



**QUEEN'S
UNIVERSITY
BELFAST**

Incremental model learning for spectroscopy-based food analysis

Diaz, K., Georgouli, K., Koidis, A., & Martinez Del Rincon, J. (2017). Incremental model learning for spectroscopy-based food analysis. *Chemometrics and Intelligent Laboratory Systems*, 167, 123-131. <https://doi.org/10.1016/j.chemolab.2017.06.002>

Published in:
Chemometrics and Intelligent Laboratory Systems

Document Version:
Peer reviewed version

Queen's University Belfast - Research Portal:
[Link to publication record in Queen's University Belfast Research Portal](#)

Publisher rights

© 2017 Elsevier B. V. This manuscript version is made available under the CC-BY-NC-ND 4.0 license <http://creativecommons.org/licenses/by-nc-nd/4.0/>, which permits distribution and reproduction for noncommercial purposes, provided the author and source are cited.

General rights

Copyright for the publications made accessible via the Queen's University Belfast Research Portal is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The Research Portal is Queen's institutional repository that provides access to Queen's research output. Every effort has been made to ensure that content in the Research Portal does not infringe any person's rights, or applicable UK laws. If you discover content in the Research Portal that you believe breaches copyright or violates any law, please contact openaccess@qub.ac.uk.

Incremental model learning for spectroscopy-based food analysis

Katerine Diaz^a, Konstantia Georgouli^b, Anastasios Koidis^b, Jesus Martinez Del Rincon^{c,*}

^a*Computer Vision Centre, Universidad Autonoma de Barcelona, Spain*

^b*Institute for Global Food Security, Queens University Belfast, UK*

^c*Institute of Electronics, Communications and Information Technology, Queens University Belfast, UK*

Abstract

In this paper we propose the use of incremental learning for creating and improving multivariate analysis models in the field of chemometrics of spectral data. As main advantages, our proposed incremental subspace-based learning allows creating models faster, progressively improving previously created models and sharing them between laboratories and institutions without requiring transferring or disclosing individual spectra samples. In particular, our approach allows to improve the generalization and adaptability of previously generated models with a few new spectral samples to be applicable to real-world situations. The potential of our approach is demonstrated using vegetable oil type identification based on spectroscopic data as case study. Results show how incremental models maintain the accuracy of batch learning methodologies while reducing their computational cost and handicaps.

Keywords: Incremental model learning, IGDCV technique, Subspace based learning, Identification, Vegetable oils, FT-IR spectroscopy

2010 MSC: 00-01, 99-00

*Corresponding author

Email addresses: `k Diaz@cvc.uab.es` (Katerine Diaz), `kgeorgouli01@qub.ac.uk` (Konstantia Georgouli), `t.koidis@qub.ac.uk` (Anastasios Koidis), `j.martinez-del-rincon@qub.ac.uk` (Jesus Martinez Del Rincon)

1. Introduction

In the last decade the use of chemometrics in food analysis is steadily growing. This is caused because the output of most analytical methods is nowadays multivariate data matrices (spectroscopic, chromatographic/mass spectrometry data, isotopic, sensorial, etc) which cannot be manually analysed and demand appropriate chemometric analysis in order to process and capture the most important and relevant information in the data. Selection of multivariate methods (e.g. classification methods) however is often limited to a set of well known standard methods (e.g. PLS-DA and SIMCA classification methods) and researchers are faced with some persisting problem with the chemometric models that they generate [1].

Among these problems that must be addressed, the generality of the models created to new conditions is the most important one. While extensive research has been done to create models under controlled conditions, for a small problem or dataset, the applicability of those models in real world -e.g. in food testing in the food industry or in routine analysis in a regulated testing laboratory- is very scarce. This is due to the overfitting of the model to the calibration set when only one instrument, one analytical laboratory or, in general, one set of assumptions are taken into consideration to create the models. Thus, when these models are tested in other slightly different conditions, they report much lower performances than the expected one. Recalibrating or recreating similar models to work in those situations may be an extremely arduous task, with a similar time and effort scale to the design, and tuning of the first model.

To avoid a full recalibration, model updating and calibration transfer techniques have been proposed to cover the transfer of multivariate classification models between different spectrometers [2, 3], temperatures [3, 4], harvesting seasons [4] and even different geographical regions [5]. Calibration transfer techniques [2] allow mapping the new spectra to the primary model spectra domain by calculating a transformation matrix from one domain to the other. [Different calibration transfer techniques have been recently explored in chemical sensor](#)

31 [arrays to overcome inherent sensor variability \[6, 7, 8\]](#). Only a small set of sam-
32 ples are required to be measured in both the primary and secondary conditions.
33 However, in many applications it is not realistic that exactly the same sample
34 can be measured, e.g. the same food sample from two different geographical lo-
35 cations. More interesting are methods based on model updating by augmenting
36 sample spectra from a new condition. While many sample would normally be
37 required to span to the new conditions [4], which amounts to a full recalibra-
38 tion, approaches based on Tikhonov regularisation (TR) [3, 5] only needs a few
39 samples to update the model. As disadvantage, TR still requires access to the
40 initial samples to recompute the updated model, with the consequent compu-
41 tational cost of involving all samples in the optimisation, and its performance
42 heavily relies on a meta-parameter that controls the balance between the initial
43 model and the augmented samples, and which can only be tuned empirically.
44 [Finally, some recursive learning approaches \[9, 10\] propose a framework where](#)
45 [both incremental and decremental stages are used to improve the initial model.](#)
46 [However, to fully exploit their potential and being able to remove old samples,](#)
47 [access to the initial samples is also required.](#)

48 Moreover, new samples are analysed on a routine basis and new data is ge-
49 nerated including cases when new component classes are needed to be created
50 (in authentication/adulteration studies, in traceability, proximate analysis pre-
51 diction etc). As a result, existing and validated models may stop being useful
52 and/or applicable. It is then necessary to retrain them. However, this requires
53 access to the original samples, which may be lost or unavailable. Similarly, if
54 an external laboratory, or other third party such as a company or an institution
55 wishes to improve an existing model, the access to the original samples may
56 be tricky or impossible, with privacy or confidentiality issues playing a role. In
57 all these previously described situations, it is clear that evolving a chemometric
58 model may be a better solution than recreating or retraining it as a full new
59 batch. This will only require access to the existing models and the new samples.
60 It will also be a more efficient manner to store the information, reducing the
61 memory and physical space required and it can potentially decrease the time to

62 create an improved model.

63 While incremental learning has been used and proposed in other fields [11,
64 12, 13, 9, 10], its intrinsic advantages have been scarcely exploited in the field
65 of food analysis and chemometrics [14, 15, 16, 17, 18, 19]. Bhattacharyya et
66 al. [14, 15] applied neural networks for identification of seven different black
67 tea classes. Their incremental approach allow to add new classes of black tea
68 to the original set. In Tudu et al. (2009) [16], the same researchers applied
69 incremental fuzzy logic to the black tea identification. Cernuda et al. [17, 18, 19]
70 proposed a flexible fuzzy inference system for the monitor of the concentration
71 of sulphuric acid (H_2SO_4), sodium sulfate (Na_2SO_4) and zinc sulfate (ZnSO_4) in
72 viscose production and in the melamine resin production process, which allows
73 online adaptation of parameters and structural changes in the model. However,
74 techniques based on neural networks and fuzzy logic are scarcely used in food
75 science, reducing the impact of these incremental approaches, and they require
76 huge amounts of calibration samples to generate the calibration models, which
77 is unlikely for most food analysis scenarios.

78 In this paper we aim to extend the use of incremental learning in the field
79 of food analysis and chemometrics. Among the variety of incremental learning
80 techniques, we have chosen subspace based learning as the family of machine
81 learning to apply due to their proved ability to evolve online [13], the ability
82 to generate efficient models using a reduced number of calibration samples,
83 and the extensive use of some of the basic subspace based methods such as
84 Principal Component Analysis (PCA), and Soft independent modelling of class
85 analogies (SIMCA)- in food science [20, 21], both for exploratory analysis [22]
86 and classification [23, 24, 25]. Thus, the present work introduces the use of an
87 incremental subspace based learning technique, called Incremental Generalized
88 Discriminative Common Vectors (IGDCV), which allows efficiently adding new
89 data samples and classes to a knowledge base. In this way, our methodology
90 is able to update the model to the new scenario without recalculating the full
91 projection or accessing the previously processed calibration data, while retaining
92 the previously acquired knowledge. Our approach is evaluated using vegetable

93 oil type identification [22, 26, 27, 28] as case study and results are compared
94 against a non incremental learning technique, i.e. an equivalent batch method.
95 Three different incremental scenarios are tested in this application area: when
96 new samples are available to improve the model, when new classes must be
97 identified by the model, and when new instruments are used in the identification
98 process.

99 **2. Incremental Learning Framework**

100 Several incremental feature extraction based on linear subspace methods
101 have been proposed and used on many practical applications. Among them, we
102 find the Incremental approaches of the PCA [29], Linear Discriminant Analysis
103 (LDA) [30] and DCV [31]. While PCA-based incremental approaches are simple
104 and versatile, they are not optimal for discrimination and classification purposes
105 since no class information is used to obtain principal components which may lead
106 to unsuited subspaces. On the contrary, LDA is a supervised technique which
107 makes use of the class information to obtain the most discriminative space by
108 maximizing the distance between classes while minimizing the distance between
109 the samples within the same class. However, LDA-based approaches cannot
110 be applied when the dimension of the sample space is larger than the number
111 of samples in the calibration set, since the within-class scatter matrix will be
112 singular. This problem is known as the *Small Sample Size* SSS problem [32],
113 and it is frequent in spectroscopic and chromatographic application, where the
114 number of variables per sample is in the order of thousands while the total
115 number of samples used for calibration rarely goes above the hundreds [22].

116 Among the approaches that have been proposed to solve the SSS problem,
117 the Generalized Discriminative Common Vectors (GDCV) has been proved [13]
118 to provide discriminative subspaces for classification regardless of the SSS as-
119 sumption. GDCV is a variation of LDA [33, 34] which introduces the idea
120 of approximate extended null and reduced range subspaces of the within-class
121 scatter matrix. Given the good performance of GDCV batch approaches, we

122 proposed the use of Incremental GDCV [13] as the base of our online learning
 123 framework for food analysis, where new information is added while retaining
 124 the previously acquired knowledge, without accessing the previously processed
 125 calibration data.

126 2.1. IGDCV

127 Formally, let the calibration set X be composed of c classes, where every
 128 class j has m_j samples. The total number of samples in the calibration set is
 129 $M = \sum_{j=1}^c m_j$. Let x_j^i be a d -dimensional column vector which denotes the i^{th}
 130 sample from the j^{th} class. The within-class scatter matrix, S_w^X , is defined as,

$$S_w^X = \sum_{j=1}^c \sum_{i=1}^{m_j} (x_j^i - \bar{x}_j)(x_j^i - \bar{x}_j)^T = X_c X_c^T \quad (1)$$

131

132 where \bar{x}_j is the average of the samples in the j^{th} class, and the centered
 133 data matrix, X_c consists of column vectors $(x_j^i - \bar{x}_j)$ for all $j = 1 \dots c$ and
 134 $i = 1 \dots m_j$.

135 The extension of the null space of S_w^X (which implies restricting the corre-
 136 sponding range space) is done from the [Eigen-Value Decomposition \(EVD\)](#) of
 137 S_w^X .

$$EVD(S_w^X) : U_r \Lambda_r U_r^T \quad (2)$$

138 where $U_r \in \mathbb{R}^{d \times r}$ are the eigenvectors associated to the nonzero eigenvalues
 139 Λ_r . The scattering added to the null space can be measured as the trace
 140 $tr(U_\alpha^T S_w^X U_\alpha)$. This quantity is up to $tr(S_w^X)$ when no directions are remo-
 141 ved, $U_\alpha = U_r$, and decreases as more and more important directions disappear
 142 from U_r . Consequently, the scattering preserved after a projection, U_α , can be
 143 written as follows

$$\alpha = 1 - \frac{tr(U_\alpha^T S_w^X U_\alpha)}{tr(S_w^X)} \quad (3)$$

144

145 The projection basis fulfilling the above conditions for a given value of α
 146 can be obtained through U_r , such that r is reassigned. The α value is the main
 147 parameter of GDCV, which can be tuned by using cross-validation over the
 148 training set. The GDCV method can be summarized as

- 149 1. Obtain U_α such that $S_w^X = U_r \Lambda_r U_r^T$. where Λ_α contains the smallest
 150 eigenvalues in Λ_r and $tr(\Lambda_\alpha) = \alpha \cdot tr(\Lambda_r)$
- 151 2. Project class means as $x_{gcv}^j = \bar{x}_j - U_\alpha U_\alpha^T \bar{x}_j$. These are the so-called
 152 generalized common vectors of each class.
- 153 3. Define $X^{com} = [x_{gcv}^1 \dots x_{gcv}^c]$ and let X_c^{com} be its centered version with
 154 regard to the mean, $\bar{x}_{gcv} = \frac{1}{c} \sum_{j=1}^c x_{gcv}^j$
- 155 4. Obtain the projection $W \in \mathfrak{R}^{d \times (c-1)}$ such that $tr(W^T X_c^{com} X_c^{com T} W)$ is
 156 maximum.

Thus, by using the projection matrix W , any sample x_i can be projected in
 the discriminative subspace $gdcv$ for an easier classification, according to

$$x_i^{gdcv} = W^T \cdot (x_i - \bar{x}_{gcv}) \quad (4)$$

157 In an incremental learning scenario, once an initial dataset X has been used
 158 to obtain U_α, Λ_α and W , a new set of sample Y will be available in a later stage
 159 to improve the learned projection. This new set of data Y may be composed
 160 of a single sample or several ones that may belong to pre existing classes or
 161 to fully new categories. In the general case, the new dataset Y consists of n_j
 162 samples from each class, resulting in a total of $N = \sum_{j=1}^c n_j$ new samples to be
 163 considered in the learning process.

164 The IGDCV method allows obtaining $U'_\alpha, \Lambda'_\alpha$ and W' corresponding to the
 165 new complete dataset, $[X \ Y]$, without having to reapply the GDCV algorithm
 166 to $[X \ Y]$. Instead, they will be obtained incrementally by adding the effect of
 167 new data, Y , into the previous solution corresponding to X , such that

$$S_w^Z = S_w^X + Y_c Y_c^T + A A^T \quad (5)$$

168 where Y_c consists of column vectors $(y_j^i - \bar{y}_j)$ for all $j = 1 \dots c$ and $i = 1 \dots n_j$.
 169 $A = [a_1 \dots a_c]$ is a matrix whose columns are the c weighted average differences
 170 given by

$$a_j = \sqrt{\frac{m_j n_j}{m_j + n_j}} (\bar{x}_j - \bar{y}_j), \quad j = 1 \dots c \quad (6)$$

171

The IGDCV algorithm is summarized as

Algorithm 1. *IGDCV Algorithm*

Parameter: α , $0 < \alpha \leq 1$

Input: $Y \in \mathbb{R}^{d \times N}$, $\{n_j\}_{j=1}^c$, $N = \sum_{j=1}^c n_j$

From previous iteration: $U_\alpha \in \mathbb{R}^{d \times r}$, $\Lambda_\alpha \in \mathbb{R}^{r \times r}$, $\bar{x}_j \in \mathbb{R}^d$, $\{m_j\}_{j=1}^c$

Output: $U'_\alpha \in \mathbb{R}^{d \times r'}$, $\Lambda'_\alpha \in \mathbb{R}^{r' \times r'}$, $\bar{x}'_j \in \mathbb{R}^d$, $\{m'_j\}_{j=1}^c$

Method:

1. Compute \bar{y}_j , Y_c , A
 2. Compute $V = \text{orth}([Y_c \ A] - U_\alpha U_\alpha^T [Y_c \ A])$
 3. Build $M_\alpha = \begin{bmatrix} \Lambda_\alpha & 0 \\ 0 & 0 \end{bmatrix} + [U_\alpha \ V]^T Y_c Y_c^T [U_\alpha \ V] + [U_\alpha \ V]^T A A^T [U_\alpha \ V]$
 4. Compute R and Λ' by eigendecomposing M_α
 5. Compute $\beta = (1 - \alpha) \frac{\text{tr}(\Lambda_r)}{\text{tr}(\Lambda')} + \alpha$
 6. Split R and Λ' in R_β and Λ_β by β
 7. Let $U'_\alpha = [U_\alpha \ V] R_\beta$ and $\Lambda'_\alpha = \Lambda_\beta$
 8. Update: $m'_j = m_j + n_j$, $j = 1, \dots, c$
 $\bar{x}'_j = (m_j \bar{x}_j + n_j \bar{y}_j) / m'_j$
 9. Project class means as $x_{gcv}^j = \bar{x}'_j - U'_\alpha U'^T_\alpha \bar{x}'_j$.
-

Figure 1: Incremental Generalized Discriminant Common Vector (IGDCV) algorithm.

172

173 If some of the data vectors in Y correspond to new classes which are not
 174 present in X , the expressions of the IGDCV algorithm are valid by extending
 175 the value of c and setting $m_j = 0$ in X for all new classes. Both if m_j or n_j are
 176 zero for any class j , the corresponding mean is undefined and the corresponding
 177 column in A , a_j , should be set to zero. If all data vectors in Y correspond to

178 new classes, then the whole matrix A is the zero matrix and can be removed
179 from all expressions.

180 The overall cost of the IGDCV is dominated by the cost of step 7 in Fig. 1,
181 $O(dr'^2)$ where r' is the expected rank of the range space preserved that heavily
182 depends on the parameter α .

183 *2.2. Classification*

184 After applying IGDCV, samples can be projected into a discriminative sub-
185 space where meaningful conclusion can be extracted, if used as exploratory
186 analysis, or an automatic classification can be achieved. The performed super-
187 vised learning ensures that the different classes to be recognized are as separate
188 as possible, making the classification problem very simple, since the complexity
189 of the problem has been moved to the previous stage. Thus, we have coupled
190 our incremental subspace learning with a k-Nearest Neighbors (kNN) classifier
191 in order to provide this functionality.

192 Two advantages are derived from the use of the KNN classifier. First, given
193 its simplicity, the performance of IGDCV will be directly reflected in the experi-
194 ments, which could otherwise be masked by a more complex classifier. Second,
195 since no calibration is required in KNN, the online learning of the classifier will
196 be automatic when the subspace is updated.

197 **3. Case of study**

198 In order to evaluate the potential and advantages of incremental learning,
199 the problem of identifying vegetable oil types using spectroscopic analysis was
200 chosen as case study. This is a relevant case of study [22, 26, 27, 28, 25] brought
201 into attention due to European Regulation 1169/2011, which requires producers
202 of foods that contain refined vegetable oil blends to label the oil types. In this
203 context, deliberate or accidental errors in the label are common, leading to con-
204 sumer misinformation [21], so automatic identification and verification of the
205 provided information is required. From an analytical point of view, testing an

206 unknown vegetable oil to identify its origin and composition is a very difficult
207 task [27, 35], but where spectroscopy -such as FTIR- and subspace-based met-
208 hods have demonstrated their capabilities [22]. However, the performed single
209 lab validation [22, 36, 37] of current approaches, which is common but undesi-
210 rable in the field, indicates that the real performance in realistic conditions may
211 be far from the reported accuracy.

212 4. Materials and methods

213 4.1. Samples

214 A data set of 630 vegetable oil samples was used in this study. Two different
215 classification problems are considered with respect to the number of classes.
216 Calibration models were developed for 6 classes and 12 classes of vegetable oils
217 (see Table 1). For the 6-class problem, the classes to be predicted are labelled
218 as PO: palm oil /palm stearin /palm olein, RS: sunflower /rapeseed oil and
219 their mixtures, PKOC: palm kernel oil /coconut oil and binary mixtures of the
220 above. For the 12-class problem, the classes are PO: palm oil /palm stearin
221 /palm olein, RO: rapeseed, SO: sunflower, PKO: palm kernel, CCO: coconut,
222 and all the binary combinations of the above oils. The 12-class model provides
223 more resolution because it clearly distinguishes between the individual botanical
224 origins, and it is therefore a more complex problem, while the 6-class model
225 groups some origins together according to their similarities. This allows us to
226 test our approach at to different levels of complexity, which are related to the
227 expected level of resolution to be detected.

228 4.2. FT-IR spectral acquisition

229 The acquisition of most FT-IR spectra samples was performed using a Nico-
230 let iS5 Thermo spectrometer (Thermo Fisher Scientific, Dublin, Ireland) equip-
231 ped with a DTGS KBr detector and a KBr beam splitter. Spectra were acquired
232 from 4000 to 550 cm^{-1} co-adding 32 interferograms at 4 cm^{-1} resolution with

Table 1: Different oil types for the 6 and 12-class problem.

	Class	Samples		Class	Samples
1	PO	104	1	PO	104
2	RS	114	2	RO	36
3	PKOC	36	3	SO	23
4	RS-PKOC	83	4	PKO	26
5	RS-PO	181	5	CCO	10
6	PO-PKOC	112	6	RO-PO	98
			7	SO-PO	83
			8	RO-PKO	51
			9	SO-PKO	32
			10	RO-SO	55
			11	PO-PKO	66
			12	PO-CCO	46

233 a diamond attenuated total reflectance (iD5 ATR) accessory. Absorbance va-
 234 lues were recorded at each spectrum point. The final sample spectrum was the
 235 average of three replicates with initial 7157 data points.

236 Through an interlaboratory experiment sixteen extra FT-IR instruments
 237 were used to acquire several extra oil spectra, as shown in Table 2. A total of
 238 nine samples including pure oils and oil admixtures were prepared in our lab and
 239 sent to each of the instruments participated to collect spectra representatives
 240 of most classes with all instruments. The acquisition parameters have been
 241 harmonized so that they are compatible with every FT-IR instrument. Linear
 242 interpolation was applied to spectra from different instruments in order to get
 243 the desirable number of variables.

244 4.3. Data pre-treatment

245 The resulting FT-IR spectral profiles underwent some typical preprocessing
 246 techniques in order to reduce or remove any random or systematic variation in
 247 the data [38]. Five steps are involved in this phase. Specifically, prior to the ap-
 248 plication of the multivariate models, Standard Normal Variate (SNV) [39], first
 249 order derivative [40], S-Golay filter [41] [polynomial order=2,frame size=9] and

Table 2: Instruments for the interlaboratory experiment. (Note: N/a - not available)

Id	Participant	FT-IR Instrument	Detector	Year	Samples
1	Our lab (Institute for global food security, QUB)	Thermo Fisher Scientific Nicolet iS5	DTGS	2012	486
2	Teagasc, Food Research Centre	Bio-Rad Excalibur FTS 3100	DTGS	2001	9
3	PerkinElmer Ltd	PerkinElmer Spectrum 2	DTGS	2012	9
4	PerkinElmer Ltd	PerkinElmer Frontier	DTGS	2013	9
5	Brennan and Co.	Bruker Alpha	DTGS	2013	9
6	Public Analyst Scientific Services	PerkinElmer Spectrum 100	LiTaO3	2007	9
7	LGC Limited (UK)	PerkinElmer Spectrum One	DTGS	2001	9
8	Premier Analytical Services (Premierfoods)	Bio-Rad Excalibur FTS300MX	DTGS	2002	9
9	Institute of Food Research (IFR)	Nicolet MagnaIR 860	DTGS	1998	9
10	Institute of Food Research (IFR)	Bio-Rad FTS6000	DTGS	1996	9
11	Institute of Food Research (IFR)	Thermo Fisher Scientific Nicolet iN10MX/iZ10	DTGS	2011	9
12	Shimadzu (Mason Technology)	Shimadzu IRA nity-1S	DLaTGS	n/a	9
13	Antech(IRE)	Thermo Fisher Scientific TruDefender FTX	DLaTGS	n/a	9
14	Agri-Food and Biosciences Institute (AFBI)	PerkinElmer Spectrum One	MIR TGS	n/a	9
15	Walloon Agricultural Research Centre (CRA-W)	Bruker Vertex 70	DLaTGS	2007	9
16	Walloon Agricultural Research Centre (CRA-W)	Bruker Vertex 70	DLaTGS	2012	9
17	Walloon Agricultural Research Centre (CRA-W)	Bruker Vertex 70	MCT	2012	9

250 Pareto scaling [42] were applied for removing the scatter, correcting the baseline,
251 smoothing the data points and scaling the data for preventing the dominance of
252 high absorbances respectively. At the end of this preprocessing procedure, the
253 irrelevant spectra area was cut out by selecting only the wavelengths between
254 654.23 and 1875.43 cm-1 and between 2520.02 and 3120.74 cm-1, corresponding
255 to relevant fatty acid involved in oil identification [22, 25]. In total, 3781 varia-
256 bles are resulted. All chemometric data preprocessing was performed by means
257 of in-house Matlab routines (The MathWorks Inc., USA).

258 5. Results

259 Using our case of study, three scenarios where the potential of incremental
260 learning is relevant will be tested. In the first scenario, an oil type identification
261 model is trained with a few calibration samples. After this initial calibration,
262 new samples for each of the oil types to identify become available and are added
263 to the model for improving the initial performance. In the second scenario, a
264 simple model is initially trained to distinguish between just two oil types, and
265 then extended to identify new oil types, up to 12. In the third scenario, the oil
266 type identification model created by a single lab and using a single spectroscopy
267 analyser is extended and enhanced to be effective when used in other laboratories
268 and instruments.

269 For comparison purposes, the batch version of IGDCV, batch GDCV, is
270 used as a baseline. By using the exact batch equivalent version, we ensure the
271 comparison is performed in the same conditions. The batch version requires to
272 recreate the model every time that several, or even one single sample is available
273 and added to the calibration set and therefore, access to the original samples
274 is always obliged. The aim is then to ensure the same or similar classification
275 performance to the batch method while reducing the computational time and
276 removing the requirement of having access to the original calibration samples
277 by the incremental approach. The α parameter was empirically optimised in
278 the range (0, 0.3] with steps of 0.01 for each scenario, so that the batch GDCV

279 provided the best accuracy result prior to any incremental step or addition of
280 any new data. Then, the same value of α is used for both GDCV and IGDCV
281 and keep constant over all the iterations. Thus, we aim to simulate a carefully
282 fined-tuned initial pre-existing model to be further evolved.

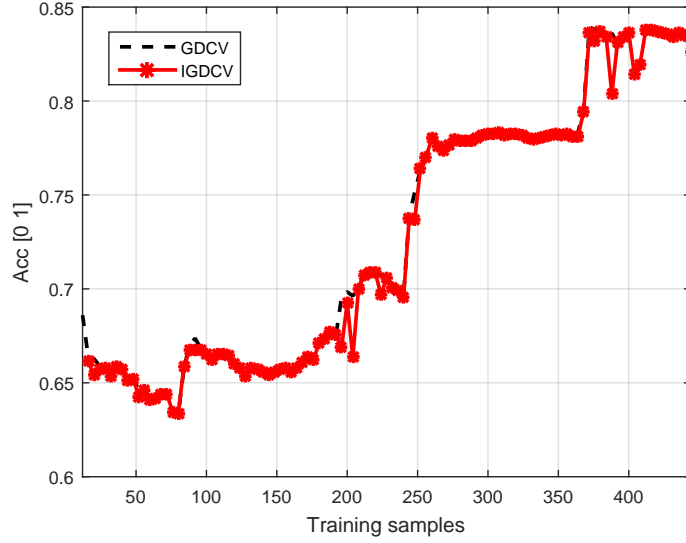
283 *5.1. First scenario: New samples*

284 In this experiment, we simulate a scenario where, for a given problem, an
285 initial dataset is captured and the corresponding model is created. Then, new
286 samples become available for calibration at different stages that can be used to
287 improve the initial model and its performance. To do so, the 6 classes dataset is
288 used. Cross validation is applied as evaluation protocol to avoid bias regarding
289 the chosen samples. Ten iterations are performed, each with a random 70/30
290 split, i.e. the dataset is divided in 70% for calibration and 30% for validation
291 in each iteration with no overlap between calibration and validation sets to
292 avoid bias in the results. Results are then averaged over the splits to generate
293 the final value. From the calibration samples, initially only 12 samples with
294 representatives of all classes are used to generate a model. Then, in incremental
295 step of 4 samples each, the model is evolved.

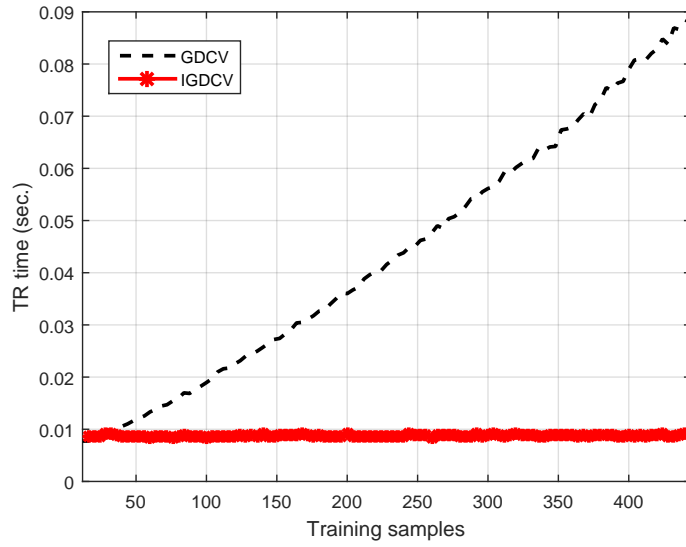
296
297 Fig. 2 shows the results of both incremental and batch methods, with the
298 preserved scattering parameter set to $\alpha = 0.13$. As expected, models perform
299 better as more calibration samples are available for learning from. Regarding
300 the incremental learning, it can be observed how the accuracy of the incremental
301 approach does not suffer, when compared with the batch algorithm, from not
302 having access to the initial samples but only to the previous model. Moreover,
303 when comparing the computational time required to generate the models (see
304 Fig. 2b), one can notice the great difference in efficiency of using an incremental
305 method regarding regenerating larger and larger models from scratch.

306 *5.2. Second scenario: New classes*

307 In this experiment, we simulate a scenario where a model has been created
308 for a simpler identification problem that is then extended to cope with a more



(a) Accuracy (ACC)



(b) Training (TR) time

Figure 2: Batch GDCV and incremental IGDCV methods regarding new samples. Scattering parameter $\alpha = 0.13$.

309 complex problem. In the initial model, only 2 different oils are expected to be
310 distinguished (Oil 1, 2) and this is incrementally evolved to identify more and
311 more classes up to the total of the 12 species.

312 Similarly to the previous scenario, cross validation is also used as evaluation
313 protocol, where 10 iterations are performed, each with a random 70/30 split,
314 i.e. the 70% of the samples from each class is used for calibration and 30% are
315 reserved for validation. Results are then averaged over the splits to generate the
316 final value and the dispersion bar. In each iteration step, all calibration samples
317 for a new class are added to the previous model.

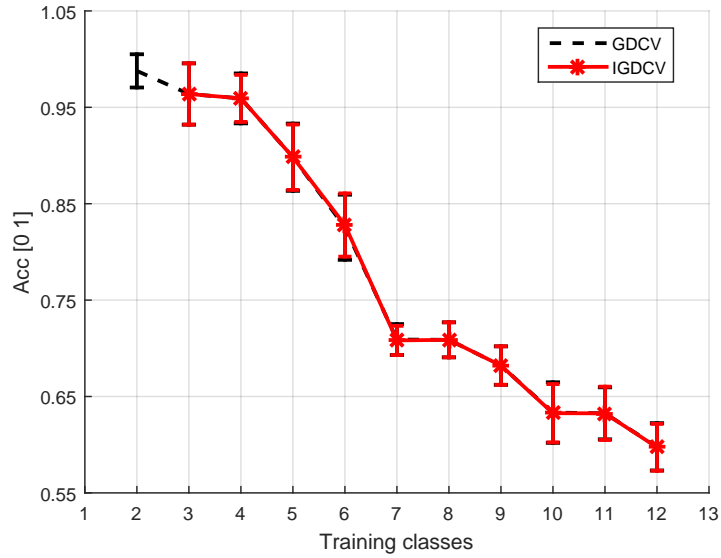
318

319 Fig. 3 shown the results of both incremental and batch method. As ex-
320 pected, the more classes must be identified, the more complex the problem and,
321 therefore, the accuracy decreases. Similarly to scenario 1, the potential of incre-
322 mental learning is stated again by conserving the accuracy of the batch approach
323 while reducing drastically the computational time and the access to the initial
324 samples.

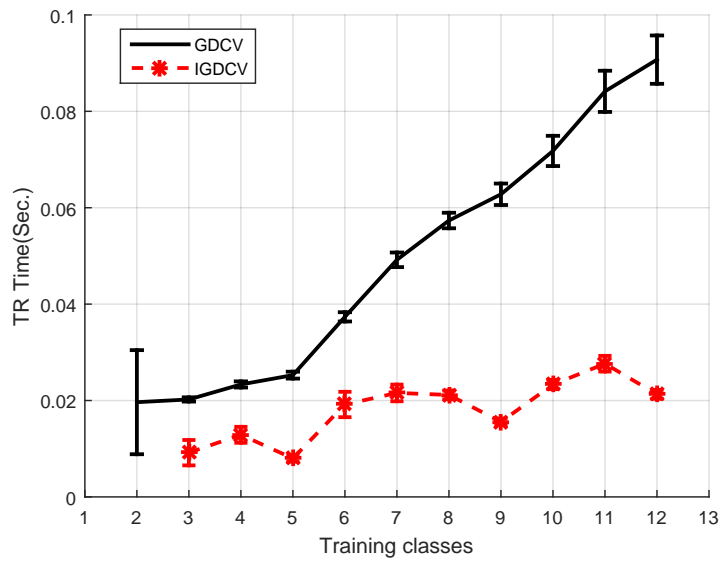
325 5.3. *Third scenario: New instruments*

326 In this scenario, we demonstrate the potential of the incremental learning to
327 generalise previously existing models so that they can then be used by others
328 laboratories using different instruments.

329 It has been shown that models created under controlled conditions, e.g. from
330 a single calibration set when only one instrument was used, perform poorly
331 when operating in real world conditions and report much lower performances
332 than what it is expected from them. This can be corroborated by generating a
333 model trained with 70% of the samples from instruments 1 (see Table 2). This
334 model is first tested with the remaining 30% of the samples belonging to the very
335 same instruments, and then tested with the samples from all other instruments.
336 Similarly to previous scenarios, cross validation is used as evaluation protocol,
337 where 10 random iterations are performed. Results for the 6 and 12 classes
338 problems are depicted in Table 3.



(a) Accuracy (ACC)



(b) Training (TR) time

Figure 3: Batch GDCV and incremental IGDCV methods regarding new classes. Scattering parameter $\alpha = 0.07$.

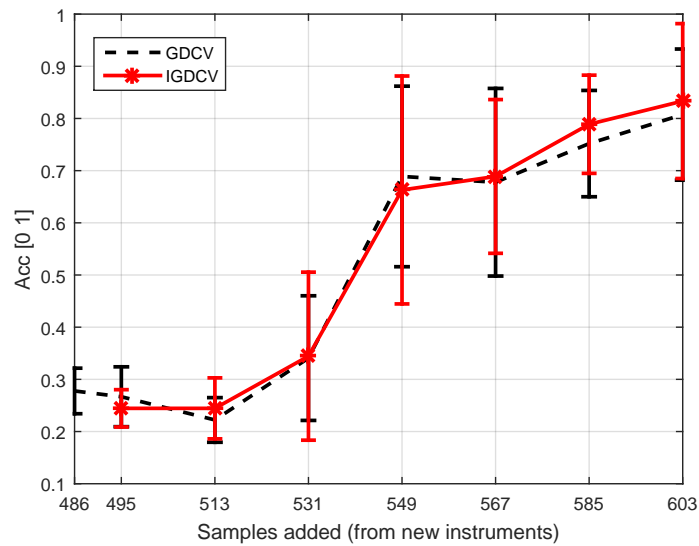
Table 3: Accuracy of GDCV model when using (2nd column) samples of the same instrument in the test set, and (3rd column) samples of different instruments in the test set to the instrument used in calibration.

Classes	Same Inst. in Test	New Inst. in Test
6	0.72 ± 0.04	0.28 ± 0.06
12	0.62 ± 0.03	0.14 ± 0.03

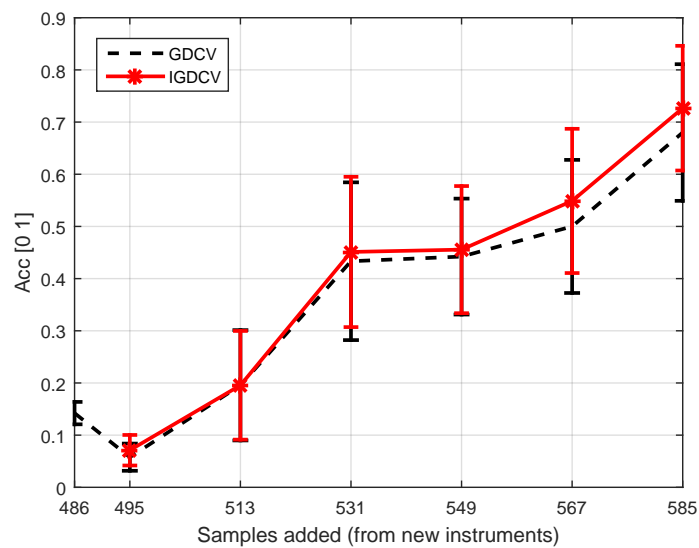
339 It can be noticed how an apparently good model, with reported accuracies
 340 60-70%, underperforms dramatically under more complicated environments or
 341 conditions. It is therefore clear the necessity of improving an existing model in
 342 order to operate more broadly.

343 We simulate this situation in this third scenario, where we evaluate the po-
 344 tential of incremental learning to improve the generality of a previously created
 345 model initially created in a single laboratory. An initial model is trained with
 346 all samples from a single instrument. Then, the samples of a new instrument
 347 are added in a first step to evolve the model, followed by incremental steps of
 348 all samples belonging to new 2 instruments in each step. Two experiments are
 349 performed, one where the model has to identify 6 classes and the other one with
 350 12 classes. Cross validation is used, repeating the experiment 10 times, where
 351 different instruments are randomly left out for the validation. In the 6 classes
 352 experiment, all samples from 3 different instruments are reserved for evaluating
 353 the system and up to 14 instruments are used in calibration. In the 12 classes
 354 experiment, all samples from 5 different instruments are kept for evaluating the
 355 system and up to 12 instruments are used in calibration.

356
 357 Fig. 4 shown the results of both incremental and batch method. It can
 358 be noticed how using more instruments and collaborating between different
 359 labs allows to radically improved the performance of a given method. Both 6
 360 and 12 class experiments behave similarly with slightly lower performance in
 361 the 12 classes due to the higher difficulty of the problem. We can see how



(a) 6 Classes



(b) 12 Classes

Figure 4: Accuracy (ACC) rate of the batch GDCV and the incremental IGDCV regarding new samples from new instruments. Scattering parameter $\alpha = 0.02$.

362 the incremental learning allows not only replicating the batch results but also
363 it improves them regarding computational time, Fig. 5. It is also important
364 to notice, how only a few samples from new instruments are needed (only 9
365 samples are available, see Table2) in our approach to improve significantly the
366 final accuracy.

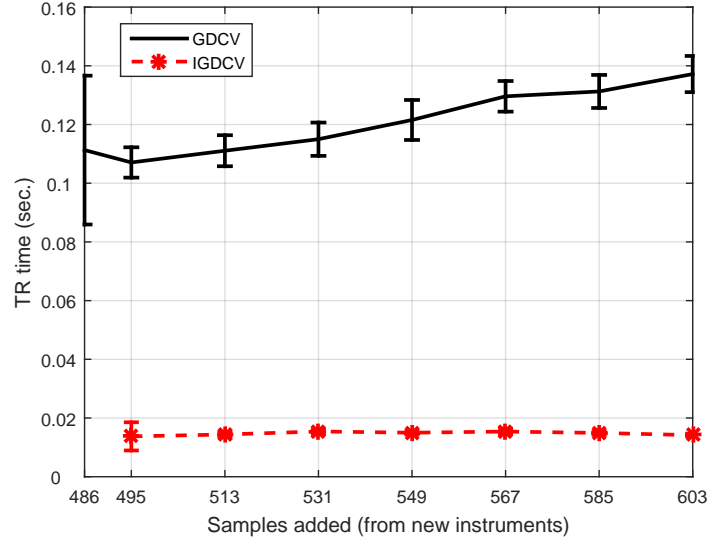
367

368 5.4. IGDCV as exploratory analysis tool

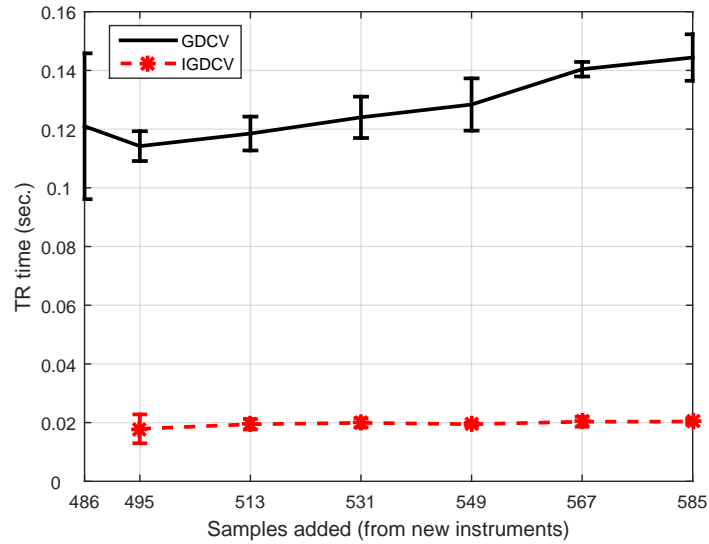
369 Apart from the benefits of using the IGDCV that were described earlier,
370 IGDCV can also be used as an exploratory analysis tool, similarly to PCA
371 [22]. In this regard, projecting the samples in the learned IGDCV can provide
372 valuable information regarding the complexity of the problem, the likelihood
373 of the model to accurately predict the correct answer and the quality of the
374 samples. Furthermore, *its incremental nature provides an extra functionality
375 not available in PCA, GDCV or other batch methods, since once a model is
376 created, a specific new sample(s) can be assessed in terms of its adequacy to be
377 included in the analysis and/or in the calibration set of the following iteration
378 of the model.*

379 Figure 6 shows the evolution of the model for the 6 class problem in the
380 first scenario, i.e. when samples are incrementally added. It can be observed
381 how, while in the first space it is not very clear what are pure or admixture oil
382 samples due to lack of data, this relationship is clearer the more online learning
383 iterations occurs and more relevant samples are added.

384 Figure 7 shows the evolution of the model for the 12 class problem in the
385 second scenario, i.e. when samples belonging to new classes are incrementally
386 added. It can be observed how the complexity of the problem grows: while in
387 the first space the 3 classes could be easily identified and separated, the space
388 is more cluttered when the number of classes increases. This visualization can
389 be used to decide which classes could not be resolved, and therefore should be
390 excluded, due to their similar properties which are translated in their overlap
391 in the space.

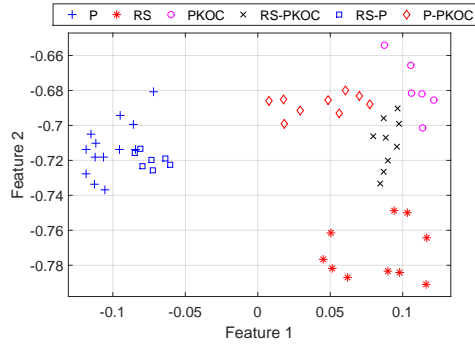


(a) 6 Classes

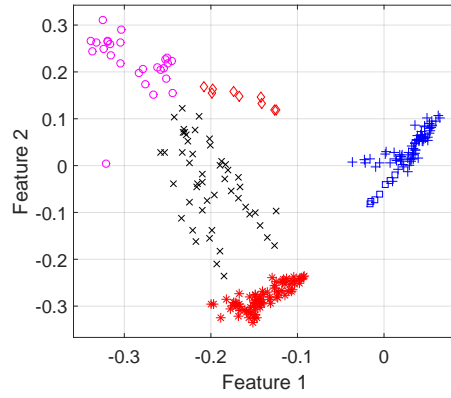


(b) 12 Classes

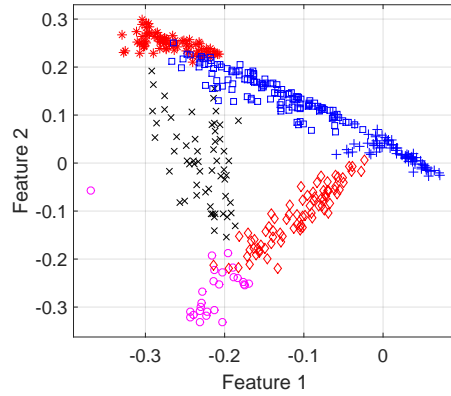
Figure 5: Training (TR) time of the batch GDCV and the incremental IGDCV regarding new samples from new instruments. Scattering parameter $\alpha = 0.02$.



(a) Initial model with 44 samples

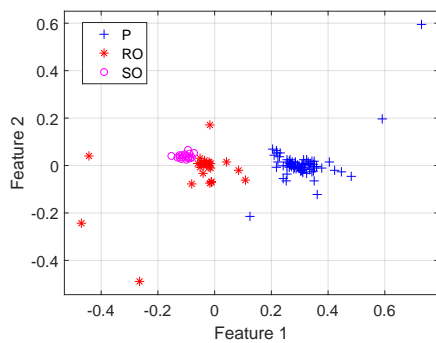


(b) After 200 samples have been incrementally learned

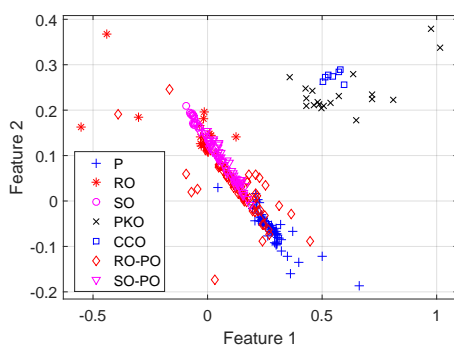


(c) After 400 samples have been incrementally learned

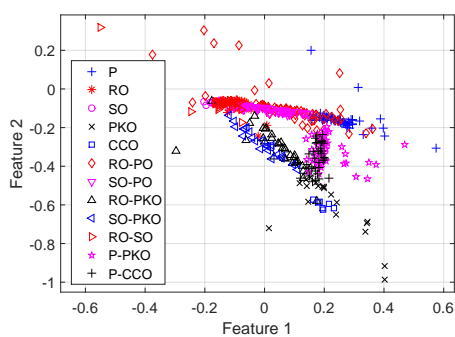
Figure 6: Samples projected into the two discriminant dimensions of the learned subspace, for the first scenario where samples are added incrementally (6-class problem)



(a) Initial model with 3 oil classes



(b) After 4 classes have been incrementally learned



(c) After 9 classes have been incrementally learned

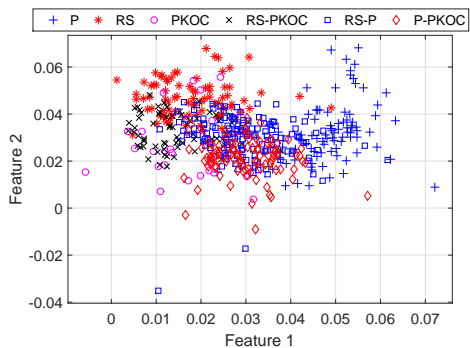
Figure 7: Samples projected into the two discriminant dimensions of the learned subspace, for the second scenario where classes are added incrementally (12-class problem)

392 Finally, Figure 8 shows the evolution of the model for the 6 class problem
393 in the third scenario, i.e. when samples belonging to new instruments are in-
394 crementally added. It can be seen how the initial model is clearly insufficient
395 to solve the problem and how adding more and more instruments seems a good
396 idea to improve discrimination between classes. It could also be used to decide
397 in which moment adding more instruments may not be convenient anymore,
398 since the subspace will not evolve further, as seen between the second and third
399 projections. Please notice how this visualization correlates with the quantitative
400 results in Figure 4, where accuracy improvement reduces after 3 iterations, i.e.
401 6 instruments.

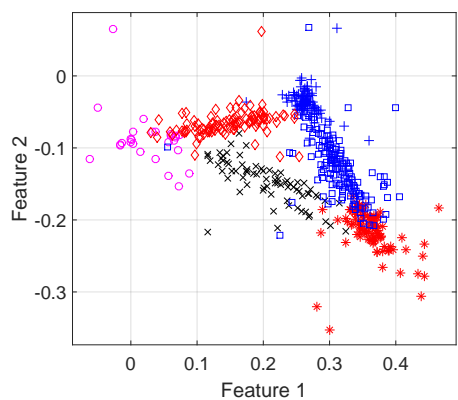
402 As can be seen, by using a incremental method for exploratory analysis,
403 relevant information is provided to food scientist such as the detection of errors
404 in the sample preparation or data generation, or the likelihood of an improved
405 model by using a new batch of samples. Furthermore, this experiments were
406 performed in a fraction of the time required by the batch method GDCV. Thus,
407 Figure 6.b) and c) were generated in 28% and 13% of the batch time respectively,
408 Figure 7.b) and c) in 40% and 22% of the batch time and Figure 8.b) and c) in
409 16% and 14% of the batch time.

410 6. Conclusion

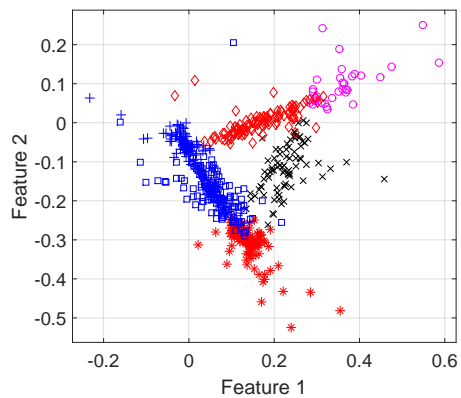
411 In this paper we apply the concept of incremental learning in food science
412 and proposed the use of a subspace based learning method, both in its incremen-
413 tal and batch method as a new chemometric analysis tool. GDCV and IGDCV
414 can be used as both classification and exploratory techniques, without some of
415 the constraints that PCA or LDA exhibits, such as requiring large number of
416 samples. The potential of incremental learning to improve and share models
417 between analytical laboratories using different acquisition equipment is demon-
418 strated through three different scenarios. By adding a very small number of
419 samples to a preexisting model, our approach allows improving significantly the
420 accuracy as well as to adapt the model to a new problem or scenario. The



(a) Initial model using samples from 2 instruments



(b) After samples from 6 instruments have been incrementally learned



(c) After samples from 10 instruments have been incrementally learned

Figure 8: Samples projected into the two discriminant dimensions of the learned subspace, for the third scenario where instruments are added incrementally (6-class problem)

421 IGCV incremental approach presented here has the advantage of maintaining
422 or improving the accuracy while reducing the computational and spatial cost,
423 and removing the hassle and privacy issues associated to share raw samples and
424 wasting time and effort reproducing the models and tuning the analytical tools.
425 As future work, we aim to extend our incremental subspace learning method to
426 other cases of studies in chemometrics as well as integrating IGDCV as part of
427 a new version of SIMCA. [We also aim to study the use of decremental learning](#)
428 [in chemometrics and add a decremental stage to our online learning framework.](#)

Acknowledgements

The authors would like to thank all the participants (research centres, public services and private food testing labs) that helped to perform the interlaboratory experiment. This research was supported with funding from The Department Learning and Employment Northern Ireland (DELNI) (PhD studentship block grant) and the Department of Environment, Food and Rural Affairs (DEFRA) of the UK (Grant no. FAO 157).

References

- [1] S. Lohumi, S. Lee, H. Lee, B.-K. Cho, A review of vibrational spectroscopic techniques for the detection of food authenticity and adulteration, *Trends in Food Science & Technology* 46 (2015) 85–98.
- [2] A. J. Myles, T. A. Zimmerman, S. D. Brown, Transfer of multivariate classification models between laboratory and process near-infrared spectrometers for the discrimination of green arabica and robusta coffee beans, *Applied spectroscopy* 60 (10) (2006) 1198–1203.
- [3] M. R. Kunz, J. H. Kalivas, E. Andries, Model updating for spectral calibration maintenance and transfer using 1-norm variants of tikhonov regularization, *Analytical chemistry* 82 (9) (2010) 3642–3649.

- [4] M. Golic, K. B. Walsh, Robustness of calibration models based on near infrared spectroscopy for the in-line grading of stonefruit for total soluble solids content, *Analytica Chimica Acta* 555 (2) (2006) 286–291.
- [5] M. R. Kunz, J. Ottaway, J. H. Kalivas, C. A. Georgiou, G. A. Mousdis, Updating a synchronous fluorescence spectroscopic virgin olive oil adulteration calibration to a new geographical region, *Journal of agricultural and food chemistry* 59 (4) (2011) 1051–1057.
- [6] Evaluation of calibration transfer strategies between metal oxide gas sensor arrays, *Procedia Engineering* 120 (2015) 261 – 264.
- [7] Calibration transfer and drift counteraction in chemical sensor arrays using direct standardization, *Sensors and Actuators B: Chemical* 236 (2016) 1044 – 1053.
- [8] L. Fernandez, S. Guney, A. Gutierrez-Galvez, S. Marco, Calibration transfer in temperature modulated gas sensor arrays, *Sensors and Actuators B: Chemical* 231 (2016) 276 – 284.
- [9] Y. Liu, N. Hu, H. Wang, P. Li, Soft chemical analyzer development using adaptive least-squares support vector regression with selective pruning and variable moving window size, *Industrial and Engineering Chemistry Research* 48 (12) (2009) 5731–5741.
- [10] Y. Liu, H. Wang, J. Yu, P. Li, Selective recursive kernel learning for on-line identification of nonlinear systems with narx form, *Journal of Process Control* 20 (2) (2010) 181–194.
- [11] D. Ross, J. Lim, R. Lin, M. Yang, Incremental learning for robust visual tracking, *IJCV* 77 (1–3) (2008) 125–141.
- [12] Y. Peng, S. Pang, G. Chen, A. Sarrafzadeh, T. Ban, D. Inoue, Chunk incremental idr/qr lda learning, in: *The 2013 International Joint Conference on Neural Networks (IJCNN)*, 2013, pp. 1–8.

- [13] K. Diaz-Chito, F. Ferri, W. Díaz-Villanueva, Incremental generalized discriminative common vectors for image classification, *IEEE Trans. Neural Networks and Learning Systems* 26 (8) (2015) 1761–1775.
- [14] N. Bhattacharyya, A. Metla, R. Bandyopadhyay, B. Tudu, A. Jana, Incremental pnn classifier for a versatile electronic nose, in: *3rd International Conference on Sensing Technology*, 2008, pp. 242–247.
- [15] B. Tudu, A. Metla, B. Das, N. Bhattacharyya, A. Jana, D. Ghosh, R. Bandyopadhyay, Towards versatile electronic nose pattern classifier for black tea quality evaluation: An incremental fuzzy approach, *IEEE Transactions on Instrumentation and Measurement* 58 (9) (2009) 3069–3078.
- [16] B. Tudu, A. Jana, A. Metla, D. Ghosh, N. Bhattacharyya, R. Bandyopadhyay, Electronic nose for black tea quality evaluation by an incremental {RBF} network, *Sensors and Actuators B: Chemical* 138 (1) (2009) 90 – 95.
- [17] C. Cernuda, E. Lughofer, L. Suppan, T. Röder, R. Schmuck, P. Hintenaus, W. Märzinger, J. Kasberger, Evolving chemometric models for predicting dynamic process parameters in viscose production, *Analytica chimica acta* 725 (2012) 22–38.
- [18] C. Cernuda, E. Lughofer, G. Mayr, T. Röder, P. Hintenaus, W. Märzinger, J. Kasberger, Incremental and decremental active learning for optimized self-adaptive calibration in viscose production, *Chemometrics and Intelligent Laboratory Systems* 138 (2014) 14–29.
- [19] C. Cernuda, E. Lughofer, P. Hintenaus, W. Märzinger, T. Reischer, M. Pawliczek, J. Kasberger, Hybrid adaptive calibration methods and ensemble strategy for prediction of cloud point in melamine resin production, *Chemometrics and Intelligent Laboratory Systems* 126 (2013) 60–75.
- [20] E. Szymanska, J. Gerretzen, J. Engel, B. Geurts, L. Blanchet, L. Buydens,

Chemometrics and qualitative analysis have a vibrant relationship, *TrAC Trends in Analytical Chemistry* 69 (2015) 34–51.

- [21] X. Zhang, X. Qi, M. Zou, F. Liu, Rapid authentication of olive oil by raman spectroscopy using principal component analysis, *Analytical Letters* 44 (2011) 2209–2220.
- [22] M. Osorio, S. Haughey, C. Elliott, A. Koidis, Identification of vegetable oil botanical speciation in refined vegetable oil blends using an innovative combination of chromatographic and spectroscopic techniques, *Food Chemistry* 189 (SI) (2015) 67–73.
- [23] S. Wold, M. Sjostrom, *Simca: A method for analyzing chemical data in terms of similarity and analogy*, *Chemometrics Theory and Application* (1977) 243–282.
- [24] M. Osorio, S. Haughey, C. Elliott, A. Koidis, Evaluation of methodologies to determine vegetable oil species present in oil mixtures: Proposition of an approach to meet the eu legislation demands for correct vegetable oils labelling, *Food Research International* 60 (2013) 66–75.
- [25] K. Georgouli, J. Martinez-Del-Rincon, A. Koidis, Continuous statistical modelling for rapid detection of adulteration of extra virgin olive oil using mid infrared and raman spectroscopic data, *Food Chemistry* 217 (2017) 735–742.
- [26] R. Maggio, L. Cerretani, E. Chiavaro, T. Kaufman, A. Bendini, A novel chemometric strategy for the estimation of extra virgin olive oil adulteration with edible oils, *Food Control* 21 (2010) 890–895.
- [27] B. Ozen, L. Mauer, Detection of hazelnut oil adulteration using ft-ir spectroscopy, *Journal of Agricultural and Food Chemistry* 50 (2002) 3898–3901.
- [28] A. Koidis, M. Osorio, Identification of oil mixtures in extracted and refined vegetable oils, *Lipid Technology* 25 (2013) 247–250.

- [29] S. Ozawa, S. Pang, N. Kasabov, A modified incremental principal component analysis for on-line learning of feature space and classifier, in: PRICAI: Trends in Artificial Intelligence, Vol. 3157, 2004, pp. 231–240.
- [30] S. Pang, S. Ozawa, N. Kasabov, Incremental linear discriminant analysis for classification of data streams, IEEE Trans. Systems, Man, and Cybernetics (Part B) 35 (5) (2005) 905–914.
- [31] K. Diaz-Chito, F. Ferri, W. Díaz-Villanueva, Image recognition through incremental discriminative common vectors, in: Advanced Concepts for Intelligent Vision Systems - ACIVS, 2010, pp. 304–311.
- [32] P. Howland, J. Wang, H. Park, Solving the small sample size problem in face recognition using generalized discriminant analysis, Pattern Recognition 39 (2) (2006) 277–287.
- [33] L.-F. Chen, H.-Y. Liao, M.-T. Ko, J.-C. Lin, G.-J. Yu, A new lda-based face recognition system which can solve the small sample size problem, Pattern Recognition 33 (10) (2000) 1713–1726.
- [34] M.-B. Zhao, Z. Zhang, T. W. S. Chow, Z. Wu, On the theoretical and computational analysis between trace ratio lda and null-space lda, in: The 2012 International Joint Conference on Neural Networks (IJCNN), 2012, pp. 1–7.
- [35] D. mejkalov, A. Piccolo, High-power gradient diffusion nmr spectroscopy for the rapid assessment of extra-virgin olive oil adulteration, Food Chemistry 118 (2010) 153–158.
- [36] A. Christy, S. Kasemsumran, Y. Du, Y. Ozaki, The detection and quantification of adulteration in olive oil by near-infrared spectroscopy and chemometrics, Analytical Sciences 20 (2004) 935–940.
- [37] E. Lopez-Diez, G. Bianchi, R. Goodacre, Rapid quantitative assessment of the adulteration of virgin olive oils with hazelnut oils using raman

- spectroscopy and chemometrics, *Journal of Agricultural and Food Chemistry* 51 (2003) 6145–6150.
- [38] O. Devos, G. Downey, L. Duponchel, Simultaneous data preprocessing and svm classification model selection based on a parallel genetic algorithm applied to spectroscopic data of olive oils, *Food Chemistry* 148 (2014) 124–130.
- [39] R. J. Barnes, M. S. Dhanoa, S. J. Lister, Standard normal variate transformation and de-trending of near-infrared diffuse reflectance spectra, *Appl. Spectrosc* 43 (5) (1989) 772–777.
- [40] B. G. Osborne, T. Fearn, P. H. Hindle, *Practical NIR spectroscopy with applications in food and beverage analysis*, Longman scientific and technical, 1993.
- [41] A. Savitzky, M. J. Golay, Smoothing and differentiation of data by simplified least squares procedures, *Analytical chemistry* 36 (8) (1964) 1627–1639.
- [42] H. C. H. R. A. Berg and, J. A. Westerhuis, A. K. Smilde, M. J. Werf, Centering, scaling, and transformations: improving the biological information content of metabolomics data, *BMC genomics* 7 (1).