

1 **Combining isotopic signatures of $n(^{87}\text{Sr})/n(^{86}\text{Sr})$ and light stable elements (C, N,**
2 **O, S) with multi-elemental profiling for the authentication of provenance of**
3 **European cereal samples**

4
5 ³Daniel GOITOM ASFAHA, ³Christophe R. QUETEL*, ⁴Freddy THOMAS, ¹Micha HORACEK,
6 ¹Bernhard WIMMER, ¹Gerhard HEISS, ²Christian DEKANT, ²Peter DETERS-ITZELSBERGER,
7 ²Stefan HOELZL, ²Susanne RUMMEL, ³Christophe BRACH-PAPA, ³Marleen VAN
8 BOCXSTAELE, ⁴Eric JAMIN, ⁵Malcolm BAXTER, ⁵Katharina HEINRICH, ⁶Daniela
9 BERTOLDI, ⁶Luana BONTEMPO, ⁶Federica CAMIN, ⁶Roberto LARCHER, ⁶Matteo PERINI, ⁵,
10 ⁷Simon KELLY, ⁸Andreas ROSSMANN, ⁹Antje SCHELLENBERG, ⁹Claus SCHLICHT, ¹⁰Heinz
11 FROESCHL, ¹¹Jurian HOOGEWERFF, ¹¹Henriette UECKERMANN

12
13 *corresponding author:

14 Christophe Quétel

15 Institute for Reference Materials and Measurements (Joint Research Centre- European Commission).

16 111 Retieseweg. B-2440 Geel, Belgium.

17 e-mail: Christophe.Quétel@ec.europa.eu

18 Tel.: +32-14-57-1658

19 Fax: +32-14-57-1863

20
21 ¹ AIT (Austrian Institute of Technology), Seibersdorf, Austria

22 ² BSPG (Bayerische Staatssammlung für Paläontologie und Geologie), München, Germany

23 ³ EC-JRC-IRMM (European Commission – Joint Research Centre – Institute for Reference Materials and
24 Measurements), Geel, Belgium

25 ⁴ Eurofins Scientific Analytics, Nantes, France

26 ⁵ FERA (Food and Environment Research Agency), York, United Kingdom

27 ⁶ IASMA (Fondazione E. Mach – Istituto Agrario di San Michele all'Adige) San Michele all' Adige, Italy

28 ⁷ IFR (Institute of Food Research), Norwich, United Kingdom

29 ⁸ IGmbH (Isolab GmbH), Schweitenkirchen, Germany

30 ⁹ LGL (Bayerisches Landesamt für Gesundheit und Lebensmittelsicherheit), Oberschleißheim, Germany

31 ¹⁰ Seibersdorf Labor GmbH, Seibersdorf, Austria

32 ¹¹ UEA (University of East Anglia), Norwich, United Kingdom

33
34 **Keywords:** cereal; authentication of geographical origin; isotopic tracers; multi-elemental profiling

35 **Abbreviations:** CRM - Certified Reference Materials; IRMS - Isotope ratio mass spectrometry;

36 MC-TIMS - Multi Collector Thermal Ionisation Mass Spectrometer; MC-ICPMS - Multi

37 Collector Inductively Coupled Plasma Mass Spectrometer; PCA - Principal Component Analysis;

38 PLS-DA - Partial Least Square Discriminant Analysis

39 **Abstract**

40 The aim of this work (from the FP6 project TRACE) was to develop methods based on the use of
41 geochemical markers for the authentication of the geographical origin of cereal samples in Europe
42 (cf. EC regulations 2081/92 and 1898/06). For the first time the potential usefulness of combining
43 $n(^{87}\text{Sr})/n(^{86}\text{Sr})$ and $\delta^{13}\text{C}$, $\delta^{15}\text{N}$, $\delta^{18}\text{O}$ and $\delta^{34}\text{S}$ isotopic signatures, alone or with key element
44 concentrations ([Na], [K], [Ca], [Cu] and [Rb], progressively identified out of 31 sets of results),
45 was investigated through multiple step multivariate statistics for more than 500 cereal samples
46 collected over 2 years from 17 sampling sites across Europe representing an extensive range of
47 geographical and environmental characteristics.

48 Both models compared involved three sample classification categories (north/south; proximity to
49 the Atlantic Ocean/to the Mediterranean Sea/to else; bed rock geologies). The first two
50 categorisations were the most efficient, particularly when using the ten variables selected together
51 (with, in some instances, element concentrations making a greater impact than the isotopic tracers).
52 Validation of models included external prediction tests on 20% of the data randomly selected and,
53 rarely done, a study on the robustness of these multivariate data treatments to uncertainties on
54 measurement results. With the models tested it was possible to individualise 15 of the sampling
55 sites.

56

57 **I. Introduction**

58 In Europe, guaranteeing the geographical origin of a food product (cf. EC regulations 2081/92 and
59 1898/06) is regarded as an assurance of quality and safety. The FP6 project TRACE funded by the
60 European Commission was launched in 2005, aiming to deliver tools “that will enhance consumer
61 confidence in the authenticity of food” (FP6-TRACE project-website, 9th April 2010). An important
62 component of this project was to study how geochemical markers and the relationships between
63 these markers could be used for the characterisation of the region of origin of food products. For
64 instance, it was shown that variations in H, C, N and S isotopic composition and element
65 concentrations are useful for the differentiation between samples coming from some different
66 regions in Europe in the case of lamb meat, honey and olive oils (Camin et al., 2007, 2010;
67 Schellenberg et al., 2010). The present study also reports research findings from the TRACE project
68 but on a broader scope level and with regards to wheat. For the first time the potential usefulness of
69 C, N, O, S and Sr isotopic signatures combined to multi-element profiles was investigated for over
70 500 samples originating from sites in Europe representing an extensive range of geographical and
71 environmental characteristics (climate, distance to sea/ocean, geology).

72 The $\delta^{13}\text{C}$ values of plant compounds depend on photosynthetic pathways (discrimination between
73 C3 and C4 plants), the plant age and level of maturation (Farquhar et al., 2003; Smith and Epstein,
74 1971). $\delta^{13}\text{C}$ is also affected, although to a lesser extent, by several environmental factors such as
75 relative humidity, temperature, amount of precipitation and water stress (O'Leary, 1995). $\delta^{15}\text{N}$
76 values depend on the botanical type of the plant and on the bacterial activity associated to its
77 growth (Farquhar et al., 2003), and can also be affected by temperature and agricultural practices
78 (e.g. type of fertilizers used) (Bateman et al., 2005; Martin and Martin, 2003). The $\delta^{18}\text{O}$ of plant
79 material reflects the water taken up (linked to temperature, latitude, elevation, distance from the sea
80 and amount of precipitation) and exchanged (evaporative and diffusion effects during transpiration),
81 and also biosynthetic pathways including the isotopic exchange between organic molecules and
82 plant water (Bréas et al., 1998). Variations in the $\delta^{34}\text{S}$ ratio in plants have been reported to be linked
83 with the soil geology (Thode, 1991), the distance from the sea and proximity to industrial activities
84 (Camin et al., 2007). The $n(^{87}\text{Sr})/n(^{86}\text{Sr})$ ratio in soils may vary depending on the [Rb]/[Sr]
85 concentration ratio in (and age of) the surrounding rock or mineral and, thus, has the potential to
86 reflect the nature of the underground geology (Capo et al., 1998). It is expected that the
87 $n(^{87}\text{Sr})/n(^{86}\text{Sr})$ in plants should also be a good indicator of the rocks and soil conditions during plant
88 growth. Similarly, it is expected that element concentrations in food commodities such as olive oils
89 (Camin et al., 2010), potatoes (Anderson et. al., 1999), garlic (Smith, 2005), orange juice (Simpkins
90 et al. 2000) and coffee (Anderson et. al., 2002) are mainly related to the geological and

91 pedoclimatic characteristics of the site of growth and plant farming practices. A discussion on
92 reasons why the elemental composition in plant may provide unique markers of the geographical
93 origin can be found in Kelly et. al. (2005).

94 The way these tracers could be used to enable the authentication of provenance of cereal samples
95 was investigated through multiple steps multivariate statistical analysis. The method validation
96 scheme we designed included a study on the robustness of these multivariate data treatments to
97 uncertainties in measurement results.

98

99

100 **II. Experimental**

101 The measurands in cereal samples are hereby referred to as *variables*. They consist of Sr isotope
102 ratios (noted $n(^{87}\text{Sr})/n(^{86}\text{Sr})$) and 31 element concentrations (noted [E]) in the bulk, as well as δ -
103 scale values for the CNOS group (noted $\delta^{13}\text{C}$, $\delta^{15}\text{N}$, $\delta^{18}\text{O}$ and $\delta^{34}\text{S}$) in the defatted fraction.

104 $\delta^{13}\text{C}$, $\delta^{15}\text{N}$, $\delta^{18}\text{O}$ and $\delta^{34}\text{S}$ measurements were carried out in seven laboratories (Eurofins Scientific
105 Analytics, Nantes, France; FERA, York, UK; AIT, Seibersdorf, Austria; LGL, Oberschleißheim,
106 Germany; Isolab, Schweitenkirchen, Germany; IASMA, San Michele all'Adige, Italy; IFR,
107 Norwich, UK). $n(^{87}\text{Sr})/n(^{86}\text{Sr})$ measurements were carried out in four laboratories (EC-JRC-IRMM,
108 Geel, Belgium; FERA; BSPG, München, Germany; UEA, Norwich, UK). Element concentration
109 measurements were carried out in four laboratories (Seibersdorf Labor GmbH, Seibersdorf, Austria;
110 FERA; IASMA; UEA). The scientists responsible for producing these experimental data are co-
111 authoring this paper.

112

113 ***a) Sampling strategy and sample collection***

114 Cereal samples were collected from multiple farms or local producers in 19 sampling sites in
115 Europe between summer 2005 and summer 2007. The names and average geographical coordinates
116 of these sampling sites are provided in Table 1.

117 Although wheat (of different varieties, including winter wheat, Durum wheat, Emmer wheat,
118 Epeautre and Spring wheat) was the main target, barley, rye, triticale, oat (all C3 plants) and four
119 corn (C4 plant) samples were also collected from areas where wheat was not produced. The
120 sampling plan was to collect four samples from five different producers or fields within each of the
121 predefined sampling sites per year, and this was repeated for a second year around the same time.
122 The ideal scheme of 40 cereal samples per sampling site could not always be strictly followed, due
123 to practical limitations, and a total of 557 cereal samples were obtained, as shown in Table 1.
124 Algarve and Barcelona samples (incl. the four corns) were excluded for this study because they

125 were purchased on markets and their claimed origin could not be verified. Samples from Galicia
126 were collected for one season only, and only three samples could be obtained from Iceland.

127

128 ***b) Sample preparation, spectrometric measurements and validation of measurement methods***

129 *i) For $\delta^{13}\text{C}$, $\delta^{15}\text{N}$, $\delta^{18}\text{O}$ and $\delta^{34}\text{S}$ results*

130 $\delta^{13}\text{C}$, $\delta^{15}\text{N}$ and $\delta^{18}\text{O}$ were measured using a range of Isotope Ratio Mass Spectrometers (Delta Plus
131 XL, Delta Plus XP, Delta V, Delta S, ThermoFinnigan, Bremen, Germany; Isoprime, AP2003,
132 Optima, GV Instruments Ltd., Manchester, UK), IRMS. They were connected to pyrolisers
133 (TC/EA, ThermoFinnigan; PyrOH and EA3000, Eurovector, GV Instruments Ltd.) for $^{18}\text{O}/^{16}\text{O}$ or
134 an elemental analysers (Flash EA 1112, 1110, 1108 ThermoFinnigan; Costech ECS4010; NA2100
135 Proteins, Carloerba, Milan, Italy; Vario EL III, Elementar Analysensysteme GmbH, Hanau,
136 Germany) for $^{13}\text{C}/^{12}\text{C}$ and $^{15}\text{N}/^{14}\text{N}$. Samples contained very low amounts of sulphur (usually less
137 than 0.1 % according to Sieper et al., 2006) and $\delta^{34}\text{S}$ values were obtained from only one
138 instrument (Vario EL III, Elementar Analysensysteme GmbH). IRMS instrumental conditions for
139 each isotope ratio were reported by the same authors in a previous paper (Camin et al., 2007).

140 The isotopic values were expressed in $\delta\text{‰}$ versus V-PDB (Vienna - Pee Dee Belemnite) for $\delta^{13}\text{C}$,
141 versus AIR for $\delta^{15}\text{N}$, versus V-SMOW (Vienna – Standard Mean Ocean Water) for $\delta^{18}\text{O}$, and
142 versus V-CDT (Vienna Canyon-Diablo-Troilite) for $\delta^{34}\text{S}$ according to the following formula: $[(R_s-$
143 $R_{std})/R_{std}] \times 1000$, where R_s is the isotope ratio measured for the sample and R_{std} is the isotope ratio
144 of the international standard. The values were calculated against in-house standards, which were
145 themselves calibrated against international reference materials: fuel oil NBS-22 (IAEA) and sugar
146 IAEA-CH-6 (IAEA) for $^{13}\text{C}/^{12}\text{C}$, USGS-40 and IAEA-N-1 (IAEA) for $^{15}\text{N}/^{14}\text{N}$, Benzoic Acid
147 IAEA-601 (IAEA) and IAEA-CH-6 (IAEA) for $^{18}\text{O}/^{16}\text{O}$ and IAEA-S-1 for $^{34}\text{S}/^{32}\text{S}$.

148 The measurement repeatability was $\leq 0.2\text{‰}$ for $\delta^{13}\text{C}$, $\leq 0.5\text{‰}$ for $\delta^{18}\text{O}$, $\leq 0.1\text{‰}$ for $\delta^{15}\text{N}$ and \leq
149 0.3‰ for $\delta^{34}\text{S}$. An inter comparison was organized between the seven laboratories in charge of
150 IRMS measurements for validation purposes. The test material for the comparison (wheat flour)
151 was also used periodically in each laboratory as quality control, to monitor deviations over time.
152 The dispersions of results between laboratories during the comparison were 0.2‰ for $\delta^{13}\text{C}$, 0.9‰
153 for $\delta^{18}\text{O}$, 0.8‰ for $\delta^{15}\text{N}$ and 1.2‰ for $\delta^{34}\text{S}$.

154 *ii) For $n(^{87}\text{Sr})/n(^{86}\text{Sr})$ results*

155 The instrumentation and laboratories involved in the Sr isotopic measurements were the same as
156 those of the study on mineral water samples published recently (Voerkelius et al., 2010). There
157 were one Multi Collector Thermal Ionisation Mass Spectrometer (MC-TIMS – MAT 261/262,
158 Thermo Finnigan) and three Multi Collector Inductively Coupled Plasma Mass Spectrometers (MC-

159 ICPMS – Nu Plasma, Nu Instruments; Axiom, ex- VG Instruments; Isoprobe, ex-Micromass). MC-
160 TIMS measurements were carried out on tungsten single filaments. For MC-ICPMS measurements,
161 typically, the instrumentation was equipped with a minicyclonic jacketed cinnabar spray chamber
162 or a porous membrane based desolvation unit, a 0.2 ml min⁻¹ microconcentric nebuliser and Ni or Pt
163 sampler and skimmer cones.

164 Sample uptake was 0.1 - 0.3 g for MC-TIMS measurements, and 0.4 - 1 g for MC-ICPMS
165 measurements. Samples were dry-ashed for MC-TIMS measurements, or acid digested in a
166 microwave oven for MC-ICPMS measurements. Residuals were dissolved in concentrated HNO₃ to
167 allow separation of Sr from other elements such as Ca, Ba or Rb by ion chromatography on a Sr
168 specific crown ether resin (Sr-spec®). After separation the Sr concentration in solution was 10-500
169 ng g⁻¹ in about 2.5 g of 2-3% HNO₃.

170 The same standard operating procedure (SOP, cf. details in report of FP6-TRACE Deliverable
171 D15.9, 2009) applied to all laboratories to allow for comparison of measurement results. This SOP
172 included recommendations on Rb/Sr separations, corrections for procedural blanks, the prescription
173 of the $n(^{86}\text{Sr})/n(^{88}\text{Sr})$ and $n(^{87}\text{Rb})/n(^{85}\text{Rb})$ reference values used for corrections (0.1194 ± 0 and
174 0.38565 ± 0.00030 , respectively). It also required running regularly (every 1 to 4 samples
175 maximum) experiments on the NIST-987 Sr isotopic CRM for the purpose of contributing to the
176 validation of each methods developed. The relative dispersion of results between laboratories for
177 identical cereal test materials during four comparisons was up to approximately 0.9‰ depending on
178 physical (flour/coarse grains, grains with or without husk etc.) and chemical ([Rb] and [Sr]/[Rb]
179 concentration ratio) characteristics of samples.

180 *iii) For [E] results*

181 Measurements of element concentrations were performed on four quadrupole ICP-MS instruments
182 (two Agilent 7500ce, Tokyo, Japan; one Elan 6100 and one Elan 6000, Perkin Elmer Sciex,
183 Toronto, Canada). The first two were fitted with concentric Micromist nebulisers (Glass Expansion,
184 Melbourne, Australia) and water-cooled Scott double-pass spray chambers, the third with a cross-
185 flow nebuliser and the fourth with a concentric nebuliser and a cyclonic water-cooled spray-
186 chamber. In all cases peristaltic pumps were used to regulate sample flow rates. RF power settings
187 ranged from 1100 to 1500W on the four instruments. Pt cones were used on the Elan 6000, while Ni
188 cones were used on the other instruments.

189 Microwave digestion was employed by three of the laboratories, digesting between 0.3 and 0.5 g
190 cereal with distilled nitric acid, either with or without addition of ultrapure hydrochloric acid. The
191 fourth laboratory used ultraviolet digestion of 0.5 g sample material with 4 ml nitric acid and 1 ml

192 distilled hydrochloric acid. All labware involved was cleaned beforehand with solutions of between
193 5% and 10% nitric acid. Results were corrected for cereal moisture content.

194 NIST1567a Wheat Flour was used for method validation purposes and measured with all sample
195 batches. There was agreement within uncertainty boundaries with certified values for all elements
196 of interest in this study.

197
198 ***c) Data set description and the multivariate data analysis tools applied***

199 In our study the $n(^{87}\text{Sr})/n(^{86}\text{Sr})$ and the CNOS isotopic signatures were treated alone and in
200 combination with [E] results. In the first case it was necessary to exclude an additional 13 samples
201 (of 32) from Jura Krakowska and 2 samples (of 9) from Jylland for which four out of five of the
202 isotopic data were not reported. Thus, the overall cereal data set consisted of 512 samples from 17
203 sampling sites. When isotopic and [E] results were combined all 527 samples could be considered
204 (e.g. including the 15 samples excluded in the previous case).

205 Prior to the application of multivariate statistical calculations some data pre-treatment was required
206 (Eriksson et al., 2001; Esbensen, 2006). For each variable unit variance scaling and mean-centering
207 of data was applied (division by the standard deviation followed by subtraction of the mean value).
208 In addition, [Na] and [Rb] data required log-transformation in order to normalise the distributions.

209 In this study *Principal Component Analysis* (PCA) and *Partial Least Square Discriminant Analysis*
210 (PLS-DA) modelling tools were used. PCA provides a more comprehensive overview of all data by
211 producing a few orthogonal (uncorrelated) *principal components* (PCs) which extract the main
212 information about the data set (Eriksson et al., 2001; Esbensen, 2006). PCA is normally applied at
213 early stage of multivariate data analysis as “exploratory tool”. PLS-DA modelling is more suitable
214 for smaller numbers of defined classes and maximises the separation between them.

215 PCA and PLS-DA offer a number of useful parameters and diagnostic tools expressed graphically
216 and numerically. These include PC-score plots, $R^2X(\text{cum})$, $R^2Y(\text{cum})$, $Q^2Y(\text{cum})$ and VIP
217 (*Variable Importance in the Projection*). The R^2 parameters are a quantitative measure of the
218 “goodness of fit” of a given model (R^2X for the “predictors”, and R^2Y for the “responses”),
219 whereas the Q^2 parameters indicate the “goodness of prediction” (predictive ability) of a given
220 model. There exist some rules to help identifying the best balance between the predictive power and
221 a reasonable level of fitness of the model. According to Eriksson et al. (2001) “generally, a $Q^2 >$
222 0.5 is regarded as good and a $Q^2 > 0.9$ as excellent, but these guidelines are of course heavily
223 application dependent”, and “differences between R^2 and Q^2 larger than $0.2-0.3$ indicate the
224 presence of many irrelevant model terms or few outlying data points”. For PLS-DA the relative
225 importance of variables is illustrated by VIP plots. Variables with a VIP value greater than unity

226 play the most important roles for the discrimination of the classes. For more elaborated
227 explanations and mathematical expressions related to these diagnostic tools and parameters readers
228 are referred to Eriksson et al. (2001).

229 In this work, all the multivariate computations were carried out using SIMCA-P Ver. 12 (Umetrics
230 AB, Sweden) software. The Statistica Ver. 7 (Statsoft Inc., USA) software was used to produce the
231 box-whisker plots for the variables.

232 Our validation scheme included *external prediction* tests on 20% of the data assigned as
233 independent data sets (randomly selected). These data were compared (on class membership)
234 against prediction results obtained from PLS-DA models built with the remaining 80% of the data
235 set. Models were considered satisfactory when external prediction rates were $\geq 70\%$.

236

237

238 **III Results and discussion**

239 To have an overview of all 527 samples from 17 sites with 36 descriptor variables, a PCA was
240 applied to the data set (see Figure ES1 provided as electronic supplementary material).

241 There were no clear trends or patterns to distinguish samples of different sampling sites. Thus we
242 decided to apply PLS-DA, based on supervised classification, and to compare two models, using
243 the five isotopic variables alone and in combination with [E] results.

244 The classification categories we chose (Table 1) were according to *latitude* (two classes: north and
245 south of Europe), *proximity to a marine environment* (three classes: Atlantic, Mediterranean and
246 Inland) and *geology of the underlying bed rock* (four classes: Shale, Acid magmatic, Limestone and
247 Basaltic). The way we established these classification categories and the results obtained are
248 discussed in the subsequent sections.

249 The five isotopic variables ($n(^{87}\text{Sr})/n(^{86}\text{Sr})$, $\delta^{13}\text{C}$, $\delta^{15}\text{N}$, $\delta^{18}\text{O}$ and $\delta^{34}\text{S}$) were systematically
250 considered as they were found alternately significant depending on the classes considered for the
251 data, as shown later. Preliminary investigations (successive PLS-DA and VIP plots) were carried
252 out to evaluate the significance of the 31 element concentration variables available. From the first
253 list of nine apparently influential concentration variables identified, four ([Co], [Ga], [Cd] and [Cs])
254 were further eliminated due to low concentration in the majority of the samples and rather poor
255 analytical figures of merit (measurement reproducibility amongst the laboratories for cereal samples
256 was, from the TRACE quality assurance report (2009), $\leq 20\%$ for [Cd], 20% to 50% for [Co] and
257 over 50% for [Ga] and [Cs]). That left only [Na], [K], [Ca], [Cu] and [Rb]. Hereafter, the two
258 methods of modelling will be referred to as the “5-variable” and “10-variable” (the 5 isotopic
259 markers and [Na], [K], [Ca], [Cu] and [Rb]) models, respectively.

260

261 **a) PLS discriminant analysis models**

262 *i) Classification according to latitude (north and south classes)*

263 The motivation was the assumption of a relationship between climatic conditions and the latitudinal
264 geographic position of each site where the cereal samples were harvested. After several model
265 refining attempts, the optimal boundary line for the discrimination between north and south samples
266 was found to be around latitude line 47° N through “central” Europe. In both cases PLS-DA
267 computations with this classification resulted in two optimal significant PCs, and relatively high
268 values of *goodness of prediction* Q2(cum) (figures not between brackets in Table 2). As shown in
269 the score plots for PC1 and PC2 (Figure 1), the two classes were well sorted mainly in the direction
270 of PC1. The dashed line was added to delineate the two groups visually.

271 It can be seen on VIP plots (see Figure ES2 provided as electronic supplementary material) that the
272 first two most influential variables for a north/south differentiation were $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$, which is
273 coherent with both being potential indicators of climatic conditions (Heaton et al., 2008; Kelly et
274 al., 2002, 2005; Rossmann et al., 2000; Suzuki et al., 2008). More remarkably [Cu] played a
275 significant discriminatory role in the case of the *10-variables* model. As shown by box-whisker
276 plots (see Figure ES3 provided as electronic supplementary material) [Cu] was rather higher for
277 samples collected in the south than in the north. It is reasonable to assume that the amount of
278 copper in a cereal plant is proportional to the amount of Cu available from the soil where it was
279 grown. What we observed for our cereal data might be explained by the highest Cu concentrations
280 in topsoil and subsoil in Europe reported by Foregs maps (Salminen, 2005) for southern regions,
281 around the Mediterranean basin and within the Iberian peninsula, and along the west coasts of
282 France, UK and Norway. The distribution of copper in subsoil is mainly related to regional and
283 local geology, and to mineralisation. According to Salminen (2005) the distribution of copper in
284 topsoil may also be influenced by anthropogenic contamination (pollution by agricultural sewage
285 enriched in Cu or the use of copper sulphate as a fungicide in fruit cultivation and vineyards).

286 PLS-DA model fittings for external prediction tests showed satisfactory results with Q2(cum)
287 values similar to those of PLS-DA models generated using the entire data set. The external
288 prediction rates described in column 3 and 4 of Table 3 (figures not between brackets, on 103 and
289 106 samples, respectively) were 100 % for 9-10 out of 17 sampling sites. The inclusion of [E] data
290 did not make a significant difference overall: there were improvements for samples from
291 Fraenkische Alb, Allgaeu, Chalkidiki, and Galicia, and degradations for samples from Muehlviertel,
292 Limousin, Firenze, Jura Krakowska and Cornwall. The external prediction rates were less than 70

293 % only for Galicia, Chalkidiki (the *5-variables* model only) and Firenze (both models). These
294 results indicate that both models proposed are robust/stable for the north/south discrimination.

295 *ii) Classification according to proximity to a marine environment (Atlantic,*
296 *Mediterranean and Inland classes)*

297 We then examined whether the proximity to oceanic/sea conditions (via sea-spray deposition,
298 unusual historical sedimentation conditions, and so on) could be characterised.

299 Three classes of samples were considered: “Atlantic” and “Mediterranean” for sampling sites less
300 than 100 km away from the respective coasts, and “Inland” for the remaining samples (cf. Table 1).

301 PLS-DA with the *5-* and *10-variables* models resulted in three and four PCs, respectively. As
302 illustrated in Figure 2 (score plots with the first two PCs) the choice of these three classes
303 (delineated with dotted lines) for sorting samples was relevant. There was also a significant
304 improvement of the modelisation when incorporating the five [E] variables, as confirmed by the
305 increase of Q2(cum) values in Table 2 (from 0.49 to 0.66).

306 VIP values for both models are presented in Figure ES2. $\delta^{34}\text{S}$, $\delta^{18}\text{O}$ and $\delta^{13}\text{C}$ were the most
307 discriminatory isotopic variables but, with the *10-variables* model, only after [Na] and just before
308 [K].

309 Furthermore, box-whisker plots in Figure ES3 show that “Atlantic” class cereal samples can be
310 distinguished from others as they exhibit, globally, higher values of [Na], $\delta^{34}\text{S}$ and [K] (associated
311 with the emission of sea-spray and the deposition of sea-salt). Cereal samples close to the
312 Mediterranean Sea can also be discriminated but essentially on the basis of tracers of climatic
313 conditions (highest $\delta^{13}\text{C}$ and $\delta^{18}\text{O}$ values, globally), thus in this case a north/south issue in line with
314 conclusions from the previous discussion.

315 External prediction tests were run in the same way as described before, and the new PLS-DA model
316 parameters (PCs, R2, Q2) were almost identical to the original ones (Table 2).

317 The results in column 5 and 6 in Table 3 (figures not between brackets), with the *5-* or the *10-*
318 *variables* models respectively, indicate 100% success rate for 5 and 9 sampling sites, and < 70%
319 success rate for 4 sampling sites in both cases. For Jylland, Galway and Iceland though, the number
320 of samples available for the external prediction tests was scarce and the results obtained (100% or
321 0%) must be interpreted cautiously. When considering sampling sites with $\geq 70\%$ success rate (and
322 apart from the 3 sites with ≤ 3 samples), the number of samples predicted correctly is either
323 identical for both modelling approaches (for 7 sites) or better by 11% to 29% with the *10-variables*
324 modelling (for Marchfeld, Gaeuboden, Allgaeu and Jura Krakowska). Therefore it was possible to
325 conclude, first, that the proximity (or not) to a marine environment could be correctly predicted for
326 a very large majority of cereal samples and, second, that the combination with [E] variables (and

327 particularly [Na] and [K]) significantly improved this prediction ability as compared to the use of
328 isotopic variables alone. The sites with < 70% success rates to the external prediction tests for both
329 types of models (also apart from the 3 sites with ≤ 3 samples) were Chalkidiki, Firenze and Galicia,
330 similarly to what had been observed previously with the “Latitude” classification (*5-variables*
331 model only).

332 *iii) Classification according to bedrock geology (Acid Magmatic, Shale, Limestone*
333 *and Basaltic classes)*

334 For the third test we sorted cereal samples into four classes (“Shale/mudstone/clay/loess” incl.
335 sandstones and other clastic sediments, “Acid Magmatic”, “Limestone” and “Basaltic”) according
336 to the bedrock geology of their respective sampling sites, described in Table 1.

337 Application of PLS-DA to these four classes produced three and five PCs when using the *5-* and *10-*
338 *variables* models, respectively. Q2(cum) values (Table 2) indicated that the predictive ability of
339 these models had deteriorated in comparison to the previous two classifications. Score plots with
340 the first two PCs are shown in Figure 3. The “Acid Magmatic” class was consistently better sorted
341 than the other 3 classes for both modelling approaches. With only 3 samples the “Basaltic” class
342 was not discriminated at all, although there was a slight difference whether the *5-variables* or the
343 *10-variables* model was considered.

344 VIP plots for the PLS-DA on bed rock geology classes (Figure ES2) show [Rb], $n(^{87}\text{Sr})/n(^{86}\text{Sr})$ and
345 $\delta^{15}\text{N}$ as being the most influential variables. Since the VIP score of [Ca] and $\delta^{34}\text{S}$ were ≥ 1 within
346 confidence intervals, these variables may also be considered relevant for this classification. As
347 explained earlier $n(^{87}\text{Sr})/n(^{86}\text{Sr})$ values change depending on the [Rb]/[Sr] concentration ratio in the
348 surrounding geological structure. Thus the major role played by [Rb] and $n(^{87}\text{Sr})/n(^{86}\text{Sr})$ in this
349 classification is logical. Box-whisker plots in Figure ES3 show that with higher values globally for
350 both variables the “Acid Magmatic” class differentiates from the other two major classes. The fact
351 that [Rb] was significantly more influential than $n(^{87}\text{Sr})/n(^{86}\text{Sr})$ is noteworthy. Furthermore, $\delta^{15}\text{N}$
352 values were the highest globally for the “Shale/mudstone/clay/loess” class. $\delta^{15}\text{N}$ values may change
353 depending on mineral origin (from -6 ‰ to +6 ‰) or organic origin (from +1 ‰ to 37 ‰) of the
354 fertilisers (Bateman and Kelly, 2007). The trend observed here might indicate agricultural practices
355 specific to farming territories associated with this geological class, although this statement is
356 speculative since its verification was beyond the scope of this study.

357 As expected, PLS-DA computations for external predictions tests showed no improvements of the
358 Q2(cum) values. According to results in column 7 and 8 in Table 3 the origin of samples was
359 predicted 100% correctly only for Marchfeld (both models), for Limousin (*5-variables* model) and
360 for Orkney and Carpentras (*10-variables* model). External prediction rates for Sicily (both models),

361 Firenze (*5-variables* model) and Allgaeu (*10-variables* model) were also satisfactory ($\geq 70\%$), and
362 5 of these correctly predicted sites belonged to the “Shale/mudstone/clay/loess” class. The
363 capability to identify correctly samples from Firenze was also a positive result (not the case with
364 previous two classifications). Globally, external prediction rates were again slightly better when
365 taking into account the five [E] variables than with the five isotopic variables alone.
366 Overall, samples from Galicia and Chalkidiki were the only ones (also with those from Firenze in
367 the case of the *10-variables* approach) that could not be identified with the combination of
368 classification categories investigated. The outcome for the “Bed rock geology” sorting was less
369 useful than for the “Latitude” and “Marine/Inland” sorting. Several possible reasons can be
370 envisaged to explain this difference. Geological backgrounds are often not homogeneous (for
371 Galicia and Chalkidiki in particular it is more a patchwork than a uniform system, as visible from
372 maps by Asch, 2005). Besides, there is not necessarily a straightforward relationship between bed
373 rock geologies and the compositions of the mineral fractions available for collection by plants
374 growing on the soil surface. Understanding the intimate interactions between soil and plants was
375 outside the scope of TRACE but this supplementary dimension could be a meaningful inclusion in a
376 follow-up project.

377

378 ***b) Robustness to uncertainties on measurement results as additional model validation***

379 The validation approach described in the previous sections for our models was based on the
380 calculation of Q2(cum) values and the run of external prediction tests. We also investigated a much
381 more unusual way of validating these models by examining their robustness to inter-comparison
382 variability. PLS-DA based models assume exactness of the input data and these mathematical tools
383 are not designed to handle measurement uncertainties. We introduced changes to our data set based
384 on the maximum dispersions of results ‘*d*’ observed during our inter-comparisons (0.2‰ for $\delta^{13}\text{C}$,
385 0.9‰ for $\delta^{18}\text{O}$, 0.8‰ for $\delta^{15}\text{N}$, 1.2‰ for $\delta^{34}\text{S}$ and, in relative terms, 0.9‰ for $n(^{87}\text{Sr})/n(^{86}\text{Sr})$ and
386 40% for [E]). A component ‘*c*’ was added to all original observations corresponding to the product
387 of ‘*d*’ and a randomly generated ‘*r*’ ranging from -1 to +1 (equ. 1).

388

$$389 \text{Equation 1} \quad c = r*d$$

390

391 A new series of PLS-DA models for each classification category were produced from the simulated
392 data set. The new prediction results (values between brackets in Tables 2 and 3) showed that such
393 combinations of multivariate analysis were globally robust to the fluctuations imposed to the
394 original data set although some degradation was observed. Q2(cum) values decreased by no more

395 than 0.05. Trends for the external prediction rates remained also quite similar. Transition from \geq
396 70% to $< 70\%$ was observed with the *10-variables* “Latitude” model for Chalkidiki and Galicia,
397 with the *10-variables* “Marine/Inland” model for Jura Krakowska, with the *5-variables* “bed rock
398 geology” model for Firenze and Sicily and for the *10-variables* “bed rock geology” model for
399 Algaeu and Carpentras. Transition from $< 70\%$ to $\geq 70\%$ was observed with the *5-variables* “bed
400 rock geology” model for Carpentras, with the *10-variables* “bed rock geology” model for Firenze
401 and Gauboden. Overall, external prediction results for all the *10-* and *5- variables* models differed
402 for not more than 11 samples (over 103 or 106 test samples) compared to results from tests with the
403 original data.

404 These simulations with these dispersion ranges illustrate the potential sensitivity of the prediction
405 tools investigated to the quality of the experimental data, and provide an indication of what could
406 be the maximum experimental uncertainty tolerable for the models proposed to work.

407

408

409 **IV Conclusions**

410 The sequential approach described in this study was successful. Grouping sampling sites
411 successively according to the *latitude* (north and south classes) and the *proximity to a marine*
412 *environment* (Atlantic, Mediterranean and Inland classes) was particularly efficient, using ten
413 carefully chosen variables (five isotopic tracers combined with [Na], [K], [Ca], [Cu] and [Rb]
414 concentrations). Grouping sampling sites according to the *geology of the underlying bed rock*
415 (Shale, Acid magmatic, Limestone and Basaltic classes) was less useful, probably due to the
416 simplicity of the classification used and the fact that bed rock geologies alone may not be adequate
417 to predict soil chemistry. Moreover, this combination of classification categories allowed the
418 identification of 12 unique and generic ‘identities’ (column 6 in Table 1), thus providing multiple
419 ways of describing geographical locations based on easy to use principles.

420 Another important lesson learned was that in some instances element concentrations made a greater
421 impact as variables than the five isotopic tracers ([Na] versus $\delta^{34}\text{S}$ as a proxy to a marine
422 environment and to the Atlantic ocean in particular; [Rb] versus $n(^{86}\text{Sr})/n(^{88}\text{Sr})$ as a proxy to a
423 certain bed rock geology).

424 Although our validation scheme was quite extensive, it is natural to examine the question of the
425 domains of applicability. Clearly, it cannot be claimed that the models proposed will always be able
426 to predict/identify correctly the geographical origin of all possible cereal sample collected in
427 Europe. The 17 sampling sites chosen do not cover all the variability on this continent in term of
428 geological backgrounds, soil characteristics and climatic conditions. However, a study like this one,

429 because of its unique size (number of sampling sites, number of samples per site, number of
430 variables assessed per sample), allows the design of a tool applicable to a much wider range of
431 cereal samples in Europe than any previous comparable studies, and proposes a proof of concept
432 applicable to other types of grains and food products than only cereals.

433 Our results also demonstrate the feasibility, the potential interest and also the limitations of such a
434 large size study at the scale of a continent. Given the variety of parameters investigated and the
435 great number of samples involved, this kind of project requires the collaboration of many
436 organisations. A tight coordination between partners, including vast efforts in the field of data
437 quality assurance, is mandatory. The quality of experimental data used for modelling purposes
438 cannot be considered better than the difference between measurement results (on similar samples)
439 observed for the different project partners. And our simulations have shown that these multivariate
440 data treatments are not insensitive to such levels of uncertainties.

441

442 **References**

443

- 444 1. Anderson, K. A., Magnuson, B. A., Tschirgi, M. L., Smith, B., 1999. Determining the
445 Geographic Origin of Potatoes with Trace Metal Analysis Using Statistical and Neural Network
446 Classifiers. *Journal of Agricultural and Food Chemistry*, 47, 1568-1575.
- 447 2. Anderson, K. A., Smith, B., 2002. Chemical Profiling to Differentiate Geographic Growing
448 Origins of Coffee. *Journal of Agricultural and Food Chemistry*, 50, 2068-2075.
- 449 3. Asch, S., IGME-map. *IGME 5000*, 2005. *The 1:5 Million International Geological Map of*
450 *Europe and Adjacent Areas*, BGR, Hannover.
- 451 4. Bateman, A. S., Kelly, S. D., Jickells, T. D., 2005. Nitrogen Isotope Relationships between Crops
452 and Fertilizer: Implications for Using Nitrogen Isotope Analysis as an Indicator of Agricultural
453 Regime. *Journal of Agricultural and Food Chemistry*, 53, 5760-5765.
- 454 5. Bateman, A. S., Kelly, S. D., 2007. Fertilizer nitrogen isotope signatures. *Isotopes in*
455 *Environmental Health Studies*, 43, 237-247.
- 456 6. Bréas, O., Guillou, C., Reniero, F., Sada, E., Angerosa, F., 1998. Oxygen-18 measurement by
457 continuous flow pyrolysis/isotope ratio mass spectrometry of vegetable oils. *Rapid*
458 *Communications in Mass Spectrometry*, 12, 188-192.
- 459 7. Camin, F., Bontempo, L., Heinrich, K., Horacek, M., Kelly, S. D., Schlicht, C., Thomas, F.,
460 Monahan, F. J., Hoogewerff, J., Rossmann, A., 2007. Multi-element (H,C,N,S) stable isotope
461 characteristics of lamb meat from different European regions. *Analytical and Bioanalytical*
462 *Chemistry*, 389, 309-320.
- 463 8. Camin, F., Larcher, R., Nicolini, G., Bontempo, L., Bertoldi, D., Perini, M., Schlicht, C.,
464 Schellenberg, A., Thomas, F., Heinrich, K., Voerkelius, S., Horacek, M., Ueckermann, H.,
465 Froeschl, H., Wimmer, B., Heiss, G., Baxter, M., Rossmann, A., Hoogewerff, J., 2010. Isotopic and
466 elemental data for tracing the origin of European olive oils. *Journal of Agricultural and Food*
467 *Chemistry*, 58, 570-577.
- 468 9. Capo, R. C., Stewart, B. W., Chadwick, O. A., 1998. Strontium isotopes as tracers of ecosystem
469 processes: theory and methods. *Geoderma*, 82, 197-225.
- 470 10. Eriksson, L., Johansson, E., Kettaneh-Wold, N., Wold, S., 2001. *Multi- and Megavariate Data*
471 *Analysis: Principles and Applications*. Umetrics AB: Umeå.
- 472 11. Esbensen, K. H., 2006. *Multivariate Data Analysis - In Practice. An Introduction to*
473 *Multivariate Data Analysis and Experimental Design*. Fifth ed., CAMO Software, Oslo.
- 474 12. Farquhar, G. D., Ehleringer, J. R., Hubick, K. T., 2003. Carbon Isotope Discrimination and
475 Photosynthesis. *Annual Review of Plant Physiology and Plant Molecular Biology*, 40, 503-537.

- 476 13. FP6-TRACE project-website, 9th April, 2010. <http://www.trace.eu.org/>.
- 477 14. Heaton, K., Kelly, S. D., Hoogewerff, J., Woolfe, M., 2008. Verifying the geographical origin
478 of beef: The application of multi-element isotope and trace element analysis. *Food Chemistry*, 107,
479 506-515.
- 480 15. Kelly, S., Baxter, M., Chapman, S., Rhodes, C., Dennis, J., Brereton, P., 2002. The application
481 of isotopic and elemental analysis to determine the geographical origin of premium long grain rice.
482 *European Food Research and Technology*, 214, 72-78.
- 483 16. Kelly, S., Heaton, K., Hoogewerff, J., 2005. Tracing the geographical origin of food: The
484 application of multi-element and multi-isotope analysis. *Trends in Food Science & Technology*, 16,
485 555-567.
- 486 17. Martin, G. J., Martin, M. L., 2003. Climatic significance of isotope ratios. *Phytochemistry*
487 *Reviews*, 2, 179-190.
- 488 18. O'Leary, M. H., 1995. Environmental effects on carbon fractionation in terrestrial plants. In
489 *Stable Isotopes in the Biosphere*, Wada, E.; Yoneyama, T.; Minigawa, M.; Ando, T.; Fry, B. D.,
490 Eds. Kyoto University Press, Kyoto,; pp 78–91.
- 491 19. TRACE quality assurance report, 2009. In *Report of the FP6-TRACE Deliverable D15.9 “Final*
492 *WP15 QA Report”*; 16 September, 2009.
- 493 20. Rossmann, A., Haberhauer, G., Hölzl, S., Horn, P., Pichlmayer, F., Voerkelius, S., 2000. The
494 potential of multielement stable isotope analysis for regional origin assignment of butter. *European*
495 *Food Research and Technology*, 211, 32-40.
- 496 21. Salminen, R. (ed.). 2005. *Geochemical Atlas of Europe. Part 1: Background Information,*
497 *Methodology and Maps*. Geological Survey of Finland: Espoo.
- 498 22. Schellenberg, A., Chmielus, S., Schlicht, C., Camin, F., Perini, M.; Bontempo, L., Heinrich, K.,
499 Kelly, S. D., Rossmann, A., Thomas, F., Jamin, E., Horacek, M., 2010. Multielement stable isotope
500 ratios (H, C, N, S) of honey from different European regions. *Food Chemistry* 121, 770–777.
- 501 23. Sieper, H.-P., Kupka, H.-J., Williams, T., Rossmann, A., Rummel, S., Tanz, N., Schmidt, H.-L.,
502 2006. A measuring system for the fast simultaneous isotope ratio and elemental analysis of carbon,
503 hydrogen, nitrogen and sulfur in food commodities and other biological material. *Rapid*
504 *Communications in Mass Spectrometry*, 20, 2521-2527.
- 505 24. Simpkins, W. A., Louie, H., Wu, M., Harrison, M., Goldberg, D., 2000. Trace elements in
506 Australian orange juice and other products. *Food Chemistry*, 71, 423-433
- 507 25. Smith, B. N., Epstein, S., 1971. Two Categories of ¹³C/¹²C Ratios for Higher Plants. *Plant*
508 *Physiology*, 47, 380-384.

- 509 26. Smith, R. G., 2005. Determination of the Country of Origin of Garlic (*Allium sativum*) Using
510 Trace Metal Profiling. *Journal of Agricultural and Food Chemistry*, 53, 4041-4045
- 511 27. Suzuki, Y., Chikaraishi, Y., Ogawa, N. O., Ohkouchi, N., Korenaga, T., 2008. Geographical
512 origin of polished rice based on multiple element and stable isotope analyses. *Food Chemistry*, 109,
513 470-475.
- 514 28. Thode, H., 1991, Sulfur isotopes in nature and the environment: an overview. In *Stable*
515 *Isotopes. Natural and anthropogenic sulfur in the environment (SCOPE 43)*, Krouse, H.; Grinenko,
516 V., Eds. Wiley, Chichester,; pp. 2-25.
- 517 29. Voerkelius, S., Lorenz, G. D., Rummel, S., Quézel, C. R., Heiss, G., Baxter, M., Brach-Papa, C.,
518 Deters-Itzelsberger, P., Hoelzl, S., Hoogewerff, J., Ponzevera, E., Van Bockstaele, M.,
519 Ueckermann, H., 2010. Strontium isotopic signatures of natural mineral waters, the reference to a
520 simple geological map and its potential for authentication of food. *Food Chemistry* 118, 933-940.
521

522 **Tables**

523

524 Table 1

Site Name (country code)	Average GPS coordinates (deg., min. and sec.)	samples	Latitude	Marine/Inland	Bed rock geology	Combined 'identity'
Marchfeld (AT)	N48 13 43 E16 49 30	41	N	I	Sh	1
Gauboden (DE)	N48 49 05 E12 34 19	34	N	I		
Allgaeu (DE)	N48 03 37 E10 36 07	38	N	I		
Jylland (DK)	N56 18 43 E09 59 01	7/9*	N	A		3
Orkney (GB)	N58 58 07 W02 56 14	24	N	A		
Firenze (IT)	N43 57 29 E11 18 58	46	S	M		AM
Sicily (IT)	N37 45 14 E14 36 00	40	S	M		
Muehlviertel (AT)	N48 27 54 E14 04 29	25	N	I	4	
Limousin (FR)	N45 59 01 E02 12 56	40	S	I	5	
Cornwall (GB)	N50 04 33 W05 40 32	40	N	A	6	
Galicia (ES)	N43 03 34 W08 04 44	20	S	A	7	
Chalkidiki (GR)	N40 22 27 E23 36 41	40	S	M	8	
Fraenkische Alb (DE)	N49 57 04 E11 06 35	40	N	I	L	
Jura Krakowska (PL)	N50 10 00 E19 45 03	19/32*	N	I		9
Galway (IE)	N53 09 06 W08 56 52	15	N	A		10
Carpentras (FR)	N44 12 00 E05 19 20	40	S	M		11
Iceland (IS)	N63 32 24 W19 39 31	3	N	A	B	12

525

* *Only with the 10-variables model.*

526

Table 1. Sampling sites, average GPS coordinates, number of samples per site and classification categories: Latitude (North, N, and South, S), Marine/Inland (Inland, I; Atlantic, A; and Mediterranean, M) and Bed rock geology (Shale/mudstone/clay/loess incl. sandstone and other clastic sediments, Sh; Acid Magmatic, AM; Limestone, L; and Basaltic, B)

527

528

529

530 Table 2

Classification	Models	PCs	R2X(cum)	R2Y(cum)	Q2(cum)	R2Y(cum)-Q2(cum)
Latitude	5V	2 – (2)	0.52 – (0.51)	0.59 – (0.57)	0.58 – (0.57)	0.00 (0.01)
	10V	2 – (2)	0.33 – (0.32)	0.62 – (0.60)	0.62 – (0.59)	0.01 (0.01)
Marine/Inland	5V	3 – (3)	0.72 – (0.71)	0.50 – (0.46)	0.49 – (0.45)	0.00 (0.00)
	10V	4 – (4)	0.66 – (0.60)	0.67 – (0.65)	0.66 – (0.64)	0.01 (0.01)
Bed rock geology	5V	3 – (3)	0.73 – (0.73)	0.27 – (0.24)	0.26 – (0.23)	0.01 (0.01)
	10V	5 – (4)	0.77 – (0.56)	0.37 – (0.33)	0.35 – (0.30)	0.02 (0.04)

531

Table 2. Main PLS-DA parameters estimated for the 5-variables model (5V) and the 10-variables model (10V). PCs is the number of principal components. Figures between brackets correspond to results obtained for the simulated data set (addition to all original data of component corresponding to twice the value of the stated standard uncertainties multiplied by a randomly generated number ranging from -1 to +1)

532

533

534

535

536

537

538

539

540 Table 3

Sample site	n	Latitude		Marine/Inland		Bed rock geology	
		5-variables model	10-variables model	5-variables model	10-variables model	5-variables model	10-variables model
Marchfeld	8	100 - (100)	100 - (100)	88 - (100)	100 - (100)	100 - (88)	100 - (100)
Muehlviertel	5	100 - (100)	80 - (80)	100 - (100)	100 - (100)	40 - (20)	60 - (40)
Fraenkische Alb	8	75 - (75)	100 - (88)	75 - (75)	75 - (88)	0 - (0)	13 - (13)
Gauboden	7	100 - (100)	100 - (100)	71 - (86)	100 - (100)	14 - (0)	57 - (71)
Allgaeu	8	88 - (88)	100 - (100)	75 - (75)	100 - (88)	38 - (38)	88 - (63)
Jylland	1	100 - (100)	100 - (100)	100 - (100)	0 - (0)	0 - (0)	0 - (0)
Carpentras	8	100 - (100)	100 - (100)	100 - (100)	100 - (100)	63 - (75)	100 - (63)
Limousin	8	100 - (100)	88 - (88)	88 - (88)	88 - (88)	100 - (100)	88 - (88)
Chalkidiki	8	63 - (50)	75 - (63)	25 - (38)	13 - (25)	50 - (38)	50 - (50)
Galway	3	100 - (100)	100 - (100)	0 - (0)	100 - (100)	0 - (0)	0 - (0)
Firenze	9	56 - (56)	44 - (44)	44 - (44)	44 - (44)	78 - (67)	56 - (78)
Sicily	8	88 - (100)	88 - (88)	88 - (75)	88 - (88)	75 - (63)	88 - (75)
Jura Krakowska	7* / 4	75 - (75)	71 - (71)	75 - (75)	86 - (57)	0 - (0)	29 - (29)
Galicia	4	50 - (50)	75 - (50)	0 - (0)	50 - (25)	0 - (0)	50 - (25)
Cornwall	8	100 - (100)	75 - (75)	100 - (88)	100 - (100)	25 - (25)	63 - (63)
Orkney	5	100 - (100)	100 - (100)	100 - (100)	100 - (100)	80 - (80)	100 - (100)
Iceland	1	100 - (100)	100 - (100)	0 - (0)	100 - (100)	0 - (0)	0 - (0)

* Only with the 10-variables model

541

542

543

544

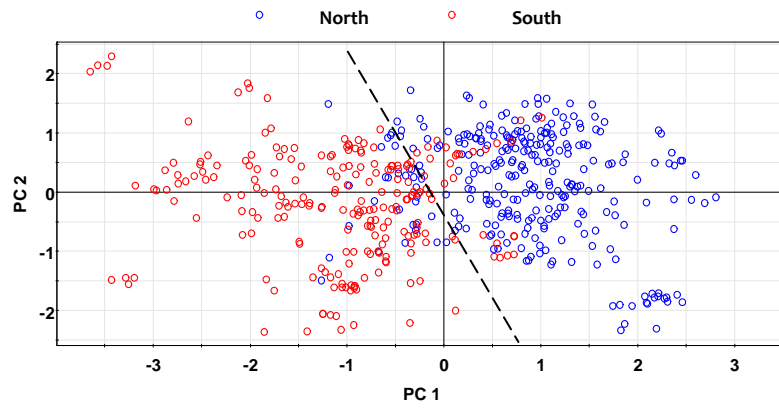
545

Table 3. External prediction results (%) for “Latitude”, “Marine/Inland” and “Bed rock geology” classification categories using 103 or 106 samples (for 5- and 10-variables models, respectively) as test data sets. External prediction results between brackets are for the data set simulated to estimate the impact of measurement uncertainty in these models.

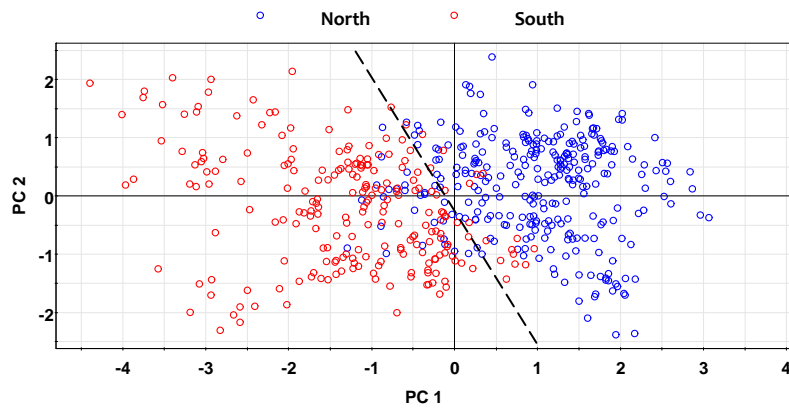
546 **Figures**

547

548 Figure 1



(a) 5-variables model ($n=512$)



(b) 10-variables model ($n=527$)

Figure 1. "Latitude" classification (classes North and South) PLS-DA score plots for the 5-variables (a) and the 10-variables (b) models

549

550

551

552

553

554

555

556

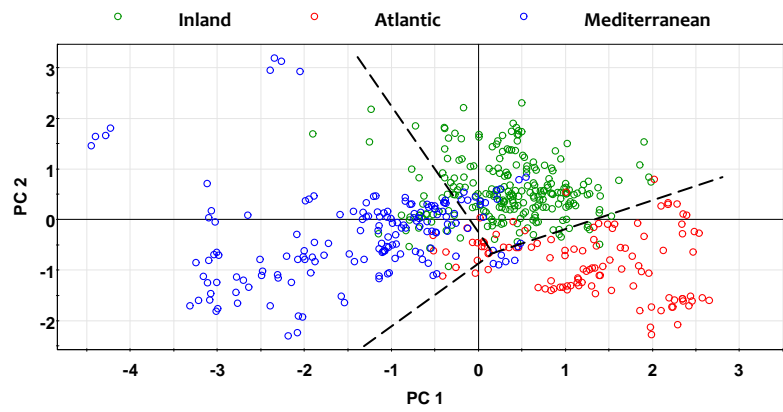
557

558

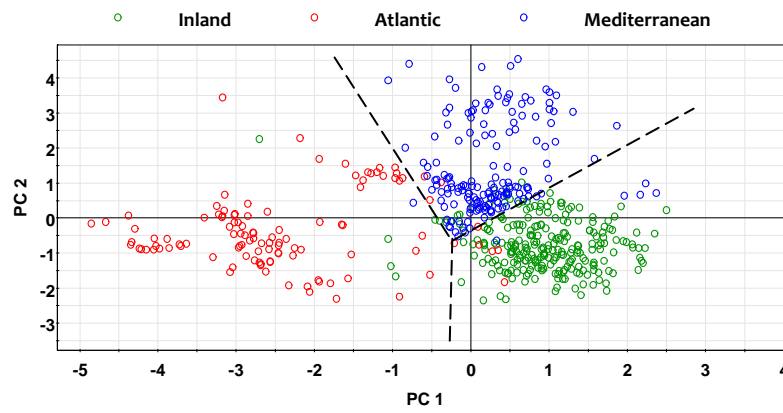
559

560

561



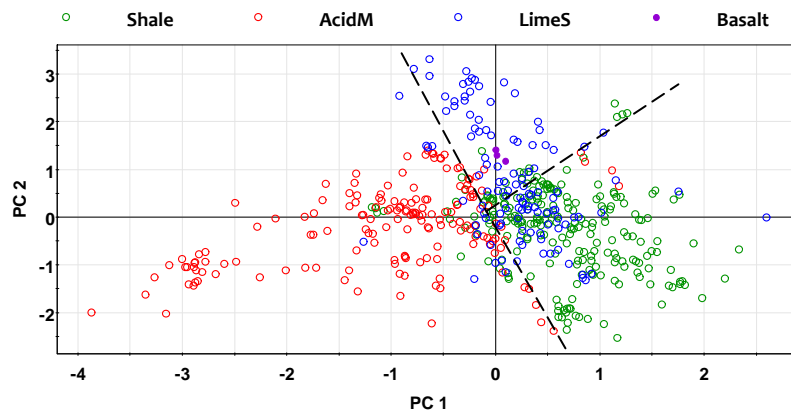
(a) 5-variables model (n=512)



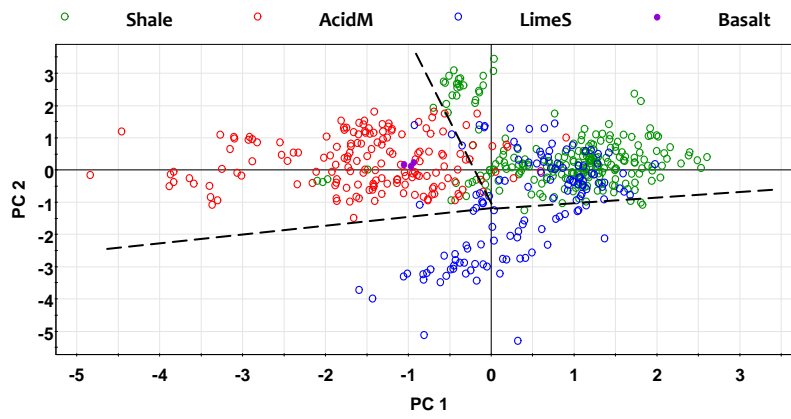
(b) 10-variables model (n=527)

Figure 2. “Marine/Inland” classification (classes Inland, Atlantic and Mediterranean) PLS-DA score plots for the 5-variables (a) and the 10-variables (b) models

563
564
565
566
567
568
569
570
571
572
573
574
575
576



(a) 5-variables model (n=512)



(b) 10-variables model (n=527)

Figure 3. “Bed rock geology” classification (classes Shale, Acid Magmatic, Limestone and Basaltic) PLS-DA score plots for the 5-variables (a) and the 10-variables (b) models