 484 [38] Z. Movasaghi, S. Rehman, I. U. Rehman *Appl Spectrosc Rev*. **2007**, *42*, 493-541.
- 485 [39] L. Zanotti, E. Bignotti, S. Calza, E. Bandiera, G. Ruggeri, C. Galli, G. Tognon, M. Ragnoli, C.
486 Romani, R. A. Tassi *Clin Chem Lab Med*. **2012**, *50*, 2189-2198.
- 487 [40] P. B. Panici, G. Scambia, G. Baiocchi, L. Perrone, S. Greggi, F. Battaglia, S. Mancuso *Gynecol*
488 *Obstet Invest*. **1989**, *27*, 208-212.
- 489 [41] Y. Ueda, T. Enomoto, T. Kimura, T. Miyatake, K. Yoshino, M. Fujita, T. Kimura *Cancers*. **2010**, *2*,
490 1312.
- 491 [42] G. Farias-Eisner, F. Su, T. Robbins, J. Kotlerman, S. Reddy, R. Farias-Eisner *Am J Obstet Gynecol*.
492 **2010**, *202*, 73.e71-75.
- 493 [43] C. Bettegowda, M. Sausen, R. J. Leary, I. Kinde, Y. Wang, N. Agrawal, B. R. Bartlett, H. Wang, B.
494 Luber, R. M. Alani *Sci Transl Med*. **2014**, *6*, 224ra224-224ra224.
- 495 [44] Y. I. Elshimali, H. Khaddour, M. Sarkissyan, Y. Wu, J. V. Vadgama *Int J Mol Sci*. **2013**, *14*, 18925-
496 18958.
- 497 [45] L. A. Diaz Jr, A. Bardelli *J Clin Oncol*. **2014**, *32*, 579-586.
- 498 [46] A. Vallée, M. Marcq, A. Bizieux, C. E. Kouri, H. Lacroix, J. Bennouna, J.-Y. Douillard, M. G. Denis
499 *Lung Cancer*. **2013**, *82*, 373-374.
- 500 [47] S. El Messaoudi, F. Rolet, F. Mouliere, A. R. Thierry *Clinica Chimica Acta*. **2013**, *424*, 222-230.
- 501 [48] S. E. Glassford, B. Byrne, S. G. Kazarian *Biochim Biophys Acta - Proteins Proteom*. **2013**, *1834*,
502 2849-2858.
- 503 [49] K. Ataka, S. T. Stripp, J. Heberle *Biochim Biophys Acta – Biomembranes*. **2013**, *1828*, 2283-2293.
- 504 [50] J.-Y. Xu, T.-W. Chen, W.-J. Bao, K. Wang, X.-H. Xia *Langmuir*. **2012**, *28*, 17564-17570.
- 505 [51] R. Adato, H. Altug *Nat Commun*. **2013**, *4*, 2154.