# Determination of the Geographical Origin of Hazelnuts (*Corylus avellana* L.) by Near-Infrared Spectroscopy (NIR) and a Low-Level Fusion with Nuclear Magnetic Resonance (NMR)

Navid Shakiba[a,b], Annika Gerdes[a,b], Nathalie Holz[a], Sören Wenck[b], René Bachmann[c], Tobias Schneider[a], Stephan Seifert[b], Markus Fischer[b], Thomas Hackl[a,b,*]

[a]Institute of Organic Chemistry, University of Hamburg, Martin-Luther-King-Platz 6, 20146 Hamburg, Germany,

[b]HAMBURG SCHOOL OF FOOD SCIENCE - Institute of Food Chemistry, University of Hamburg, Grindelallee 117, 20146 Hamburg, Germany

[c]Landeslabor Schleswig-Holstein, Max-Eyth-Straße 5, 24537 Neumünster, Germany

*Corresponding author: Tel.: +49-40 42838-2804; E-Mail: Thomas.Hackl@chemie.uni-hamburg.de

**ABSTRACT**

Fourier-transform near-infrared (FT-NIR) spectroscopy was used to determine the geographical origin of 233 hazelnut samples of various varieties from five different countries (Germany, France, Georgia, Italy, Turkey). The experimental determination of the geographical origin of hazelnuts is important, because there are usually large price differences between the producer countries and thus a risk of food fraud that should not be underestimated. The present work is a feasibility study using a low-cost method, as high-field NMR and UPLC-QTOF-MS have already been used for this question. Sample sets were split with repeated nested cross validation and an ensemble of discriminant classifiers with random subspaces was used to build the classification models. By using a preprocessing strategy consisting of multiplicative scatter correction, bucketing and the mean averaging of five measured spectra per sample, a test accuracy of $90.6 \pm 3.9\%$ was achieved, which rivals results obtained with much more expensive infrastructure. The application of the feature selection approach surrogate minimal depth showed that the successful classification is mainly caused by protein signals. In addition, a low-level data fusion of the NIR and NMR data was performed to assess how well the two methods complement each other. The data fusion was compared to a complementary approach, where the classification results based on the individual NIR and NMR models were jointly examined. The data fusion performed better than the individual methods with a test accuracy of $96.6 \pm 2.8\%$. A comparison of the outliers in all classification models shows conspicuities in always the same samples, indicating that robust classification models are obtained.

**KEYWORDS**

Geographical Origin, NIR, NMR, Data Fusion, Hazelnut, Feature selection

2

## 1. INTRODUCTION

Hazelnuts (*Corylus avellana* L.) are a globally traded food with a production volume of approximately 1,125,000 t in 2019.[1] Turkey is the main producing country with a volume of 776,000 t in 2019, representing 69% of the world production. Other major producing countries are Italy, Azerbaijan, the USA, Chile, China and Georgia, although producer prices vary widely in some cases. For example, the price for a ton of hazelnuts from Georgia in 2019 was only 1550 USD/t, but in Italy it was 3600 USD/t, as these hazelnuts are considered to be of particularly high quality. Such a wide price range is bound to provide a financial incentive for food fraud, where hazelnuts from a low-price producing country are falsely declared with a different origin to increase profits.

Bachmann *et al.* (2018) and Klockmann *et al.* (2016) explored the issue of determining the geographical origin of hazelnuts using high-resolution instrumentation, [1]H NMR spectroscopy and ultraperformance liquid chromatography quadrupole time-of-flight mass spectrometry (UPLC-QTOF-MS) in combination with chemometric evaluation strategies.[2,3] These studies showed that it is possible to distinguish the origin of hazelnuts using metabolomics approaches. However, these tools are quite expensive and require a high level of scientific expertise, which is limiting, especially for smaller laboratories and small and medium-sized food companies.

Fourier-transform near-infrared spectroscopy offers a cost-effective way to determine the geographical origin of food, as has already been shown with various foods such as pistachio, wheat, almonds and walnuts.[4–7] In addition, NIR can be used for a wide range of food-related issues in the food sector, e.g. quality control of olive oil, determination of storage time of pork and identifying oxidation of vegetable oils.[8–10] Other advantages of NIR spectroscopy are the absence of hazardous chemicals, the non-destructive nature, a fast measurement time and the fact that no extraction is required. Two Italian research groups have already used NIR spectroscopy to distinguish 'Nocciola Romana', which carries a Protected Designation of Origin (PDO), from other hazelnuts; however, a holistic comparison of several countries of origin has not yet taken place.[11,12] Both studies examined whole, shelled hazelnuts in order to develop a non-destructive rapid method. Based on a comparison of different preparation techniques for NIR measurement to determine the geographical origin of almonds, we decided to analyze the samples after homogenization and freeze-drying, as we expected this approach to provide a higher information content and a better representation of the sample populations.[13] One aim of this study is therefore to investigate the ability of NIR spectroscopy to determine the geographical origin of hazelnut samples, as there is a need for such a low-cost analytical

3

71   method. This could be used in industry for incoming goods inspection. To establish such a

72   method, we compared various preprocessing strategies and classification approaches. In

73   addition, surrogate minimal depth (SMD) was applied, a random forest based approach for

74   feature selection and relation analysis that has already been used to study other vibrational

75   spectroscopic data.[14–16]

76   The newly acquired NIR and the already existing NMR data were selected for low-level fusion

77   due to their one-dimensionality and potentially complementary nature. Low-level fusion

78   involves concatenation of the datasets with or without prior preprocessing methods.[17] In the

79   case of hazelnuts as a matrix, the NIR captures mainly non-specific information on groups of

80   substances with high concentration, e. g. lipids, carbohydrates and proteins, while the $^1$H NMR

81   measurement of the polar extract provides more specific information on substances such as

82   organic acids, amino acids and specific carbohydrates. To the best of our knowledge, this is the

83   first publication on the experimental determination of the geographical origin of food by

84   combining NIR and high-field NMR data in a multiclass model. The aim of the low-level data

85   fusion approach of the NIR and NMR data is to obtain a statistical model that is better than the

86   individual methods.

87

88

4

## 2. Materials and Methods

*2.1. Hazelnut Samples*

In a previous study by our group 262 raw hazelnut samples were used for the determination of the geographical origin by means of $^1$H NMR.[2] Authentic reference material was provided by partners, distributors and suppliers. Of these, only 233 samples could be used for NIR analysis, as the sample material for some samples has already been used up. Samples from a total of five countries were analyzed, with several samples coming from economically important growing regions. In this study, 27 German samples, 116 French samples, 15 Georgian samples, 37 Italian samples and 38 Turkish samples were used. The same samples were also taken for the $^1$H NMR analysis and the low-level data fusion. More detailed information on origin and variety are given in the Supporting Information (Table S1).

*2.2. Sample Treatment*

All hazelnut samples were treated according to Bachmann et al.[2] Hazelnut samples were frozen in liquid nitrogen before they were homogenized with a Grindomix GM 300 knife mill and dry ice was added. After evaporation of the dry ice, the samples could be directly used for NMR analysis. To prepare the samples for NIR measurement, the homogenized samples were then freeze-dried for 48 hours.

*2.3. NIR spectroscopy*

1.250 g ($\pm$ 0.005 g) of the ground and freeze-dried hazelnut samples were thawed at 22 °C ($\pm$ 2 °C) in closed glass vials (52.0 mm x 22 mm x 1.2 mm, Nipro Diagnostics Germany GmbH, Ratingen, Germany) preceding NIR measurement.

The NIR measurements were performed on a TANGO FT-NIR spectrometer (Bruker Optics, Bremen, Germany) equipped with an integrating sphere. Spectra were recorded in reflectance mode at room temperature (22 $\pm$ 2 °C), with the wavenumber range set to 11546-3949 cm$^{-1}$ collecting 50 scans at a resolution of 2 cm$^{-1}$. Each sample was analyzed five times by shaking the lyophilisate in the glass vial between measurements.

*2.4. NMR spectroscopy*

The NMR spectra used for the data fusion were acquired by Bachmann *et al.* in an earlier study on a Bruker Avance III 400 MHz spectrometer (Bruker Biospin, Rheinstetten, Germany) operating at 400.13 MHz with the noesygppr1d pulse program.[2]

5

119 *2.5. NIR spectra preprocessing*

120 All preprocessing techniques were performed using MATLAB R2020b (The MathWorks Inc.,
121 Natick, MA, USA). Multiplicative scatter correction (MSC) was applied to ensure a good
122 comparability between samples. MSC is a commonly used preprocessing step to normalize the
123 data and remove artifacts from the samples by using the mean spectrum of the available
124 data.[18] These artifacts are mostly due to differences in particle size of the powdered sample,
125 which leads to non-uniform scattering effects.[19] After MSC, either no derivative, the first
126 derivative or the second derivative was applied to the spectra. The approaches that used a
127 derivative also utilized a Savitzky-Golay smoothing filter with a window size of 11 and a
128 polynomial order of 2 to minimize the negative effects of a derivative on the signal-to-noise
129 ratio.[18] Next, variable reduction was achieved by calculating the mean of five adjacent
130 features into one bucket, leading to a reduction from 3720 variables (spectral range: 11538-
131 3949 cm$^{-1}$) to 744 NIR-buckets. Finally, the mean or median of the five measured spectra per
132 sample was determined. The classification results of the different preprocessing strategies were
133 then compared.

134 *2.6. NMR spectra preprocessing*

135 The NMR data were processed with Topspin 3.2 (Bruker Biospins, Rheinstetten, Germany). A
136 Fourier transformation with a line broadening factor of 0.3 was applied on the FIDs, then
137 baseline corrected and phased. Integrals of signals and regions from the NMR spectra were
138 determined manually in AMIX 3.9.14 (Bruker Biospins, Rheinstetten, Germany) as variable
139 sized buckets and normalized to total intensity by scaling. A total of 222 NMR-buckets were
140 defined for each sample. The mean and median of the triplicate measurement was
141 determined.[2]

142 *2.7. NIR-NMR low-level fusion*

143 For the low-level fusion, the 744 NIR buckets of the best performing model were concatenated
144 with the 222 NMR buckets, once mean and once median averaged.[17] Autoscaling was used
145 as a scaling method for the fusion data.[20,21]

146 *2.8. Multivariate data analysis*

147 Multivariate data analysis was performed with MATLAB R2020b including the Classification
148 Learner app (The MathWorks Inc., Natick, MA, USA). Samples were split into five equal parts
149 with a stratified nested cross-validation stratified by geographic origin.[22] The internal
150 validation was also split fivefold to avoid overfitting during model training. Repeated nested

6

151  cross-validation (RNCV) was iterated five times to obtain an average result, resulting in 25

152  different corresponding training and test sets, with each sample being part of 20 training and 5

153  test sets, as different sample splits can lead to large differences in model accuracy. The

154  Classification Learner app was used to determine which classifiers would be suitable for the

155  NIR, NMR and low-level fusion data, and then the chosen classifier was used for automatic

156  model training and subsequent validation. The ensemble of discriminant classifiers using

157  random subspaces was the best performing method and was trained with 372 subspace

158  dimensions and 30 learning cycles.[23] The test accuracies given are the mean of the test

159  accuracies of all sample splits from the RNCV. The macro-$F_1$ score is calculated as the

160  arithmetic mean of $F_1$ score of the five classes, formed from the harmonic mean of the class-

161  wise precision and sensitivity. In addition, Fleiss' kappa was calculated to determine the degree

162  of agreement of the classifiers in each model.[24]

163  *2.9. Feature selection and relation analysis with surrogate minimal depth (SMD)*

164  The software R in version 3.6.3 and the R package SurrogateMinimalDepth in version 0.2.0

165  (https://github.com/StephanSeifert/SurrogateMinimalDepth) were utilized for feature selection

166  with the parameters ntree = 10000, mtry = 143, min.node.size = 1 and s = 149. In order to

167  compensate for the class imbalance, case.weights were chosen accordingly meaning that

168  samples from rare classes were sampled more frequently for training. Subsequently, the relation

169  parameter mean adjusted agreement of the selected features was determined and depicted in a

170  heatmap generated by the R package pheatmap in version 1.0.12. For the random forest

171  classification, the R package ranger was applied with the above described parameters.[25] Since

172  random forests provides an internal validation, no cross validation scheme had to be applied

173  and all of the samples were utilized simultaneously.

174

## 3. RESULTS AND DISCUSSION

*3.1. NIR-Spectroscopy*

Hazelnuts are rich in fat (~61%), carbohydrates (~17%), protein (~15%) and have a water content of ~5%.[26] The fatty acid profile is dominated by oleic (72.8-83.5%), linoleic (7.6-16.6%) and palmitic (4.1-6.8%) acid and is similar to that of olive oil.[27] A model NIR spectrum of a hazelnut sample is shown in the Supporting Information (S2). The NIR spectrum shows strong similarities to those of other species of nuts, because of akin nutrient composition.[6,7] Due to the broad absorption and the overlapping of the signals of the individual substances of the complex matrix, no peaks in the spectra can be clearly assigned to specific metabolites. Instead, signals and regions in the spectra can usually be assigned to different molecular vibrations caused mainly by compound classes of macronutrients. The peak at 8550 cm$^{-1}$ can be assigned to HC=CH (C-H second overtone) caused by unsaturated fatty acids. Other signals that can be associated with lipids are the second overtone of C-H (C-H, C-H$_2$, C-H$_3$) stretching vibrations between 8500-8000 cm$^{-1}$, the first overtone of C-H between 5900-5600 cm$^{-1}$ and the combination bands of the methylenic CH$_2$ between 4500-4000 cm$^{-1}$.[11,12] The first overtone of N-H and O-H of proteins can be observed in the region between 7100-6100 cm$^{-1}$.[28] Another region related to proteins is between 4900-4600 cm$^{-1}$, caused by the combination band of peptide bonds.[28]

Principal component analysis (PCA) is arguably the most widely used unsupervised method for reducing the complexity of metabolomics data while preserving variance as much as possible and revealing underlying class information.[29] The limitations of PCA as an exploratory method are that underlying patterns cannot be uncovered if the intragroup variance of the sample groups is greater than the intergroup variance.[30] The advantages of PCA include an initial unbiased look at the data to examine the extent to which samples are similar within and outside their groups and to identify potential outliers.[31] Figure 1A shows the PCA scores plot of the unprocessed samples, where the first principal component (PC) contains 85.0% of the variance and 8.3% the PC 2. The plot shows a cluster of all samples with no outliers or clear group separation. Nevertheless, the different groups show similarities. The French samples are in the center-left of the plot, while the German samples are below and the Georgian samples are above. The Italian and Turkish samples mostly scatter from the center to the right side of the plot along the first principal component (PC1). The PCA scores plot of the preprocessed data is shown in Figure 1B with PC 1 accounting for 60.2% of the variance and PC 2 representing 23.7%. The plot shows a coherent cluster with no outliers, but with less clear

8

spatial allocations of the different groups. As expected, PCA cannot identify separate groups with respect to the origin of the samples. Hence, supervised multivariate analysis was performed to determine the geographical origin of the hazelnut samples.
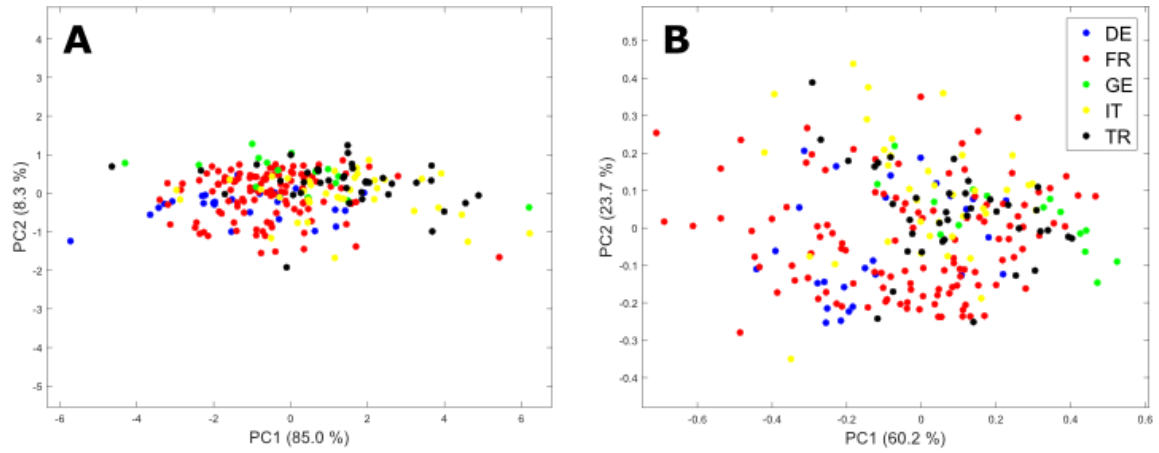


**Figure 1. (A)** PCA scores plot of unprocessed, mean averaged NIR spectra. **(B)** PCA scores plot of NIR spectra after MSC, bucketing and mean averaging.

**Table 1.** Test accuracy, macro-F1 score and Fleiss' kappa coefficient of the different preprocessing strategies of the NIR analysis attained by using an ensemble of discriminant classifiers using random subspaces.

| Strategy | Preprocessing | Test Accuracy | Macro-$F_1$ | Fleiss' Kappa |
|---|---|---|---|---|
| NIR-I | MSC – Mean | 89.5 ± 4.3% | 85.9% | 87.4% |
| NIR-II | MSC – Median | 84.2 ± 5.3% | 80.1% | 78.5% |
| NIR-III | MSC – Bucketing – Mean | 90.6 ± 3.9% | 88.1% | 88.7% |
| NIR-IV | MSC – Bucketing – Median | 81.3 ± 5.2% | 76.7% | 77.8% |
| NIR-V | MSC – 1. Derivative – Smoothing – Bucketing – Mean | 76.6 ± 5.3% | 70.5% | 73.7% |
| NIR-VI | MSC – 1. Derivative – Smoothing – Bucketing – Median | 75.1 ± 5.3% | 68.5% | 75.8% |
| NIR-VII | MSC – 2. Derivative – Smoothing – Bucketing – Mean | 67.6 ± 4.8% | 57.9% | 67.7% |
| NIR-VIII | MSC – 2. Derivative – Smoothing – Bucketing – Median | 68.6 ± 4.8% | 56.9% | 70.9% |
| NIR-IX | Cut – MSC – Bucketing – Mean | 83.6 ± 3.8% | 79.6% | 77.7% |
| NIR-X | Cut – MSC – Bucketing – Median | 75.2 ± 6.0% | 70.6% | 66.4% |
| NIR-SMD | MSC – Bucketing – Mean – SMD feature selection | 86.4 ± 4.7% | 82.1% | 78.1% |

For the application of supervised approaches different preprocessing strategies and different classifiers were applied and different parameters to assess their performance were utilized. The test accuracy is the most common measure for machine learning models. The macro-$F_1$ score

219    provides information about the mean class-wise precision and sensitivity of a sample split.[32]

220    Fleiss' kappa is a measure of inter-rater reliability for determining the homogeneity of the

221    RNCV's ratings of the samples, regardless of whether they were correctly allocated or not.[33]

222    The algorithm used for multivariate analysis was an ensemble of discriminant classifiers using

223    random subspaces.[23] This algorithm showed the best results for strategies I-IV. However,

224    strategies IV-VIII showed slightly better results with other classifiers (data not shown), but for

225    clarity and comparability the same classifier was used for each strategy and the results are

226    shown in Table 1. The mean spectra of preprocessing strategies I, V and VII are depicted in the

227    Supporting Information (S3).

228    The adverse effects on the signal-to-noise ratio due to smoothing and the use of the derivative

229    of the spectra is reflected in the relatively poor results of preprocessing strategies V-VIII.

230    Strategies VII and VIII, which use the second derivative, have a test accuracy of $67.6 \pm 4.8\%$

231    and $68.6 \pm 4.8\%$, respectively. Including the information from the confusion matrices shows an

232    even worse picture. Macro-$F_1$ scores of 57.9% and 56.9% for strategies VII and VIII,

233    respectively, show a more dramatic decline to their test accuracies compared to other

234    preprocessing strategies. This is probably due to the fact that the French sample group has the

235    highest number of samples and a large proportion of the samples from other countries are

236    misclassified as French. The strategies with the first derivative already show significantly better

237    results with a test accuracy of $76.6 \pm 5.3\%$ for strategy V and $75.1 \pm 5.3\%$ for strategy VI. The

238    macro-$F_1$ scores are also much closer to test accuracies. Since additive effects are observed in

239    the spectrum, the use of the first derivative is reasonable in theory. In the practice of this study,

240    however, the negative effects of the derivative on the signal-to-noise ratio may have led to

241    poorer predictive performance of the models. The NIR preprocessing strategies IX and X cut

242    off the spectrum above the wavenumber of 9000 $cm^{-1}$. This is a common preprocessing step as

243    this region is usually not very information-rich and contains mainly bands from the third and

244    the fourth overtone vibrations.[28,34] Both strategies lead to a decrease in test accuracy of 7.0%

245    and 6.1% compared to strategies III and IV, which contain features from the whole spectrum.

246    This suggests that the bands in this region are relevant for the research question. Although all

247    strategies forgoing the derivative show good model performance, the results of strategies I and

248    III show the highest test accuracy of $89.5 \pm 4.3\%$ and $90.6 \pm 3.9\%$, a macro-$F_1$ score of 85.9%

249    and 88.1% and a Fleiss' kappa coefficient of 87.4% and 88.7%. NIR strategy I only used MSC

250    and mean averaging as preprocessing steps, while NIR-III also used bucketing of the variables.

251    Although strategy I shows similar results to those of strategy III, this preprocessing approach

252    is not pursued further because bucketing ensures a more robust model, reduces the risk of

10

overfitting and requires less computing time. In summary, the best preprocessing method is one of the simplest. Forgoing a part of the spectra and using derivatives resulted in a loss of information and thus lower classification accuracies. Consistent with all preprocessing strategies except those using the second derivative is that mean averaging lead to a higher test accuracy. The advantage of the median is its robustness and protection against outliers but using the mean average can improve the spectral resolution, leading to better classification results.

**Predicted Class**

|  | DE | FR | GE | IT | TR | Sensi-tivity |
|---|---|---|---|---|---|---|
| DE | 20.8 | 5 |  | 0.4 | 0.8 | 77.0% |
| FR | 1.6 | 111.8 | 1 | 1.6 |  | 96.4% |
| GE |  | 1.2 | 11.8 | 0.4 | 1.6 | 78.7% |
| IT |  | 5.8 | 0.8 | 30.4 |  | 82.2% |
| TR |  | 1.6 |  | 0.2 | 36.2 | 95.3% |
| Precision | 92.9% | 89.2% | 86.8% | 92.1% | 93.8% | 90.6% |

*(Row labels DE, FR, GE, IT, TR under "True Class")*

**Figure 2.** Confusion matrix of NIR preprocessing strategy III. The values given correspond to the mean of the five runs of the RNCV. Mean classification accuracy (90.6%), precision and sensitivity scores of the classes are also given. Confusion matrices of the other NIR preprocessing strategies are in the Supporting Information (S5). [Color]

The good classification results show the impact of geographical influences on the macronutrient profile of hazelnut, despite variable factors such as post-harvest processing and different harvest years, making NIR spectroscopy well suited for determining the geographical origin. An external classification accuracy of $90.6 \pm 3.9\%$ for a multiclass model with five classes is impressive in the field of geographical origin determination of food using NIR (figure 2). To put this result in relation to other publications in this field: A classification model determining the geographical origin of walnuts from seven countries achieved a classification accuracy of $77.0 \pm 1.6\%$ using a linear discriminant analysis.[7] Another study that investigated the geographical origin of almonds obtained a classification accuracy of $80.3 \pm 1.5\%$ when

11

272 comparing six countries of origin with a support vector machine model.[6] Less complex
273 models than that one, comparing only two sample groups of hazelnuts, have been developed by
274 Moscetti et al. (2014) and Biancolillo et al. (2018).[11,12] Moscetti et al. (2014) compared
275 'Nocciola Romana' hazelnuts, which have a Protected Designation of Origin (PDO) indication,
276 with hazelnuts of the 'Tonda di Giffoni' and 'Barrettona' cultivars and reported a classification
277 accuracy of 95.5% using a support vector machine algorithm.[35] Biancolillo et al. (2018)
278 investigated a similar question by comparing the 'Nocciola Romana' PDO with 'other'
279 hazelnuts originating from Italy or the USA, resulting in a correct classification rate of 93.9%
280 for 'Nocciola Romana' and 95.1% for 'others' by partial least square discriminant analysis.

281 In order to identify features that are responsible for this successful classification and to analyze
282 their relationships, the feature selection approach surrogate minimal depth (SMD) was applied
283 to the NIR data of preprocessing strategy III. Unlike other feature selection techniques, SMD
284 does not evaluate the importance of the features individually, but by including their relations
285 with each other.[14] SMD selected 245 of 744 buckets, and the high number of selected features
286 can be explained by the fact that the bands in the NIR spectrum are quite broad and many
287 features belong to the same signal. To obtain a more comprehensive interpretation of the
288 important features, the mean adjusted agreement, a relation parameter that takes into account
289 the mutual association to the result, is also obtained by SMD. The results of this relation analysis
290 are presented in a heat map (Figure 3A) and in a spectrum colored according to the respective
291 clusters of the relation analysis (Figure 3B). The heat map shows six distinct clusters that mainly
292 contain neighboring buckets confirming the conclusion previously drawn from the high number
293 of selected features. Somewhat surprisingly, four of the clusters are located in the wavenumber
294 range between 11300-8700 cm$^{-1}$ and contain only low intensity signals. Cluster 1 (red), cluster
295 2 (blue) and cluster 3 (purple), which show moderate to strong relations to each other, are in
296 the region between 11500-10200 cm$^{-1}$, which can be assigned to the third overtone of C-H from
297 methyl and methylene. However, cluster 3 (purple) contains features in the range of 10600-
298 10200 cm$^{-1}$, where also bands of the N-H stretch-second overtone are found. The importance
299 of the spectral regions represented by the clusters 1-4 is also confirmed by the results of the
300 NIR preprocessing strategy IX, which did not include the buckets over 9000 cm$^{-1}$ and gave a
301 lower classification accuracy of $83.6 \pm 3.8\%$ (Table 1). However, the fact that this accuracy is
302 still quite high shows that the features in clusters 5 (green) and 6 (gray) are even more important
303 for classification. Moreover, these clusters show very interesting relations between the distinct
304 regions 7000-6000 cm$^{-1}$, 4800-4700 cm$^{-1}$ and 4400-4300 cm$^{-1}$. Signals in the region between
305 7000-6000 cm$^{-1}$ can be assigned to the first overtone of N-H from the peptide bond and side

12

chains of amino acids as well as the first overtone of O-H.[28,36] N-H combination from proteins, C-H/C=O lipid associated and O-H combination bands are in the region between 4800-4700 cm$^{-1}$.[28] Signals in the region between 4410-4390 cm$^{-1}$ can be ascribed to C=O and N-H in α-helix and β-sheet structures in peptides.[28] Since all of these related bands can be assigned to functional groups of proteins, we can conclude that the successful classification is caused by different protein compositions of the hazelnut samples.



**Figure 3.** Results of the SMD feature selection and relation analysis: A Heatmap of the relation parameter mean adjusted agreement of the 245 selected features (A), as well as an example NIR hazelnut spectrum with buckets colored according to the respective associated clusters (B) are shown.

13

316  It is also worth noting that the random forest analysis on which SMD is based on has a much
317  lower classification accuracy of 60.1% compared to the previous analysis (Supporting
318  Information, Figure S5.11). The main difficulty of the random forest model was to distinguish
319  between similar groups of Germany and France, and Turkey and Italy. To assess how much of
320  the relevant information is contained in the significant buckets, we repeated the RNCV and
321  classification using only these features resulting in a classification accuracy of 86.4 ± 4.7%. As
322  this is only 4.2% less than the model with the best performance using all features, this shows
323  that the 245 selected buckets carry the most important information for classification and that
324  the SMD performs well for feature selection, even though the random forest analysis shows
325  comparatively poor results for classification. The averaged NIR spectra of samples from the
326  five countries of origin (Supporting Information, Figure S4) were overlaid to check for
327  differences in the spectra. In agreement with the feature selection results, the spectral regions
328  11500-10000 $cm^{-1}$, 7000-6000 $cm^{-1}$ and 5100-4500 $cm^{-1}$ show the largest differences. The
329  region between 4400-4000 $cm^{-1}$ also shows differences, but the signals in this region were
330  susceptible to matrix effects, as evident by the inconsistencies of this region across the five
331  measurements of a single sample.

332  *3.2. NIR-NMR-Data Fusion*

333  In the original publication, which used [1]H NMR spectroscopy to determine the geographical
334  origin of hazelnuts based on the polar metabolome, 262 hazelnut samples were divided into a
335  training set containing two-thirds of the samples (172) and a test set with one-third of the
336  samples (90).[2] Subsequently, several classification algorithms were trained with the training
337  set and the test set was used to externally evaluate model performance. As with the results of
338  the NIR spectroscopy, an ensemble of discriminant classifiers with a random subspace
339  algorithm showed the best results. This is quite interesting because despite different observed
340  features, both datasets appear to have similar underlying structures in the processed data. The
341  model achieved a cross validation accuracy of 91% for the training set and an accuracy of 96%
342  for the test set. Due to the use of a smaller number of samples and a stratified repeated nested
343  cross validation to capture the variance of different sample splits, the NMR classification results
344  were recalculated using the same external splits into training and test sets. The NMR data were
345  Fourier transformed, baseline corrected, phased and 222 regions were defined as variable sized
346  buckets, which were normalized to total intensity by scaling.[2] The NMR measurements were
347  performed in triplicate, so the median and mean were compared for averaging. The
348  classification results of the NMR and Fusion preprocessing strategies are shown in Table 2.

| Strategy | Preprocessing | Test Accuracy | Macro-$F_1$ | Fleiss' Kappa |
|---|---|---|---|---|
| NMR-I | NMR – Mean | 94.3 ± 3.2% | 93.6% | 93.2% |
| NMR-II | NMR – Median | 95.1 ± 3.0% | 94.4% | 94.2% |
| Fusion-I | NIR-III + NMR-Mean Fusion | 96.1 ± 2.7% | 95.3% | 95.1% |
| Fusion-II | NIR-III + NMR-Median Fusion | 96.1 ± 3.2% | 95.4% | 94.0% |
| Fusion-III | NIR-III + NMR-Mean Fusion – Autoscaled | 96.6 ± 2.8% | 96.0% | 96.4% |
| Fusion-IV | NIR-III + NMR-Median Fusion – Autoscaled | 96.0 ± 2.7% | 95.4% | 95.8% |

The results of the mean and median were quite similar. Median averaging of the NMR data yielded a test accuracy of 95.1 ± 3.0%, a macro-$F_1$ score of 94.4% and a Fleiss' Kappa coefficient of 94.2%. NMR-I, using the mean of the NMR data, achieved a 0.8% lower classification rate and similarly low values for the macro-$F_1$ score and the Fleiss' Kappa coefficient. Due to the similar results of the two NMR spectroscopy strategies, both datasets were tested in a low-level data fusion with the NIR spectroscopy strategy III data. In this case, all buckets used for NIR, and NMR analysis were combined in a matrix resulting in 966 features. The dataset was then scaled using autoscaling, resulting in a standard deviation of one for each feature.[37] Fusion-I combined the data of NIR-III and NMR-I, while Fusion-II used the median averaged NMR buckets. Both fusions yielded similar results, with a test accuracy of 96.1% and only small differences in the other measures. Test accuracy of Fusion-I is 1% higher compared to NMR-II, which only uses the median averaged NMR buckets, indicating a better classification performance of the model. Fusion-III and Fusion-IV subjected the datasets from Fusion-I and -II to autoscaling. Fusion-IV used the fused dataset of NIR strategy III and the median averaged NMR buckets leading to a test accuracy of 96.0 ± 2.7% thus showing almost the same results as Fusion-I and -II. Fusion-III instead used mean averaged NMR buckets in the fusion and yielded a test accuracy of 96.6 ± 2.8%, a macro-F1 score of 96.7% and a Fleiss' kappa coefficient of 96.4%, thus showing the best results for each measure out of all examined models. Compared to the NMR-II approach, test accuracy increased by 1.5%. Such an increase is quite large considering that the statistical measures of all models are higher than 90%. A further comparison of NMR-II (Supporting Information, Figure S6.2) and Fusion-III (Figure 4) shows improvements in the classification of the German, French and Italian samples, while the accuracy for Georgian and Turkish samples remains the same.

15

**Predicted Class**

| | DE | FR | GE | IT | TR | Sensi-tivity |
|---|---|---|---|---|---|---|
| **DE** | 23.8 | 2.8 | | 0.4 | | 88.1% |
| **FR** | 0.6 | 114.4 | 0.4 | 0.4 | 0.2 | 98.6% |
| **GE** | | | 15 | | | 100% |
| **IT** | | 0.8 | 0.8 | 35.4 | | 95.7% |
| **TR** | | 1 | 0.2 | 0.2 | 36.6 | 96.3% |
| **Precision** | 97.5% | 96.1% | 91.5% | 97.3% | 99.5% | 96.6% |

*(Row labels DE, FR, GE, IT, TR are under the "True Class" axis label.)*

**Figure 4.** Confusion matrix of Fusion-III. The values correspond to the mean of the five runs of the RNCV and the mean classification accuracy (96.6%), precision and sensitivity scores of the classes are also shown. Confusion matrices of the NMR strategies and the other data fusions are shown in the Supporting Information (S6).

The individual allocations of the classification model of Fusion-III (Supporting Information, Table S7) were examined to obtain information about problematic samples. In total, 13 samples were misclassified at least once, and four samples were misclassified in each split of RNCV. These were two German samples, one from Bavaria and one from Rhineland-Palatinate, one French sample and one Turkish sample. These samples were also misclassified at least once in the individual models of NIR-III and NMR-II, but only the Turkish sample was falsely classified in all RNCV splits. One of the misclassified German samples was from Rhineland-Palatinate, which was the only sample from this federal state. It was classified three times as a French sample and twice as an Italian sample. The other German sample is of mixed variety from the municipality of Aiglsbach in Bavaria, which has always been misclassified as French. As there are five samples from Aiglsbach, the reason for these misclassifications is probably not the lack of enough samples for training but the individual composition of this sample. For a similar reason, a Turkish mixed variety sample from the Düzce region was always misclassified as French even though there are 13 samples with the same characteristics that were not misclassified once. The models did not show clear results regarding the French sample, which was misclassified in all five splits. It was classified twice as Georgian and Italian sample and once as Turkish sample. Again, these misclassifications are probably due to the individual

16

composition of this sample, as there were seven other samples of the Pauetet variety from the Midi-Pyrénées region, which were always correctly classified.

For comparison, a complementary approach, i. e. cross-checking the NIR and NMR classification results with respect to misclassified samples, was performed using the best performing models based on NIR-III and NMR-II (Supporting Information, Table S7). The results of the two models show a large overlap in the classification of the samples. 192 of the 233 samples analyzed were correctly classified by both models. A total of 41 samples were misclassified at least once in the RNCV by the NIR approach and 21 samples by the NMR approach. Again, the models overlap, as 11 samples were misclassified at least once by both models. This demonstrates that both classification models perform well and misclassify similar samples. However, it is also shown that the models also provide complementary information. In total, only two samples were incorrectly classified in all five splits, all as French, of the RNCV and in both models indicating the conservative nature of this approach. One of these samples is a Turkish sample of a mixed variety from the province of Düzce and the other is of the 'Tonda di Giffoni' variety from the Campania region in Italy, which has a Protected Geographical Indication (PGI). The Turkish sample is the same one that Fusion-III misclassified. Of all 233 samples, eight are of the 'Tonda di Giffoni' variety from France and three from Italy. This suggests that the metabolome of the misclassified Italian 'Tondi di Giffoni' sample may be more influenced by the cultivar than by environmental factors, so more samples of this cultivar from Italy are needed to adequately train the models. Of the 11 samples that were misclassified at least once in the complementary approach, eight were also misclassified at least once in Fusion-III, indicating the similarity of the results. Another factor to consider at this point is the fact that these samples were obtained from suppliers and in principle there is the possibility of a mix-up, even if it is very unlikely.

The question remains whether a data fusion should be used when determining the geographical origin of hazelnuts. If both methods have already been used, data fusion may even give better results than the methods on their own. In this case, the fusion and subsequent autoscaling of the NIR dataset, which used MSC, bucketing and mean, with the NMR dataset, which used the mean of the buckets, gave the best individual model for the hazelnut geographical origin question. However, the complementary approach of analyzing samples sequentially using NIR and NMR spectroscopy proves to be a more conservative and reliable method. This method would also be suitable for transfer to industry, where NIR analysis is used as a level 1 analysis and conspicuous samples are subsequently analyzed in an analytical laboratory using NMR as a level 2 analysis.

17

**Declaration of Competing Interest**

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

18

462 **Appendix: Supporting Information**

463 Supporting Information of this article can be found online at

## 4. References

[1] Food and Agriculture Organization of the United Nations, Production of Crops, 2021. http://www.fao.org/faostat/en/#data (accessed 25 February 2021).

[2] R. Bachmann, S. Klockmann, J. Haerdter, M. Fischer, T. Hackl, 1H NMR Spectroscopy for Determination of the Geographical Origin of Hazelnuts, J. Agric. Food Chem. 66 (2018) 11873–11879. https://doi.org/10.1021/acs.jafc.8b03724.

[3] S. Klockmann, E. Reiner, R. Bachmann, T. Hackl, M. Fischer, Food Fingerprinting: Metabolomic Approaches for Geographical Origin Discrimination of Hazelnuts (Corylus avellana) by UPLC-QTOF-MS, J. Agric. Food Chem. 64 (2016) 9253–9262. https://doi.org/10.1021/acs.jafc.6b04433.

[4] R. Vitale, M. Bevilacqua, R. Bucci, A.D. Magrì, A.L. Magrì, F. Marini, A rapid and non-invasive method for authenticating the origin of pistachio samples by NIR spectroscopy and chemometrics, Chemometrics and Intelligent Laboratory Systems 121 (2013) 90–99. https://doi.org/10.1016/j.chemolab.2012.11.019.

[5] H. Zhao, B. Guo, Y. Wei, B. Zhang, Near infrared reflectance spectroscopy for determination of the geographical origin of wheat, Food Chem. 138 (2013) 1902–1907. https://doi.org/10.1016/j.foodchem.2012.11.037.

[6] M. Arndt, M. Rurik, A. Drees, C. Ahlers, S. Feldmann, O. Kohlbacher, M. Fischer, Food authentication: Determination of the geographical origin of almonds (Prunus dulcis Mill.) via near-infrared spectroscopy, Microchemical Journal 160 (2021) 105702. https://doi.org/10.1016/j.microc.2020.105702.

[7] M. Arndt, A. Drees, C. Ahlers, M. Fischer, Determination of the Geographical Origin of Walnuts (Juglans regia L.) Using Near-Infrared Spectroscopy and Chemometrics, Foods 9 (2020). https://doi.org/10.3390/foods9121860.

[8] R.J. Mailer, Rapid evaluation of olive oil quality by NIR reflectance spectroscopy, Journal of the American Oil Chemists' Society 81 (2004) 823–827. https://doi.org/10.1007/s11746-004-0986-4.

[9] Q. Chen, J. Cai, X. Wan, J. Zhao, Application of linear/non-linear classification algorithms in discrimination of pork storage time using Fourier transform near infrared (FT-NIR) spectroscopy, LWT - Food Science and Technology 44 (2011) 2053–2058. https://doi.org/10.1016/j.lwt.2011.05.015.

[10] G. Yildiz, R.L. Wehling, S.L. Cuppett, Method for Determining Oxidation of Vegetable Oils by Near-Infrared Spectroscopy, Journal of the American Oil Chemists' Society 78 (2001) 495–502. https://doi.org/10.1007/s11746-001-0292-1.

20

498 [11] R. Moscetti, E. Radicetti, D. Monarca, M. Cecchini, R. Massantini, Near infrared
499 spectroscopy is suitable for the classification of hazelnuts according to Protected
500 Designation of Origin, J. Sci. Food Agric. 95 (2015) 2619–2625.
501 https://doi.org/10.1002/jsfa.6992.

502 [12] A. Biancolillo, S. de Luca, S. Bassi, L. Roudier, R. Bucci, A.D. Magrì, F. Marini,
503 Authentication of an Italian PDO hazelnut ("Nocciola Romana") by NIR spectroscopy,
504 Environ. Sci. Pollut. Res. Int. 25 (2018) 28780–28786. https://doi.org/10.1007/s11356-
505 018-1755-2.

506 [13] M. Arndt, M. Rurik, A. Drees, K. Bigdowski, O. Kohlbacher, M. Fischer, Comparison of
507 different sample preparation techniques for NIR screening and their influence on the
508 geographical origin determination of almonds (Prunus dulcis MILL.), Food Control 115
509 (2020) 107302. https://doi.org/10.1016/j.foodcont.2020.107302.

510 [14] S. Seifert, S. Gundlach, S. Szymczak, Surrogate minimal depth as an importance measure
511 for variables in random forests, Bioinformatics 35 (2019) 3663–3671.
512 https://doi.org/10.1093/bioinformatics/btz149.

513 [15] S. Seifert, Application of random forest based approaches to surface-enhanced Raman
514 scattering data, Sci. Rep. 10 (2020) 5436. https://doi.org/10.1038/s41598-020-62338-8.

515 [16] V. Živanović, S. Seifert, D. Drescher, P. Schrade, S. Werner, P. Guttmann, G.P. Szekeres,
516 S. Bachmann, G. Schneider, C. Arenz, J. Kneipp, Optical Nanosensing of Lipid
517 Accumulation due to Enzyme Inhibition in Live Cells, ACS Nano 13 (2019) 9363–9375.
518 https://doi.org/10.1021/acsnano.9b04001.

519 [17] E. Borràs, J. Ferré, R. Boqué, M. Mestres, L. Aceña, O. Busto, Data fusion methodologies
520 for food and beverage authentication and quality assessment - a review, Anal. Chim. Acta
521 891 (2015) 1–14. https://doi.org/10.1016/j.aca.2015.04.042.

522 [18] Å. Rinnan, F. van den Berg, S.B. Engelsen, Review of the most common pre-processing
523 techniques for near-infrared spectra, TrAC Trends in Analytical Chemistry 28 (2009)
524 1201–1222. https://doi.org/10.1016/j.trac.2009.07.007.

525 [19] T. Isaksson, T. Næs, The Effect of Multiplicative Scatter Correction (MSC) and Linearity
526 Improvement in NIR Spectroscopy, Applied Spectroscopy 42 (1988) 1273–1284.
527 https://doi.org/10.1366/0003702884429869.

528 [20] J.E. Jackson, A User's Guide To Principal Components, John Wiley & Sons, 1991.

529 [21] J. Meurs, scaledata, 2021. https://github.com/jorismeurs/scaledata (accessed 28 January
530 2021).

531   [22] S. Watermann, C. Schmitt, T. Schneider, T. Hackl, Comparison of Regular, Pure Shift, and

532        Fast 2D NMR Experiments for Determination of the Geographical Origin of Walnuts,

533        Metabolites 11 (2021). https://doi.org/10.3390/metabo11010039.

534   [23] T.K. Ho, The random subspace method for constructing decision forests, IEEE Trans.

535        Pattern Anal. Machine Intell. 20 (1998) 832–844. https://doi.org/10.1109/34.709601.

536   [24] J.L. Fleiss, B. Levin, M.C. Paik, Statistical Methods for Rates and Proportions, third ed.,

537        Wiley, 2003.

538   [25] M.N. Wright, A. Ziegler, ranger A Fast Implementation of Random Forests for High

539        Dimensional Data in C++ and R, J. Stat. Soft. 77 (2017).

540        https://doi.org/10.18637/jss.v077.i01.

541   [26] U.S. Department of Agriculture, USDA Food and Nutrient Database for Dietary Studies

542        2017-2018. http://www.ars.usda.gov/nea/bhnrc/fsrg (accessed 25 February 2021).

543   [27] P.L. Benitez-Sánchez, M. Len-Camacho, R. Aparicio, A comprehensive study of hazelnut

544        oil composition with comparisons to other vegetable oils, particularly olive oil, European

545        Food Research and Technology 218 (2003) 13–19. https://doi.org/10.1007/s00217-003-

546        0766-4.

547   [28] J. Workman, L. Weyer, Practical guide and spectral atlas for interpretive near-infrared

548        spectroscopy, second ed., CRC Press, Boca Raton, FL, 2012.

549   [29] B. Worley, R. Powers, Multivariate Analysis in Metabolomics, Curr. Metabolomics 1

550        (2013) 92–107. https://doi.org/10.2174/2213235X11301010092.

551   [30] S. Guo, P. Rösch, J. Popp, T. Bocklitz, Modified PCA and PLS: Towards a better

552        classification in Raman spectroscopy–based biological applications, Journal of

553        Chemometrics 34 (2020). https://doi.org/10.1002/cem.3202.

554   [31] F. Gharibnezhad, L.E. Mujica, J. Rodellar, Applying robust variant of Principal

555        Component Analysis as a damage detector in the presence of outliers, Mechanical Systems

556        and Signal Processing 50-51 (2015) 467–479.

557        https://doi.org/10.1016/j.ymssp.2014.05.032.

558   [32] A. Tharwat, Classification assessment methods, ACI 17 (2021) 168–192.

559        https://doi.org/10.1016/j.aci.2018.08.003.

560   [33] T.R. Nichols, P.M. Wisner, G. Cripe, L. Gulabchand, Putting the Kappa Statistic to Use,

561        Qual Assur J 13 (2010) 57–61. https://doi.org/10.1002/qaj.481.

562   [34] T. Segelke, S. Schelm, C. Ahlers, M. Fischer, Food Authentication: Truffle (Tuber spp.)

563        Species Differentiation by FT-NIR and Chemometrics, Foods 9 (2020).

564        https://doi.org/10.3390/foods9070922.

565 [35] European Commission, Regulation (EC) No 510/2006 'Nocciola Romana' PDO, Official
566      Journal of the European Union (2008).

567 [36] E.W. Ciurczak, B. Igne, J. Workman, D.A. Burns (Eds.), Handbook of near-infrared
568      analysis, CRC Press/Taylor & Francis Group, Boca Raton, 2021.

569 [37] R.A. van den Berg, H.C.J. Hoefsloot, J.A. Westerhuis, A.K. Smilde, M.J. van der Werf,
570      Centering, scaling, and transformations: improving the biological information content of
571      metabolomics data, BMC Genomics 7 (2006) 142. https://doi.org/10.1186/1471-2164-7-
572      142.