

1 Estimation of Composition of Quinoa (*Chenopodium quinoa* Willd.)

2 Grains by Near-Infrared Transmission Spectroscopy

3
4
5 Christian Encina-Zelada^{1,2,3}, Vasco Cadavez¹, Jorge Pereda², Luz Gómez-Pando⁴, Bettit
6 Salvá-Ruíz², José A. Teixeira³, Martha Ibañez⁴, Kristian H. Liland⁵, Ursula Gonzales-
7 Barron^{1*}

8
9 ¹CIMO Mountain Research Centre, School of Agriculture, Polytechnic Institute of Braganza,
10 Portugal.

11 ²Department of Food Technology, Faculty of Food Industries, National Agricultural
12 University La Molina, Lima, Peru.

13 ³Department of Biological Engineering, School of Engineering, University of Minho,
14 Portugal.

15 ⁴Cereals and Andean Crops Programme, Faculty of Agronomy, National Agricultural
16 University La Molina, Lima, Peru.

17 ⁵Nofima AS – Norwegian Institute of Food, Fisheries and Aquaculture Research,
18 Osloveien 1, N-1430, Ås, Norway

19
20
21
22 *Corresponding author: Ursula A. Gonzales-Barron; Phone: +351 273 303 325; E-mail:
23 ubarron@ipb.pt; Mailing address: School of Agriculture, Polytechnic Institute of Braganza
24 Campus de Santa Apolónia, Apartado 1172, 5301-854 Portugal

25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

ABSTRACT

The aim of this study was to develop robust chemometric models for the routine determination of dietary constituents of quinoa (*Chenopodium quinoa* Willd.) using Near-Infrared Transmission (NIT) spectroscopy. Spectra of quinoa grains of 77 cultivars were acquired while dietary constituents were determined by reference methods. Spectra were subjected to multiplicative scatter correction (MSC) or extended multiplicative signal correction (EMSC), and were (or not) treated by Savitzky-Golay (SG) filters. Latent variables were extracted by partial least squares regression (PLSR) or canonical powered partial least squares (CPPLS) algorithms, and the accuracy and predictability of all modelling strategies were compared. Smoothing the spectra improved the accuracy of the models for fat (root mean square error of cross-validation, RMSECV: 0.319 – 0.327%), ashes (RMSECV: 0.224 – 0.230%), and particularly for protein (RMSECV: 0.518 – 0.564%) and carbohydrates (RMSECV: 0.542 – 0.559%), while enhancing the prediction performance, particularly, for fat (root mean square error of prediction, RMSEP: 0.248 – 0.335%) and ashes (RMSEP: 0.137 – 0.191%). Although the highest predictability was achieved for ashes (SG-filtered EMSC/PLSR: bootstrapped 90% confidence interval for RMSEP: [0.376 – 0.512]) and carbohydrates (SG-filtered MSC/CPPLS: 90% CI RMSEP: [0.651 – 0.901]), precision was acceptable for protein (SG-filtered MSC/CPPLS: 90% CI RMSEP: [0.650 – 0.852]), fat (SG-filtered EMSC/CPPLS: 90% CI RMSEP: [0.478 – 0.654]) and moisture (non-filtered EMSC/PLSR: 90% CI RMSEP: [0.658 – 0.833]).

Keywords: Canonical, partial least squares, chemometrics, scatter correction, Savitzky-Golay

50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74

1. Introduction

Quinoa (*Chenopodium quinoa* Willd.) is a pseudocereal originating from the surroundings of the Titicaca Lake (Peru and Bolivia), which has been cultivated for centuries in the Andean countries. Quinoa is known as a pseudo-cereal because its seeds are used as cereal grains; although its nutritional quality is superior to that of the common cereals (Vega-Gálvez et al., 2010; Jancurová, Minarovicová, & Dandar, 2009).

Near infrared transmission (NIT) spectroscopy can presently provide rapid and accurate analysis of starch, moisture, protein, and oil contents in whole kernel cereals (Büchman, Josefsson & Cowe, 2001; Miralbés, 2004; and Pojić, Mastilović, Pestorić, & Radusin, 2008). However, when analysing intact samples by diffuse reflectance or transmittance spectroscopy, uncontrolled variations in light scattering are often a dominating artifact that complicates subsequent chemometric modelling (Panero, Panero, Panero, & Silva, 2013). This undesired scattering variation is due to uncontrolled physical variations of the samples, such as particle size and shape, sample packing, surface and orientation of the particles (Cantor, Hoag, Ellison, Khan, & Lyon, 2011). In order to minimise the multiplicative interference of scatter and particle size for the construction of robust models, NIT spectra are subjected to processing techniques for signal correction (i.e., multiplicative scatter correction and extended multiplicative signal correction) and noise removal (i.e., Savitzky-Golay derivatives).

75 Processed spectroscopy data matrices are then related with physicochemical data using
76 multivariate calibration methods (Ferreira, Pallone, & Poppi, 2015). Partial least squares
77 regression (PLSR) is currently considered as one of the most robust multivariate regression
78 techniques as it is associated with prediction errors that are lower than those of the principal
79 component analysis (Wold, Martens, & Wold 1983; Moghimi, Aghkhani, Sazgarnia, &
80 Sarmad, 2010). Recently, a generalisation of PLSR has been proposed that incorporates
81 discrete and continuous responses, additional measurements, and individual weighting of
82 observations. The technique is known as Canonical Powered Partial Least Squares (CPPLS)
83 because the optimal latent variables are found by combining PLS methodology and canonical
84 correlation analysis (Indahl, Liland, & Næs, 2009; Mevik, Wehrens, & Liland, 2015). Thus,
85 the objective of this study was three-fold: (i) to assess the feasibility of accurately quantifying
86 dietary constituents of quinoa (moisture, protein, fat, ashes and carbohydrates) whole grains
87 by NIT spectroscopy; (ii) to compare the robustness and prediction capability of the PLSR
88 and CPPLS multivariate models after scatter correction of the spectra; and (iii) to assess to
89 what extent smoothing filters applied to scatter-corrected spectra can further improve the
90 performance of the PLSR and CPPLS algorithms.

91

92 **2. Methodology**

93

94 **2.1 Samples and proximate composition analysis**

95

96 The samples utilised in this study were quinoa (*Chenopodium quinoa* Willd.) whole grains of
97 orange, beige, black and yellow colour, corresponding to 77 different cultivars. They were all
98 harvested in Peru at the National Agricultural University La Molina (Lima) and the Regional
99 Development Centre – Highland (Junin), between 2010 and 2012. Moisture, protein, fat and

100 ashes contents were determined in triplicate using the reference methods 925.10, 920.87
101 (conversion factor of 6.25), 923.05 and 923.03, respectively, as described by the Association
102 of Official Agricultural Chemists (AOAC, 2000). Total carbohydrate content was calculated
103 by difference as: 100 - (weight in grams [protein + fat + water + ashes] in 100 g of quinoa).
104 Proteins, fat, ashes and carbohydrate contents were then converted into dry basis (db).

105

106 **2.2 Near-infrared transmission (NIT) spectra acquisition**

107

108 NIT spectra were acquired by placing the whole grains directly in an Infratec 1241 grain
109 analyser (Module Foss Tecator, Denmark), using 60-mm quartz cuvettes, and scanning the
110 region 850-1048 nm (wavenumber range of 11765 – 9524 cm^{-1}). The spectra were recorded at
111 scanning step intervals of 2 nm to give 100 data points per sample. A total of 10 frequency
112 scans were performed per sample, and carefully assessed for consistency. Raw spectral data
113 (i.e., a vector of 100 data points per sample) were linked to the chemical analyses data on a
114 spreadsheet. To correct for the non-linearity in the measure of transmittance (T), T was
115 transformed into absorbance (A) by taking the base 10 logarithm of the reciprocal of the
116 transmittance values ($A = \log 1/T$).

117

118 **2.3 NIT spectral pre-processing**

119

120 To minimise the multiplicative effects of light scattering, spectra were subjected to
121 multiplicative scatter correction (MSC) or extended multiplicative signal correction (EMSC).
122 MSC is a transformation method used to compensate for additive and multiplicative effects in
123 spectral data (Maleki, Mouazen, Ramon, & De Baerdemaeker, 2007). Both EMSC and MSC
124 attempt to separate physical light scattering effects from chemical (vibrational) light

125 absorbance, yet EMSC is a modification of the standard MSC which adds polynomials to the
126 correction model in addition to the constant baseline effect and reference scaling of MSC
127 (Martens & Stark, 1991; Panero et al., 2013). The basic EMSC with polynomials of degree 2
128 was applied. For each of the dietary constituents analysed, PLSR and CPPLS multivariate
129 models were then fitted to the MSC- or EMSC- pre-processed spectra; thereby producing four
130 treatments (MSC/PLSR, EMSC/PLSR, MSC/CPPLS and EMSC/CPPLS) which were
131 compared in terms of predictability.

132

133 In addition, Savitzky-Golay (SG) derivative filters (Savitzky and Golay, 1964) were applied
134 after correcting spectra for scattering (MSC or EMSC) to assess whether the predictive
135 performance of the PLSR and CPPLS models could be further enhanced. SG smoothing
136 performs a piece-wise polynomial fitting with specified polynomial degree (p), window
137 length (w), and derivative order (m) to the spectrum. Thus, SG filters produced by all possible
138 combinations of $m=\{1, 2\}$, $p=\{2, 3, 4\}$ and $w=\{3, 5, 7, 9, 11\}$ were applied to each of the
139 MSC and EMSC scatter-corrected spectra.

140

141 **2.4 Chemometric multivariate data analysis**

142

143 The extraction of information from quinoa grain's pre-processed spectra to estimate moisture,
144 protein, fat, ashes and carbohydrates contents was performed by the PLSR and CPPLS
145 chemometric algorithms. For the CPPLS models estimating moisture content, the additional
146 variables were protein, fat, ashes and quinoa cultivar. For the estimation of protein by CPPLS,
147 the additional variables were moisture, fat, ashes and cultivar; whereas for the estimation of
148 fat, the additional variables were moisture, protein and ashes. The additional variables for
149 ashes content CPPLS models were moisture, fat and quinoa cultivar, while those for

150 carbohydrates content were moisture, ashes and fat. Selection of the additional variables for
151 each dietary constituent's CPPLS model was carried out by trial and error.

152

153 As a first step, the full data set was divided into a subset for calibration (~80% data, 62
154 samples) and the remaining ~20% (15 samples) for prediction or validation, by means of
155 random split stratified by cultivar. PLSR and CPPLS were fitted separately to MSC and
156 EMSC scatter-corrected spectra with and without SG filters. The performance of the different
157 models (a model is defined as a combination of a pre-processing filter and a chemometric
158 multivariate algorithm) was determined by *cross-validation* as an *internal calibration* method
159 using the calibration data set. In our case, the leave-one-out (LOO) method was used. Briefly,
160 in the LOO method, each sample is removed one at a time from the calibration set, a new
161 calibration performed and a prediction score calculated for the sample removed. This
162 procedure is repeated until every sample has been left out once. The performance of the
163 model was assessed by the root mean square error of cross-validation (RMSECV), which is
164 deemed as the best single estimate of the prediction capability of the model (González-
165 Martín, Moncada, Fischer, & Escuredo 2014; Mevik & Wehrens, 2007). Then, the optimal
166 number of components of a model was selected at the first RMSECV local minimum, rather
167 than the absolute minimum (to avoid overfitting). For such a number of components, the root
168 mean square error of calibration (RMSEC) was computed. In addition, the coefficients of
169 correlation between reference values and values fitted by cross-validation (R_{CV}) and the
170 calibration model (R_C) were computed.

171

172 Following completion of the calibration, models were validated using the prediction data set.
173 Model performance was evaluated by obtaining the root mean square error of prediction
174 (RMSEP) and the coefficient of correlation (R_P) between reference values and those predicted

175 by the model. For each of the four treatments (i.e., MSC/PLSR, EMSC/PLSR, MSC/CPPLS
176 and EMSC/CPPLS), the SG filters leading to the highest accuracy were identified. To assess
177 the best model(s) for each dietary constituent, the model had to present not only a low RMSE
178 but also a high R. The entire NIT spectra analysis was conducted using the “pls” (Mevik et
179 al., 2015), “emsc” (Liland, 2016) and the “prospectr” (Stevens & Ramirez-López, 2013)
180 packages implemented in the R software version 3.2.5 (R Core Team, 2016).

181

182 **3. Results and Discussion**

183

184 **3.1 Proximate composition analysis of quinoa**

185

186 The values reported in this study for fat (5.35 – 7.78% db) and ashes (2.51 – 4.11% db; Table
187 1) were comparable to those reported by Repo-Carrasco-Valencia, Hellström, Pihlava, &
188 Mattila (2010) for six ecotypes of similar Peruvian quinoa (fat: 4.36-7.59% db, and ashes:
189 2.57-3.44% db). However, they found considerably higher protein content (12.55-16.08% db)
190 and lower carbohydrates content (67.13-77.02% db) than those found in this report (8.33 –
191 11.38% db; and 78.48 – 82.89% db, respectively). Analysing quinoa samples from Peru,
192 Bolivia and Brazil, Ferreira et al. (2015) encountered substantially higher fat (6.19 – 15.52%
193 db) and ashes (3.07 – 9.15% db) contents than those of our study. The variation in ashes are
194 influenced by the dependence of the mineral content on type of soil and fertiliser application.
195 Moisture is the compound most variable among published studies (from 8.26-11.51% in
196 Repo-Carrasco-Valencia et al. (2010) up to 25.66 – 33.16% in Ferreira et al. (2015)) because
197 it depends upon drying and storage of seeds. The standard deviations suggest that sufficient
198 variation in the dietary compounds existed among the quinoa cultivars to develop
199 chemometric models.

200

201 **3.2 Pre-processing methods for signal correction and smoothing of quinoa's NIT spectra**

202

203 The first step of signal pre-treatment is crucial as redundant information should be removed
204 from the spectra. With corrected spectra, the repeatability and reproducibility of the
205 chemometric multivariate model can be increased (Stevens & Ramirez-Lopez, 2013). In the
206 first instance, the transmittance spectra of the quinoa grains without any processing pointed to
207 the occurrence of multiplicative scaling effects (Figure 1, top left), which were still present
208 when spectra were transformed into absorbance (Figure 1, top right). Such transformation is
209 needed to move signal processing to a domain where Beer-Lambert's law applies and additive
210 effects of compounds are linear. Light scattering, one of the main causes of multiplicative
211 scale effects (i.e., scale differences) in spectral data, was corrected by both methods, MSC
212 (Figure 1, bottom left) and EMSC (Figure 1, bottom right), although the application of EMSC
213 yielded a better signal correction. Whereas MSC was developed to remove both scaling
214 effects (a multiplicative factor) and baseline shift effects (an additive factor), EMSC was
215 designed to allow the separation of multiplicative physical effects (path length, light
216 scattering, etc.) from additive chemical effects (absorbance of analytes and interferants) and
217 additive physical effects (temperature shifts, baseline variations, etc.) (Panero et al., 2013).
218 Hence, additive effects, chemical and/or physical, must have been also present in the raw
219 spectra.

220

221 In general, when SG first (SG1) and second (SG2) derivative filters were applied to either the
222 MSC- or the EMSC-corrected spectra, the peaks below and above the baseline were
223 emphasised. It was not unexpected that EMSC+SG pre-processing (Figure 2, bottom)
224 produced cleaner signals than MSC+SG pre-processing (Figure 2, top), as EMSC yielded a

225 better correction for light scattering and additive effects than MSC. However, whether the
226 application of SG1 or SG2 pre-processing smoothing filter produces better signals should be
227 determined by the resulting predictive capacity of the chemometric models.

228

229 **3.3 Comparisons between scatter correction methods and multivariate algorithms**

230

231 For moisture, protein and ashes contents, regardless of the chemometric algorithm used (i.e.,
232 PLSR or CPPLS), the application of EMSC to the spectra produced lower errors (i.e.,
233 RMSECV) by up to ~4.8% in the case of protein, than those produced by MSC treatments
234 (Table 2). Comparing EMSC and MSC performance, Panero et al. (2013) similarly found
235 lower RMSEC and RMSEP values when applying the former scatter correction method on
236 marzipan spectra for NIR determination of moisture. Correspondingly, for moisture, protein
237 and ashes contents, correcting the signal scatter by EMSC led to higher R_{cv} values (range of
238 0.572 – 0.769) than those produced by the simpler MSC (0.564 – 0.742; Table 2).
239 Considering that the models fitted to EMSC-processed spectra consistently led to fewer
240 optimal components (3 – 7) than those fitted to MSC-processed spectra (4 – 8), it can be
241 stated that EMSC, with their resulting lower cross-validation errors and higher cross-
242 validation correlation coefficients, had a tendency to produce more robust models than MSC
243 for the NIT determination of moisture, protein and ashes. Nevertheless, in the cases of fat and
244 carbohydrates, irrespective of the algorithm used for model calibration, the behaviour was the
245 opposite; this is, MSC-treated spectra yielded more robust chemometric models – as implied
246 by their lower RMSECV and higher R_{cv} – than the EMSC-treated spectra did, although with
247 at most one more component (Table 2). For fat and carbohydrates, EMSC may have overfitted
248 the baseline such that chemical information was discarded along with the scatter correction.

249

250 The multivariate regression methods also affected the accuracy of prediction for the models.
251 In the analyses of all dietary components, the CPPLS algorithm led invariably to a selection
252 of fewer optimal components (3-5) than PLSR (6-8). This was an anticipated outcome since
253 CPPLS was developed as a compression method for the extraction of more predictive
254 information in the first few components than ordinary PLSR (Indahl et al., 2009). For this
255 reason, within each dietary constituent, the models with the combination CPPLS/EMSC
256 yielded the lowest optimal number of components (3-4) while the combination PLSR/MSC
257 yielded the highest optimal number of components (7-8). For instance, for the protein
258 constituent, the 8 optimal latent variables in the combination PLSR/MSC was brought down
259 to 3 in the combination CPPLS/EMSC. In all dietary constituents – except fat – there was a
260 clear effect of the multivariate regression on the RMSEC and RMSEP values, being the
261 CPPLS algorithm associated to higher errors (Table 2).

262

263 With the exception of carbohydrates, when the quinoa grains' spectra were MSC scatter-
264 corrected, the use of the PLSR or CPPLS algorithm produced very similar cross-validation
265 errors (RMSECV) for the estimation of moisture (0.575; 0.579%), protein (0.614; 0.613%),
266 fat (0.326; 0.325%) and ashes (0.231; 0.233%). However, the effect of the regression
267 algorithm on RMSECV values became more noticeable when spectra were pre-processed by
268 EMSC for the chemometric models determining moisture (RMSECV: 0.566; 0.578%) and
269 carbohydrates (0.620; 0.638%). When applied to EMSC-treated spectra, the PLSR algorithm
270 produced more accurate models – lower RMSECV in all dietary constituents – than those
271 produced by CPPLS. Even for moisture, protein and ashes, the PLSR/EMSC treatment
272 yielded the highest R_{cv} and R_c values among the four treatments. This may arise from the
273 higher optimal number of components consistently picked by the PLSR algorithm (Table 2).

274

275 Earlier, Ferreira et al. (2015) proposed a series of chemometric models to estimate the
276 proximate composition of quinoa from Fourier transform near-infrared (FTIR) spectra. In
277 order to contrast the accuracy of our models with their FTIR models, the coefficient of
278 variation ($CV=RMSECV/\text{mean}$) was calculated as a common metric for comparison since it is
279 a dimensionless number less sensitive to difference in means. The chemometric models
280 presented in this study were more accurate than those obtained in Ferreira et al. (2005), as
281 indicated by the considerably lower CV of our models for moisture (5.3 – 5.5% as opposed to
282 5.9%), protein (5.8 – 6.2% as opposed to 14.9%), fat (4.9 – 5.2% as opposed to 11.7%),
283 carbohydrates (0.73 – 0.79% as opposed to 7.0%) and ashes (7.0 – 7.4% as opposed to
284 15.5%). Similarly, the external validation CV ($RMSEP/\text{mean}$) obtained from our models for
285 protein (5.5 – 6.4%) and fat (5.6 – 4.1%) were far lower than those reported by González-
286 Martín et al. (2013) (10.4% and 8.3%, respectively). Nonetheless, when contrasting the
287 estimates of correlation between the reference and the spectral methods, the R_{CV} (0.56 – 0.77)
288 and R_C (0.51 – 0.83; Table 2) found in our models were, as a whole, lower than those reported
289 by both González-Martín et al. (2013) (R_{CV} : 0.89 – 0.96) and Ferreira et al. (2015) (R_C : 0.86 –
290 0.91). The lower correlation coefficients encountered in this study may have been a
291 manifestation of our effort to avoid overfitting by consistently selecting the number of latent
292 variables that minimise RMSECV. Moreover, by definition, the coefficient of determination
293 tends to decrease when the range of the dependent variable is lower. The ranges of protein
294 (8.33 – 11.4% db), fat (5.35 – 7.78%), carbohydrates (78.5 – 82.9%) and ashes (2.51 –
295 4.11%) assayed from our quinoa samples were narrow in comparison to those from the quinoa
296 samples surveyed in Ferreira et al. (2015) (protein: 11.4 – 36%, fat: 6.19 – 15.52%,
297 carbohydrates: 43.6 – 76.4% and ashes: 3.07 – 9.15%).

298

299 **3.4 Influence of SG derivative filters on robustness of chemometric models**

300

301 Table 3 compiles the SG combinations (m, p, w) leading to the highest predictability within
302 each of the four treatments (i.e., MSC/PLSR, EMSC/PLSR, MSC/CPPLS and
303 EMSC/CPPLS). Although for protein, the same SG filter type (m=1, p=2, w=9) produced the
304 best model's accuracy in the four treatments, this did not necessarily hold for the other dietary
305 constituents (Table 3).

306

307 Regardless of the signal correction method and the multivariate algorithm used, SG filtering
308 of quinoa's spectra improved the accuracy of the chemometric models, yet to different
309 degrees: the reduction in RMSECV and RMSEC in the models for moisture (reduction by 1.3
310 – 2.6% and 8 – 14%, respectively), fat (1.5 – 5.3% and 0.4 – 1.1%) and ashes (2.1 – 2.2% and
311 2.1 – 10.6%) were all slight in comparison to the considerable reduction in those statistics in
312 the models for protein (8.0 – 11.9% and 20.5 – 28.5%) and carbohydrates (8.9 – 12.4% and
313 24.2 – 35.0%). Similarly, SG-filtering improved the correlation statistics of calibration: as
314 before, the increase in R_{CV} and R_C values was slight in the models for moisture (increase by
315 2.6 – 5.2% and 0 – 6.4%, respectively), fat (1.4 – 5.0% and 0 – 0.5%) and ashes (0 – 1.8%
316 and 1.1 – 7.1%), whereas the improvement was substantial in the models for protein (13.9 –
317 17.3% and 15.6 – 42.2%) and carbohydrates (8.0 – 14.5% and 10.8 – 33%) (percentual
318 differences not shown but calculated from Table 2 and 3).

319

320 The improved RMSECV, RMSEC, R_{CV} and R_C statistics from the models with SG filters for
321 protein and carbohydrates, may be associated to the fact that, for protein and carbohydrates,
322 filtering the spectra led to a higher number of optimal components in the MSC/PLSR (from 8
323 to 12, and 7 to 12, respectively), EMSC/PLSR (6 to 10, and 7 to 10), MSC/CPPLS (4 to 8,
324 and 4 to 10) and EMSC/CPPLS (3 to 6, and 3 to 8) models. Due to the higher number of

325 components extracted from the SG spectra, the fitting capacity of the protein and
326 carbohydrates models was improved; although the CPPLS algorithm performed better than
327 the PLSR algorithm in the prediction of the test data – as suggested by the differences in
328 RMSEP and R_P . Filtering the spectra with SG largely enhanced the predictive capacity of the
329 models for fat (RMSEP decreased by 1.0 – 20.4%, and R_P increased by 1.8 – 24.7%) and
330 ashes (RMSEP decreased by 0.0 – 30.8%, and R_P increased by 0.0 – 32.3%), while, as
331 mentioned before, filtering enhanced the prediction performance of the models for protein
332 (RMSEP decreased by 15.8%, and R_P increased by 19.8%), and carbohydrates (RMSEP
333 decreased by 24.8%, and R_P increased by 30.6%) only when CPPLS was used. In the
334 particular case of moisture, only the treatment MSC/CPPLS produced better predictions when
335 spectra were SG-filtered (RMSEP decreased by 10.4%, and R_P increased by 14.1%).

336

337 **3.5 Validated chemometric models for quinoa's dietary constituents**

338

339 Taking the four treatments together (Table 3), the models estimating ashes and carbohydrates
340 presented generally the highest predictive capacity, as deduced from the ranges of R_{CV} (0.744
341 – 0.761; and 0.750 – 0.767, respectively) and R_P (0.847 – 0.925; and 0.728 – 0.807,
342 respectively). However, the models for protein (R_{CV} : 0.651 – 0.717; R_P : 0.625 – 0.760) and
343 fat (R_{CV} : 0.716 – 0.732 ; R_P : 0.565 – 0.804) were of slightly lower predictive performance,
344 while the models for moisture (R_{CV} : 0.504 – 0.611; R_P : 0.441 – 0.539) were of fair
345 predictability.

346

347 Considering that a good model should bear low values of RMSECV and RMSEP, and high
348 values of R_{CV} and R_P , the final model for each quinoa's constituent was selected among those
349 presented in Table 2 and 3. For the moisture response, little-to-no gain in prediction

350 performance was attained by SG-filtering the spectra with the many combinations tested.
351 Thus, for this variable, the best model was achieved using a non-filtered spectra treated by
352 MSC and extracting 8 PLSR components, which rendered a prediction CV (RMSEP/mean) of
353 5.60% and an R_P of 0.596 (other statistics for this model pointed out in bold in Table 2). For
354 the other dietary constituents, better performance was achieved using SG-filtered spectra of
355 window size 9 and first derivative, except for the fat variable which used second derivative.
356 For the NIT determination of ashes, the PLSR algorithm also produced the best model when
357 fitted to EMSC-treated spectra. The 5 optimal latent variables extracted yielded on the test
358 data a CV of 4.38% and R_P of 0.925. For the protein, fat and carbohydrates variables, the
359 CPPLS multivariate algorithm performed better: whilst the best predictability of protein
360 (CV=5.35% and $R_P=0.760$) was achieved by extracting 8 components from MSC-treated
361 spectra, the best model for carbohydrates was produced by extracting 10 components from
362 MSC-treated spectra (CV=0.80% and $R_P=0.807$). With a CV=3.79% and $R_P=0.804$, fat could
363 be estimated by a CPPLS model produced from a EMSC-treated spectra with only 3 latent
364 variables.

365

366 Finally, in order to further characterise the prediction performance of each of the final
367 models, uncertainty about the correlation coefficient of prediction (R_P) was built by
368 bootstrapping. At each of the 1000 iterations, a new 80% calibration/20% validation data
369 partition was randomly obtained, the chosen model was fitted to the calibration data with the
370 pre-determined number of components, and R_P was extracted from the test data. The
371 histograms of R_P built for each of the final models (Figure 3, left) show that the NIT model
372 for estimating ashes had the lowest uncertainty (i.e., narrow spread) about R_P , and therefore
373 was the most robust chemometric model. The wider spread of the R_P histogram for moisture
374 corroborated that, among the five dietary constituents studied, the model for moisture

375 presented the lowest precision. The degree of fitting and predictability of the final models can
376 be appreciated from the scatter plots between the reference values and those fitted (Figure 3,
377 middle) and predicted (Figure 3, right) from the NIT calibration models. The best agreement
378 between observed and predicted values was observed for ashes and carbohydrates; although,
379 as a whole, the degree of dispersion in the predictions is acceptable, bearing in mind that
380 chemical analyses also have associated errors.

381

382 **4. Conclusions**

383

384 Regardless of the multivariate algorithm used, light scattering correction of quinoa grains'
385 NIT spectra by EMSC consistently led to proximate composition models of better cross-
386 validation statistics – except for fat and carbohydrates – than those produced by MSC-treated
387 spectra. Both EMSC, as opposed to MSC; and CPPLS, as opposed to PLSR, led to fewer
388 optimal components. When spectra were treated by different types of SG filters, the optimal
389 latent variables reduced correspondingly in each of the four treatments (i.e., MSC/PLSR,
390 EMSC/PLSR, MSC/CPPLS, EMSC/CPPLS), except for the models predicting protein and
391 carbohydrates, in which the behaviour was the opposite. In addition, smoothing the quinoa's
392 spectra enhanced the accuracy of the models for fat, ashes, and particularly for protein and
393 carbohydrates, while improving also the prediction performance, particularly, for fat and
394 ashes determination. Although the most robust models could be developed for ashes (SG-
395 filtered EMSC/PLSR: 90% confidence interval for RMSEP [0.376 – 0.512] as determined by
396 bootstrap) and carbohydrates (SG-filtered MSC/CPPLS: 90% CI RMSEP: [0.651 – 0.901]),
397 the predictability was still acceptable for the other dietary constituents; namely, protein (SG-
398 filtered MSC/CPPLS: 90% CI RMSEP: [0.650 – 0.852]), fat (SG-filtered EMSC/CPPLS:
399 90% CI RMSEP: [0.478 – 0.654]) and moisture (non-filtered EMSC/PLSR: 90% CI RMSEP:

400 [0.658 – 0.833]). Thus, in this study, satisfactory predictions of the dietary constituents of
401 quinoa grains could be achieved by using NIT technology. The main advantages of the
402 technique are the rapid determination for routine analysis, the reduced costs and absence of
403 sample preparation and waste generation.

404

405 **Acknowledgments**

406

407 Mr. Encina-Zelada acknowledges the financial aid provided by the Peruvian National
408 Programme of Scholarships and Student Loans (PRONABEC) in the mode of PhD grants
409 (Presidente De La República-183308). Dr. Gonzales-Barron wishes to acknowledge the
410 financial support provided by the Portuguese Foundation for Science and Technology (FCT)
411 through the award of a five-year Investigator Fellowship (IF) in the mode of Development
412 Grants (IF/00570).

413

414 **References**

415

- 416 1. AOAC. (2000). Official methods of analysis of the Association of Analytical
417 Chemists International. In W. Horwitz (Eds.), 17th ed. AOAC International, Gaithersburg,
418 MD, USA.
- 419 2. Büchman, N. B., Josefsson, H., & Cowe, I. A. (2001). Performance of European
420 artificial neural network (ANN) calibrations for moisture and protein in cereals using the
421 Danish near infrared transmission (NIT) network. *Cereal Chemistry*, 78 (5), 572-577.
- 422 3. Cantor, S. L., Hoag, S. W., Ellison, C. D., Khan, M. A., & Lyon, R. C. (2011). NIR
423 spectroscopy applications in the development of a compacted multiparticulate system for

- 424 modified release. *Journal of the American Association of Pharmaceutical Scientists*, 12 (1),
425 262-278.
- 426 4. Ferreira, D. S., Pallone, J. A. L., & Poppi, R. J. (2015). Direct analysis of the main
427 chemical constituents in *Chenopodium quinoa* grain using Fourier transform near-infrared
428 spectroscopy. *Food Control*, 48, 91-95.
- 429 5. González-Martín, M. I., Moncada, G. W., Fischer, S., & Escuredo, O. (2014).
430 Chemical characteristics and mineral composition of quinoa by near-infrared spectroscopy.
431 *Journal of the Science of Food and Agriculture*, 94 (5), 876–881.
- 432 6. Indahl, U. G., Liland, K. H., & Næs, T. (2009). Canonical partial least squares - a
433 unified PLS approach to classification and regression problems. *Journal of Chemometrics*, 23,
434 495–504.
- 435 7. Jancurová, M., Minarovicová, L., & Dandar, A. (2009). Quinoa - a Review. *Czech*
436 *Journal of Food Sciences*, 27 (2), 71-79.
- 437 8. Liland, K. H. (2016). Extended Multiplicative Signal Correction. Package “EMSC”.
438 Date 2016-04-24. Repository CRAN. Available online at: [https://cran.r-](https://cran.r-project.org/web/packages/EMSC/index.html)
439 [project.org/web/packages/EMSC/index.html](https://cran.r-project.org/web/packages/EMSC/index.html) (Accessed: 16.05.2016).
- 440 9. Maleki, M. R., Mouazen, A. M., Ramon, H., & De Baerdemaeker, J. (2007).
441 Multiplicative scatter correction during on-line measurement with near infrared spectroscopy.
442 *Biosystems Engineering*, 96 (3), 427-433.
- 443 10. Martens, H., & Stark, E. (1991). Extended multiplicative signal orrection and spectral
444 interference subtraction: new preprocessing methods for near infrared spectroscopy. *Journal*
445 *of Pharmaceutical and Biomedical Analysis*, 9 (8), 625-635.
- 446 11. Mevik, B. H., & Wehrens, R. (2007). The pls package: principal component and
447 partial least squares regression in R. *Journal of Statistical Software*, 18 (2), 1-24.

- 448 12. Mevik, B. H., Wehrens, R., & Liland, K. H. (2015). Pls: Partial Least Squares and
449 Principal Component Regression. R package version 2.5-0. Available online at: [https://cran.r-](https://cran.r-project.org/web/packages/pls/)
450 [project.org/web/packages/pls/](https://cran.r-project.org/web/packages/pls/) (Accessed: 16.05.2016).
- 451 13. Miralbés, C. (2004). Quality control in the milling industry using near infrared
452 transmittance spectroscopy. *Food Chemistry*, *88* (4), 621-628.
- 453 14. Moghimi, A., Aghkhani, M. H., Sazgarnia, A., & Sarmad, M. (2010). Vis/NIR
454 spectroscopy and chemometrics for the prediction of soluble solids content and acidity (pH)
455 of kiwifruit. *Biosystems Engineering*, *106* (3), 295-302.
- 456 15. Panero, P. S., Panero, F. S., Panero, J. S., & Silva, H. E. B. (2013). Application of
457 extended multiplicative signal correction to short-wavelength near infrared spectra of
458 moisture in marzipan. *Journal of Data Analysis and Information Processing*, *1* (3), 30-34.
- 459 16. Pojić, M., Mastilović, J., Pestorić, M., & Radusin, T. (2008). The ensuring of
460 measurements for cereal quality determination. *Food Processing, Quality and Safety*, *35* (1),
461 11-18.
- 462 17. R Core Team. (2016). R: A language and environment for statistical computing. R
463 Foundation for Statistical Computing, Vienna, Austria. Available online at: [http://www.R-](http://www.R-project.org/)
464 [project.org/](http://www.R-project.org/) (Accessed: 04.02.2016).
- 465 18. Repo-Carrasco-Valencia, R., Hellström, J. K., Pihlava, J. M., & Mattila, P. H. (2010).
466 Flavonoids and other phenolic compounds in Andean indigenous grains: Quinoa
467 (*Chenopodium quinoa*), kañiwa (*Chenopodium pallidicaule*) and kiwicha (*Amaranthus*
468 *caudatus*). *Food Chemistry*, *120* (1), 128-133.
- 469 19. Savitzky, A., & Golay, M. (1964). Smoothing and differentiation of data by simplified
470 least squares procedures. *Analytical Chemistry*, *36*, 1627-1639.

- 471 20. Stevens, A., & Ramirez-Lopez, L. (2013). An introduction to the prospectr package.
472 Vignette R package version 0.1.3. Available online
473 at: <https://github.com/antoinestevens/prospectr> (Accessed: 16.05.2016).
- 474 21. Vega-Gálvez, A., Miranda, M., Vergara, J., Uribe, E., Puente, L., & Martínez, E.
475 (2010). Nutrition facts and functional potential of quinoa (*Chenopodium quinoa* willd.), an
476 ancient Andean grain: a review. *Journal of the Science of Food and Agriculture*, *90* (15),
477 2541–2547.
- 478 22. Wold, H., Martens, H., & Wold, S. (1983). The multivariate calibration method in
479 chemistry solved by the PLS method. In A. Ruhe, & B. Kågström (Eds.), *Proceedings of the*
480 *Conference of Matrix Pencils, Lecture Notes in Mathematics* (pp. 286–293). Springer Verlag:
481 Heidelberg.

482

483

484 TABLE CAPTIONS

485

486 Table 1. Summary statistics of the major dietary compounds of quinoa samples in % dry
487 basis, except for moisture (% wet basis)

488

489 Table 2. Accuracy of prediction of NIT chemometric models for quinoa constituents defined
490 by signal correction type (MSC: multiplicative scatter correction, or EMSC: extended
491 multiplicative signal correction) and multivariate algorithm (PLSR: partial least squares
492 regression, or CPPLS: canonical powered partial least squares), as measured by the root mean
493 square errors of cross-validation (RMSECV), calibration (RMSEC) and prediction (RMSEP),
494 and the coefficients of correlation between reference values and those estimated by cross-

495 validation (R_{CV}), calibration (R_C) and prediction (R_P), all of them computed at the minimum
496 number of components

497
498 Table 3. Effect of the best Savitzky-Golay smoothing filter (m: derivative order, p:
499 polynomial order and w:window size) on the accuracy of prediction of NIT chemometric
500 models for quinoa constituents defined by signal correction type (MSC: multiplicative scatter
501 correction, or EMSC: extended multiplicative signal correction) and multivariate algorithm
502 (PLSR: partial least squares regression, or CPPLS: canonical powered partial least squares),
503 as measured by the root mean square errors of cross-validation (RMSECV), calibration
504 (RMSEC) and prediction (RMSEP), and the coefficients of correlation between reference
505 values and those estimated by cross-validation (R_{CV}), calibration (R_C) and prediction (R_P), all
506 of them computed at the minimum number of components

507
508

509 FIGURE CAPTIONS

510
511 Figure 1. Untransformed or raw near-infrared transmittance spectra of quinoa whole grains
512 (top left), spectra transformed into absorbance (top right), and absorbance spectra corrected
513 for scattering applying multiplicative scatter correction (MSC; bottom left) or extended
514 multiplicative signal correction (EMSC; bottom right)

515
516 Figure 2. Effects of applying Savitzky-Golay first- (SG1; left) and second-derivative (SG2;
517 right) with polynomial degree 3 and window size 5 to quinoa grains spectra previously
518 corrected by multiplicative scatter correction (MSC; top) or extended multiplicative signal
519 correction (EMSC; bottom)

520

521 Figure 3. Prediction performance of NIT chemometric models for moisture, protein, fat, ashes
522 and carbohydrates contents in quinoa grains, as evaluated by the uncertainty about the
523 correlation coefficient of prediction (R_p) built by bootstrapping (left), and the scatter plots
524 between chemical reference values and those fitted to the calibration data set (middle) and
525 predicted using the validation data set (right)

526